



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Sede Amministrativa: Università degli Studi di Padova

Dipartimento di Ingegneria Industriale

CORSO DI DOTTORATO DI RICERCA IN: INGEGNERIA INDUSTRIALE

CURRICOLO: INGEGNERIA CHIMICA E AMBIENTALE

CICLO: XXXI

**PHARMACEUTICAL DEVELOPMENT AND MANUFACTURING
IN A QUALITY BY DESIGN PERSPECTIVE:
METHODOLOGIES FOR DESIGN SPACE DESCRIPTION**

Coordinatore: Ch.mo Prof. Paolo Colombo

Supervisore: Ch.mo Prof. Massimiliano Barolo

Dottorando: Gabriele Bano

Foreword

The fulfillment of the research results included in this Dissertation involved the intellectual and financial support of many people and institutions, to whom the author is very grateful.

Most of the research activity that led to the results reported in this Dissertation has been carried out at CAPE-Lab, Computer-Aided Process Engineering Laboratory, at the Department of Industrial Engineering of the University of Padova (Italy), under the supervision of prof. Massimiliano Barolo and prof. Fabrizio Bezzo. Part of the work was carried out at Rutgers University, Piscataway, NJ (USA) during a 6-month stay under the supervision of prof. Marianthi Ierapetritou, and part represents a collaboration with Dr. Emanuele Tomba from GlaxoSmithKline Biologicals SA- Siena (Italy).

Financial support has been provided by the University of Padova and by *Fondazione Ing. Aldo Gini*, Padova (Italy).

All the material reported in this Dissertation is original, unless explicit references to studies carried out by other people are indicated. In the following, a list of publications stemmed from this project is reported.

CONTRIBUTIONS IN INTERNATIONAL JOURNALS

Bano, G., Wang, Z., Facco, P., Bezzo, F., Barolo, M., Ierapetritou, M. (2018) A novel and systematic approach to identify the design space of pharmaceutical processes. *Comput. Chem. Eng.* **115**, 309-322.

Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018) Probabilistic design space determination in pharmaceutical development: a Bayesian/latent variable approach. *AIChE J.* **64**, 2438- 2449.

Bano, G., Meneghetti, N., Facco, P., Bezzo, F., Barolo M. (2017) Uncertainty back-propagation in PLS model inversion for design space determination in pharmaceutical product development. *Comput. Chem. Eng.* **101**, 110-124.

CONTRIBUTIONS SUBMITTED TO INTERNATIONAL JOURNALS

Bano, G., Facco, P., Ierapetritou, M., Bezzo, F., Barolo, M. (2018) Design space maintenance by online model adaptation in pharmaceutical manufacturing. Submitted to *Comput. Chem. Eng.*

CONTRIBUTIONS IN INTERNATIONAL JOURNALS (in preparation)

Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018) Design space description in pharmaceutical development and manufacturing: a review. *In preparation*.

CONTRIBUTIONS IN PEER-REVIEWED CONFERENCE PROCEEDINGS

Bano, G., Wang, Z., Facco, P., Bezzo, F., Barolo, M., Ierapetritou, M. (2018). Dimensionality reduction in feasibility analysis by latent variable modeling. In: *Computer-Aided Chemical Engineering 44, Proc. of the 13th International Symposium on Process Systems Engineering – PSE 2018, July 1-5 2018* (M.R. Eden, M.G. Ierapetritou, G.P. Towler, Eds.), Elsevier, Amsterdam (The Netherlands), 1477-1482.

CONFERENCE PRESENTATIONS

Bano, G., Wang Z., Facco P., Bezzo F., Barolo M., Ierapetritou M. (2018). Dimensionality reduction in feasibility analysis by latent variable modeling. Oral presentation at: *PSE 2018, July 1-5 2018* (San Diego; CA- U.S.A).

Bano, G., Facco, P., Bezzo, F., Barolo, M (2017) Handling parametric and measurement uncertainty for design space determination in pharmaceutical product development: a Bayesian approach. Oral presentation at: *World Congress of Chemical Engineering WCCE10*. October 1-5, Barcelona (Spain).

Bano, G., Meneghetti, N., Facco, P., Bezzo, F., Barolo, M. (2017) Determinazione dello spazio di progetto di un nuovo prodotto farmaceutico: caratterizzazione dell'incertezza. Oral presentation at: *Convegno GRICU 2017: gli orizzonti dell'ingegneria chimica*. September 12-14, Anacapri (NA; Italy).

CONTRIBUTIONS IN CONFERENCE PROCEEDINGS

Facco, P., Bano, G., Bezzo, F., Barolo M. (2017) Multivariate statistical approaches to aid pharmaceutical processes: product development and process monitoring in a Quality-by-design perspective. EuroPACT 2017. May 10-12, Potsdam (Germany).

Abstract

In the last decade, the pharmaceutical industry has been experiencing a period of drastic change in the way new products and processes are being conceived, due to the introduction of the *Quality by design* (QbD) initiative put forth by the pharmaceutical regulatory agencies (such as the Food and Drug Administration (FDA) and the European Medicines Agency (EMA)).

One of the most important aspects introduced in the QbD framework is that of design space (DS) of a pharmaceutical product, defined as “the multidimensional combination and interaction of input variables (e.g. material attributes) and process parameters that have been demonstrated to provide assurance of quality”. The identification of the DS represents a key advantage for pharmaceutical companies, since once the DS has been approved by the regulatory agency, movements within the DS do not constitute a manufacturing change and therefore do not require any further regulatory post-approval. This translates into an enhanced flexibility during process operation, with significant advantages in terms of productivity and process economics.

Mathematical modeling, both first-principles and data-driven, has proven to be a valuable tool to assist a DS identification exercise. The development of advanced mathematical techniques for the determination and maintenance of a design space, as well as the quantification of the uncertainty associated with its identification, is a research area that has gained increasing attention during the last years.

The objective of this Dissertation is to develop novel methodologies to assist the (i) determination of the design space of a new pharmaceutical product, (ii) quantify the assurance of quality for a new pharmaceutical product as advocated by the regulatory agencies, (iii) adapt and maintain a design space during plant operation, and (iv) design optimal experiments for the calibration of first-principles mathematical models to be used for design space identification.

With respect to the issue of **design space determination**, a methodology is proposed that combines surrogate-based feasibility analysis and latent-variable modeling for the identification of the design space of a new pharmaceutical product. Projection onto latent structures (PLS) is exploited to obtain a latent representation of the space identified by the model inputs (i.e. raw material properties and process parameters) and surrogate-based feasibility is then used to reconstruct the boundary of the DS on this latent representation, with significant reduction of the overall computational burden. The final result is a compact representation of the DS that can be easily expressed in terms of the original physically-relevant input variables (process

parameters and raw material properties) and can then be easily interpreted by industrial practitioners.

As regards the **quantification of “assurance” of quality**, two novel methodologies are proposed to account for the two most common sources of model uncertainty (structural and parametric) in the model-based identification of the DS of a new pharmaceutical product.

The first methodology is specifically suited for the quantification of assurance of quality when a PLS model is to be used for DS identification. Two frequentist analytical models are proposed to back-propagate the uncertainty from the quality attributes of the final product to the space identified by the set of raw material properties and process parameters of the manufacturing process. It is shown how these models can be used to identify a subset of input combinations (i.e., raw material properties and process parameters) within which the DS is expected to lie with a given degree of confidence. It is also shown how this reduced space of input combinations (called experiment space) can be used to tailor an experimental campaign for the final assessment of the DS, with a significant reduction of the experimental effort required with respect to a non-tailored experimental campaign. The validity of the proposed methodology is tested on granulation and roll compaction processes, involving both simulated and experimental data.

The second methodology proposes a joint Bayesian/latent-variable approach, and the assurance of quality is quantified in terms of the probability that the final product will meet its specifications. In this context, the DS is defined in a probabilistic framework as the set of input combinations that guarantee that the probability that the product will meet its quality specifications is greater than a predefined threshold value. Bayesian multivariate linear regression is coupled with latent-variable modeling in order to obtain a computationally friendly implementation of this probabilistic DS. Specifically, PLS is exploited to reduce the computational burden for the discretization of the input domain and to give a compact representation of the DS. On the other hand, Bayesian multivariate linear regression is used to compute the probability that the product will meet the desired quality for each of the discretization points of the input domain. The ability of the methodology to give a scientifically-driven representation of the probabilistic DS is proved with three case studies involving literature experimental data of pharmaceutical unit operations.

With respect to the issue of the **maintenance of a design space**, a methodology is proposed to adapt in real time a model-based representation of a design space during plant operation in the presence of process-model mismatch.

Based on the availability of a first-principles model (FPM) or semi-empirical model for the manufacturing process, together with measurements from plant sensors, the methodology jointly exploits (i) a dynamic state estimator and (ii) feasibility analysis to perform a risk-based

online maintenance of the DS. The state estimator is deployed to obtain an up-to-date FPM by adjusting in real-time a small subset of the model parameters. Feasibility analysis and surrogate-based feasibility analysis are used to update the DS in real-time by exploiting the up-to-date FPM returned by the state estimator. The effectiveness of the methodology is shown with two simulated case studies, namely the roll compaction of microcrystalline cellulose and the penicillin fermentation in a pilot scale bioreactor.

As regards the **design of optimal experiments** for the calibration of mathematical models for DS identification, a model-based design of experiments (MBDoe) approach is presented for an industrial freeze-drying process. A preliminary analysis is performed to choose the most suitable process model between different model alternatives and to test the structural consistency of the chosen model. A new experiment is then designed based on this model using MBDoe techniques, in order to increase the precision of the estimates of the most influential model parameters. The results of the MBDoe activity are then tested both *in silico* and on the real equipment.

Riassunto

Negli ultimi anni, le modalità con le quali le aziende farmaceutiche sviluppano e successivamente producono nuovi prodotti sta subendo un cambiamento radicale, a causa dell'introduzione dell'iniziativa *Quality by Design* (QbD) da parte delle agenzie di regolamentazione, come l'americana Food and Drug Administration (FDA) o l'agenzia europea del farmaco (European Medicines Agency; EMA).

Lo scopo dell'iniziativa è quello di promuovere l'adozione di un approccio scientifico e sistematico nelle fasi di (i) sviluppo di prodotto e di processo, (ii) trasferimento di processo da scala laboratorio a scala industriale e (iii) monitoraggio e controllo di processo. Il fine ultimo dell'iniziativa è quello di stimolare l'utilizzo di tecniche scientifiche rigorose all'interno di un contesto industriale, come quello farmaceutico, tipicamente orientato all'empirismo e poco propenso al cambiamento (in larga parte a causa degli stringenti vincoli regolatori imposti in passato).

Uno dei concetti più importanti introdotti nell'ambito dell'iniziativa QbD è quello di spazio di progetto (*design space*; DS) di un nuovo prodotto farmaceutico. Lo spazio di progetto è definito come l'insieme delle combinazioni delle proprietà delle materie prime (ad esempio, la composizione del principio attivo) e dei parametri di processo (portate, temperature, etc...) per le quali la qualità finale del prodotto risulta essere garantita. Le aziende interessate possono (in forma volontaria) identificare lo spazio di progetto di un nuovo prodotto e sottoporlo all'approvazione delle agenzie regolatorie: se il DS viene approvato, il processo può essere esercito, senza dover richiedere una nuova approvazione all'ente regolatore, anche cambiando le condizioni operative e/o le proprietà delle materie prime, purché la loro combinazione rimanga all'interno dello spazio di progetto. Ciò comporta un notevole vantaggio in termini di flessibilità nell'esercizio del processo produttivo, consentendo alle aziende di ottimizzare la marcia dell'impianto con conseguente notevole riduzione dei costi di esercizio.

Al fine di ottenere l'approvazione del DS, le agenzie regolatorie richiedono che vengano rispettati due concetti chiave, ovvero che l'azienda: (i) *dimostri* con un metodo scientifico rigoroso la procedura utilizzata per l'identificazione dello spazio di progetto; (ii) *quantifichi* il grado di certezza (e quindi, la probabilità) che la qualità del prodotto venga garantita qualora il processo venga esercito all'interno dello spazio di progetto.

Da un punto di vista ingegneristico, entrambe le richieste possono essere affrontate attraverso l'utilizzo congiunto di modellazione matematica e osservazioni sperimentali. Le tecniche modellistiche possono essere basate esclusivamente su dati storici di laboratorio o di impianto (modelli cosiddetti *data-driven*), su una conoscenza approfondita dei fenomeni fisici e chimici che avvengono nel processo (modelli a principi primi o *knowledge-driven*) o su una

combinazione di entrambi (modelli semi-empirici). In tutti i suddetti casi, vi è la necessità di sviluppare metodologie sistematiche che, basandosi sulla tipologia di modello prescelta, consentano di soddisfare entrambe le richieste degli enti regolatori.

Nonostante il numero di contributi scientifici in questo ambito sia aumentato in modo significativo negli ultimi anni, permane la necessità di sviluppare e integrare le tecniche esistenti per sviluppare metodologie che consentano di risolvere quattro aspetti ancora irrisolti, e cioè:

1. identificare con il minor carico computazionale possibile lo spazio di progetto, e ottenerne una rappresentazione immediata che possa essere interpretata anche dal personale non esperto in modellazione matematica;
2. fornire una metrica dell'incertezza associata alla predizione, mediante modello, dello spazio di progetto;
3. adattare la predizione dello spazio di progetto durante la marcia dell'impianto, per far fronte a mancanze strutturali del modello e/o deviazioni parametriche dello stesso;
4. progettare esperimenti ottimali per l'identificazione di modelli a principi primi al fine di individuare rapidamente lo spazio di progetto di nuovi prodotti.

L'obiettivo di questa Dissertazione è proporre tecniche di modellazione avanzata per affrontare le quattro problematiche sopra descritte. L'approccio scientifico adottato si basa sia sullo sviluppo di nuove procedure metodologiche, sia sull'integrazione e combinazione di tecniche note.

La problematica riguardante **la quantificazione dell'incertezza** associata alla predizione del DS di un nuovo prodotto farmaceutico viene affrontata nei Capitoli 3 e 4.

In particolare, nel Capitolo 3 viene presentata una metodologia per quantificare l'incertezza associata ad una predizione del DS qualora un modello di proiezione su strutture latenti (PLS: *projection onto latent structures*) venga utilizzato a tal scopo. La metodologia si basa sullo sviluppo di due nuovi modelli di propagazione dell'incertezza dallo spazio della qualità di prodotto (cioè dallo spazio delle uscite del modello) allo spazio delle proprietà delle materie prime e parametri di processo (cioè lo spazio degli ingressi del modello). I due modelli vengono derivati in modo analitico sfruttando nozioni di statistica frequentista, e successivamente integrati in un contesto di inversione di modelli a variabili latenti. La procedura consente di individuare un sottospazio del dominio degli ingressi, chiamato spazio degli esperimenti (ES; *experiment space*) nel quale è possibile garantire, con un livello di fiducia preassegnato, che lo spazio di progetto risieda. Da un punto di vista pratico, l'individuazione dello spazio degli esperimenti può essere sfruttata per condurre una campagna sperimentale mirata in una zona ristretta del dominio degli ingressi, al fine di ottenere più rapidamente l'identificazione definitiva del DS del nuovo prodotto. Ciò comporta un vantaggio in termini di riduzione dei tempi e dei costi della campagna sperimentale, con conseguente guadagno in termini di

competitività da parte dell'azienda che ne voglia trarre beneficio. L'efficacia della metodologia proposta viene verificata su diversi casi studio riguardanti alcune operazioni unitarie tipiche dell'industria farmaceutica (ad esempio, granulazione a secco/a umido), utilizzando sia dati simulati al calcolatore che dati sperimentali di letteratura.

Nel Capitolo 4 viene presentata una metodologia che, attraverso la combinazione di tecniche di statistica Bayesiana e di regressione multivariata, consente di ottenere una metrica (ossia una valutazione quantitativa) dell'incertezza associata alla predizione di uno spazio di progetto. Tale metrica viene definita come la probabilità (intesa in termini Bayesiani) che il prodotto soddisfi le specifiche di qualità suddette. In quest'ottica, lo spazio di progetto viene identificato come l'insieme di combinazioni di caratteristiche delle materie prime e parametri di processo per i quali la *probabilità* che il prodotto sia in specifica supera un certo valore di soglia definito dall'utente e consistente con le richieste dell'ente regolatorio.

La metodologia proposta si basa su due idee chiave: (i) l'utilizzo di un modello PLS per la riduzione delle dimensioni dello spazio degli ingressi, con conseguente riduzione del carico computazionale e ottenimento di una rappresentazione compatta di tale spazio; (ii) l'utilizzo di metodi Monte Carlo basati su catene di Markov per la ricostruzione delle funzioni di densità di probabilità (PDF: *probability density function*) delle uscite del modello per diversi valori degli ingressi. Il risultato finale è una rappresentazione compatta (in molte situazioni, bidimensionale) dello spazio di progetto *probabilistico*, garantendo un pieno adempimento alle richieste degli enti regolatori. L'efficacia della metodologia viene verificata anche in questo caso su diversi casi di studio relativi ad operazioni unitarie tipiche dell'industria farmaceutica, coinvolgenti sia dati simulati al calcolatore che dati sperimentali di letteratura.

Per quanto riguarda la problematica dell'**identificazione dello spazio di progetto** di un nuovo prodotto farmaceutico, nel Capitolo 5 viene presentata una metodologia che combina un modello PLS e una tecnica denominata analisi di fattibilità basata su surrogati (*surrogate-based feasibility analysis*) per identificare il DS. Il modello PLS viene utilizzato per ridurre la dimensione dello spazio degli ingressi mediante l'introduzione di un numero ridotto di variabili latenti, cioè combinazioni lineari degli ingressi reali. L'analisi di fattibilità basata su una funzione base di tipo radiale (*radial-basis function*; RBF) viene successivamente utilizzata per identificare lo spazio di progetto all'interno dello spazio latente. La riduzione della dimensione dello spazio degli ingressi viene utilizzata per rendere l'analisi di fattibilità possibile da un punto di vista computazionale, con un significativo vantaggio sia in termini di durata totale della simulazione, sia in termini di precisione ottenibile per l'identificazione del DS. L'efficacia della metodologia proposta viene dimostrata sia su una linea di produzione continua (coinvolgente sei operazioni unitarie) di compresse, sia su esempi relativi a dati simulati al calcolatore.

Con riferimento alla problematica **dell'adattamento in linea della predizione dello spazio di progetto**, nel Capitolo 6 viene proposta una metodologia che consente di correggere in linea la predizione del DS sulla base delle misure raccolte durante la marcia dell'impianto. La metodologia sfrutta in maniera congiunta uno stimatore dinamico di stato e l'analisi di fattibilità descritta al Capitolo 5 per adattare in tempo reale la predizione del DS in presenza di un disallineamento fra le predizioni di modello e il comportamento dell'impianto (*process/model mismatch*). Basandosi sulle misure in linea disponibili, lo stimatore dinamico viene utilizzato per correggere in tempo reale alcuni dei parametri critici del modello e per ottenere una stima accurata dello stato del sistema. L'analisi di fattibilità (basata o meno su modelli surrogati) viene successivamente utilizzata per predire lo spazio di progetto sulla base della correzione apportata al modello da parte dello stimatore. Il DS adattativo può quindi essere utilizzato per ottimizzare in tempo reale il processo produttivo, supportare in tempo reale il processo decisionale (ad esempio, la chiusura o apertura manuale di determinate valvole da parte degli operatori), o fungere da base per la definizione dei set-point del sistema di controllo implementato nell'impianto. La procedura proposta è stata verificata in silico su un processo continuo di granulazione di un formulato e su un processo di produzione di penicillina in un bioreattore. Entrambi i casi studio analizzati dimostrano l'affidabilità della metodologia.

Per quanto concerne **la progettazione di esperimenti ottimali per l'identificazione di modelli a principi primi**, da usarsi successivamente per identificare il DS di un nuovo prodotto, nel Capitolo 7 viene presentata un'applicazione relativa ad un processo industriale di liofilizzazione. Dopo una fase preliminare di identificazione del modello a principi primi da utilizzare per la descrizione del processo, viene eseguita una verifica della consistenza strutturale di tale modello e vengono identificati i parametri che maggiormente influiscono sulle sue predizioni. Sulla base di tale analisi preliminare, viene progettato un nuovo esperimento attraverso l'utilizzo di tecniche di progettazione ottimale degli esperimenti basata su modello (*model-based design of experiments*; MBD_{oE}). L'esperimento viene progettato in modo tale da massimizzare l'informazione che può essere estratta dai dati sperimentali ai fini della calibrazione di modello. L'esperimento progettato è stato condotto nell'impianto industriale, e i dati ottenuti hanno consentito di migliorare significativamente la precisione e l'accuratezza della stima dei parametri del modello.

Table of contents

CHAPTER 1 – MOTIVATION AND STATE OF THE ART	1
1.1 PHARMACEUTICAL INDUSTRY: FACTS AND FIGURES	1
1.1.1 ECONOMIC OUTLOOK	1
1.1.2 PHARMACEUTICAL R&D: PAST, PRESENT & FUTURE	4
1.1.3. PHARMACEUTICAL MANUFACTURING: PAST, PRESENT & FUTURE	9
1.1.4 REGULATORY ASPECTS: PAST, PRESENT & FUTURE.....	12
1.1.4 REGULATORY ASPECTS: PAST, PRESENT & FUTURE.....	14
1.2 QUALITY BY DESIGN PARADIGMS: REGULATORY OVERVIEW	17
1.2.1 PHARMACEUTICAL DEVELOPMENT	19
1.2.1.1 Step (a): identification of the Quality Target Product Profile (QTPP)	20
1.2.1.2 Step (b): identification of the product Critical Quality Attributes (CQAs)	20
1.2.1.3 Step (c): product and process design and understanding	21
1.2.1.4 Step (d): design space (DS) description.....	23
1.2.1.5 Step (e): definition of a control strategy	25
1.2.2 TECHNOLOGY TRANSFER	26
1.2.3 COMMERCIAL MANUFACTURING & CONTINUAL PROCESS IMPROVEMENT	26
1.2.4 TECHNICAL AND ECONOMIC BENEFITS OF QBD: A CRITICAL REVIEW.....	27
1.3 MATHEMATICAL MODELLING FOR QBD IMPLEMENTATION	29
1.4 DESIGN SPACE APPLICATIONS THROUGHOUT DRUG LIFECYCLE	32
1.4.1 DESIGN SPACE DESCRIPTION AND UNCERTAINTY QUANTIFICATION	33
1.4.1.1 Design space description strategies for data-driven models	33
1.4.1.2 Design space description strategies for knowledge-driven models.....	36
1.4.1.3 Design space description strategies for hybrid models	38
1.4.2 DESIGN SPACE PROPAGATION FROM INDIVIDUAL UNITS TO AN ENTIRE PROCESS	38
1.4.3 DESIGN SPACE SCALE-UP AND SCALE-OUT	39
1.4.4 DESIGN SPACE MAINTENANCE	41
1.5 OBJECTIVES OF THE RESEARCH.....	42
1.6 DISSERTATION ROADMAP	45
CHAPTER 2 - METHODS	49
2.1 PROJECTION ON LATENT STRUCTURES (PLS).....	50

2.1.1 PRE- AND POST- PLS MODELING ACTIVITIES.....	51
2.1.1.1 Data pretreatment	52
2.1.1.2 Selection of the optimal number of LVs	52
2.1.1.3 PLS model diagnostics	53
2.2 PREDICTION UNCERTAINTY IN PLS MODELING.....	54
2.2.1 OLS-TYPE METHODS	55
2.2.2 LINEARIZATION-BASED METHODS.....	57
2.2.3 RE-SAMPLING BASED METHODS.....	57
2.2.4 THE UNSCRAMBLER METHOD.....	59
2.2.5 ESTIMATION OF σ^2 AND DEGREES OF FREEDOM	60
2.2.5.1 Summary	61
2.3 KEY CONCEPTS OF BAYESIAN STATISTICS	63
2.3.1 PRIOR DISTRIBUTION $p(\boldsymbol{\theta})$	64
2.3.2 POSTERIOR DISTRIBUTION $p(\boldsymbol{\theta} \mathbf{y})$	65
2.3.3 POSTERIOR PREDICTIVE DISTRIBUTION OF A NEW RESPONSE \mathbf{y}_p	65
2.3.4 CREDIBLE REGION VS CONFIDENCE REGION	66
2.4 BAYESIAN MULTIVARIATE LINEAR REGRESSION.....	67
2.4.1 LIKELIHOOD FUNCTION IN MULTIVARIATE LINEAR REGRESSION	68
2.4.2 PRIOR DISTRIBUTIONS	68
2.4.3 POSTERIOR DISTRIBUTION	70
2.4.4 POSTERIOR PREDICTIVE DISTRIBUTION OF A NEW RESPONSE.....	70
2.5 SURROGATE-BASED FEASIBILITY ANALYSIS	72
2.5.1 SURROGATE MODELS	72
2.5.1.1 Kriging surrogate	73
2.5.1.2 Radial-basis function (RBF) surrogate	73
2.5.2 ADAPTIVE SAMPLING	74
2.5.3 SURROGATE ACCURACY	75
2.6 STATE ESTIMATION ALGORITHMS	77
2.6.1 EXTENDED KALMAN FILTER (EKF).....	77
2.6.1.1 Prediction step	77
2.6.1.2 Measurement update step	78
2.6.2 UNSCENTED KALMAN FILTER (UKF)	79
2.6.2.1 Prediction step	79
2.6.2.2 Measurement update step	81
2.6.3 ENSEMBLE KALMAN FILTER (EnKF)	82
2.6.3.1 Prediction step	82

2.6.3.2 Measurement update step	83
2.7 STATE ESTIMATION AND ONLINE MODEL RECALIBRATION.....	83
CHAPTER 3 – UNCERTAINTY BACK-PROPAGATION IN PLS MODELING FOR DESIGN SPACE DETERMINATION	85
3.1 INTRODUCTION.....	85
3.2 METHODS	87
3.2.1 PROJECTION TO LATENT STRUCTURES (PLS) FOR PHARMACEUTICAL PRODUCT DEVELOPMENT	87
3.2.2 FORMULATION OF PREDICTION UNCERTAINTY.....	90
3.2.3 BACK-PROPAGATION OF UNCERTAINTY IN PLS MODEL INVERSION.....	91
3.3 UNCERTAINTY BACK-PROPAGATION AND EXPERIMENT SPACE DETERMINATION	95
3.3.1 PROPOSED METHODOLOGY	95
3.4 CASE STUDIES	101
3.4.1 CASE STUDY #1: MATHEMATICAL EXAMPLE	101
3.4.2 CASE STUDY #2: ROLL COMPACTION OF AN INTERMEDIATE DRUG LOAD FORMULATION..	101
3.4.3 CASE STUDY #3: HIGH-SHEAR WET GRANULATION OF A PHARMACEUTICAL BLEND	102
3.4.4 CASE STUDY #4: DRY GRANULATION BY ROLLER COMPACTION.....	103
3.4.5 CASE STUDY #5: WET GRANULATION	103
3.5 RESULTS FOR CASE STUDY #1	103
3.5.1 QUALITY TARGET DESCRIBED BY AN EQUALITY CONSTRAINT	103
3.5.2 QUALITY TARGET DESCRIBED BY AN INEQUALITY CONSTRAINT.....	108
3.6 RESULTS FOR CASE STUDY #2	111
3.7 RESULTS FOR CASE STUDY #3	113
3.8 RESULTS FOR CASE STUDY #4	113
3.9 RESULTS FOR CASE STUDY #5	114
3.10 CONCLUSIONS.....	115
CHAPTER 4 – A BAYESIAN/LATENT VARIABLE APPROACH FOR DESIGN SPACE DETERMINATION.....	117
4.1 INTRODUCTION.....	117
4.2 REVIEW OF PARTIAL LEAST-SQUARES (PLS) REGRESSION.....	119
4.3 BAYESIAN DESIGN SPACE: MATHEMATICAL FORMULATION.....	120
4.4 BAYESIAN IDENTIFICATION OF THE DESIGN SPACE	123
4.4.1 PROPOSED METHODOLOGY	124
4.4.2 SELECTION OF THE OPTIMAL NUMBER OF SAMPLES	128
4.5 CASE STUDIES	130

4.5.1 CASE STUDY #1: DRY GRANULATION OF A PHARMACEUTICAL BLEND BY ROLLER COMPACTION	130
4.5.2 CASE STUDY #2: HIGH-SHEAR WET GRANULATION OF A PHARMACEUTICAL BLEND	131
4.5.3 CASE STUDY #3: ROLL COMPACTION OF AN INTERMEDIATE DRUG LOAD FORMULATION ..	132
4.6 RESULTS	133
4.6.1 RESULTS FOR CASE STUDY #1	133
4.6.2 RESULTS FOR CASE STUDY #2	135
4.6.3 RESULTS FOR CASE STUDY #3	138
4.7 DISCUSSION	139
4.8 CONCLUSIONS	140
CHAPTER 5 – FEASIBILITY ANALYSIS AND LATENT VARIABLE MODELING FOR DESIGN SPACE DETERMINATION	143
5.1 INTRODUCTION	143
5.2 MATHEMATICAL BACKGROUND	145
5.2.1 PARTIAL LEAST SQUARES REGRESSION (PLS)	146
5.2.2 FEASIBILITY ANALYSIS	147
5.2.3 SURROGATE-BASED FEASIBILITY ANALYSIS	147
5.2.4 RADIAL BASIS FUNCTION (RBF) SURROGATE MODEL	148
5.2.5 ADAPTIVE SAMPLING	149
5.2.6 SURROGATE ACCURACY	150
5.3 FEASIBILITY ANALYSIS AND REDUCTION OF INPUT SPACE DIMENSIONALITY	151
5.3.1 PROPOSED METHODOLOGY	151
5.4 CASE STUDIES	158
5.4.1 CASE STUDY #1: LOW DIMENSIONAL TEST PROBLEM	159
5.4.2 CASE STUDY #2: HIGH DIMENSIONAL TEST PROBLEM	160
5.4.3 CASE STUDY #3: SIMULATION OF CONTINUOUS DIRECT COMPACTION FOR PHARMACEUTICAL PRODUCTION	161
5.5 RESULTS	163
5.5.1 RESULTS FOR CASE STUDY #1	163
5.5.2 RESULTS FOR CASE STUDY #2	165
5.5.3 RESULTS FOR CASE STUDY #3	167
5.6 CONCLUSIONS	170
CHAPTER 6 – DESIGN SPACE MAINTENANCE BY ONLINE MODEL ADAPTATION IN PHARMACEUTICAL MANUFACTURING	173
6.1 INTRODUCTION	173
6.2 METHODS	175
6.2.1 STATE ESTIMATOR FOR DAE SYSTEMS	175

6.2.2 FEASIBILITY ANALYSIS FOR DYNAMICAL SYSTEMS.....	177
6.3 FEASIBILITY ANALYSIS FOR DESIGN SPACE DETERMINATION	179
6.3.1 FEASIBLE REGION AND DESIGN SPACE.....	179
6.3.2 LIMITATIONS OF EXISTING MODEL-BASED DESIGN SPACE DESCRIPTION METHODOLOGIES	180
6.4 DESIGN SPACE MAINTENANCE BY ONLINE MODEL ADAPTATION.....	182
6.4.1. PROPOSED METHODOLOGY.....	183
6.5 CASE STUDIES	185
6.5.1 CASE STUDY #1: ROLLER COMPACTION OF MICROCRYSTALLINE CELLULOSE	186
6.5.2 CASE STUDY #2: PENICILLIN FERMENTATION	188
6.6 RESULTS AND DISCUSSION	190
6.6.1 RESULTS AND DISCUSSION FOR CASE STUDY #1	190
6.6.1.1 SIMULATION SET-UP	190
6.6.1.2. Offline model-based DS identification (step #1).....	190
6.6.1.3 State estimation for online model adaptation (steps #2 and #3)	193
6.6.1.4 Online model-based DS maintenance (step #4).....	194
6.6.2 RESULTS AND DISCUSSION FOR CASE STUDY #2	197
6.6.2.1 Simulation set-up	197
6.6.2.2 Offline model-based DS identification (step #1)	198
6.6.2.3 Identification of uncertain model parameters (step #2).....	200
6.6.2.4 State estimation for online model adaptation (step #3).....	201
6.6.2.5 Online model-based DS maintenance (step #4)	202
6.7. CONCLUSIONS	204

CHAPTER 7 – MODEL-BASED DESIGN OF EXPERIMENTS IN PHARMACEUTICAL FREEZE DRYING.....

7.1 FREEZE DRYING: PROCESS DESCRIPTION.....	208
7.1.1 PRIMARY DRYING	210
7.2 MATHEMATICAL MODELING OF PRIMARY DRYING.....	212
7.3 MODEL-BASED DESIGN OF EXPERIMENTS (MBDOE): BRIEF OVERVIEW	217
7.3.1 DESIGN OF AN OPTIMAL EXPERIMENT	217
7.3.2 ADDITIONAL MBDoE ACTIVITIES.....	219
7.4 PRELIMINARY ASSESSMENT OF FREEZE-DRYING MODEL PERFORMANCE	219
7.4.1 DESCRIPTION OF THE HISTORICAL DATASETS.....	220
7.4.1.1 Calibration datasets	220
7.4.1.2 Validation datasets	223
7.4.2 STEP #1: GRAVIMETRIC ESTIMATION OF K_v	225

7.4.2.1 Results for dataset A	226
7.4.3. STEP #2: SENSITIVITY ANALYSIS FOR K_v	230
7.5.4 STEP #3: SENSITIVITY ANALYSIS FOR R_p	232
7.4.5 STEP #4: PARAMETER ESTIMATION	235
7.4.6 STEP #5: MODEL VALIDATION	238
7.5 MODEL-BASED DESIGN OF EXPERIMENTS FOR PRIMARY DRYING	240
7.5.1 PROBLEM STATEMENT	241
7.5.2 DESIGNED EXPERIMENT: “IN SILICO” RESULTS	242
7.6 FINAL REMARKS AND FUTURE WORK	244
APPENDIX A	255
APPENDIX B	259
REFERENCES	263
ACKNOWLEDGMENTS	289

List of symbols

Acronyms

AR	=	acceptance region
CI	=	confidence interval
CMA	=	critical material attribute
CP	=	coverage probability
CPP	=	critical process parameter
CPPD	=	conditional posterior predictive distribution
CQA	=	critical quality attributes
CR	=	confidence region
DAE	=	differential- algebraic equation
df	=	degrees of freedom
DoE	=	design of experiments
DQ	=	direct quadrature
DS	=	design space
DSp	=	design space projection
EI	=	expected improvement
EMA	=	European Medicines Agency
ES	=	experiment space
ESS	=	experiment space shrinkage
FDA	=	Food and Drug Administration
FPM	=	first-principles model
ICH	=	international conference on harmonization for technical requirements for registration of pharmaceuticals for human use
KPI	=	key performance indicator
KS	=	knowledge space
KSA	=	Kennard-Stone's algorithm
LSEF	=	least-squares error functional
LV	=	latent variable
LVM	=	latent variable model

MCMC	=	Markov-chain Monte-Carlo
NIPALS	=	noniterative partial least squares
ODE	=	ordinary differential equation
PCA	=	principal component analysis
PLS	=	partial least-squares
PPD	=	posterior predictive distribution
PSD	=	particle size distribution
QbD	=	quality by design
RBF	=	radial-basis function
R&D	=	research and development
SEC	=	standard error of calibration
SPE	=	squared prediction error
SQP	=	sequential quadratic programming

Chapter 1

Motivation and state of the art*

The objective of this Chapter is to provide an overview of the background and motivations of this Dissertation. First, facts and figures on the current state and future trends of the pharmaceutical industry are presented. Then, the Quality by design (QbD) initiative that has been introduced in the pharmaceutical context is analyzed from a regulatory perspective, with particular focus on the concept of design space (DS) of a new drug. The main challenges that, from a process systems engineering perspective, need to be addressed for the identification and maintenance of a design space are then pinpointed. Finally, the state-of-the art on the ways these challenges have, up to now, being addressed in the scientific literature is presented and critically discussed.

1.1 Pharmaceutical industry: facts and figures

This section provides facts and figures on the state of the pharmaceutical industry, with particular focus on four different thematic areas: *i*) economic evolution; *ii*) pharmaceutical *research and development*; *iii*) pharmaceutical *manufacturing*; *iv*) role of the regulatory agencies.

1.1.1 Economic outlook

Whilst the economic outlook is characterized by positive expectations on the healthcare market growth and global medicine expenditure, the pharmaceutical industry is facing unprecedented challenges due to a large number of patent expirations, lower market growth in the so-called *pharmerging* countries¹ and increased pricing and market pressures. A recent market report proposed by the QuantileIMS Institute (QuantileIMS Institute, 2016) forecasts that the total expenditure for drugs will reach \$1.5 trillion by 2021, with an increase of 33% with respect to

* Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018) Design space description in pharmaceutical development and manufacturing: a review. *In preparation*.

¹ The countries are: China, Brazil, India, Russia, Mexico, Turkey, Poland, Saudi Arabia, Indonesia, Egypt, Philippines, Pakistan, Vietnam, Bangladesh, Argentina, Algeria, Colombia, South Africa, Chile, Nigeria, Kazakhstan.

the global medicine expenditure of 2016, but with a moderate decrease on the compound annual growth rate² (CAGR), that is expected to range between 4% and 7 % with respect to the nearly 9% annual growth rate of the years 2014 and 2015. A more conservative market report that has been recently proposed by EvaluatePharma[®] (EvaluatePharma, 2018) suggests that the above figures should be revised to \$ 1.2 trillion by 2024, with an average annual growth rate of 6.4%, as a consequence of the sales losses due to the increased genericization of drugs and due to a slower performance of the market of biosimilars. Similar figures are reported in the economic analysis suggested by Deloitte (Deloitte, 2018). As reported in Table 1.1, the US market is expected to remain the world's largest pharmaceutical market, but the annual growth rate is forecast to decrease significantly from the 12% reached in 2015 to 6-7% in the next five years. On the other hand, China is expected to continue playing a key role as the second largest market, with a robust average annual growth rate of 12%. As for Europe, the situation is expected to be less thriving than in the rest of the highly industrialized countries, with a forecast annual growth rate of 1-4 %, mainly due to the weak economic growth in the region. Finally, the growth of the pharmaceutical market is expected to proceed at a slower pace in the *pharmerging* countries (excluding China), from an average of 7% in the year 2012-2016 to an average of 4% for the years 2016-2021.

Table 1.1 Top 20 countries ranking forecast by 2021 based on the pharmaceutical market size (100 = U.S. market size). Adapted from *QuantileIMS Institute, 2016*.

Rank #	Country	Market size index (100 = U.S.) in constant USD
1	U.S.	100
2	China	25
3	Japan	14
4	Germany	8
5	Brazil	6
6	U.K.	6
7	Italy	5
8	France	5
9	India	5
10	Spain	4
11	Canada	4
12	South Korea	2
13	Russia	2
14	Turkey	2
15	Australia	2
16	Mexico	2
17	Saudi Arabia	1
18	Poland	1
19	Argentina	1
20	Egypt	1

² CAGR is the rate of return of an investment that would be obtained by investing the profits at the end of each year.

The general consensus on the significant increase of the global medicine expenditure for the next years is supported by the fact that the global population is expected to increase by an average 1.24% by 2030, with an increase of people aged 65-80 from 22% in 2000 to 28% in 2030 (QuantileIMS Institute, 2016). This aspect, together with an easier availability of drugs to more people due to the increasing urbanization and growing of the middle class (Rauschnabel, 2018), represents a key driver for the expansion of the healthcare market in the next years.

The top 20 pharmaceutical companies are expected to keep their leading position during the next years, but a reduction from 63.2% (in 2017) to 51.9% (in 2024) of the total market share of these companies is expected to happen in favor of the other market players. Table 1.2 gives an overview of the top 20 pharma companies and their forecast prescription drug sales for the period 2017-2024 (EvaluatePharma, 2018). The same table suggests that the next years will be characterized by an increased competition between the different pharma players, and therefore innovation and differentiation will play a key role for the maintenance of the market shares.

Table 1.2 Worldwide (WW) prescription drug sales and market share for the top 20 pharma companies for the period 2017-2024. Adapted from EvaluatePharma, 2018.

Rank #	Company	Prescription sales (\$ bn)		Market share (%)		Rank change (+/-)
		2017	2024	2017	2024	
1	Novartis	41.9	53.2	5.3	4.4	+1
2	Pfizer	45.4	51.2	5.8	4.3	-1
3	Roche	41.7	50.6	5.3	4.2	+0
4	Johnson & Johnson	34.4	47.4	4.4	3.9	+1
5	Sanofi	34.1	44.2	4.3	3.7	+1
6	GlaxoSmithKline	28.7	38.4	3.6	3.2	+1
7	Merck & co.	35.4	38.0	4.5	3.2	-3
8	AbbVie	27.7	37.2	3.5	3.1	+0
9	AstraZeneca	19.8	31.7	2.5	2.6	+2
10	Bristol-Myers Squibb	19.3	28.7	2.4	2.4	+2
11	Amgen	21.8	24.8	2.8	2.1	-1
12	Novo Nordisk	17.0	24.6	2.2	2.0	+4
13	Celgene	12.9	23.7	1.6	2.0	+8
14	Eli Lilly	18.5	22.2	2.3	1.8	-1
15	Bayer	17.7	19.7	2.2	1.6	+0
16	Gilead Sciences	25.7	19.0	3.3	1.6	-7
17	Boehringer Ingelheim	14.3	18.3	1.8	1.5	+2
18	Shire	14.4	17.7	1.8	1.5	+0
19	Takeda	13.3	17.0	1.7	1.4	+1
20	Allergan	14.9	16.8	1.9	1.4	-3
Total top 20		498.8	624.7	63.2	51.9	
Other		290.0	578.8	36.8	48.1	
Total		788.8	1203.5	100.0	100.0	

In terms of employment, the pharmaceutical industry is expected to maintain a leading role in the manufacturing sector, thanks to a steady employment growth both in high- and low-income countries. A recent review from the International Federation of Pharmaceutical Manufacturers & Associations (IFPMA, 2017) reports that the pharmaceutical industry employed around 5.1 million people worldwide in 2014, with an increase of 1.5 million with respect to the 2006 and with steady expectations on the annual growth for the period 2014–2021. Moreover, due to its peculiar feature of employing high-skilled workers with high qualifications, the pharmaceutical sector induces the creation of many other indirect jobs within the manufacturing sector. The qualified training and the direct exposure to innovative processes and new technologies that pharmaceutical employees experience also represents an asset for the entire workforce. In fact, it has been proven (WifOR, 2016) that many pharmaceutical workers decide to exploit the knowledge gained during their working experience to start new companies or to transfer their skills to other manufacturing sectors, thus fostering economic development.

Different positions and solutions have been proposed to boost the economic performance of the pharmaceutical industry with respect to the forecast figures described above. However, the most important analysts (QuantileIMS Institute, 2016; Deloitte, 2018; EvaluatePharma, 2018; KPMG, 2018; McKinsey, 2018; PwC, 2018;) seem to agree on the following key actions that pharmaceutical companies should implement to increase their competitiveness:

1. Increase R&D productivity in terms of discovery and launch of new molecular entities (NMEs) and biologicals;
2. Increase the automation of the manufacturing processes, in line with the Industry 4.0 paradigm (Kagermann *et al.*, 2011) that has been recently gained attention in the manufacturing sector;
3. Exploit *big data* methodologies to extract information on the patient profiles, and therefore design new drugs that are perfectly suited to the patient needs;
4. Optimize the supply chain and consider strategic clinical, regulatory and scientific partnerships.

In view of the above, it is clear that the introduction of innovative solutions and technologies to assist all the different stages of pharmaceutical production (i.e. drug discovery, product and process development, product manufacturing) will represent a competitive advantage for companies in an era of patent expiration and strict governments' budget constraints.

1.1.2 Pharmaceutical R&D: past, present & future

The discovery, test and approval of new drugs represents the milestone upon which pharmaceutical companies can diversify and strengthen their assets. Research and development

(R&D) is the activity that is deemed to support the discovery, launch and approval of NMEs and new biological entities (commonly referred to as biologicals)³ in the market.

The process that leads to the launch of a new drug on the market is typically very long and is comprised by several steps, as schematically shown in Fig. 1.1.

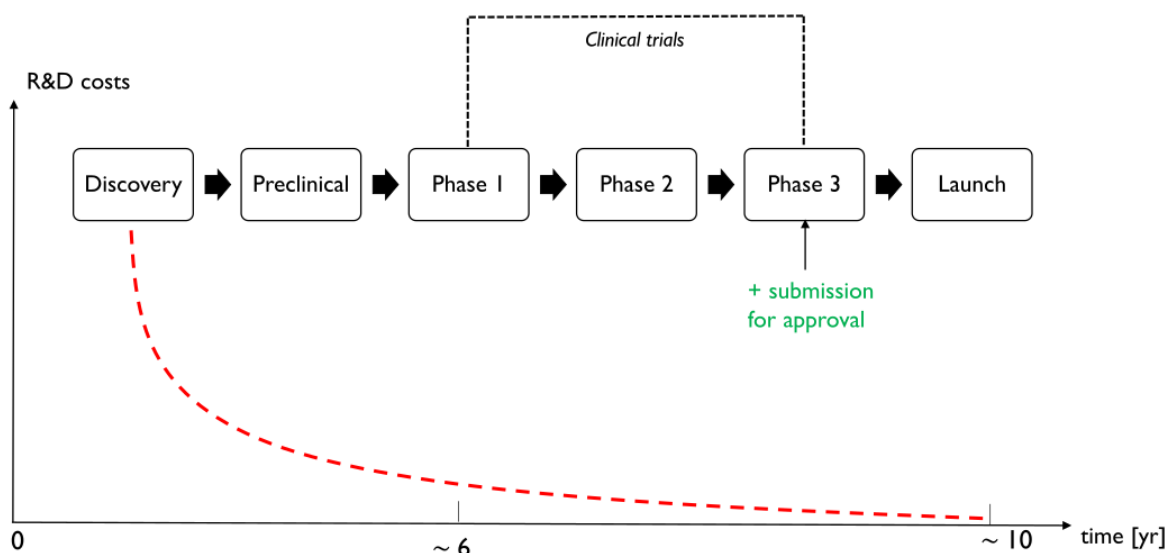


Figure 1.1 Pre-launch stages for a new pharmaceutical drug and qualitative profile of R&D costs during the pre-launch stage.

Research deals with the first two steps of the process, which are:

- 1) Drug discovery. Chemical or biological compounds that exhibit the potential to treat new or existing conditions are initially screened. A promising compound is, on average, found every 5,000-10,000 candidates (IFPMA, 2017). This step may take up to 6 years and the chance of success are, as explained above, less than 0.01% (EFPIA, 2017; IFPMA, 2017).
- 2) Pre-clinical trials. At this stage, the safety and efficacy of the compound is tested *before* testing its performance on humans. Different types of pre-clinical tests exist: these include, but are not limited to, pharmacodynamics and pharmacokinetic tests, toxicity tests, ADME (adsorption, distribution, metabolism and excretion) tests, and are typically performed both *in silico* and *in vivo*. It is estimated that a top-20 pharmaceutical company has, on average, around 250 compounds under pre-clinical tests per year (IFPMA, 2017; EvaluatePharma, 2018) and the completion of this step may take up to 2 years.

³ New molecular entities are drugs that contain an active moiety that has not been previously approved by the regulatory agency. They are typically new *chemical* entities. Biologicals represent biologic compounds or vaccines that have not been previously approved by the regulatory agency.

The budget required to perform the two research tasks described above is typically very large, and the return on investment is very low, due to the low chance of success. It is estimated (EFPIA, 2017, Deloitte, 2017) that around 21% of the overall budget for the launch of a new drug is used in these two stages (Fig. 1.2).

Development refers to the activities that are performed after the new drug successfully passes the pre-clinical trials and is characterized by a series of clinical trials (i.e. on humans). These clinical trials are divided into three different phases, which need to subsequently obtain positive results (i.e., the drug needs to pass Phase I in order to enter Phase II, and so on). The phases are:

- 3) Phase I. This phase is intended as the first stage for the assessment of the efficacy and safety of the new drug on humans. It involves a small number of healthy volunteers (typically in the range 20-100) and the rate of success is, on average, around 65% (IFPMA, 2017). The main target of this step is to assess a suitable dose (or range of doses) for the new drug to be effective.
- 4) Phase II. This phase is intended to assess if the drug has any biological activities or side effects on humans. The number of volunteers involved is larger (100-500) and the rate of success for this phase drastically decreases to around 40% (IFPMA, 2017). Most of the new drugs fail during this phase, where it is typically discovered that the drug does not work as planned or may exhibit toxicity effects (Deloitte, 2017).
- 5) Phase III. This phase involves a much larger number of volunteers (1000-5000) and is aimed at assessing the real effectiveness of the drug with respect to the current treatments. Due to the large number of people involved, this phase represents a costly and time-consuming step, but the rate of success are typically higher than for phase II (an average success rate is around 50%).

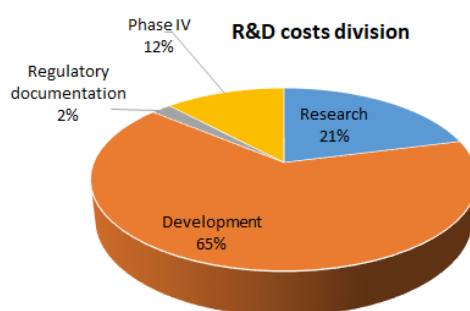


Figure 1.2 R&D costs division in the pharmaceutical industry (IFPMA, 2017).

The overall *development* process requires 6 to 8 years (IFPMA, 2017, Deloitte, 2017, EvaluatePharma, 2018) and up to 65% of the overall budget for the launch of the new drug (Fig. 1.2). Therefore, the overall *research and development* process is deemed to absorb more than 85% of the total budget for the development and launch of a new molecular entity.

If the new drug successfully passes all the phases of the clinical trials, the following step is the submission for approval to the regulatory agency (green text in Fig. 1.1).

This step may take up to 2 years and absorbs around 2% of the total budget. If successful, the drug is ready for its launch on the market. However, periodic clinical trials are required after the launch in order to review the safety and efficacy of the drug. This step is typically referred to as Phase IV clinical trials, and may absorb up to 12% of the total budget (Deloitte, 2017). A qualitative trend of the R&D costs (thus excluding the approval and Phase IV costs) during the development of a new drug is shown in Fig. 1.1 (red dashed line).

From the above considerations, it can be concluded that R&D costs represent a significant portion of the overall costs that pharmaceutical companies are forced to face (~ 20%; Deloitte, 2017) and are characterized by low chances of success and continuously decreasing returns on investments (Deloitte, 2017; IFPMA, 2017; McKinsey, 2017; EvaluatePharma, 2018). Deloitte (2017) estimates that the overall cost for the development of a new drug (from discovery to launch) has reached, on average, the all-time peak of 165 billion dollars in 2017, and is expected to increase during the next years. EvaluatePharma® (EvaluatePharma, 2018) estimates an average growth rate for the worldwide R&D costs in the pharmaceutical industry of 3.2% for the period 2018-2024. The past and the expected trend for the worldwide R&D spend is reported in Table 1.3.

Table 1.3. *Worldwide pharma R&D spend and growth rate with respect to the previous year in the period 2010-2024. Adapted from EvaluatePharma, 2018.*

Year	R&D spend [\$bn]	R&D spend growth with respect to the previous year [%]
2010	129	(-)
2011	137	+6.2
2012	136	-0.4
2013	138	+1.7
2014	144	+4.4
2015	149	+3.5
2016	159	+6.3
2017	165	+3.9
2018	172	+4.1
2019	177 ^(*)	+3.1
2020	183 ^(*)	+3.0
2021	188 ^(*)	+3.1
2022	194 ^(*)	+3.0
2023	199 ^(*)	+2.5
2024	204 ^(*)	+2.6

(*) = forecast

If, from one side, the global R&D spend is expected to increase nonstop during the next years, the expectations on the number of NMEs and biologicals that can potentially be approved by the regulatory agencies seem promising. In fact, after a very poor performance in 2016, the

number of approved NMEs and biologicals has significantly increased in 2017, and historically high levels are expected to be reached in the next six years. Table 1.4 shows the number of NMEs and biologicals that have been approved by the Food and Drug Administration (FDA) in the period 2002-2017. The number of drugs that are expected to be launched in the market for the next six years is forecast to be, on average, 45 per year (QuantileIMS, 2016; EvaluatePharma, 2018) with a peak of 63 new approved drugs forecast for 2024. It is worth noticing that these forecasts are subject to a high degree of uncertainty, since regulatory requirements are expected to become stricter year after year, and the pre-launch stage is characterized by unpredictable events that may lead to a regulatory rejection.

Table 1.4. *FDA new drug approvals (NMEs + biologicals). Adapted from EvaluatePharma, 2018.*

Year	# of NMEs approved	# of biologicals approved	Total NMEs+biologicals approved
2002	17	9	26
2003	21	14	35
2004	31	7	38
2005	18	10	28
2006	18	11	29
2007	16	9	25
2008	21	10	31
2009	20	15	35
2010	15	11	26
2011	24	11	35
2012	34	10	44
2013	25	10	35
2014	30	21	51
2015	33	23	56
2016	15	12	27
2017	36	19	55

In summary, pharmaceutical R&D is inherently a high-risk, high-reward endeavor. However, the steep increase of the costs for drug development that is expected to continue during the next years is building up a remarkable pressure on R&D productivity. Pharmaceutical companies are forced to deal with an increasing competition on the global market, and new technologies that can speed-up the R&D process are becoming essential tools to cope with this problem. Although it is impossible to find a common agreement between the different analysts on the tools that should be used to increase R&D productivity, some recognized solutions (Deloitte, 2017; McKinsey, 2018; PwC, 2018) that can be proposed are:

- the use of artificial intelligence (AI) techniques to assist and speed-up the drug discovery stage. AI algorithms can be used to analyze large amount of data from different sources, such as pre-clinical or clinical-trials of already existing compounds, in order to get new insights that may help researchers to identify new patients' needs and therefore new drug discoveries. Moreover, AI algorithms can be used to identify

potential candidates for clinical trials through targeted advertising (thanks to huge amount of data available, for example, from social media) and to assist a better understanding of the patients' needs.

- Use of real world evidence (RWE)⁴ to obtain a better understanding of rare diseases, to identify new potential diseases, to assist clinical trials or to expedite the enrolment of volunteers.
- Use of *in silico* tools (i.e., model-based tools) to assist: *i*) the product development stage; *ii*) the process (or “recipe”) development stage, i.e. the way the new drug is designed to be manufactured. In more detail, modeling tools are expected to predict clinical outcomes, inform clinical trials design and support evidence of safety and efficacy of the new drug during the product development stage. On the other hands, *in silico* simulations are expected to assist the conceptual development, the design and the optimization of the manufacturing process for the production of the new drug.
- Use of robotic and cognitive automation to automatize certain aspects of the R&D process, such as the design and monitoring of the clinical trials, or organize and speed-up the large amount of documentation that needs to be filled in during the drug development process and in the submission for approval by the regulatory agency.

Most of the techniques described above have long been known and used within the process systems engineering (PSE) community, which therefore can give a key contribution to boost pharma R&D productivity.

1.1.3. Pharmaceutical manufacturing: past, present & future

As described in the previous section, pharmaceutical R&D (which not only includes product development, but also process development) is worldwide recognized as a time-consuming, highly specialized cutting-edge activity, due to the strict requirements that new products need to fulfill before being launched to the market. However, on the other side, pharmaceutical manufacturing has been traditionally characterized by strict and experience-based procedures, that led to the adoption of several process malpractices and high percentages of rejected products during the production campaigns. The main reason for this lack of innovation is to be found in the strict regulatory environment within which companies were forced to operate. The regulatory agencies, in fact, by addressing all the focus towards the achievement of the desired drug safety and efficacy, were unintentionally stimulating companies to invest their money on obtaining new patents, rather than on innovating and revamping their manufacturing processes. The final result was that, in terms of manufacturing procedures and facilities, the

⁴ RWE is the evidence obtained from observational data that are not obtained during clinical trials and that are typically stored in electronic health records, billing activities databases, registries, mobile devices and so on.

pharmaceutical industry was (and still is) lagging behind several other industries, despite its cutting-edge R&D activity and its role for the society.

A turning point in the relationship between pharmaceutical companies and regulatory agencies was first introduced in 2004 (FDA, 2004a; FDA, 2004b) by the American FDA. FDA was the first agency to acknowledge that the rigid regulatory environment was the main restraint for companies to not invest on innovation and on the adoption of cutting-edge technologies in product manufacturing. Several initiatives (FDA, 2004a; FDA, 2004b; FDA, 2004c; FDA, 2006) culminated in the introduction of the Quality by Design (QbD) framework (as opposed to Quality by Testing), whose objective is to favor a flexible and efficient environment within which high quality products can be manufactured without extensive regulatory oversight. QbD promotes the adoption of systematic and science-based (as opposed to experience-based) tools to assist: *i*) product and process design (i.e., product and process development); *ii*) technology transfer; *iii*) product manufacturing. A detailed description of all these aspects is given in § 1.2. Regarding product manufacturing, the introduction of QbD represented (and is still representing) the main driving force for the modernization of the practices and technologies that are adopted in pharmaceutical manufacturing. In terms of product manufacturing, the early stage of QbD (2009-2011) was mainly focused on the implementation of some of the new ideas introduced within this new framework (adoption of process analytical technologies (PAT); exploitation of real-time release testing and in/on-line quality testing to pre-existing and well-established manufacturing processes and procedures (Reklaitis, 2017). These processes are traditionally operated batchwise and involve multiple units and stages that are typically followed in sequence according to strict procedures. In the following years (2012-2015), as a result of the adoption of new paradigms to assist the product and process development stage (culminating in the definition of the design space of the product to be manufactured), particular effort was put towards the development of efficient control strategies and monitoring tools to guarantee process operation within the boundary of the design space. The focus was, again, mainly towards traditional technologies operated batchwise, and still represents the most important target that pharmaceutical companies are addressing to improve their process operation⁵ (Reklaitis et al., 2017).

Continuous manufacturing (CM) of pharmaceuticals has gained attention from the scientific community, the industrial practitioners and the regulators (Lee et al., 2015; Fisher et al., 2016; Kopcha, 2017; Nasr, 2017; Yu and Kopcha, 2017). It is worth noticing that the possibility of adopting continuous operation for some of the units involved during drug manufacturing was first proposed in the 90s (Paul, 1990) and kept being proposed throughout the years (Dienenmann, 2000). Potentially, CM could enable a much faster and more efficient (in terms

⁵ It is estimated that more than 95% of the unit operations involved in the pharmaceutical industry are operate batchwise or semi-batchwise, and, despite the increasing attention towards continuous manufacturing, this percentage is expected to remain stable for the next 5-8 years (Troup and Georgakis, 2013; Reklaitis *et al.*, 2017).

of reduction of products that need to be discarded due to off-specifications) operation with respect to batch manufacturing, as a result of the tighter handling of product specifications that can be obtained with this technology. Moreover, it would enable a substantial reduction of the energy needs of the manufacturing line and would involve a sensible reduction of the risk of human error due to its high degree of automation. However, some challenges still need to be faced for a practical implementation of CM, e.g., high equipment and installation costs, the costs of developing and implementing the advanced process control strategies that are needed to guarantee product specifications, and the strategies to be adopted in the presence of a product recall. These drawbacks, together with the strict regulatory environment within which companies were forced to operate, represented the main reasons that kept CM as a potential, but not applicable alternative to the traditional batch approach in the early 2000s. The recent launch of the QbD initiative has revived the interest towards this technology, even though the lack of a regulatory framework on this topic is still limiting a systematic approach to its implementation. To have an idea of the gap that currently exists between the scientific literature (where the number of publications on CM has exponentially grown in the past 5 years) and the practical implementation of CM (Collins, 2018), it is worth considering that the only marketed products that have been obtained through CM are the two reported by the FDA (Kopcha, 2017) and a new drug that has been recently approved (Collins, 2018).

Interestingly, the opinion of industrial practitioners on the real benefits of CM is not always as optimistic as the one of many regulators and academics (Reklaitis, 2017). For example, Collins (2018) highlights the fact that the costs for implementing CM are comparable with the costs required to upgrade (i.e., to adopt the technologies encouraged within the QbD framework) the existing batch operations. However, while CM is characterized by a total lack of experience and a scarce (or inexistent) regulatory framework, batch manufacturing is well-established and is characterized by a wide and detailed regulation. Therefore, it is worth questioning the real advantage of investing in CM with respect to investing in the upgrade of the existing batch operations.

Despite the different views on the benefits that CM may bring to the pharmaceutical industry, it is commonly accepted that batch manufacturing will still be the preferred route towards drug production, while CM will keep gaining increasing attention especially for the production of the new drugs that will be launched in the market (Reklaitis, 2017; Collins, 2018). Independently of the type of operation adopted, the next years are expected to bring an increased level of digitalization within the manufacturing processes, in line with the Industry 4.0 revolution (Kagermann *et al.*, 2011) that has recently involved other manufacturing sectors. The digitalization of drug manufacturing is expected to involve an extensive use of process simulators, advanced process control and monitoring techniques, supply chain automation and *big data* technologies to extract information from the huge amount of data that are collected throughout the product lifecycle (Herwig *et al.*, 2017; Ding, 2018; Rauschnabel, 2018). A

futuristic overview of the challenges that must be faced towards the implementation of the Industry 4.0 paradigm in pharma is given by Herwig *et al.* (2017). As for R&D (i.e., product and process development), process systems engineering tools are expected to give a key contribution to the solution of these challenges (Reklaitis, 2017).

1.1.4 Regulatory aspects: past, present & future

As briefly mentioned in the previous section, the role of the regulatory agencies is of paramount importance for protecting and improving human health, as well as for influencing the decisions and strategies that pharmaceutical companies adopt to keep their competitiveness in the global market. These agencies include the already cited FDA, the European Medicines Agency (EMA) and many other regulatory bodies that are spread all over the world (a full list can be found in Global Regulatory Authority Websites, 2018). Some of the activities carried out by these agencies are (Collins, 2018): *i*) the review of information for new drug submissions from pharmaceutical companies; *ii*) the review of updated process information throughout the lifecycle of medicines; *iii*) the regular inspection of manufacturing sites around the world and *iv*) the continuous assurance of compliance of quality. Some of these agencies share alignment on the most important regulations and quality expectations for the medicines that are introduced in the market through the *International Council on Harmonization (ICH)*⁶, while others have specific expectations that are not always shared by the other regulatory bodies. Despite this complex regulatory framework, several common ideas can be found between the different regulatory entities, and the discussion proposed in this Dissertation will mainly focus on the ideas and guidelines enforced within the context of ICH.

The need to establish common regulations for good practices in pharmaceutical development and manufacturing was largely driven by toxicological safety-related disasters such as the elixir of sulfanilamide incident in 1932 (Immel, 2001) or the sulfathiazole disaster in 1941 (Swann, 1999). In the former case, diethylene glycol was used within the elixir, without testing its toxicity (which is now well-known) to humans. The latter incident, instead, refers to the commercialization by the Winthrop Chemical Company of New York of sulfathiazole tablets contaminated with phenobarbital, which caused hundreds of deaths and thousands of injuries, and that was proven to be due to a manufacturing control failure and serious firm's irregularities. These two disasters pushed FDA to conceive (and finalize in 1987) the current Good Manufacturing Practices (cGMP), that were intended to guarantee the safety and efficacy of drugs through the use of methods, facilities and controls for the manufacturing, processing, packing or holding of these products (Immel, 2001; Woodcock, 2013; Collins, 2018).

Based on the context in which these cGMPs were first conceived (i.e. as a reaction to the incidents described above), it is easy to understand why the main focus of these documents was

⁶ ICH is composed by the regulatory authorities of Europe, U.S. and Japan, together with industrial experts.

addressed towards ensuring the toxicological safety of the drug, rather than promoting the adoption of innovative methodologies and technologies for drug development and manufacturing. Indeed, interpretations of these cGMPs varied amongst companies (Collins, 2018), and the two main strategies adopted by the regulatory bodies to ensure compliance with the cGMPs were: *i*) the manufacturing site inspection approach (Junod, 2004), and *ii*) the preparation of guidance documents, which contain nonbinding recommendations and invite the companies to engage in a proactive discussion with the regulatory agency. These guidance documents were not intended to substitute the cGMPs (whose adherence is mandatory), but rather as a series of suggestions that, if followed, would have helped the company to comply with them.

The guidance documents' approach immediately received an incredibly positive feedback from companies. Since these documents contain the current regulatory agency's thinking on the topic, many manufacturers started to interpret their recommendations as compulsory (Collins, 2018), and they are nowadays recognized as the fastest and most efficient way to understand the position of the regulatory agencies on a given topic. Given the success of this formula, and based on the increased awareness of the regulatory agencies to foster innovation within the pharmaceutical context (§ 1.1.3), a series of important initiatives was put forth by the American FDA in the early 2000s. First, in 2004, the "*Pharmaceutical cGMPs for the 21st century – A risk-based approach*" was introduced (FDA, 2004b). The targets of this new set of cGMPs were, among others, to foster the adoption of new technological advances in pharmaceutical manufacturing, introduce quality system and risk management approaches for product development and manufacturing and, in a nutshell, promote the adoption of science-based approaches and state-of-the-art technologies within the existing pharmaceutical framework. Many other initiatives followed this original document (FDA, 2004c; FDA, 2006) and the formalization of this new approach for pharmaceutical development and manufacturing was finalized through a series of ICH guidelines (ICH 2006, 2009a, 2009b, 2011, 2012) that represent the foundation of the broad framework which is now known as Quality by Design (QbD). The introduction of the QbD initiative aimed at replacing the traditional procedural and inspection-based approach to assure compliance (i.e., product quality) with a new scientific and risk-based approach, characterized by a greater understanding of the product and process considered, as well as the sources of variability that can potentially affect product quality. The fundamental idea behind QbD is that product quality should be "built" into the product since its conception and design, rather than verified after its manufacturing. This requires a mechanistic understanding of the relationship between product quality and all the sources of variability that may have an impact on it.

The launch of QbD represented a turning point for the pharmaceutical community, and the practical implementation of this paradigm still remains the core of the current pharmaceutical innovation (that not only involves the R&D process, but also the manufacturing stage, as

explained in § 1.1.3). From their side, the regulatory agencies regularly release documents (which are, mainly, guidelines or drafts of guidelines, in view of the considerations described above) aimed at describing their current thinking on specific topics. A brief record of the main ICH guidelines that have been proposed throughout the years and that contain the main ideas of the QbD paradigm are reported in Table 1.5. As can be noticed, the release of these guidelines is still an ongoing process (two final guidelines and one draft have been released in 2018), and represent a reference for the practical implementation of the QbD paradigm. Moreover, in order to centralize the regulatory activities in terms of review, policy, research and science activities, the FDA has recently created the Office for Pharmaceutical Quality (OPQ) (Yu and Woodstock, 2015), whose aim is to uniform and ease any communications between the regulatory bodies and the manufacturers. A detailed description of the activities carried out by this “central” regulatory body is reported by Yu and Woodstock (2015).

1.1.4 Regulatory aspects: past, present & future

As briefly mentioned in the previous section, the role of the regulatory agencies is of paramount importance for protecting and improving human health, as well as for influencing the decisions and strategies that pharmaceutical companies adopt to keep their competitiveness in the global market. These agencies include the already cited FDA, the European Medicines Agency (EMA) and many other regulatory bodies that are spread all over the world (a full list can be found in Global Regulatory Authority Websites, 2018). Some of the activities carried out by these agencies are (Collins, 2018): *i*) the review of information for new drug submissions from pharmaceutical companies; *ii*) the review of updated process information throughout the lifecycle of medicines; *iii*) the regular inspection of manufacturing sites around the world and *iv*) the continuous assurance of compliance of quality. Some of these agencies share alignment on the most important regulations and quality expectations for the medicines that are introduced in the market through the *International Council on Harmonization (ICH)*⁷, while others have specific expectations that are not always shared by the other regulatory bodies. Despite this complex regulatory framework, several common ideas can be found between the different regulatory entities, and the discussion proposed in this Dissertation will mainly focus on the ideas and guidelines enforced within the context of ICH.

The need to establish common regulations for good practices in pharmaceutical development and manufacturing was largely driven by toxicological safety-related disasters such as the elixir of sulfanilamide incident in 1932 (Immel, 2001) or the sulfathiazole disaster in 1941 (Swann, 1999). In the former case, diethylene glycol was used within the elixir, without testing its toxicity (that is now well-known) to humans.

⁷ ICH is composed by the regulatory authorities of Europe, U.S. and Japan, together with industrial experts.

Table 1.5. Main ICH guidelines describing the key ideas of the QbD paradigm (adapted and updated from Tomba, 2013).

Date	Document	Reference	Title	Type
06/01/2006	ICH Q9	ICH (2006)	Q9- Quality risk management (QRM) system.	Final guidance
04/07/2009	ICH Q10	ICH(2009a)	Q10 – Pharmaceutical quality system	Final guidance
11/20/2009	ICH Q8(R2)	ICH(2009b)	Q8(R2)- Pharmaceutical development	Final guidance
11/01/2011	ICH Q&A	ICH (2011)	Q8, Q9, Q10 Questions and Answers	Final guidance
07/25/2012	ICH Q&A from training sessions	ICH (2012a)	Q8, Q9, Q10 points to consider	Final guidance
11/09/2012	ICH Q11	ICH (2012b)	Q11- Development and manufacture of drug substances	Final guidance
09/30/2016	ICH Q7	ICH (2016)	Q7- Good manufacturing practice guidance for active pharmaceutical ingredients. Guidance for industry	Final guidance
02/23/2018	ICH Q11 Q&A	ICH (2018a)	Q11- Development of manufacture of drug substances – Questions and answers	Final guidance
04/19/2018	ICH Q7 Q&A	ICH(2018b)	Q7- Good manufacturing practice guidance for active pharmaceutical industry. Guidance for industry- Questions and answers	Final guidance
05/30/2018	ICH Q12	ICH(2018c)	Q12- Technical and regulatory considerations for pharmaceutical product lifecycle management. Core guideline for industry	Draft guidance

The latter incident, instead, refers to the commercialization by the Winthrop Chemical Company of New York of sulfathiazole tablets contaminated with phenobarbital, which caused hundreds of deaths and thousands of injuries, and that was proven to be due to a manufacturing control failure and serious firm's irregularities. These two disasters pushed FDA to conceive, and after few years (specifically, in 1987) finalize the current Good Manufacturing Practices (cGMP), that were intended to guarantee the safety and efficacy of drugs through the use of methods, facilities and controls for the manufacturing, processing, packing or holding of these products (Immel, 2001; Woodcock, 2013; Collins, 2018). Based on the context in which these

cGMPs were first conceived (i.e. as a reaction to the incidents described above), it is easy to understand why the main focus of these documents was addressed towards ensuring the toxicological safety of the drug, rather than promoting the adoption of innovative methodologies and technologies for drug development and manufacturing. Indeed, interpretations of these cGMPs varied amongst companies (Collins, 2018), and the two main strategies adopted by the regulatory bodies to ensure compliance with the cGMPs were: *i*) the manufacturing site inspection approach (Junod, 2004) and *ii*) the preparation of guidance documents, which contain nonbinding recommendations and invite the companies to engage in a proactive discussion with the regulatory agency. These guidance documents were not intended to substitute the cGMPs (whose adherence is mandatory), but rather as a series of suggestions that, if followed, would have helped the company to comply with them.

The guidance documents' approach immediately received an incredibly positive feedback from companies. Since these documents contain the current regulatory agency's thinking on the topic, many manufacturers started to interpret their recommendations as compulsory (Collins, 2018), and they are nowadays recognized as the fastest and most efficient way to understand the position of the regulatory agencies on a given topic. Given the success of this formula, and based on the increased awareness of the regulatory agencies to foster innovation within the pharmaceutical context (§ 1.1.3), a series of important initiatives was put forth by the American FDA in the early 2000s. First, in 2004, the "*Pharmaceutical cGMPs for the 21st century – A risk-based approach*" was introduced (FDA, 2004b). The targets of this new set of cGMPs were, among others, to foster the adoption of new technological advances in pharmaceutical manufacturing, introduce quality system and risk management approaches for product development and manufacturing and, in a nutshell, promote the adoption of science-based approaches and state-of-the-art technologies within the existing pharmaceutical framework. Many other initiatives followed this original document (FDA, 2004c; FDA, 2006) and the formalization of this new approach for pharmaceutical development and manufacturing was finalized through a series of ICH guidelines (ICH 2006, 2009a, 2009b, 2011, 2012) which represent the foundation of the broad framework which is now known as *Quality by Design* (QbD). The introduction of the QbD initiative aimed at replacing the traditional procedural and inspection-based approach to assure compliance (i.e. product quality) with a new scientific and risk-based approach, characterized by a greater understanding of the product and process considered, as well as the sources of variability that can potentially affect product quality. The fundamental idea behind QbD is that product quality should be "built" into the product since its conception and design, rather than verified after its manufacturing. This requires a mechanistic understanding of the relationship between product quality and all the sources of variability that may have an impact on it.

The launch of QbD represented a turning point for the pharmaceutical community, and the practical implementation of this paradigm still remains the core of the current pharmaceutical

innovation (that not only involves the R&D process, but also the manufacturing stage, as explained in § 1.1.3). From their side, the regulatory agencies regularly release documents (which are, mainly, guidelines or drafts of guidelines, in view of the considerations described above) aimed at describing their current thinking on specific topics. A brief record of the main ICH guidelines that have been proposed throughout the years and that contain the main ideas of the QbD paradigm are reported in Table 1.5. As can be noticed, the release of these guidelines is still an ongoing process (two final guidelines and one draft have been released in 2018) and represent a reference for the practical implementation of the QbD paradigm. Moreover, in order to centralize the regulatory activities in terms of review, policy, research and science activities, FDA has recently created few years ago the Office for Pharmaceutical Quality (OPQ) (Yu and Woodstock, 2015), whose aim is to uniform and ease any communications between the regulatory bodies and the manufacturers. A detailed description of the activities carried out by this “central” regulatory body is reported in Yu and Woodstock, 2015.

1.2 Quality by design paradigms: regulatory overview

The concepts embedded within the Quality by Design initiative can be understood through a rigorous interpretation of the guidelines that have been issued, and are still being issued, by the regulatory agencies. From a general perspective, the interpretation of QbD proposed by the regulators can be schematically summarized as a sequence of 5 activities (Yu *et al.*, 2014; Yu *et al.*, 2015):

- 1) identification of a quality target product profile (QTPP), from which the critical quality attributes (CQAs) of the drug under development can be assessed;
- 2) *product* design and understanding, including the identification of the raw material properties that are critical for product quality (often defined as critical material attributes; CMA);
- 3) *process* design and understanding, including the identification of critical process parameters (CPPs) and the derivation of a mechanistic relationship between the raw material properties and CPPs with the product CQAs;
- 4) definition of a control strategy, in terms of identifying the specifications for the API, excipient and the final drug, and in terms of developing a strategy to control each step of the manufacturing process;
- 5) exploitation of process capability and adoption of a continual process improvement strategy.

Table 1.5. Brief record of the main ICH guidelines that describe the key ideas of the QbD paradigm. Adapted and updated from Tomba, 2013.

Date	Document	Reference	Title	Type
06/01/2006	ICH Q9	ICH (2006)	Q9- Quality risk management (QRM) system.	Final guidance
04/07/2009	ICH Q10	ICH(2009a)	Q10 – Pharmaceutical quality system	Final guidance
11/20/2009	ICH Q8(R2)	ICH(2009b)	Q8(R2)- Pharmaceutical development	Final guidance
11/01/2011	ICH Q&A	ICH (2011)	Q8, Q9, Q10 Questions and Answers	Final guidance
07/25/2012	ICH Q&A from training sessions	ICH (2012a)	Q8, Q9, Q10 points to consider	Final guidance
11/09/2012	ICH Q11	ICH (2012b)	Q11- Development and manufacture of drug substances	Final guidance
09/30/2016	ICH Q7	ICH (2016)	Q7- Good manufacturing practice guidance for active pharmaceutical ingredients. Guidance for industry	Final guidance
02/23/2018	ICH Q11 Q&A	ICH (2018a)	Q11- Development of manufacture of drug substances – Questions and answers	Final guidance
04/19/2018	ICH Q7 Q&A	ICH(2018b)	Q7- Good manufacturing practice guidance for active pharmaceutical industry. Guidance for industry- Questions and answers	Final guidance
05/30/2018	ICH Q12	ICH(2018c)	Q12- Technical and regulatory considerations for pharmaceutical product lifecycle management. Core guideline for industry	Draft guidance

Based on the definition proposed in § 1.1.2 and § 1.1.3, steps 1-4 are typical R&D activities, while steps 5 is related to the manufacturing stage⁸. This concept has been eventually formalized in the ICH-Q10 guideline (ICH, 2009a), where the activities that characterize the lifecycle of a pharmaceutical product are classified as follows:

- a) *Pharmaceutical development*, that collects steps 1-3 of the previous classification;

⁸ Note that the definition of the control strategy, even if practically tested and implemented during the manufacturing stage, may be conceived during the process development stage (i.e., step 3).

- b) *Technology transfer*, whose goal is to transfer product and process knowledge from development to manufacturing (*scale-up*) or from different manufacturing sites (*scale-out*);
- c) *Commercial manufacturing* (or product manufacturing), which collects steps 4 and 5 of the previous classification.

Moreover, a new step (called *product discontinuation*) has been added to the previous classification. The goal of this step is to develop strategies for an effective management of the terminal stage of the product lifecycle (e.g. product recall from the market). The key activities and the key terminology that is adopted in the three steps (a,b,c) described above is briefly reviewed in the following sections.

1.2.1 Pharmaceutical development

The definition of pharmaceutical development and the steps that should be implemented during this activity according the QbD framework are thoroughly described in the regulatory guideline ICH Q8(R2) (ICH, 2009b). ICH defines pharmaceutical development as the activity “whose aim is to design a quality product and its manufacturing process to consistently deliver the intended performance of the product” (ICH, 2009b). The elements that, according to the QbD framework, pharmaceutical development should include are schematically shown in Fig. 1.3.

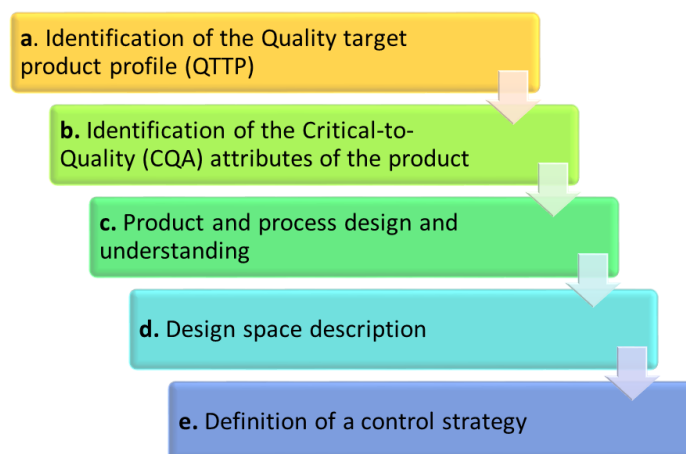


Figure 1.3 Elements that should be considered in pharmaceutical development according to ICH (ICH, 2009b).

Steps (a), (b), (c), (e) are considered minimum requirements, while step (d) represents the core step to implement an enhanced QbD approach.

1.2.1.1 Step (a): identification of the Quality Target Product Profile (QTPP)

Before going through the different steps of pharmaceutical development, the definition of quality for a pharmaceutical drug must be given. According to the ICH Q6A guideline, quality is defined as the suitability of either a drug substance or drug product for its intended use (ICH, 1999). A revised definition of quality is proposed in the ICH Q9 guideline, where quality is defined as the degree to which a set of inherent properties of a product fulfills requirements (ICH, 2006). A prospective summary of the quality characteristics of the product that ideally should be achieved to ensure the desired safety and efficacy is defined as the Quality Target Product Profile (QTPP). The identification of the QTPP is a key step of pharmaceutical development and directly influences the success of all the remaining activities. Raw and coworkers (Raw *et al.*, 2011) analyzed the impact that a wrong identification of the QTPP can have on the subsequent development stages, emphasizing that this step is often underestimated by companies. Following the regulatory guidelines (ICH, 2009b), the definition of the QTPP should include the following aspects:

- intended use in clinical setting, route of administration, dosage form, delivery systems;
- dosage strength(s);
- container closure system;
- therapeutic moiety release or delivery and attributes affecting pharmacokinetic characteristics (e.g. dissolution, aerodynamic performance) appropriate to the drug product dosage form being developed;
- drug product quality criteria (e.g. sterility, purity, stability and drug release) appropriate for the intended marketed product.

Once identified, the QTPP forms the basis for the identification of the critical quality attributes (CQAs) of the product.

1.2.1.2 Step (b): identification of the product Critical Quality Attributes (CQAs)

A critical quality attribute (CQA) is defined by ICH (ICH, 2009b) as a physical, chemical, biological or microbiological property or characteristic of an output material, including a finished drug product, that should be within an appropriate limit, range, distribution to ensure the desired product quality. According to this definition, the CQAs may not only be related to the drug product, but also to the intermediates (i.e., in-process materials). CQAs should be derived from the QTPP identified at the previous step, and criticality should be primarily based upon the severity of harm to the patient should the product fall outside the acceptable range for that attribute (Yu *et al.*, 2014). Examples of CQAs could include all those aspects affecting product purity, drug release and stability, or specific properties such as aerodynamic properties for inhaled products, sterility for parenterals and so on. The CQAs of intermediates could include specific properties (e.g. particle size distribution, bulk density etc...) that may potentially affect product CQAs. It is worth noticing that, in several situations, the CQAs of the

intermediates are improperly referred to as “product” CQAs. For the sake of simplicity, in this Dissertation, this terminology will be adopted.

1.2.1.3 Step (c): product and process design and understanding

Steps (c) and (d) of the methodology presented in (§ 1.2) can be classified according to two different and sequential activities: *i) product design and understanding* and *ii) process design and understanding*.

Product design is concerned with the assessment of the ability of the product to meet the desired QTPP and maintain the desired QTPP over the product shelf life. *Product understanding* is related to the identification of the raw material properties that can potentially affect the product CQAs, and on the ability to link these material properties to the product CQAs (ICH, 2012b; Yu *et al.*, 2014). *Process understanding* is related to the identification of the critical process parameters (CPPs) that can potentially affect the product CQAs, and on the development of a link between the CPPs and the CQAs. Finally, *process design* is concerned with the design of a manufacturing process that can consistently guarantee to produce the desired product and to maintain the desired product CQAs throughout the entire production campaign.

With respect to product design, the key elements that should be taken into account at this stage are (ICH, 2012b; Yu *et al.*, 2014):

- physical, chemical and biological characterization of the API;
- identification and selection of the excipient type and grade, and understanding of its potential variability;
- understanding of the interaction between drug and excipients.

Examples of physical properties of the API that may be considered are solubility, dissolution rate, stability and many others. Chemical properties may include, for example, chemical stability in solid state; examples of biological properties are bioavailability and membrane permeability. The final target of this step is to design, through an open-ended activity of formulation optimization, a product that meet the desired QTPP and that can guarantee to preserve this QTPP throughout its entire lifecycle.

With respect to product understanding, the identification of the raw material properties (critical material attributes; CMAs) that affect the product CQAs is the first step that must be followed. Following the regulatory parlance, a CMA is a physical, chemical, biological or microbiological property or characteristic of an *input* material that should be within an appropriate limit, range or distribution to assure product quality (ICH, 2009b). It is worth noticing the similarity of this definition with that of CQA (§ 1.2.1.2), but the difference is that, while CQAs are related to *output* materials (such as intermediates and the drug product), CMAs are related to *input* materials (such as the drug substance and the excipients). Moreover, the CQA of an intermediate may become a CMA of that same intermediate for a downstream manufacturing step (Yu *et al.*, 2014). Given that investigating all the material attributes that may potentially

affect the product CQAs is not practically feasible (due to the large number of attributes to be investigated), it is recommended (ICH, 2009b; Yu *et al.*, 2014) to follow a risk-based approach based on prior knowledge and scientific understanding to identify the CMAs and relate them to the product CQAs. This activity can include the following steps (Yu *et al.*, 2014; Yu *et al.*, 2017):

- qualitative determination of all possible known input material attributes that may affect the product CQAs;
- identification of potentially high-risk attributes through risk assessment methodologies and scientific knowledge;
- design and execution of experiments (eventually designed with *design of experiments*; DoE) based on predefined levels or ranges for the attributes identified at the previous step;
- analysis of the experimental data with data-driven or (whenever possible) mechanistic modeling to assess the criticality of each attribute.

With respect to the process design activity, the process that should be used to manufacture the drug product must be conceived at this stage. The target of process design is to identify all the unit operations and manufacturing stages (along with their operating mode, i.e. batch or continuous) that should be adopted for the manufacturing of the drug product. The design of the manufacturing process should not only guarantee to obtain the desired product characteristics, but also that these characteristics can be achieved throughout the entire production campaign. The regulatory guidelines (ICH, 2009b) explicitly require a detailed description of the unit operations and equipment involved in the manufacturing process when applying for a new drug approval, emphasizing that every choice should be justified *scientifically*, and not based only on experience (e.g., from previous products). Moreover, scientific proof of the appropriateness of the equipment and unit configurations chosen should be given.

With respect to process understanding, the key elements that should be followed are intrinsically similar to those of product understanding (Yu *et al.*, 2014):

- qualitative determination of all possible known process parameters that may impact on the product CQAs;
- identification of potentially high-risk process parameters through risk assessment methodologies and scientific knowledge;
- design and execution of experiments (eventually designed) based on predefined levels or ranges for the parameters identified at the previous step;
- analysis of the experimental data with data-driven or mechanistic modeling to assess if a given process parameter is critical.

The main target of this stage is the identification of the critical process parameters (CPPs) that have an impact on the product CQAs during manufacturing and the identification of a link

between the CPPs and the CQAs. These can be obtained with a combination of experimentation and modeling (data-driven or mechanistic), as will be discussed in § 1.3. In many situations, the models that link the CMAs with the CQAs obtain during product understanding are integrated with the models that link the CPPs and CQAs obtained at this stage, thus obtaining an “integrated” model that related the CMAs and CPPs with the product CQAs. Once this “link” is obtained, the next step is the identification of the design space of the drug product considered.

1.2.1.4 Step (d): design space (DS) description

The concept of design space (DS) of a pharmaceutical process represents one of the fundamental paradigms of the QbD initiative. According to the regulatory documents (ICH, 2009b), the DS is defined as “the multidimensional combination and interaction of input variables (e.g. material attributes) and process parameters that have been demonstrated to provide assurance of quality”. Based on the definitions given in § 1.2.1.3, the definition of DS can also be reformulated as the multivariate space identified by the CMAs and CPPs that guarantee to obtain the desired product CQAs.

The establishment of the design space is to be considered as the ultimate result of the product and process understanding activities, and it represents the most important evidence of the change in the relationship between the regulatory bodies and pharmaceutical companies introduced within the QbD initiative. In fact, if the design space is approved by the regulatory agency, the process can be run anywhere within the design space without requiring any regulatory post-approval endorsement. This translates into an enhanced process flexibility, in contrast with the traditional approach where every change in the process had to be communicated to the regulatory body for assessment and approval. The DS therefore represents a “window” of operating conditions within which the company can optimize the manufacturing process, which can increase process profitability.

Five remarks are worth being discussed with respect to the regulatory definition of DS and its implications. Additional (technical) comments will be provided in § 1.3.

- Remark #1: the identification of the design space, even if strongly encouraged by the regulatory bodies, is a non-binding activity that can be performed on a volunteer basis. The regulatory guidelines, in fact, explicitly state that the design space “can be proposed by the applicant and is subject to regulatory assessment and approval” (ICH, 2009b).
- Remark #2: in the definition given by the regulatory agencies, the multivariate nature of the design space is explicitly emphasized. In mathematical terms, this means that the design space is given by the manifold spanned by the vector field of each relevant input parameter (CPPs and CMAs) and that the correlation between these inputs should be taken into account when defining the DS. Accordingly, univariate experiments (i.e., experiments performed by changing one input factor at a time) cannot constitute the basis for the definition of a DS. The definition of the DS as a combination of proven

acceptable ranges⁹ is therefore not permitted (ICH, 2009b), since it lacks of understanding of the possible interactions between the input parameters. Interestingly, it should be noted that the regulatory documents are not entirely consistent in this respect. In fact, while stating that the DS cannot be expressed in terms of PARs, they also state that “a design space can be expressed in terms of ranges of material attributes and process parameters, or through more complex mathematical relationships” (ICH, 2009b). Clearly, the first part of this statement is not consistent with the previous one or, provided that those attributes are obtained as a result of a multivariate analysis (i.e., accounting for input correlation), is to be considered at least misleading.

- **Remark #3:** the regulatory documents leave the freedom to establish independent design spaces for each single unit operation of the manufacturing line, or a single design space that spans multiple operations (ICH, 2009b). From the one side, the establishment of independent design spaces for single unit operations is typically simpler to develop; however, on the other side, the definition of a single design space for the overall manufacturing line can provide an enhanced operational flexibility (ICH, 2009b). In practical terms, this poses two technical problems: *i*) how to derive the design space of an entire manufacturing line given the independent design spaces of the single unit operations of that line (if the first option is adopted); *ii*) how to deal with the complexity and high dimensionality of a single DS for a whole manufacturing process (if the second option is adopted).
- **Remark #4:** when a design space is to be established, the scale at which the DS has been developed should be stated and the relevance of the proposed design space with respect to the plant manufacturing process should be described. In this respect, the regulatory guidelines provide qualitative recommendations without giving substantial details. The regulatory documents leave the freedom to develop the DS at any scale the manufacturer may consider relevant (ICH, 2009b). However, if the DS developed at a small scale (e.g., laboratory scale) is claimed to be applicable to the manufacturing process, the rationale behind the *scale-up* of the DS that has been adopted should be described (ICH, 2009b). In this regard, the regulatory guidelines suggest to describe the design space in terms of relevant scale-independent parameters, and to give a thorough description of the models adopted for scaling (ICH, 2009b).
- **Remark #5:** from a general perspective, it should be emphasized that the regulatory guidelines only provide general indications on how to establish the design space, but no specific technical recommendations are given. This leaves the freedom to the manufacturer to develop and propose the rationale adopted for the description of the DS. If, from one side, this represents a clear advantage for the applicants, since it

⁹ A proven acceptable range (PAR) is defined as a “characterised range of a process parameter for which operation within this range, while keeping the others parameters constant, will result in producing a material meeting relevant quality criteria”.

provides greater flexibility on the choice of the experimental/modeling strategy to be adopted, from the other side it paves the way for different interpretations of the regulatory documents, thus compromising the development of a common implementation framework. This problem has been recently highlighted by industrial practitioners (Watson *et al.*, 2018) and will be discussed in more detail in § 1.4.4..

1.2.1.5 Step (e): definition of a control strategy

The definition of control strategy given by the regulatory documents is “a planned set of controls, derived from current product and process understanding that ensures process performance and product quality” (ICH, 2009b). This definition is not the same definition that is typically adopted in common engineering understanding, where the term control is specifically related to the concept of process control. In fact, the definition of control strategy proposed by the regulatory agencies involves three different levels of controls (Yu *et al.*, 2014), namely:

- Level 1. This type of control coincides with the engineering definition of process control, i.e., it is related to the monitoring of CMAs and automatic adjustment of CPPs to consistently obtain the desired product CQAs. In other terms, Level 1 control is related to the exploitation of (feedback/feedforward/advanced) process control strategies to guarantee that desired product quality characteristics (set-points) can be obtained by appropriate manipulation of the process parameters (i.e., CPPs), even in the presence of disturbances (i.e., raw material properties fluctuations).
- Level 2. This type of control is related to the implementation of non-automatic decisions (i.e., decisions implemented manually by operators) within the boundaries of the design space, with the aim of reducing end-product testing.
- Level 3. This type of control is the one that has been traditionally implemented in the pharmaceutical industry. It relies on extensive end-product testing and tight handling of process parameters. In practical terms, this type of control is equivalent to an *a posteriori* product quality verification (quality by testing), and has nothing to do with the engineering definition of control.

To give a simple engineering interpretation of three different levels of control proposed by the regulatory bodies, a rough classification is that Level 1 control coincides with the concept of closed-loop process control (i.e., feedback, feedforward, or other advanced process control techniques), Level 2 with the concept of open-loop process control (i.e., manual control), while Level 3 is equivalent to not having process control at all. The regulatory documents encourage the use of a combination of Level 1 and Level 2 controls for a correct implementation of the QbD framework (ICH, 2009b; Yu *et al.*, 2014). It should be noted that regulatory documents are not entirely clear on the relationship existing between control strategy and design space. In fact, while in one point of the body of the ICH Q8(R2) guideline the identification of the design

space and the definition of a control strategy are presented as two sequential activities (ICH, 2009b), in other sections of the same guideline it is stated that “(...) an appropriate control strategy can, for example, include a proposal for a design space(s) and/or real-time release testing”. It is not entirely clear, therefore, if the control strategy should guarantee that the process can be operated within the boundary of the DS previously established, or if the identification of the DS can be interpreted as part of the control strategy itself.

1.2.2 Technology transfer

Technology transfer includes all the activities required to transfer the manufacturing process from the laboratory scale to the commercial scale (possibly through a pilot plant). These activities include not only the transfer of equipment and process operating conditions, but also all the subsidiary technologies tested and implemented at the laboratory scale (e.g., sensors, analyzers, control instrumentation etc...). Within these activities, the *scale-up* of the design space and control strategy discussed in § 1.2 represent key steps to guarantee that the desired product quality can be obtained also in the commercial scale plant. According to the regulatory recommendations, technology transfer is an activity that should be performed using a risk-based approach (ICH, 2006) and should pose the basis for a continual process improvement strategy. Despite simple general recommendations on how technology transfer should be implemented, it is worth noticing that the regulatory documents do not provide technical recommendations on how this step should be implemented. Therefore, manufacturers are left with the freedom to prove the effectiveness of the scale-up procedures adopted (Yu *et al.*, 2014). This has the advantage of giving great flexibility on how to face this problem, but at the same time leaves companies with several different options and interpretations (Watson *et al.*, 2018).

1.2.3 Commercial manufacturing & continual process improvement

Commercial manufacturing is the culmination of the product/process design and technology transfer activities, and is meant to obtain a mass production of the desired drug. As briefly discussed in §1.1.3, pharmaceutical processes typically involve several operating units that may be operated batchwise or continuously. Despite very few entire manufacturing lines are operated continuously, the presence of single units operated continuously has long been considered as a standard practice in pharmaceutical manufacturing (Garcia-Munoz *et al.*, 2018). The regulatory documents explicitly mention two elements that should be considered when mass production of the drug has been reached: process capability and continual process improvement (ICH, 2006; ICH, 2009b, Yu *et al.*, 2014).

Process capability is defined as “a measure of the inherent variability of a stable process that is in a state of statistical control in relation to the established acceptance criteria” (ICH, 2006; Yu *et al.*, 2014). In practical terms, the definition of process capability given by the regulators is

equivalent to the concept of process *monitoring*, to be intended as the early identification of potential sources of common-cause variation within the process. Once these common sources of variation have been identified, mitigation actions should be taken *via* the control strategy to guarantee that the process will give the desired product quality (i.e., process *control*).

The regulatory concept of *continual process improvement* refers to a set of activities that the company can carry out to enhance the ability of the process to maintain product quality (ICH, 2006; Yu *et al.*, 2014). These activities may include, but are not limited to (Yu *et al.*, 2014):

- measuring key aspects of the current process and collect relevant data;
- analyzing the data to investigate and verify cause and effect relationships;
- improving or optimizing the current process based upon data analysis ;
- continuously *monitoring* the process and implementing control systems such as statistical process control.

The regulatory agencies encourage manufacturers to exploit the additional process knowledge that can be obtained during manufacturing (e.g., through sensors, analyzers etc...) to improve their processes (ICH, 2009a; ICH, 2009b). In this respect, a specific mention is given to the periodic maintenance of design spaces obtained using mathematical models. The ICH Q8(R2) document, in fact, explicitly states that “(...) for certain design spaces using mathematical models, periodic maintenance could be useful to ensure the model’s performance...” and that “expansion, reduction or redefinition of the design space could be desired upon gaining additional process knowledge” (ICH, 2009b). Interestingly, while specific mention to the possibility of updating the design space during the product lifecycle is given, no indication on *how* this should be done and communicated to the regulatory bodies are given.

Solicited by manufacturers (Herwig *et al.*, 2017; Watson *et al.*, 2018), ICH has recently released a draft for a new guidance, the Q12 guidance on *Technical and regulatory considerations for pharmaceutical product lifecycle management*, that aims at removing the “(...) technical and regulatory gaps that limit the full realization of more flexible regulatory approaches to post-approval changes as describe in ICH Q8(R2) and Q10 Annex 1...” (ICH, 2018). This document is still a draft guidance and is currently under public evaluation for final approval. Once finalized, the guidance will represent the regulatory reference for a common implementation of product lifecycle approaches, including design space updates and control strategy updates.

1.2.4 Technical and economic benefits of QbD: a critical review

Since the launch of the QbD initiative, several surveys and reviews have been published to understand the reaction of the pharmaceutical companies to this new paradigm, both in technical and economic terms.

IBM (IBM, 2005) was the first one to forecast the impact of QbD from an economic point of view. As an example, for a drug with US\$ 1 billion peak annual sales, it was estimated that practical implementation of the QbD paradigm could have generated an extra US\$ 1.6 billion

over the entire drug lifecycle (IBM, 2005). Despite this rather optimistic forecast, the opinion of market observers started becoming more prudent when the QbD implementation process started taking place within industry. In 2009 a private survey of the FDA, which was then made publicly available by McKinsey (McKinsey, 2009), revealed that, despite the potential of QbD to promote innovative technological advancements, many industrial companies were still skeptics on its real economic benefits (in terms of return on investment). FDA identified four different challenges occurring within companies and six internal challenges occurring within the regulatory bodies. The four challenges involving companies were briefly summarized as follows (McKinsey, 2009):

- disconnection between cross functional areas of the same company, e.g. R&D and manufacturing;
- lack of belief in business case, i.e., many companies proved to be skeptic about the timing and the investments required to practically implement QbD;
- lack of suitable technological equipment for a correct implementation of QbD;
- alignment with third parties (e.g., suppliers) on QbD implementation.

As per the regulatory bodies, the following criticalities were emphasized:

- presence of a small portion of people *within* FDA not supportive towards the real efficacy of QbD;
- lack of technical and specific guidance on how to implement *in practice* the principles of QbD;
- misalignment between international regulatory bodies;
- lack of preparation on the topic of some regulators that could be deduced during on-site inspections;
- scarce collaboration between companies and regulators;
- scarce ability of the regulatory to inspire confidence on the actual benefits of QbD implementation.

The regulatory agencies promoted different initiatives to face both the internal challenges and the challenges that were directly investing pharmaceutical companies. Positive results were obtained and the QbD paradigms started to spread at a very fast pace within the community. The increasing interest of companies towards a systematic implementation of the QbD approach was first confirmed by an industrial survey¹⁰ conducted by Kourti and coworkers (Kourti and Davis, 2012), where the benefits of QbD in terms of improved product and process understanding and improved robustness in manufacturing were extensively acknowledged. A later survey (Cook *et al.*, 2014) conducted by the American Association of Pharmaceutical Scientists (AAPS), involving 149 pharmaceutical scientists¹¹, reported that the majority of

¹⁰ The survey involved 12 different pharma companies, including bio-tech companies.

¹¹ The composition of the scientists was as follows: 88% from industry, 7% from academia, 4% from regulatory bodies, 1% others.

respondents (54 to 76%) acknowledged an extensive use of QbD tools and elements for product and process development. Moreover, respondents acknowledged the positive benefits of QbD in terms of the impact it can have on the patient's health (78%), as well as on internal processes such as knowledge management (85%), decision-making (79%) and lean manufacture (71%) (Cook *et al.*, 2014). However, scientists seemed to be more skeptics (more than 50%) on the economic benefits (in terms of return on investments) of the QbD implementation. Very recent surveys on this topic (Chatfield, 2017; Lundsberg-Nielsen and Bruce, 2017) confirm that QbD has now reached a mature and generalized level of implementation between companies, and the economic benefits that can be obtained are starting getting acknowledged. However, critical gaps still exist between the regulatory bodies and pharmaceutical companies on the technical details for a correct implementation between the QbD (Herwig *et al.*, 2017; Watson *et al.*, 2018; Collins, 2018). Process systems engineering tools are expected to play a key role to reduce this gap (Herwig *et al.*, 2017; Reklaitis, 2017) and to enhance the profitability of a correct QbD implementation.

1.3 Mathematical modelling

1.3 Mathematical modelling for QbD implementation

As extensively discussed in the previous sections, the most important target of the QbD paradigm is to promote the adoption of rigorous scientific tools to assist the different stages (pharmaceutical development, technology transfer, commercial manufacturing) of the lifecycle of a new pharmaceutical drug. One of the most important tools that can be used to achieve this purpose is mathematical modelling (García-Muñoz and Oksanen, 2010).

The importance of mathematical modelling for a correct implementation of the QbD paradigm is extensively emphasized in the regulatory documents (FDA, 2004c; ICH, 2009a; ICH, 2009b, ICH, 2012a). Since the final scope of every development/manufacturing activity is the achievement of the desired product quality, models are classified by regulators according to their impact (low/medium/high) in assuring this target (ICH, 2012a). Following the regulatory parlance (ICH, 2012a), low-impact models are defined as models that are typically used in product/process development (e.g., for formulation optimization). Medium-impact models are models that can be useful in assuring quality of the product but are not the sole indicators of product quality (e.g., most design space models). High-impact models are models whose predictions are significant indicators of product quality (e.g., a chemometric model for product assay or a surrogate model for dissolution).

From a PSE perspective, a much more useful classification of the models that can be exploited in pharmaceutical development/manufacturing is based on the following two criteria (Bonvin *et al.*, 2017):

- model *type*, i.e., knowledge-driven, hybrid or data-driven;

- model *scope*, i.e. based on the final purpose of implementation (e.g., models for design space determination).

From a general perspective, a model is characterized by three elements: equations, variables and parameters. The model *type* describes the amount of knowledge embedded within the model. In this regard, models can be classified as knowledge-driven (also known as *first-principles* or *mechanistic* or *deterministic*), hybrid (also known as *semi-empirical* or *gray-box*) and data-driven (also known as *data-based* or *empirical* or *black-box*).

Knowledge-driven models are based on fundamental knowledge of the underlying physical phenomena that govern the system under investigation. For these models, the *equations* describe in mathematical terms the physical laws (e.g., mass and energy balances, heat/mass transfer mechanisms) representing the system; the *variables* represent the system states; *parameters* inform on how the mathematical description given by the equations should be tuned to match the actual system behavior. Data-driven models do not embed any knowledge on the physical mechanisms involved in the system. For these models, the *equations* simply represent a convenient representation of the dataset(s) collected for the system under investigation; *variables* collect the inputs and outputs of the dataset; *parameters* inform on how equations should be tuned to match the available data. Semi-empirical models represent an intermediate situation between knowledge-driven and data-driven models, i.e. they combine fundamental knowledge on the system to describe certain phenomena with empirical reasoning to describe others phenomena.

The model *scope* describes the purpose of implementation (i.e., the application) of the model. Models can be exploited to assist the three stages of the drug lifecycle (pharmaceutical development, technology transfer, commercial manufacturing and continual process improvement). In pharmaceutical development, models can be used to assist all the sequential activities discussed in § 1.2.1 (QTTP and CQAs identification; product and process understanding; product and process design, including design space identification and definition of a control strategy) with the purpose of accelerating the launch of new products in the market. With respect to technology transfer, models can be used to assist both the *scale-up* and *scale-out* of the manufacturing process (including its design space and its control strategy), with the purpose of facilitating its transfer between different scales or between different sites. Finally, with respect to commercial manufacturing and continual process improvement, models can be used to assist product quality monitoring and control, as well as to enhance process productivity (i.e., *via* process optimization) and to assist process intensification.

The increasing interest of the pharmaceutical community towards mathematical modelling has been thoroughly reviewed by Troup and Georgakis (2013), Rogers and Ierapetritou (2015) and, very recently, by Reklaitis *et al.* (2017).

In 2013, Troup and Georgakis (2013) presented the results of a survey involving 21 professionals from worldwide top-pharma companies on the use of mathematical modelling for

process analytics, process monitoring, plant-wide information system, unit operation modeling, quality control and process optimization. The results showed that, with respect to process analytics and monitoring, the use of multivariate chemometric models and multivariate statistical process control represented common practices in the industrial context (according to 67% of the respondents), as well as the use of statistical multivariate tool to analyze historical process data. With respect to process modelling and optimization, more than one third of the respondents revealed that data-driven approaches were adopted to model between 80-100% of the unit operations of the manufacturing lines, while the other two thirds revealed that data-driven approaches were exploited to model *at least* 60% of the unit operations. However, the use of first principles models were acknowledged for at least 10% of the unit operations, mainly involving secondary manufacturing activities. Moreover, the advantages (possibility to perform extrapolation, wider applicability range, possibility to include product physical properties) and disadvantages (costs and complexity of model development, computational burden) of first-principles models with respect to data-driven models were stressed out by the respondents.

In 2015, a review by Rogers and Ierapetritou (2015) revealed the increasing interest of the pharmaceutical community towards the use of hybrid as well as first-principles models. This interest was partially driven by the increasing attention towards continuous manufacturing in the pharmaceutical community (§ 1.1.3). Very recently, Reklaitis *et al.* (2017) revealed that data-driven, hybrid and first-principles models are nowadays jointly exploited by pharmaceutical companies to tackle different problems of the drug lifecycle. Specifically, while data-driven approaches still represent the most adopted tools for process analytics, process monitoring and process control, hybrid as well as first-principles models are starting being exploited much more frequently to assist the different phases of pharmaceutical development. Moreover, due to the recent advancements in computational power and technology, it is expected that *online* applications of these model will play a key role in the next years (Reklaitis *et al.*, 2017).

Within the different scopes of mathematical modeling in the pharmaceutical context, the identification of the design space of a new pharmaceutical product probably represents the most important one (García-Muñoz and Oksanen, 2010). For example, almost 70% of the respondents of the survey of Troup and Georgakis (2013) reported the use of multivariate approaches for DS identification, while Reklaitis *et al.* (2017) suggested that this activity is nowadays common routine in pharmaceutical R&D. Since the launch of the QbD initiative, different modelling strategies have been proposed to tackle this issue. The aim of the next sections is to propose a critical review of the most important contributions on this topic, from both industry and academia.

1.4 Design space applications throughout drug lifecycle

Considering the three different activities (pharmaceutical development, technology transfer, commercial manufacturing and continual process improvement) of a drug lifecycle, different applications related to the concept of design space can be identified. These possible problems are briefly summarized in Fig. 1.3.

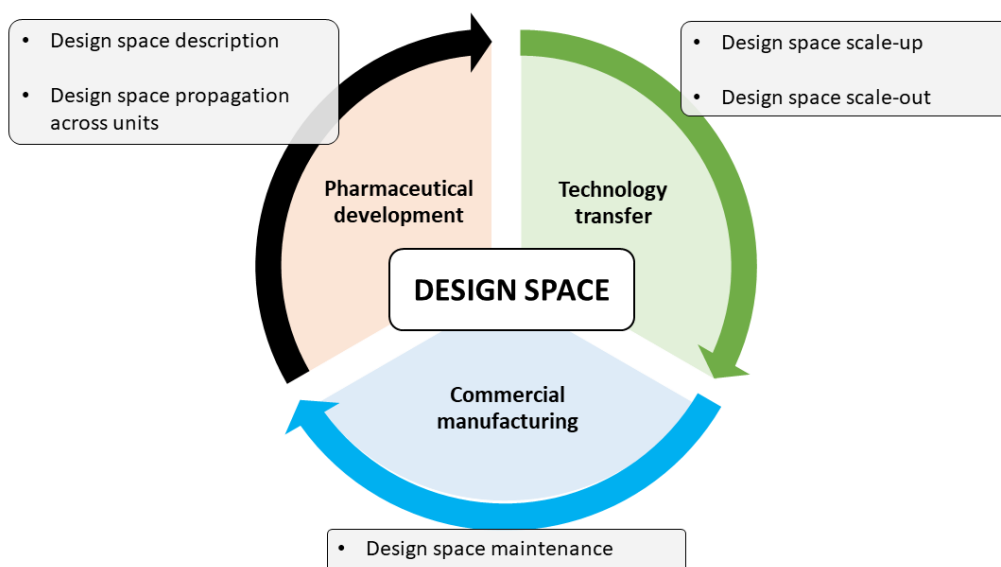


Figure 1.4 Different problems involved in a design space description during pharmaceutical development, technology transfer, and commercial manufacturing and continual process improvement.

With respect to pharmaceutical development, typical issues to be tackled are:

- design space description and uncertainty quantification;
- design space *propagation* from individual units to the entire process.

Design space *description* concerns the identification of the subset of raw material properties and CPPs that allows obtaining the desired quality characteristics for a new drug under development. This activity should be complemented with analysis of the uncertainty associate to the proposed description of the DS.

Design space *propagation* refers to the problem of defining the DS of the overall manufacturing line, based on the design spaces of the individual units arranged in the line. As discussed in §1.2.1.4, DSs are typically determined for single unit operations, but the determination of a single DS for the entire manufacturing process is always desirable to increase process flexibility.

With respect to technology transfer activities, issues related to the DS are:

- design space *scale-up*;
- design space *scale-out*.

Design space *scale-up* refers to the transfer of the DS obtained at the laboratory scale during product/process development to the commercial manufacturing scale by preserving its validity¹². On the other side, design space *scale-out* refers to the transfer of the design space between different sites producing the same product but implementing different process configurations.

With respect to commercial manufacturing and continual process improvement, the main issue related to the concept of DS is that of design space maintenance. Design space maintenance is the continuous verification and, possibly, update of the design space of the manufacturing process during plant operation, in order to preserve its validity throughout the entire production campaign.

Several contributions on the use of mathematical modeling to assist each of the activities described above have been presented in the literature. A short review of these contributions is presented in the following sections.

1.4.1 Design space description and uncertainty quantification

Exploitation of mathematical modelling for DS description is very frequent (Reklaitis, 2017). From a general perspective, DS description, as well as quantification of uncertainty in the description, require two elements:

- a) a model to relate the raw material properties and CPPs to the product CQAs;
- b) a methodology that, given the desired quality specifications and the model of point (a), returns the model-based prediction of the DS, including the uncertainty associated with this prediction.

With respect to the first point, several modelling strategies have been proposed in the literature to relate the CMAs and CPPs with the product CQAs for a new drug. A thorough review on this topic can be found in the work of Tomba (2012) and Reklaitis (2017). With respect to the second point, a summary of the different approaches (related to knowledge-driven, data-driven and hybrid models respectively) is shown in Fig.1.5.

1.4.1.1 Design space description strategies for data-driven models

Once a model to relate the CMAs and CPPs to the product CQAs has been selected and validated against experimental evidence, the next step is to determine the DS with a given mathematical technique and appropriate uncertainty metrics. As summarized in Fig.1.5, alternative approaches have been proposed to perform this task for knowledge-driven models, data-driven models and hybrid models.

¹² This means that the subset of CPPs and raw material properties that are derived from the laboratory-scale DS *must* guarantee the desired product quality on the commercial manufacturing scale.

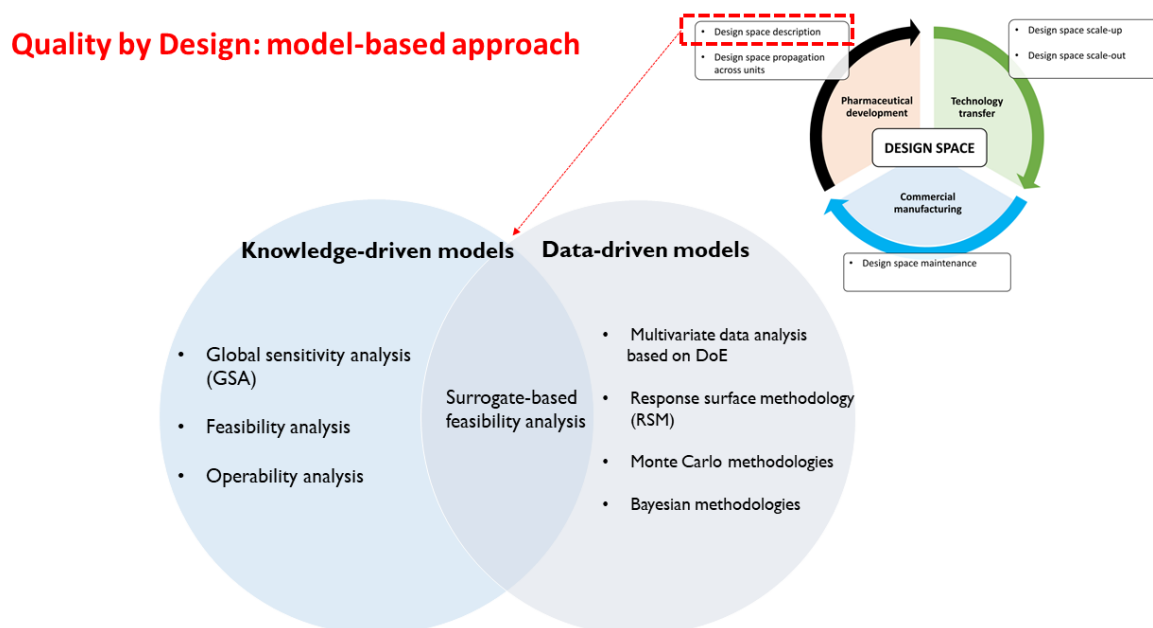


Figure 1.5 Available mathematical techniques to assist DS description.

In the context of data-driven models, multivariate data analysis based on DoE data is one of the oldest strategies that have been adopted to assist DS description. Once one selected multivariate model (e.g., multivariate linear regression) is built upon DoE-based data, this analysis simply consists of describing the DS in mathematical terms (*via* mathematical inversion of the original model) or, in most situations, in graphical terms (typically with contour surface plots that span the entire input domain). A detailed review of the different contributions on this topic has been presented by Lepore and Spavins (2008) and has been recently updated by Reklaitis *et al.* (2017).

Response surface methodology (RSM; Box and Wilson, 1951) represents a class of statistical methods that consist of three iterative steps:

1. selection of a *model* to relate CMAs and CPPs to product CQAs within the ones described above (typically, a polynomial regression model);
2. design (using standard DoE techniques) and execution of experiments in order to tune the parameters of the model until the model representativeness is satisfactory (if no satisfactory prediction fidelity is obtained by tuning the model parameters, the model structure can be changed);
3. exploitation of optimization methods to find the values of the independent variables (i.e., CMAs and CPPs) that produce the desired values of the responses (i.e., product CQAs). This step is equivalent to identify the boundary of the design space within the region of the input domain within which the model is deemed to be representative (also known as *knowledge space*).

Therefore, with respect to simple DoE-based multivariate analysis, RSMs have the peculiarity to require an iterative process that includes an optimization strategy. RSM has for long represented one of the most used techniques for DS determination in an industrial context (Lepore and Spavins, 2008; Garcia-Munoz *et al.*, 2015). However, due to the increasing attention towards hybrid and knowledge-driven models, its use has seen a slight decrease during the last years (Garcia-Munoz *et al.*, 2015). Nonetheless, several applications of RSM for DS determination have been reported. Examples include the use of RSM for the determination of the design space of fluid bed granulators (Zacour *et al.*, 2012), continuous mixers (Boukouvala *et al.*, 2010), wet granulators (Huang *et al.*, 2009), direct compression (Charoo *et al.*, 2012) or entire manufacturing lines (am Ende *et al.*, 2007). Several other applications have been recently reviewed by Reklaitis *et al.* (2017). The main advantage of these techniques is that they are very easy to use thanks to the use of specific commercial software; on the other hand, they are not able to handle model parameter uncertainty in a rigorous manner (Peterson, 2004).

Latent variable model inversion is related to the inversion of the LVMs in order to identify the subset of input combinations (CPPs and CMAs) that allows obtaining the desired product CQAs. The advantage of this approach is that the DS can be determined and represented on the *latent* space, which has usually a much smaller dimension than the original input space, thus obtaining a compact representation of the DS. The theoretical framework for LVM inversion has been originally proposed by Kourti (2006), and then implemented in practice by Garcia-Munoz *et al.* (2010) and eventually refined by Tomba *et al.*, (2012). Examples on the use of LVM inversion for DS description of tablet manufacturing lines (Liu *et al.*, 2011; Yacoub and MacGregor, 2011), granulation and roll compaction (Tomba *et al.*, 2013; Facco *et al.*, 2015) have been reported. A detailed description of these methods will be given in Chapter 2.

Monte Carlo methodologies perform a random sampling of the input domain (using one selected sampling strategy) and compute the values of the model outputs for each of the sampling points selected. The single values of the model outputs obtained for each sampling point are then used to obtain approximate distributions of the model outputs, where “distribution” is to be intended according to its frequentist meaning, not the Bayesian one. These distributions are then used to build appropriate confidence intervals (or regions, if the number of product CQAs is greater than one). The subset of input combinations that guarantee to obtain the product CQAs within these confidence regions is then identified as the DS of the product considered.

Applications of this approach with standard multivariate linear regression (Gujra *et al.*, 2007), stepwise multivariate linear regression (Debrus *et al.*, 2011) and polynomial regression models (Kauffman and Geoffroy, 2008) have been reported. Additional applications have been reviewed by Reklaitis *et al.* (2017).

Bayesian methodologies combine fundamentals of Bayesian statistics with one of the aforementioned data-driven models (in most situations, multivariate linear regression models)

to return a *probabilistic* representation of the DS of a new pharmaceutical product. The most important difference of these methods with respect, for example, to standard Monte Carlo methods is that model inputs and outputs (as well as model parameters) are not treated as *deterministic* variables (i.e., described by single numerical values), but as *random* variables (i.e., described by probability distributions). The probabilistic DS is obtained with these approaches by first building the posterior predictive distribution of the product CQAs, and then by identifying the subset of the input domain that guarantee to obtain the desired product quality characteristics with a probability greater than a predefined threshold value. The pioneer on the use of Bayesian techniques for DS determination is Peterson (2004; 2008), and some applications of these techniques for DS determination have been reported (Stockdale *et al.*, 2009; Peterson, 2009; Peterson, 2010). The main disadvantage of these methodologies is that they require expensive Markov Chain Monte Carlo (MCMC) simulations, which can become prohibitive when the number of input factors is large.

1.4.1.2 Design space description strategies for knowledge-driven models

In the context of *knowledge-driven* models, the mathematical techniques that have been proposed to assist a DS determination are:

- Global sensitivity analysis (GSA);
- feasibility analysis;
- operability analysis.

Due to the limited diffusion of knowledge-driven models for the description of pharmaceutical unit operations, the contributions that can be found in the open literature are fewer than for data-driven models. However, increasing attention has been directed towards these methodologies recently.

Global sensitivity analysis (GSA; Saltelli *et al.*, 2008) collects a series of simulation tools whose aim is to evaluate how the variability on the model outputs can be apportioned to the different model inputs. In pharmaceutical development, this can be used to:

- a) determine the CPPs and CMAs that are most influential with respect to the product CQAs;
- b) determine the probability distributions of the model outputs by spanning the entire range of variability of the model inputs;
- c) determine the subset of the input domain that allows satisfying the desired specifications on the product CQAs with a given confidence (i.e. the DS).

The three activities (a)-(b)-(c) can be performed separately or (more often) in a sequential manner. Different categories of GSA methods have been proposed, namely: *i*) screening methods; *ii*) regression-based methods; *iii*) variance-based methods; *iv*) metamodel-based methods. A thorough review on these techniques can be found in the study of Iooss and

Lemaître (2015). In the context of DS determination, the two approaches that are mostly adopted are screening methods and variance-based methods.

Screening methods (Saltelli *et al.*, 2008) consider wide ranges of variation for the model inputs and compute a global sensitivity metric for a given output as an average of local measures. The most important method belonging to this category is the Morris method (Morris, 1991). A recent application of this method for the determination of the most influential CPPs of a continuous pharmaceutical manufacturing line and subsequent DS determination has been proposed by Wang *et al.* (2017).

Variance-based methods (Helton *et al.*, 2003) decompose the variance of each single model output into several components, including the contributions of the single input factors as well as their interaction. The most famous methods belonging to this category are the so called Sobol's methods (Sobol, 1993; Sobol, 2001), which exploit a Monte-Carlo approach to compute the variance of each output and to decompose it according to the contributions of the different input factors. The use of these methods in pharmaceutical manufacturing has been strongly supported by their recent implementation in advanced process modelling environment such as gPROMS® (Process Systems Enterprise, 2014). Recent applications of Sobol's methods in pharmaceutical manufacturing have been proposed by Wang *et al.* (2017) and Garcia-Munoz *et al.* (2018), where their ability to return to rapidly screen the entire input domain and identify the boundary of the DS for a large number of input factors has been proven.

While GSA solves the DS determination problem as a “forward” problem, (i.e., by building the posterior distributions of the product CQAs through extensive sampling within the input domain and performing one simulation for each input sample), a technique that solves the DS determination problem as an “inverse” problem is feasibility analysis.

Feasibility analysis was first proposed for the design of general chemical processes (Halemane and Grossman, 1983) and that only recently has been applied for the determination of the DS of pharmaceutical processes. Differently from GSA, feasibility analysis solves the DS determination problem as an “inverse” problem: assigned the specifications (constraints) for the product CQAs, the portion of the input domain that allows satisfying those constraints (i.e., the DS) is obtained by solving a bi-level nonlinear programming optimization problem. The solution of the bi-level optimization problem requires the definition of a function, called *feasibility function*, that collects the maximum value of all the quality constraints imposed on product. Mathematical details on this technique, together with recent advances in this area, can be found in the recent work of Bhosekar and Ierapetritou (2017). Feasibility analysis has been proposed to assist DS description for single unit operations such as roller compaction (Banjeree *et al.*, 2010) granulation (Rogers and Ierapetritou, 2015; Wang and Ierapetritou, 2017) as well as for the determination of the DS of an entire continuous manufacturing line (Wang *et al.*, 2017b). The main advantage of this technique is that multivariate constraints on the product CQAs can be handled in straightforward fashion; the main disadvantage is that *all* the sources

of uncertainty (e.g., on model parameters) as well as external disturbances that may affect the DS prediction must be determined and modeled *a priori*. Additionally, this technique suffers from the curse of dimensionality, i.e., obtaining a solution within a reasonable amount of time can become prohibitive when the number of input factors is large.

Operability analysis (Vinson and Georgakis, 2000) is a mathematical technique that aims at understanding whether all the points in the desirable output ranges (i.e., product quality) are operable in the presence of expected disturbances and with the available ranges of the inputs (i.e., CPPs and CMAs). In the pharmaceutical context, the definition of DS is equivalent to the definition of “operable space” used by this technique. Description of the DS with operability analysis, although conceptually very similar to feasibility analysis, presents the substantial difference of not requiring the solution of a nonlinear optimization problem. Examples of applications for pharmaceutical unit operations have been proposed by Uztürk and Georgakis (2002) and by Georgakis (2017). The use of this technique for DS determination, however, is still very limited and mainly focused on bio-pharmaceutical operations.

1.4.1.3 Design space description strategies for hybrid models

In the context of hybrid models, mention should be given to surrogate-based feasibility analysis (Banjeree *et al.*, 2010; Rogers and Ierapetritou, 2015), which is a particular type of feasibility analysis that is applied with computationally expensive models characterized by black-box constraints on the product CQAs (i.e., hybrid models). The idea behind this technique is to build a computationally cheap approximation of the original hybrid model, which is called a surrogate, and to compute the DS of the original model based on this surrogate. Different types of surrogate can be used (Banjeree *et al.*, 2010), including kriging surrogates, radial basis functions (RBFs) and simple polynomial interpolants (Rogers *et al.*, 2015). Applications of surrogate-based feasibility analysis for DS determination have been proposed recently (Wang and Ierapetritou, 2017a; Wang *et al.*, 2017). A detailed description of this approach will be provided in § 5.2.

1.4.2 Design space propagation from individual units to an entire process

A pharmaceutical manufacturing process is typically composed by several processing units. As described in the regulatory documents (ICH, 2009), during a DS description it would always be desirable obtaining the DS for the *entire* manufacturing process rather than for the individual units only. Some attempts of defining the DS for entire manufacturing lines have been proposed recently (Wang and Ierapetritou, 2017b). However, in many practical situations, the determination of the DS is performed for each individual unit of the manufacturing line, due to the complexity (both in terms of model building effort and computational power required) of dealing with integrated flowsheet models. Strategies to describe the DS of the entire process based on the DSs of the individual unit operations would therefore be useful to fit the regulatory

guidelines. However, the contributions that can be found in the literature on this topic are very limited. From a modelling perspective, the determination of the DS of an entire process based on the DSs of the individual processing units can be interpreted as a problem of how input variability propagates across the entire manufacturing line (the output of the first unit becomes the input of the second one, and so on) till the final product. Attempts of studying this variability propagation have been proposed by Rogers and Ierapetritou (2015) and Wang and Ierapetritou (2017a) in the context of feasibility analysis, and by Metta *et al.* (2018) in the context of a DEM- PBM.

1.4.3 Design space scale-up and scale-out

The design space of a new pharmaceutical product is typically determined and validated at laboratory scale during the drug development stage. However, its final use is intended for the commercial manufacturing scale. The activity that is deemed to transfer the DS obtained at the laboratory scale to the commercial scale is defined as *DS scale-up*. From an industrial perspective, the advantage of developing a DS at laboratory scale and scaling it up to the commercial manufacturing scale is that the costs for the experimental assessment of the DS are much lower at the laboratory scale. Therefore, the main target of the DS scale-up activity is to reduce to a minimum the number of experimental runs that must be performed on the commercial-scale plant for the final assessment of the DS. In this regard, model-based approaches play a key role to assist a DS scale-up exercise.

Despite several contributions on the scale-up of pharmaceutical unit operations/processes have been published, the number of studies on DS scale-up is very limited. In the context of data-driven approaches, the available contributions are limited to the exploitation of *historical* plant data of products similar to the one under development to perform the DS scale-up activity. Latent variable models are exploited to this purpose; in particular, Garcia-Munoz *et al.* (2004, 2005) proposed a new PLS-based technique, called joint-Y PLS (JY-PLS), specifically designed to relate two datasets of the same process at two different scales, in order to identify the relationships between variables at multiple scales and identify similarities that can be exploited for DS scale-up. Liu *et al.* (2011b) exploited this technique to assist the DS scale-up for a roller compaction process, while Garcia-Munoz *et al.*, (2009) used JY-PLS to understand the effect of raw material properties on product CQAs on different scales in order to guide a DS scale-up exercise.

In the context of knowledge-driven models, a larger number of contributions can be found, and this is mostly related to the fact that knowledge-driven models require fewer plant data. These studies can be classified according to two categories:

- a) contributions based on the derivation of scale-up relationships from CFD or DEM simulations;

b) contributions based on the derivation of scale-independent formulations for the DS.

In the first category, scale-up relationships are derived from complex CFD/DEM simulations and used to scale up the DS from the laboratory scale to the production scale. Examples of this approach include applications to fluid bed drying (Parker *et al.*, 2013; Chen *et al.*, 2018), freeze drying (Zhu *et al.*, 2018), mixing (Lindenberg and Mazzotti, 2009), film and pan coating (Pandey *et al.*, 2013; Chen *et al.*, 2017) and blending (Horibe *et al.*, 2018). As an example, Pandey *et al.* (2013) developed a DEM simulation to study the effect of changing operating conditions and equipment scale on the particle motion during pan coating, and based on the DEM results developed a simple scale-up relationship for the DS of the pan coater. Zhu *et al.* (2018) followed a similar procedure for the scale-up of a freeze drying cycle, but they exploited intensive CFD simulations. Horibe *et al.* (2018) used DEM simulations to derive scale-up relationships for pharmaceutical blenders.

In the second category, a mechanistic model is typically used to determine the DS at the laboratory scale (i.e., the DS prediction is validated against experimental data obtained at the laboratory scale). Then, the model is re-formulated using scale-independent parameters (typically, dimensionless numbers) that allow transferring the DS to the larger scale. It is worth noticing that this approach is also advocated by the regulatory agencies (§ 1.1.2.). Reported applications of this approach are very limited and mostly refer to the scale-up of the DS for pharmaceutical freeze-drying processes. For example, Fissore and Barresi (2011b) and Pisano *et al.* (2013) derived a scale-independent formulation for a simple mechanistic model of the primary drying stage of a freeze-drying process, and used this formulation to derive the DS of the process at commercial scale. A very recent work of Yoshino *et al.* (2018) use a similar approach to assist the scale-up of a film coating process, while Garcia-Munoz *et al.* (2015) exploit this methodology to scale-up the DS for the adsorption of a drug substance in a packed-bed reactor.

Scale-up is not the only technology transfer activity that involves a DS description. In fact, when the product is already being manufactured in the commercial plant, there may be the need to:

- transfer the manufacturing of the given product to a different site, with a different process configuration;
- use the same process to produce a different product.

Both activities are referred to as *scale-out* activities (Reklaitis, 2017). However, whereas the problem of product transfer has been addressed by several authors (Jaeckle and MacGregor, 2000; Garcia Munoz *et al.*, 2005; Tomba *et al.*, 2013), the problem of assisting the *scale-out* of a DS has not been addressed so far in the open literature, even if its importance has been emphasized (Reklaitis, 2017).

1.4.4 Design space maintenance

If the DS for a given product is known, the input materials properties and process operating conditions lie within the DS. However, if the DS was described using a model of the process, an issue arises on whether the model originally used to simulate the process (hence, to describe the DS) is still valid under the current process state. Process/model mismatch may arise from uncertainty in the evaluation of some model parameters (e.g., a heat exchange coefficient may change due to fouling), or insufficient/inappropriate description of the underlying physical phenomena driving the process (e.g., environmental factors not accounted for). The fact that, during the plant lifecycle, the model may not be able to accurately reproduce the actual plant behavior, questions the appropriateness of the operating conditions selected by using the DS defined using that model.

When a model is used to assist DS description in a lab-scale plant, the different sources of uncertainty (e.g., parametric uncertainty) and structural mismatch (e.g., environmental factors) that can potentially affect the prediction fidelity of the model at the plant scale should be included and accounted for in the DS description. This requires knowing *a priori* all the uncertainties and disturbances that may affect plant operation. However, in practical situations, modelling all these possible sources of uncertainty and disturbances is never possible. Different phenomena may occur during plant operation (e.g., parameters drifts or unmodeled phenomena that could not be observed at the laboratory scale, such as the effect of upstream and downstream units) that can impact on the prediction fidelity of the model, hence on the model-based DS representation. In some cases, some of these phenomena can be observed and accounted for during the DS scale-up activity. However, there is no guarantee that plant operation will remain stable over long periods, hence that the initial model-based DS description will preserve its validity throughout the entire production campaign.

Following the regulatory parlance, the activity that is deemed to keep constantly updated the DS of a pharmaceutical product throughout its entire lifecycle is named DS *maintenance* (ICH, 2009). even though no indications are given by the regulators on how to perform this activity. From a modelling perspective, despite the importance of performing a regular model maintenance during pharmaceutical manufacturing is not new (Wise and Roginski, 2015; Flaten, 2018), applications are limited to the maintenance of multivariate data-driven models (e.g., PCA or PLS models) for process monitoring purposes. Contributions on how to perform a DS maintenance (with either data-driven or knowledge-driven models) are lacking. The relevance of developing such contributions within the framework of the Industry 4.0 initiative has been emphasized by Herwig *et al.* (2017), and the recent launch of the Q12 regulatory document (ICH, 2018) is expected to push the scientific community to tackle this topic in the near future (Herwig *et al.*, 2017).

1.5 Objectives of the research

Despite the increasing attention towards mathematical modelling for DS description, several issues still need to be addressed to obtain a general modelling framework for practical implementation of the regulatory guidelines. These issues can be briefly summarized into four points.

1. A key requirement advocated by the regulatory agencies for a correct description of the design space is the quantification of the *assurance* of quality for the final product. From a modeling perspective, this translates into a rigorous quantification of the *uncertainty* associated with the model-based prediction of the DS. However, the majority of the contributions presented in the literature focus on the adoption of *deterministic* approaches for DS identification rather than *probabilistic* ones: whereas the former do not consider model uncertainty, the latter do.
2. The few probabilistic approaches presented in the literature mainly focus on situations where the product quality can be described in terms of univariate quality specifications, while in most practical situations the quality of the product is expressed in terms of multivariate specifications.
3. In the context of first-principles and semi-empirical modeling, advanced techniques are required in order to handle several inputs (process parameters and raw material properties) in a computationally efficient way. Most of the proposed techniques still rely on computationally expensive simulations which, in some cases, are not even able to return a DS description with the desired level of accuracy. Moreover, providing a compact and easy-to-interpret representation the DS when a large number of inputs are involved is typically impossible, thus preventing easy interpretation of the results by regulators and non-practitioners. There is therefore the need to develop computationally efficient methodologies that allow obtaining user-friendly and ready-to-use representations of the design space.
4. While several modeling strategies have been proposed for the identification of the DS at the product and process development stage, very few studies address how to regularly assess and possibly adjust its model-based representation during product manufacturing. In fact, the prediction fidelity of the model adopted at the product/process development stage may deteriorate during plant operation, thus affecting the validity of the prediction of the DS obtained at the beginning of the product lifecycle. In other terms, several phenomena may arise during plant operation (e.g. parameter drifts, un-modeled physical phenomena such as environmental effects, effect of upstream and downstream units, etc...) that may partially or totally invalidate the prediction of the DS obtained at the product and process development stage. The importance of exploiting the additional process knowledge that can be gained during

plant operation for reduction, expansion or redefinition of the design space is clearly stated in the regulatory documents (ICH, 2009) and is part of a continual process improvement strategy that pharmaceutical companies should adopt for a correct implementation of the QbD framework. The scientific contributions on this topic are still very limited and tailored to very specific applications, while the development of general methodologies is still an open research area.

5. The identification of the DS by means of first-principles or semi-empirical models require accurate and precise estimation of the model parameters, in order to obtain a high prediction fidelity and therefore a reliable estimation of the DS. The parameter estimation activity is performed by challenging the model predictions against the experimental observations that can be collected during an experimental campaign on the real equipment. When the number of parameters is large or the parameters are strongly correlated with each other, obtaining an accurate and precise parameter can be difficult and require a large number of experiments. It is therefore desirable to design new experiments in order to *maximize* the information that can be extracted from these experiments for a precise estimation of the model parameters. A class of techniques that can be exploited to this purpose are model-based design of experiments (MBDoe) techniques. While MBDoe techniques have recently gained increasing attention in many process engineering applications, their use in pharmaceutical product and process development contexts is still very limited. The introduction of MBDoe to assist pharmaceutical product and process development is therefore still an open research area.

In view of the above, the objective of this Dissertation is to propose novel and systematic methodologies to tackle the research issues previously discussed. Specifically, the innovative contributions that can be found in this Dissertation are the following.

- **Quantification of the “assurance” of quality for a new pharmaceutical product** as advocated by the regulatory agencies, with particular focus on data-driven modeling strategies such as latent variable modeling (LVM) and multivariate linear regression. Two situations are addressed, namely: *i*) the situation where a projection onto latent structures (PLS) model is to be used for design space identification; *ii*) the situation where a general multivariate linear regression model is to be used for DS identification. With respect to the former case, the inversion of latent-variable models, such as PLS models, has proved to be an effective tool to assist the determination of the DS of a new pharmaceutical product. A challenging issue in PLS model inversion is to describe how the uncertainty on the model outputs (product quality) relates to the uncertainty on the model inputs (raw material properties and process parameters). This in turn translates

into a quantitative inclusion of the concept of assurance of quality for the final product in the model-based representation of the design space. The objective of the first part of this Dissertation is to develop a methodology to relate the uncertainty on the output of a PLS model to the uncertainty on the model inputs. The expected result is the identification of a portion of the original input domain within which the DS of the product under development is expected to lie with a confidence equal to or greater than an assigned threshold. This narrow region of the input domain is expected to be used to assist a targeted experimental campaign for the final assessment of the DS, with significant reduction of the time and costs of the experimentation.

With respect to the second aspect, the objective is to exploit elements of Bayesian statistics in order to quantify the probability that a pharmaceutical product will meet its quality specifications during manufacturing. The focus is restricted on situations where a multivariate linear regression model is used to relate the process parameters and raw material properties with the product quality attributes. The final target is to obtain a quantitative metric, i.e. the Bayesian *probability* of the product to meet its quality specifications, that completely addresses the concept of assurance of quality advocated by the regulatory agencies. The DS that can be obtained within this Bayesian framework can therefore be defined as a *probabilistic* DS, and it can be claimed that its derivation is entirely consistent with the regulatory guidelines.

- **Design space description using semi-empirical or first-principles models**, with particular focus on the reduction of the overall computational burden and the development of low-dimensional and easy to interpret representations of the design space. The research effort will be targeted on a mathematical technique that has been recently proven to be potentially very effective for DS determination with complex mathematical models. This technique is surrogate-based feasibility analysis and, although being very efficient in handling complex multivariate quality specifications and complex process models, it suffers from one important drawback, i.e. the curse of dimensionality. In other terms, when the number of input factors is large, this technique shows severe limitations in obtaining an accurate prediction of the design space within a reasonable computational time. A methodology will be developed in this Dissertation to solve this issue and an application to a continuous manufacturing line of a pharmaceutical tablet will be presented.
- **Maintenance of a model-based design space during plant operation**, based on the concept of continual process improvement and lifecycle management advocated by the regulatory agencies. The target in this respect is to develop an automatic methodology to obtain an accurate real-time representation of the design space as process operation

progresses, by exploiting the additional process knowledge that can be captured during product manufacturing. The methodology presented in this Dissertation will use concepts of model adaptation and state estimation to reach this target, and concepts of feasibility analysis and surrogate-based feasibility analysis to assist the real time DS identification.

- **Design of optimal experiments for model calibration and design space description**, with specific focus on industrial freeze-drying processes. The main objective is to show how the exploitation of model-based design of experiments (MBD_{oE}) can be used to assist the design of informative experiments for the identification of complex mechanistic models of freeze-drying processes, and how these models can be used to assist DS description. The final target is to reduce the number and cost of the experiments to be performed for model identification, thus speeding up the process development stage and boosting its profitability.

Simulated, experimental and industrial case studies will be presented throughout the Dissertation to prove the effectiveness of the proposed methodologies. A schematic roadmap for the interpretation of the Dissertation is presented in the next section.

1.6 Dissertation roadmap

The Dissertation is organized following the four research objectives presented in the previous section. A schematic roadmap of the Dissertation is shown in Fig. 1.5. After a brief review of the mathematical tools exploited throughout the Dissertation (Chapter 2), the central chapters (Chapters 3,4,5,6) collect four innovative methodologies to tackle the first three research objectives presented above. The nature of these chapters is purely *methodological*, and the effectiveness of the innovative modeling strategies proposed are tested with industrially-relevant case studies. The final chapter (Chapter 7) is focused on an industrial *application* of a well-established methodology (model-based design of experiments; MBD_{oE}): its main innovation is therefore related to the area of application of the methodology, rather than the methodology itself.

With respect to the issue of quantification of assurance of quality, in Chapter 3 a methodology is proposed to handle the back-propagation of uncertainty from the outputs (i.e. product quality) to the inputs (i.e. raw material properties and process parameters) of a PLS model. The methodology is exploited to identify a small portion (defined as the *experiment space*; ES) of the input domain within which the design space of a new pharmaceutical product is expected to lie with a given degree of confidence. It will be shown with both simulated and experimental case studies that the identification of such experiment space allows the product developer to

tailor his/her experimental campaign on a smaller domain of input combinations, with significant reduction of the time and cost of the experiments for DS assessment.

In Chapter 4, a methodology is proposed to quantify the *probability* that a new drug will meet its quality specifications by combining latent variable modelling (specifically, a PLS model) with a Bayesian multivariate linear regression model. The set of input combinations for which this probability (to be intended in its Bayesian interpretation) is greater than an assigned threshold is identified as the *probabilistic* (or Bayesian) design space of the new product. The Bayesian probability can therefore be considered as a scientific and unambiguous metric to implement the concept of assurance of quality advocated by the regulatory agencies, and the probabilistic design space can be considered as fully compliant with the regulatory guidelines. It will be shown how PLS modelling can be coupled with Bayesian model calibration to reduce the dimensionality of the input domain, thus allowing :i) a significant reduction of the computational burden of the simulations (which require extensive Markov chain Monte-Carlo computations); ii) a compact representation of the design space that can be easily interpreted by regulators or non-practitioners. The effectiveness of the proposed approach will be shown with both simulated and experimental case studies of pharmaceutical unit operations.

In Chapter 5, a methodology to overcome the curse of dimensionality of surrogate-based feasibility analysis for design space description of a new pharmaceutical product with semi-empirical or first-principles models is presented. The ability of PLS to reduce the input space dimensionality is exploited to obtain a low-dimensional representation (i.e., a latent representation) of the input domain. An adaptive sampling feasibility analysis based on a radial-basis function (RBF) surrogate is then used to identify the boundary of the design space on the latent space. It will be shown how the proposed approach can be exploited to give an accurate and robust description of the design space by simultaneously reducing the computational burden for the feasibility analysis problem. The effectiveness of the methodology will be challenged with a complex integrated flowsheet model of a continuous manufacturing line of a pharmaceutical tablet.

Chapter 6 presents a methodology to obtain a real-time representation of the design space of a pharmaceutical process while plant operation progresses. Given a first-principles model of the process and measurements from plant sensors, the proposed approach exploits a dynamic state estimator and feasibility analysis to obtain the real-time representation of the design space. The state estimator is deployed to adapt online the model predictions with the available plant observations, while feasibility analysis and surrogate-based feasibility analysis are exploited to obtain the boundary of the design space with the up-to-date model returned by the state estimator. The ability of the methodology to timely track changes in the DS representation will be shown with two simulated case studies, the former involving an enforced parametric mismatch, the latter a structural mismatch.

Finally, Chapter 7 discusses an industrial application of MBD_oE techniques for the identification of a mechanistic model of a pharmaceutical freeze dryer. Particular focus is given to the primary drying stage of the process, which has a strong impact on the overall process efficiency due to its duration (~60-80% of the overall freeze-drying cycle). Optimal experiments are designed in order to extract the maximum information from the data for the estimation of two critical model parameters. A good improvement of the model identifiability will be shown by performing the designed experiments both *in silico* and in the real equipment.

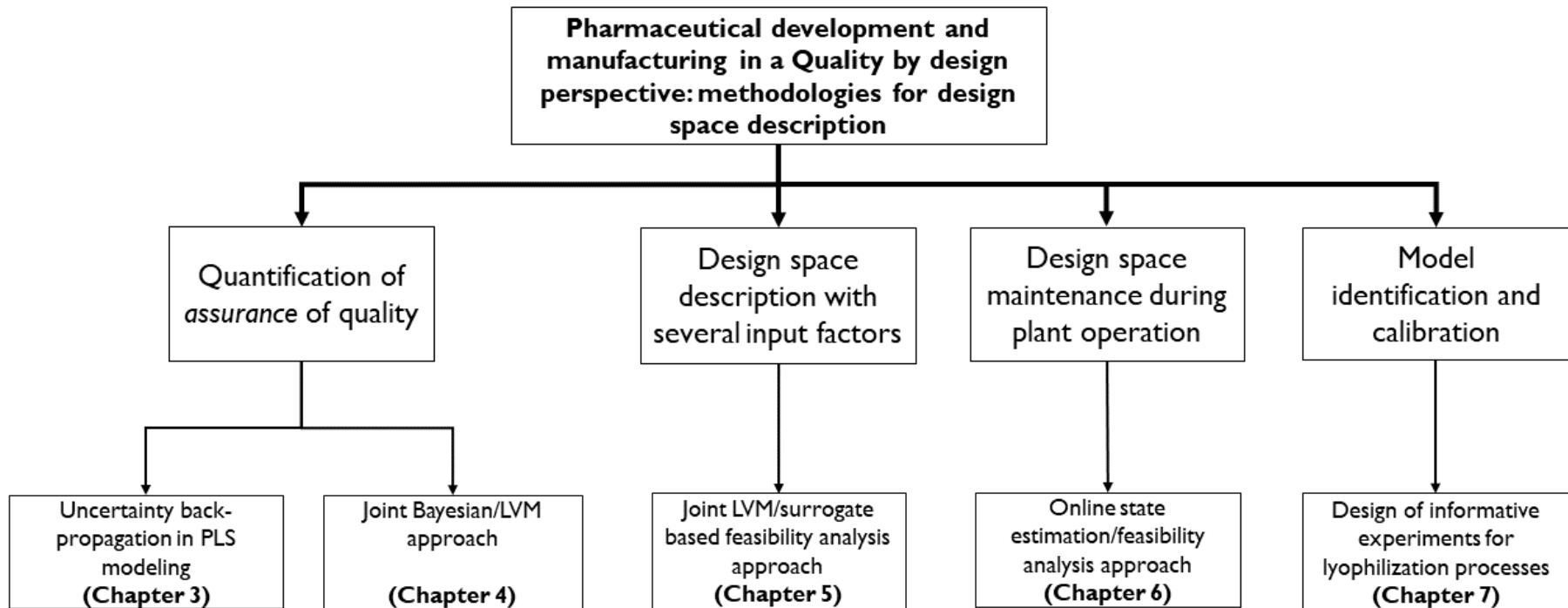


Figure 1.6. Dissertation roadmap.

Chapter 2

Methods

This chapter provides a concise overview of the mathematical techniques that are exploited in this Dissertation. The chapter is organized in four sections. In Section A, a review of projection on latent structures (PLS) for dimensionality reduction is presented. In Section B, the key concepts of Bayesian statistics and Bayesian multivariate regression are summarized. In Section C, the mathematical formulation of feasibility analysis and surrogate-based feasibility analysis is discussed. Lastly, in Section D state estimation techniques for online model adaptation and recalibration are presented.

SECTION A: DIMENSIONALITY REDUCTION TECHNIQUES

Latent variable models (LVMs) are a class of statistical models that are used for three main purposes: (i) data interpretation, (ii) dimensionality reduction and (iii) regression analysis.

Let $\mathbf{X} [N \times V]$ be a historical dataset of N observations of V *observable* variables. The aim of LVMs is to explain the correlation structure of \mathbf{X} by means of A *unobservable* variables, called latent variables (LVs), that explain the maximum multidimensional variance of the original dataset. In the presence of a strong collinearity in the original dataset, LVMs are able to capture the information retained in the matrix \mathbf{X} with a much smaller number of LVs, i.e. $A \ll V$. One of the most important LVM that can be used to interpret and reduce the dimensionality of a highly correlated dataset is principal component analysis (PCA; Jackson, 1991).

In many other applications, the historical dataset is split into a matrix of V input variables $\mathbf{X} [N \times V]$ and a matrix of M response variables $\mathbf{Y} [N \times M]$. In this context, LVMs are used to explain the joint correlation structure of \mathbf{X} and \mathbf{Y} , i.e. the LVs are selected in order to explain the maximum multidimensional variance of the input space that is mostly correlated with the output space. The identification of this correlation structure can then be used to *predict* a new response given a new set of model regressors (i.e. for regression analysis). One of the most common techniques in this area is projection on latent structures (PLS; Wold *et al.*, 1983). In the following, a brief theoretical overview of PLS and on the available methodologies to estimate prediction uncertainty in PLS modeling is presented.

2.1 Projection on latent structures (PLS)

Projection on latent structures (PLS; Wold *et al.*, 1983; Höskuldsson, 1988) is a LVM that aims at (i) explaining the joint correlation structure of the input matrix \mathbf{X} and the response matrix \mathbf{Y} and (ii) predicting a new response $\hat{\mathbf{y}}_p [1 \times M]$ given a set of new regressors \mathbf{x}_p .

PLS can be used to perform the first task or both the two tasks described above. In the former case, PLS is used as a dimensionality reduction technique. In the latter case, PLS is used as a multivariate linear regression technique. In both scenarios, PLS identifies a small set of A LVs that explain as much as possible of the joint covariance of \mathbf{X} and \mathbf{Y} . This decomposition proves to be particularly useful when a set of response variables has to be correlated with a large set of highly collinear input variables.

The structure of a PLS model can be summarized by the set of equations:

$$\mathbf{X} = \mathbf{TP}^T + \mathbf{E}_X \quad (2.1)$$

$$\mathbf{Y} = \mathbf{TQ}^T + \mathbf{E}_Y \quad (2.2)$$

$$\mathbf{T} = \mathbf{XW}^* \quad (2.3)$$

where $\mathbf{T} [N \times A]$ is the score matrix, $\mathbf{P} [V \times A]$ and $\mathbf{Q} [M \times A]$ are the \mathbf{X} and \mathbf{Y} loading matrices, \mathbf{E}_X and \mathbf{E}_Y the residuals; $\mathbf{W}^* [V \times A]$ is the weight matrix, through which the data in \mathbf{X} are projected onto the latent space to give \mathbf{T} according to Eq. (2.3).

The weight vector $\mathbf{w}_1 [V \times 1]$ related to the first LV can be computed by solving the eigenvector decomposition problem for the matrix $\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}$:

$$\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1 \quad (2.4)$$

where λ_1 is the eigenvalue associated with the first LV. Eq. (2.4) is equivalent to solving the optimization problem:

$$\max_{\mathbf{w}_1} (\mathbf{w}_1^T \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}_1) \quad (2.5)$$

$$s.t \quad \mathbf{w}_1^T \mathbf{w}_1 = 1. \quad (2.6)$$

Given \mathbf{w}_1 , the score vector $\mathbf{t}_1 [N \times 1]$ related to the first LV can be computed according to the following expression:

$$\mathbf{t}_1 = \mathbf{X} \mathbf{w}_1. \quad (2.7)$$

The weights related to the subsequent LVs can be computed according to the optimization problem (2.5), but using the deflated input and output matrices \mathbf{X}_{a+1} and \mathbf{Y}_{a+1} respectively. These matrices are defined for $a = 1, 2, \dots, A - 1$ as:

$$\mathbf{X}_{a+1} = \left(\mathbf{I}_N - \frac{\mathbf{t}_a \mathbf{t}_a^T}{\mathbf{t}_a^T \mathbf{t}_a} \right) \mathbf{X}_i \quad (2.8)$$

$$\mathbf{Y}_{a+1} = \left(\mathbf{I}_N - \frac{\mathbf{t}_a \mathbf{t}_a^T}{\mathbf{t}_a^T \mathbf{t}_a} \right) \mathbf{Y}_i \quad (2.9)$$

where \mathbf{I}_N is the $N \times N$ identity matrix. The corresponding loading vectors for the a -th LV related to X and Y respectively are given by:

$$\mathbf{p}_a^T = \frac{\mathbf{t}_a^T \mathbf{X}_a}{\mathbf{t}_a^T \mathbf{t}_a} \quad (2.10)$$

$$\mathbf{q}_a^T = \frac{\mathbf{t}_a^T \mathbf{Y}_a}{\mathbf{t}_a^T \mathbf{t}_a}. \quad (2.11)$$

The weight matrix \mathbf{W}^* that appears in Eq. (2.3) is related to the matrix \mathbf{W} that collects the weight vectors $\mathbf{w}_i, i = 1, \dots, V$ according to the following relationship:

$$\mathbf{W}^* = \mathbf{W}(\mathbf{P}^T \mathbf{W})^{-1}. \quad (2.12)$$

It is worth noticing that the \mathbf{X} -score matrix of Eq. (2.2) is sometimes substituted by the \mathbf{Y} -loading matrix $\mathbf{U} [N \times A]$. The two score matrices are related to each other by a linear relationship (called *inner relation*; Geladi and Kowalski, 1986), thus making the two formulations mathematically equivalent.

The solution of the eigenvector problem (2.4) is not straightforward from an algebraic point of view (e.g. via singular value decomposition (SVD) techniques) and is particularly cumbersome when dealing with missing data in the original dataset. Therefore, several iterative algorithms have been proposed in the literature for a computationally efficient implementation of PLS. The two most famous algorithms are the NIPALS algorithm (Wold *et al.*, 1983) and the SIMPLS algorithm (de Jong, 1993). A thorough description of these algorithms can be found in the cited references.

2.1.1 Pre- and post- PLS modeling activities

When a PLS model is used to describe a historical dataset $[\mathbf{X}; \mathbf{Y}]$, three activities must be performed at different stages of the modeling procedure, namely:

1. data pretreatment (before building the PLS model);
2. selection of the optimal number of LVs (to finalize the PLS model structure);
3. PLS model diagnostics (to validate the performance of the PLS model).

A point-to-point description of these activities is described in the following sections.

2.1.1.1 Data pretreatment

Data pretreatment is a preliminary step that is performed in order to avoid misleading results from the PLS model. This activity needs to be done to account for the different physical nature and dimensions of the variables of the original dataset.

Different data pretreatment techniques exist (Eriksson *et al.*, 2006), including mean-centering, auto-scaling, filtering and denoising. When the variables of the historical dataset are physically different, as it is often the case with process variables, the most suitable data pretreatment technique is auto-scaling.

Auto-scaling consists of subtracting to each column of the original dataset [\mathbf{X} ; \mathbf{Y}] the mean value of that column and then dividing each element for the standard deviation of the column. In other terms, each column is mean-centered and then normalized with its standard deviation. This allows obtaining unbiased directions of maximum variability from the PLS model that are not affected by the mean values of the original variables. Moreover, the scaling operation (i.e. normalization with the standard deviation) has the advantage of partially linearizing the original dataset.

All the studies presented in this Dissertation assume that the historical data have been auto-scaled according to the procedure described above.

2.1.1.2 Selection of the optimal number of LVs

A key step of the PLS model building activity is the selection of the number of LVs (i.e. the dimensionality of the latent space). From a general perspective, the selection of the number of LVs depends on several factors and should be suited to the final application of the PLS model. In most situations, it may be desirable to select a number of LVs that explain a big portion of the variability of both the input and output datasets. However, there may be situations where only one of the two aforementioned requirements needs to be fulfilled.

Different methods have been proposed in the literature to select the appropriate number of LVs.

The most famous are:

- the scree test (Jackson *et al.*, 1991);
- the eigenvalue-greater-than-one rule (Mardia *et al.*, 1979);
- cross-validation (Wold, 1978).

The scree test is a graphical procedure that monitors a metric (the explained variance R^2 of the input and output calibration datasets) and assumes that the optimal number of LVs is the one that yields to a “stabilization” to the metric profile. The underlying assumption behind this

approach is that the amount of explained variability of the original dataset reaches a “steady-state” after a certain number of LVs.

The eigenvalue-greater-than-one rule is a method that discards all the latent variables whose associate eigenvalue is smaller than one. In fact, if data are auto-scaled, the eigenvalue associated with the a -th LV roughly describes the number of original variables that are captured by the given LV. Therefore, all the latent variables that describe less than one original variable are discarded and not considered in the PLS model.

Cross-validation is a technique that selects the optimal number of LVs by minimizing the error of the PLS model in reconstructing new samples. Different cross-validation procedures are presented in the literature: the most common is the one proposed by Wold (Wold, 1978). Mathematical details can be found in the cited reference.

2.1.1.3 PLS model diagnostics

Once a PLS model has been calibrated, it is necessary to assess its performance by means of some diagnostic procedures. Three different areas can be identified: model diagnostics, sample diagnostics and variable diagnostics. All these procedures can be performed on one or both the input and output matrices.

As regard *model diagnostics*, the most important metric that is typically used is the amount of variability of the original data captured by the model, i.e. the coefficient of determination R^2 :

$$R^2 = 1 - \frac{\sum_{n=1}^N \sum_{v=1}^V (x_{n,v} - \hat{x}_{n,v})^2}{\sum_{n=1}^N \sum_{v=1}^V (x_{n,v})^2}. \quad (2.13)$$

The value $\hat{x}_{n,v}$ is the PLS reconstruction of the $x_{n,v}$ element of the original input matrix. The same definition (2.13) can be applied to the output matrix. To distinguish between the amount of \mathbf{X} - and \mathbf{y} - variability explained by the PLS model, the coefficient of determinations will be denoted as R_X^2 and R_Y^2 respectively.

As regard *sample diagnostics*, it is possible to identify potential outliers or samples of the original dataset that have a strong influence on the PLS model with two metrics, namely the Hotelling's T^2 and the squared prediction error SPE .

The Hotelling's T^2 (Hotelling, 1933) is a metric that quantifies the Mahalanobis distance from the projection of a sample on the latent space to the origin of the latent space itself. It is mainly used to assess the deviation of a given sample with respect to the average conditions of the historical dataset. The higher this metric for a given sample, the lower the adherence of the sample with respect to the calibration dataset, the higher its *leverage* with respect to the model (i.e. the higher its influence on model calibration). The Hotelling's T^2 for the n -th sample is computed according to the equation (Mardia *et al.*, 1979):

$$T_n^2 = \mathbf{t}_n^T \mathbf{\Lambda} \mathbf{t}_n = \sum_{a=1}^A \frac{t_{a,n}^2}{\lambda_a} \quad (2.14)$$

where \mathbf{t}_n is the score vector of the n -th observation, λ_a is the a -th eigenvalue and $\mathbf{\Lambda}$ is the $[A \times A]$ matrix collecting the A eigenvalues on its diagonal.

The squared prediction error (SPE) is a metric that describes the mismatch between the value of a sample \mathbf{x}_n and its model representation $\hat{\mathbf{x}}_n$. Geometrically, it represents the squared orthogonal (Euclidean) distance between the projection of the n -th observation and the latent space, and it is computed as:

$$SPE_n = (\mathbf{x}_n - \hat{\mathbf{x}}_n)^T (\mathbf{x}_n - \hat{\mathbf{x}}_n) = \mathbf{e}_n^T \mathbf{e}_n \quad (2.15)$$

where \mathbf{e}_n is the residual vector for the n -th observation. Samples with high values of SPE_n have a different correlation structure with respect to the one captured by the PLS model, and therefore are not described properly by the model.

As regard *variable diagnostics*, it is often desirable to understand which regressor variables are most influential with respect to the model responses and therefore mainly affect the PLS model. This can be quantified with a metric called VIP index (variable importance in the projection; Chong and Jun, 2005), that is given by:

$$VIP_v = \sqrt{V \frac{\sum_{a=1}^A R_{y,a}^2 (w_{v,a})^2}{\sum_{a=1}^A R_{y,a}^2}} \quad (2.16)$$

where $R_{y,a}^2$ is the amount of \mathbf{y} -variance explained by the a -th LV and $w_{v,a}$ is the weight of the v -th input variable on the a -th LV. The higher the value of VIP_v , the higher the influence of the v -th input variable on the PLS model. Typically, variables with values of VIP_v greater than one are considered as valuable predictor for the model responses (Eriksson *et al.*, 2001).

2.2 Prediction uncertainty in PLS modeling

As stated in the previous section, once a PLS model has been calibrated according to Eq. (2.1)-(2.3), a new response $\hat{\mathbf{y}}_p$ can be predicted given a set of regressors \mathbf{x}_p (i.e. the PLS model can be used for regression analysis). The model prediction $\hat{\mathbf{y}}_p$ can be obtained according to:

$$\hat{\mathbf{y}}_p = \mathbf{t}_p \mathbf{Q}^T \quad (2.17)$$

where \mathbf{t}_p is the score vector related to the new set of regressors \mathbf{x}_p :

$$\mathbf{t}_p = \frac{\mathbf{x}_p \mathbf{W}}{\mathbf{p}^T \mathbf{W}}. \quad (2.18)$$

When using a PLS model to predict a new response, prediction uncertainty must be quantified in order to assess the predictive ability of the model. Different approaches have been proposed in the literature (a thorough review can be found in the work of Zhang and Garcia-Munoz, 2012). All these methods assume that a single univariate response must be predicted (i.e. $\hat{\mathbf{y}}_p = \hat{y}_p$) and rely on the interpretation of the PLS model as a standard linear regression model (details on the derivation of this formulation can be found in Wold *et al.*, 1983):

$$y_p = \mathbf{x}_p^T \boldsymbol{\beta} + \epsilon_p. \quad (2.19)$$

The general procedure to estimate prediction uncertainty is composed by the following two steps:

1. **Step #1:** an estimation of the standard deviation of the prediction error s is obtained and the degrees of freedom df of the model are computed.
2. **Step #2:** A confidence interval for the model prediction \hat{y}_p is established assuming a t -statistics distribution with $(N - df)$ degrees of freedom and significance level α for the prediction error:

$$CI = \hat{y}_p \pm t_{\frac{\alpha}{2}, (N-df)} s. \quad (2.20)$$

The sources of uncertainty that can affect the model prediction can be classified into three categories:

1. Measurement errors in both the input and output calibration dataset;
2. uncertainty in the model parameters ;
3. structural uncertainty due to the un-modeled part of \hat{y}_p (e.g. due to nonlinearities of the original dataset that cannot be captured with the linear PLS model).

The available approaches typically consider only the second and third sources of uncertainty, i.e. they describe how the uncertainty on the model parameters propagates to the model response and accounts for structural uncertainty by assuming an error distribution for the model residuals. These uncertainty propagation models can be classified in four categories: Ordinary-Least-Squares (OLS)-type methods, linearization-based methods, resampling-based methods and the Unscrambler method. A brief review of each of these approaches is presented in the following.

2.2.1 OLS-type methods

The starting point of OLS-type methods is the OLS expression for the vector of model parameters:

$$\widehat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (2.21)$$

where \mathbf{y} is the univariate output calibration dataset. By taking the covariance of both sides of Eq. (2.19), it follows:

$$\text{cov}(\widehat{\boldsymbol{\beta}}_{OLS}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{cov}(\mathbf{u}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \quad (2.22)$$

where $\text{cov}(\mathbf{u}) = \sigma^2 \mathbf{I}_N$, with σ^2 variance of the model residuals. The combination of (2.19) and (2.22) leads to:

$$\text{var}(\hat{y}_p) = \text{var}(\mathbf{x}_p^T \widehat{\boldsymbol{\beta}}_{OLS}) + \text{var}(\epsilon_p) = \mathbf{x}_p^T \text{cov}(\widehat{\boldsymbol{\beta}}_{OLS}) \mathbf{x}_p + \sigma^2 = \sigma^2 (h_p + 1) \quad (2.23)$$

where h_p is defined as the leverage of the p observation and is given by:

$$h_p = \mathbf{x}_p^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_p. \quad (2.24)$$

The geometrical meaning of h_p corresponds to the distance of the new observation with respect to the calibration center (i.e. to the origin of the axis of the PLS model). The standard deviation of the prediction error s can thus be obtained from Eq. (2.25):

$$s = \sigma \sqrt{1 + h_p}. \quad (2.25)$$

Eq. (2.25), even though derived for an OLS estimate of the model parameters for a standard linear regression model, can be transposed to a PLS model in a straightforward manner. Mathematical details can be found in the work of Zhang and Garcia-Munoz (2009). The same formulation of Eq. (2.17) can be used to obtain the standard deviation of the prediction error, where the leverage h_p can be derived from the score vector \mathbf{t}_p of the new observation as in the following:

$$h_p = \mathbf{t}_p^T \mathbf{t}_p. \quad (2.26)$$

Eq. (2.17) can be used to build the confidence interval for the new prediction of the PLS model, according to step#2 of the methodology described above.

Faber and Kowalski (1997) proposed an extension to Eq. (2.25) that accounts for the measurement errors in both the input and response matrices. The main limitation of this approach is that it requires an estimate of the error variance of the input and output calibration datasets from replicated experiments, thus limiting its practical implementation. A thorough mathematical description of this approach can be found in the cited reference.

2.2.2 Linearization-based methods

The derivation of prediction uncertainty with these methods is obtained by first-order linearization of the nonlinear dependency of the model parameters $\hat{\boldsymbol{\beta}}$ with respect to the model response y . The dependency of $\hat{\boldsymbol{\beta}}$ on the model response y is obtained by Taylor series expansion truncated after the first term:

$$\hat{\boldsymbol{\beta}}(y) = \hat{\boldsymbol{\beta}}(y_p) + \mathbf{J}_p(y - y_p). \quad (2.27)$$

where \mathbf{J}_p is the Jacobian matrix of the derivatives of each element of $\hat{\boldsymbol{\beta}}$ with respect to y computed at $y = y_p$. By taking the covariance of both sides of Eq. (2.27), an estimate of the uncertainty on the model parameters can be obtained:

$$\text{cov}(\hat{\boldsymbol{\beta}}) \simeq \mathbf{J}_p \mathbf{J}_p^T \sigma^2. \quad (2.28)$$

By plugging Eq. (2.28) into Eq. (2.23), an estimate of the prediction uncertainty can be obtained according to:

$$\text{var}(\hat{y}_p) = \sigma^2(1 + \mathbf{x}_p^T \mathbf{J}_p \mathbf{J}_p^T \mathbf{x}_p). \quad (2.29)$$

Prediction uncertainty (2.29) can be computed once an appropriate estimate of the covariance \mathbf{J}_p is obtained. Different methods have been proposed, based on differential calculus (Phatak *et al.*, 1993) or inductive estimation algorithms (Denham, 1997 ; Seernels *et al.*, 2004). A recent work on this topic has been presented by Zhang and Fearn (2015).

2.2.3 Re-sampling based methods

The underlying idea behind these methods is to build synthetic datasets from the original calibration dataset by introducing artificial perturbations and to estimate the covariance of the model parameters based on these datasets. The two most important methods are jack-knife and bootstrap (Zhang and Garcia-Munoz, 2012).

Jack-knife consists of generating K datasets (\mathbf{X}_k, y_k) , $k = 1, 2, \dots, K$ by deleting one calibration sample at a time from the original dataset and estimating the regression coefficient $\hat{\boldsymbol{\beta}}_k$ for each reduced dataset. A set of ‘‘pseudo’’ regression coefficients $\hat{\boldsymbol{\beta}}_{ps,k}$ are then obtained according to the following equation:

$$\hat{\boldsymbol{\beta}}_{ps,k} = N\hat{\boldsymbol{\beta}} - (N - 1)\hat{\boldsymbol{\beta}}_k; \quad k = 1, \dots, K \quad (2.30)$$

where $\hat{\boldsymbol{\beta}}$ is the vector of the model parameters obtained with the entire calibration dataset. From Eq. (2.30) and after small algebraic manipulations, the covariance of $\hat{\boldsymbol{\beta}}$ can be obtained according to:

$$\text{cov}(\hat{\boldsymbol{\beta}}) = \frac{1}{N(N-1)} \sum_{i=1}^K (\hat{\boldsymbol{\beta}}_{ps,k} - \bar{\boldsymbol{\beta}}_{ps}) (\hat{\boldsymbol{\beta}}_{ps,k} - \bar{\boldsymbol{\beta}}_{ps})^T \quad (2.31)$$

where $\bar{\boldsymbol{\beta}}_{ps}$ is the vector collecting the average values of the pseudo parameters. The combination of Eq. (2.31) with Eq. (2.23) yields the estimation of prediction uncertainty using a jack-knife approach:

$$\text{var}(\hat{y}_p) = \mathbf{x}_p^T \frac{1}{N(N-1)} \sum_{i=1}^K (\hat{\boldsymbol{\beta}}_{ps,k} - \bar{\boldsymbol{\beta}}_{ps}) (\hat{\boldsymbol{\beta}}_{ps,k} - \bar{\boldsymbol{\beta}}_{ps})^T \mathbf{x}_p + \sigma^2. \quad (2.32)$$

As regard bootstrap techniques, two main approaches to estimate prediction uncertainty can be found in the literature, namely bootstrap by residuals and bootstrap by objects.

In the former approach, new residual vectors are randomly generated by drawing the original residuals with replacements and new datasets are obtained by adding these new residuals to the fitted response \hat{y}_p . The variance of the regression coefficients vector can then be derived under the assumption of independent and identically distributed re-sampled datasets. The procedure is as follows.

First, the real residuals are computed according to :

$$\epsilon_n = \frac{y_{cal,n} - \hat{y}_{cal,n}}{\left(1 - \frac{df}{N}\right)^{\frac{1}{2}}}; n = 1, 2, \dots, N \quad (2.33)$$

where $\hat{y}_{cal,n}$ is the n -th PLS fitted response for the n -th calibration sample. The new residual vectors and the corresponding new response datasets are then computed according to the set of equations:

$$\psi_n^b = \text{round} [U(0,1) \cdot N] + 1; n = 1, 2, \dots, N; b = 1, 2, \dots, B \quad (2.34)$$

$$y_n^b = \hat{y}_{cal,n} + \epsilon_{\psi_n^b}; n = 1, 2, \dots, N; b = 1, 2, \dots, B \quad (2.35)$$

where $\text{round} [\cdot]$ represents the rounding to the nearest integer number towards zero and $U(0,1)$ is a random number drawn from a uniform distribution between 0 and 1. The bootstrap procedure is repeated B times, thus generating B re-sampled datasets. The value of B should be

chosen large enough to obtain accurate estimate of the covariance of the model parameters: typical values are in the range $B = [100 \div 10000]$.

The set of B re-sampled datasets can therefore be mathematically summarized as:

$$(\mathbf{x}_n^b, \mathbf{y}_n^b) = (\mathbf{x}_{\text{cal},n}, \mathbf{y}_n^b), \quad n = 1, 2, \dots, N; \quad b = 1, 2, \dots, B \quad (2.36)$$

$$\mathbf{X}_b = (\mathbf{x}_1^b, \dots, \mathbf{x}_N^b)^T, \quad b = 1, 2, \dots, B \quad (2.37)$$

$$\mathbf{y}_b = (\mathbf{y}_1^b, \dots, \mathbf{y}_N^b)^T, \quad b = 1, 2, \dots, B \quad (2.38)$$

with $\mathbf{x}_{\text{cal},n}$ input vector for the n -th calibration sample. For each resampled dataset $[\mathbf{X}_b; \mathbf{y}_b]$, the corresponding vector of model parameters $\hat{\boldsymbol{\beta}}_b$ can be obtained by fitting the model responses of the given dataset. Once iterated for all the B datasets, the bootstrap estimated of the covariance of the vector of model parameters $\hat{\boldsymbol{\beta}}$ can be obtained according to:

$$\text{cov}(\hat{\boldsymbol{\beta}}) = \frac{1}{B-1} \sum_{b=1}^B (\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}})(\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}})^T \quad (2.39)$$

where $\bar{\boldsymbol{\beta}}$ is the average of the B regression vectors $\hat{\boldsymbol{\beta}}_b$. By plugging in Eq. (2.39) into Eq. (2.23) the estimation of prediction uncertainty can be obtained:

$$\text{var}(\hat{y}_p) = \mathbf{x}_p^T \frac{1}{B-1} \sum_{b=1}^B (\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}})(\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}})^T \mathbf{x}_p + \sigma^2. \quad (2.40)$$

Eq. (2.40) is the mathematical formulation of the variance of a new PLS response obtained with a bootstrap of residuals method.

In the bootstrap by objects approach, the procedure is the same as in the bootstrap by residuals, with the difference that the B re-sampled datasets are generated by randomly drawing sample (and not residuals) with replacements. The covariance of the model parameters vector is obtained as in Eq. (2.39), and prediction uncertainty has the exact same formulation as in Eq. (2.40).

2.2.4 The Unscrambler method

The Unscrambler formula (De Vries and Ter Braak, 1995) for the estimation of prediction uncertainty in PLS modeling is an empirical formula that is implemented in the chemometrics commercial software Unscrambler™. Its derivation has been thoroughly reviewed in the literature (De Vries and Ter Braak, 1995; Høy *et al.*, 1998). Given the formulation of the PLS prediction model (2.17), the underlying idea behind this approach is to derive two expressions for prediction uncertainty assuming first that the PLS scores are not affected by uncertainty,

and then that the model loadings are not affected by uncertainty. An average of the two expressions is then obtained leading the following estimate of prediction uncertainty:

$$\text{var}(\hat{y}_p) = V_{y_{val}} \left(1 - \frac{A+1}{N}\right) \left(h_p + \frac{V_{x_p}}{V_{x_{tot, val}}}\right). \quad (2.41)$$

where $V_{y_{val}}$ is the y -residual variance in the validation dataset, $V_{x_{tot, val}}$ is the average \mathbf{x} -residual variance in the validation dataset and V_{x_p} is the \mathbf{x} -residual variance in the new regression vector \mathbf{x}_p . Details on the mathematical derivation of Eq. (2.41) are reported in the cited references.

2.2.5 Estimation of σ^2 and degrees of freedom

The estimation of prediction uncertainty according to OLS-type methods (2.23), linearization-based methods (2.29), jack-knife approach (2.32) and bootstrap methods (2.40) requires an estimation of the standard deviation of the model residuals σ .

An unbiased estimator of σ that is typically used in the chemometrics literature is the root means squared error of calibration *RMSEC*:

$$RMSEC = \left(\frac{\sum_{n=1}^N (y_{cal,n} - \hat{y}_{cal,n})^2}{N - df} \right) \quad (2.42)$$

where $\hat{y}_{cal,n}$ is the PLS prediction of the n -th calibration object. *RMSEC* is used in place of σ to estimate prediction uncertainty according to the methods described in the previous sections. However, the use of *RMSEC* according to Eq. (2.42) poses the problem of obtaining an accurate estimation of the degrees of freedom of the model¹³. Three approaches have been proposed to estimate df (Zhang and Garcia-Munoz, 2009):

1. Naïve approach: each PLS factor (i.e. LV) is assumed to consume one degree of freedom;
2. Pseudo degrees of freedom (PDF) approach: it computes df based on the ratio of a model fit error and a predictive performance error (derivation can be found in the work of Zhang and Garcia-Munoz, 2009);
3. Generalized degrees of freedom (GDF) approach: it computes df based on the sum of the sensitivity of each fitted response to perturbations in the corresponding observed response.

In this study, the only approach that has been used to estimate df is the Naïve approach. Additional details are provided in Chapter 3.

¹³ Note that the estimation of the degrees of freedom is also intrinsically required in order to estimate prediction uncertainty in the bootstrap by objects/residuals approach. Moreover, the estimate of df is required in order to build the confidence interval for the model prediction (Step #2).

2.2.5.1 Summary

A schematic review of all the methods to estimate prediction uncertainty in PLS modeling described above is reported in Table 2.1. In the same table, the different types of model uncertainty that are taken into account by each method are also reported.

Table 2.1. Summary of the different methods to estimate prediction uncertainty in PLS modeling.

Method	Std. of prediction uncertainty $s = \sqrt{\text{var}(\hat{y}_p)}$	Measurement uncertainty on X and y?	Uncertainty on model parameters?	Structural uncertainty?
OLS-type	$s = RMSEC \sqrt{1 + h_p}$	X	✓	✓
Linearization	$s = RMSEC \sqrt{(1 + \mathbf{x}_p^T \mathbf{J}_p \mathbf{J}_p^T \mathbf{x}_p)}$	X	✓	✓
Jack-knife	$s = \sqrt{\mathbf{x}_p^T \frac{1}{N(N-1)} \sum_{i=1}^K (\hat{\boldsymbol{\beta}}_{ps,k} - \bar{\boldsymbol{\beta}}_{ps}) (\hat{\boldsymbol{\beta}}_{ps,k} - \bar{\boldsymbol{\beta}}_{ps})^T \mathbf{x}_p + RMSEC^2}$	X	✓	✓
Bootstrap residuals	$s = \sqrt{\mathbf{x}_p^T \frac{1}{B-1} \sum_{b=1}^B (\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}}) (\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}})^T \mathbf{x}_p + RMSEC^2}$	X	✓	✓
Bootstrap objects	$s = \sqrt{\mathbf{x}_p^T \frac{1}{B-1} \sum_{b=1}^B (\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}}) (\hat{\boldsymbol{\beta}}_b - \bar{\boldsymbol{\beta}})^T \mathbf{x}_p + RMSEC^2}$	X	✓	✓
Unscrambler	$s = \sqrt{V_{y, \text{val}} \left(1 - \frac{A+1}{N}\right) \left(h_p + \frac{V_{x_p}}{V_{x_{\text{tot}, \text{val}}}}\right)}$	X	✓	✓

SECTION B: BAYESIAN MODELING

Since the introduction of Bayes' theorem (Bayes, 1763), Bayesian statistical modeling has been considered as an alternative school of thought as opposed to classical frequentist statistical modeling. The two approaches rely on a completely different interpretation of the concept of probability: in frequentist statistics, probability is defined as the frequency that an event will occur over a large number of trials. On the other hand, in Bayesian statistics, probability is considered as the degree of plausibility for the occurrence of a given event.

In the process systems engineering community, Bayesian modeling has gained particular attention in the last few decades, thanks to the rapid increase of computational power for intensive numerical calculations. The ability to perform in a very short time Markov-chain Monte Carlo (MCMC) simulations, such as the ones required in Bayesian modeling, played a key role for the application of Bayesian methodologies in chemical engineering- related problems.

In this section, a brief mathematical background on Bayesian statistics is presented. A particular focus is given to Bayesian multivariate linear regression, since its practical application to problems related to the pharmaceutical industry will be presented in Chapter 4. A rigorous mathematical description of Bayesian modelling for chemical engineers can be found in the work of Lenk and DeSarbo (2000).

2.3 Key concepts of Bayesian statistics

The key assumption of Bayesian statistics is that the variables of interest are random variables described by probability density functions (PDFs). Let $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$ be a set of N observations of the multivariate responses $\mathbf{y}_m [N \times 1], m = 1, 2, \dots, M$. In the Bayesian framework, $\mathbf{y}_m, m = 1, 2, \dots, M$ are a set of random variables described by their respective probability distributions. The vector $\mathbf{y} = (y_1, y_2, \dots, y_M)$ that collects all the M response variables is therefore assumed to be a multivariate random variable with a given PDF. It is assumed that \mathbf{y} depends upon a set of parameters $\boldsymbol{\theta}$ according to a given pre-defined model (i.e. regression model, mechanistic model etc...). The PDF of \mathbf{y} , given the parameters $\boldsymbol{\theta}$ and any other pertinent information about the system I , is expressed as $p(\mathbf{y}|\boldsymbol{\theta}, I)$. I includes, for example, any relevant assumption that is made on the type and shape of the PDF of \mathbf{y} . The vector of parameters $\boldsymbol{\theta}$ is considered as a multivariate random variable with PDF $p(\boldsymbol{\theta}|\mathbf{y}, I)$. Any information known *a priori* on the values and uncertainty of the model parameters $\boldsymbol{\theta}$ is described by the PDF $p(\boldsymbol{\theta}|I)$. For sake of simplicity, I will be omitted in the following notation. Bayes' theorem can be formulated according to the following expression:

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int_{\boldsymbol{\theta}} p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}} \quad (2.43)$$

where $p(\boldsymbol{\theta})$ is called the *prior* (distribution) of $\boldsymbol{\theta}$, $p(\boldsymbol{\theta}|\mathbf{y})$ the *posterior* of $\boldsymbol{\theta}$ and $p(\mathbf{y}|\boldsymbol{\theta})$ is called the *likelihood function* and it is often written as $\mathcal{L}(\boldsymbol{\theta}|\mathbf{y})$. The prior of $\boldsymbol{\theta}$ contains all the information available a priori on $\boldsymbol{\theta}$; the likelihood function describes the plausibility of $\boldsymbol{\theta}$ given \mathbf{y} , while $p(\boldsymbol{\theta}|\mathbf{y})$ describes how $\boldsymbol{\theta}$ is distributed given the observed response data. The denominator of Eq. (2.43) is a constant normalizing factor that depends only on the data available. For this reason, Eq. (2.43) is often expressed as:

$$p(\boldsymbol{\theta}|\mathbf{y}) \sim \mathcal{L}(\boldsymbol{\theta}|\mathbf{y}) p(\boldsymbol{\theta}) . \quad (2.44)$$

Eq. (2.44) clarifies that the posterior distribution of the parameters $\boldsymbol{\theta}$, given the observed data, is obtained from the prior knowledge about $\boldsymbol{\theta}$, updated according to the likelihood function $\mathcal{L}(\boldsymbol{\theta}|\mathbf{y})$. The correct interpretation and application of Eq. (2.44) requires few comments, namely:

1. How the prior distribution $p(\boldsymbol{\theta})$ should be assigned;
2. how the final posterior $p(\boldsymbol{\theta}|\mathbf{y})$ can be computed from a practical viewpoint;
3. how the uncertainty on the prediction of a new response \hat{y}_p can be obtained in a Bayesian framework.

2.3.1 Prior distribution $p(\boldsymbol{\theta})$

The incorporation of the prior distribution $p(\boldsymbol{\theta})$ in the computation of the posterior of the model parameters is one of the key difference between Bayesian and frequentist parameter estimation. Every Bayesian computation requires an appropriate choice of the prior of $\boldsymbol{\theta}$: this is often considered as a strong argument against Bayesian modeling, since there is no general rule on the choice of $p(\boldsymbol{\theta})$. However, in many model-based engineering-related problems, the possibility of incorporating available knowledge on the system can be extremely beneficial to improve the accuracy of model predictions. In this sense, the above argument can be easily reverted and the incorporation of $p(\boldsymbol{\theta})$ in the analysis can be seen as one of the strengths of the Bayesian framework.

Although there is no general rule for the choice of $p(\boldsymbol{\theta})$, there are two different types of priors that can be used when performing a Bayesian simulation:

1. Non-informative prior distributions
2. Informative prior distributions.

Non-informative priors are used when very limited or no knowledge is available for the model parameters $\boldsymbol{\theta}$. In this case, the posterior $p(\boldsymbol{\theta}|\mathbf{y})$ is affected only by the likelihood $\mathcal{L}(\boldsymbol{\theta}|\mathbf{y})$. The typical choice of non-informative prior that can be made is a flat prior distribution for $\boldsymbol{\theta}$ over

its entire domain. However, it is worth noticing that there are other type of priors that, even though they can be classified as *informative* from a structural point of view, can be considered as *non-informative* when the values of some or all of their moments is not informative. A typical example is a Gaussian prior distribution (thus *informative* from a structural point of view) with a very large standard deviation and a random guess of its mean value. This type of prior distribution has to be considered non-informative, even though formally informative from a structural point of view. The reason why the choice of non-informative priors with informative structure can be useful in certain applications will be discussed in more detail in Chapter 4.

Informative priors can be used when general knowledge is available on θ , e.g. from previous experiments or from previous studies for the system under analysis. In this case, suitable PDF structures and values for their moments can be assigned for the prior $p(\theta)$. The definition of informative priors typically becomes more challenging when the dimensionality of θ increases (Coleman and Block, 2006). The problem can be typically solved by performing a sensitivity analysis of the posterior with respect to the conditional priors of the single parameters, and then obtaining a suitable joint prior for θ based on the results of this analysis.

2.3.2 Posterior distribution $p(\theta|\mathbf{y})$

Once the prior distribution for the model parameters has been set, the posterior $p(\theta|\mathbf{y})$ can be computed according to Bayes' theorem (2.44). The solution of (2.44) can be analytical or numerical. Obtaining an analytical expression for $p(\theta|\mathbf{y})$ is possible in very rare situations and under very strict simplifications (low dimensionality, linearity, simple priors etc...). In many engineering-related problems, the solution of (2.44) can only be obtained with numerical sampling methodologies such as Markov Chain Monte Carlo (MCMC) techniques. Different sampling algorithms have been proposed that are able to follow any joint posterior distribution of the model parameters. The most famous MCMC sampling algorithms are the Metropolis algorithm and the Metropolis-Hastings algorithm (Hastings, 1970), the modified Metropolis-Hastings algorithm and the Gibbs sampler (Geman, S. and Geman, D., 1993). The description of these algorithms is out of the scope of this Dissertation. A thorough description can be found in the cited references.

The application of these sampling strategies to obtain a discrete representation of $p(\theta|\mathbf{y})$ requires intensive computer simulations due to the (possibly) slow convergence of the Markov chains (MCs). The assessment of the convergence of the MCs and the computational requirements are thoroughly discussed in Chapter 4.

2.3.3 Posterior predictive distribution of a new response $\hat{\mathbf{y}}_p$

Once the parameters of the model have been calibrated in a Bayesian framework (i.e. their posterior distribution has been obtained), it is very often desirable to use that model to predict

a new response $\hat{\mathbf{y}}_p$ and to compute the uncertainty associated with that prediction. The strength of Bayesian modeling is that the value of a new prediction and its uncertainty are derived in a unified framework. In fact, as discussed in section 2.3, the model prediction $\hat{\mathbf{y}}_p$ is treated as random variable and it is therefore described by a PDF $p(\hat{\mathbf{y}}_p|\mathbf{y})$. The distribution $p(\hat{\mathbf{y}}_p|\mathbf{y})$ is called the *posterior predictive distribution* (PPD) of $\hat{\mathbf{y}}_p$. The PPD of $\hat{\mathbf{y}}_p$ can be formally computed as:

$$p(\hat{\mathbf{y}}_p|\mathbf{y}) = \int_{\boldsymbol{\theta}} p(\hat{\mathbf{y}}_p, \boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} p(\hat{\mathbf{y}}_p|\mathbf{y}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}. \quad (2.45)$$

Under the assumption that, for a given $\boldsymbol{\theta}$, $\hat{\mathbf{y}}_p$ and \mathbf{y} are conditionally independent, Eq. (2.45) can be simplified to:

$$p(\hat{\mathbf{y}}_p|\mathbf{y}) = \int_{\boldsymbol{\theta}} p(\hat{\mathbf{y}}_p|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}. \quad (2.46)$$

According to Eq. (2.45), the computation of the PPD of $\hat{\mathbf{y}}_p$ requires the posterior distribution $p(\boldsymbol{\theta}|\mathbf{y})$ of the model parameters, obtained according to Eq. (2.44), and the PDF $p(\hat{\mathbf{y}}_p|\boldsymbol{\theta})$, which can be obtained from the original model for the given values of $\boldsymbol{\theta}$. If analytical expressions for $p(\boldsymbol{\theta}|\mathbf{y})$ and $p(\hat{\mathbf{y}}_p|\boldsymbol{\theta})$ were available, Eq. (2.46) could be used to obtain an analytical expression for $p(\hat{\mathbf{y}}_p|\mathbf{y})$ by direct integration. However, as previously explained, the analytical expression of $p(\boldsymbol{\theta}|\mathbf{y})$ is almost never available and this PDF is usually obtained via MCMC sampling. For this reason, the PPD of $\hat{\mathbf{y}}_p$ can only be obtained by propagating the posterior PDF of the model parameters to the model responses, i.e. samples from the PPD can be obtained following the procedure:

1. Draw a sample $\boldsymbol{\theta}^{(l)}$ from the posterior distribution of the model parameters $p(\boldsymbol{\theta}|\mathbf{y})$;
2. Draw the respective model response $\hat{\mathbf{y}}_p^{(l)}$ from the model $p(\hat{\mathbf{y}}_p|\boldsymbol{\theta})$;
3. Repeat for $l = 1, 2, \dots, L$ times, with L number of samples (user-defined).

The application of the above procedure allows obtaining the PPD for a new response, with the notable drawback of increasing the total computational time (sampling from the PPDS sums up to the sampling from the posterior of the model parameters).

2.3.4 Credible region Vs Confidence region

The PPD of a new response that can be obtained from a Bayesian simulation carries information on both the predicted value of the new response vector $\hat{\mathbf{y}}_p$ and its uncertainty. In this regards, it can be said that prediction uncertainty is incorporated in a straightforward fashion in the Bayesian framework.

In many situations, it is desirable to express the uncertainty on the model predictions in terms of confidence intervals/regions around the expected value of $\hat{\mathbf{y}}_p$, similarly to what is done in classical frequentist modeling.

From a Bayesian perspective, the concepts of confidence interval and confidence region of frequentist statistics are replaced by the concepts of *credible* interval and *credible* region. Although the two definitions (confidence region/ credible region) may look similar at a first sight, they are completely different from both a methodological and significance point of view. In frequentist statistics, the predicted response $\hat{\mathbf{y}}_p$ is considered as a fixed value and the boundary of the M - dimensional confidence region as a random variable. On the other side, in Bayesian statistics, the predicted response $\hat{\mathbf{y}}_p$ is considered as a random variable and the boundary of the M - dimensional credible region as a random variable. In other terms, a frequentist $(1 - \alpha)\%$ confidence region means that if a large number of repeated samples from the original population is drawn, $(1-\alpha)\%$ of these samples would fall within the calculated confidence region, with $\alpha =$ user-defined significance level. On the other hand, a δ (%) credible interval for the random variable $\hat{\mathbf{y}}_p$ means that the probability that $\hat{\mathbf{y}}_p$ lies within the credible region is equal to δ .

From a mathematical point of view, a δ (%) credible region C_δ for the random variable $\hat{\mathbf{y}}_p$ can be defined as:

$$C_\delta: \int_{C_\delta} p(\hat{\mathbf{y}}_p | \mathbf{y}) d\hat{\mathbf{y}}_p = \delta \quad (2.47)$$

The definition of credible region (2.47) gives a measure of prediction uncertainty that can be obtained in a straightforward way for both univariate and multivariate responses. This last situation is of particular interest in engineering-related problems and gives a great advantage with respect to classical frequentist statistics, where the definition of multidimensional confidence regions for multivariate response vectors is often very complicated or obtained under strong and limiting assumptions.

2.4 Bayesian multivariate linear regression

As observed in section 2.1, many applications involve the availability of a historical dataset of model responses \mathbf{Y} [$N \times M$], that has to be correlated with a historical dataset of model inputs (or regressors) \mathbf{X} [$N \times V$]. The simplest model that can be used to correlated the input dataset with the output dasate is a linear multivariate regression model:

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E} \quad (2.48)$$

where $\mathbf{B} [V \times M]$ is the matrix of model parameters, and $\mathbf{E} [N \times M]$ is the matrix of model residuals. The model residuals are assumed to have $\mathbf{0}$ - mean and characterized by a semi-definite positive covariance matrix $\mathbf{\Sigma} [N \times M]$, equal for every model residual $\mathbf{e}_m, m = 1, 2, \dots, M$.

The Bayesian calibration of the regression model (2.48) can be obtained by sampling from the joint multivariate posterior PDF of \mathbf{B} according to Eq. (2.44). The calibrated regression model can then be used to predict a new response $\hat{\mathbf{y}}_p$ and characterize prediction uncertainty according to Eq. (2.46). The simple model structure of Eq. (2.48) allows obtaining a deeper insight on the likelihood function and posterior distribution of \mathbf{B} . The main characteristics are briefly described in the following sections.

2.4.1 Likelihood function in multivariate linear regression

The linear regression model (2.48) allows obtaining a simple analytical expression for the likelihood function $\mathcal{L}(\mathbf{B}, \mathbf{\Sigma} | \mathbf{y})$ under the assumption that the model residual are independent identically normal distributed with mean $\mathbf{0}$ and covariance $\mathbf{\Sigma}$. From the properties of normal distributions, it follows that each response vector $\mathbf{y}_m, m = 1, 2, \dots, M$ is assumed to follow a multivariate normal distribution according to:

$$\mathbf{y}_n \sim N_m(\mathbf{x}_n \mathbf{B}, \mathbf{\Sigma}), \quad n = 1, 2, \dots, N. \quad (2.49)$$

Under this assumption, the likelihood function $\mathcal{L}(\mathbf{B}, \mathbf{\Sigma} | \mathbf{y})$ can be analytically expressed as:

$$\mathcal{L}(\mathbf{B}, \mathbf{\Sigma} | \mathbf{y}) = (2\pi)^{-\frac{MN}{2}} \det(\mathbf{\Sigma})^{-\frac{N}{2}} \exp\left(-\frac{1}{2} \sum_{n=1}^N [(\mathbf{y}_n - \mathbf{x}_n \mathbf{B}) \mathbf{\Sigma}^{-1} (\mathbf{y}_n - \mathbf{x}_n \mathbf{B})^T]\right) \quad (2.50)$$

or more conveniently as:

$$\mathcal{L}(\mathbf{B}, \mathbf{\Sigma} | \mathbf{y}) \sim \det(\mathbf{\Sigma})^{-\frac{N}{2}} \exp\left(-\frac{1}{2} \text{tr} [\mathbf{\Sigma}^{-1} (\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB})]\right) \quad (2.51)$$

Expression (2.51), combined with any choice of the prior distributions of the model parameters and the covariance of the residuals, allows obtaining the joint posterior of the model parameters and the covariance of the residuals. A possible set of prior distributions that can be assigned is discussed in the next section.

2.4.2 Prior distributions

Depending on the choice of the joint prior distribution for the model parameters and the residuals covariance (informative/non-informative, as discussed in section 2.3.1), the joint posterior can be obtained by updating the given prior through the likelihood function (2.51).

As regard non-informative prior distributions, different choices have been proposed in the literature (see e.g. Geisser and Cornfield (1963)). A typical choice is that of assuming \mathbf{B} and Σ completely independent to each other and assigning a flat distribution for \mathbf{B} (i.e. a non-informative prior) and an inverse Wishart distribution for Σ (Wishart, 1958).

A common choice that is adopted is to use a structural informative prior distribution for \mathbf{B} , and assigning informative or non-informative values to the moments of this distribution according to the previous knowledge available on the model parameters. Typically, a $(V \times M)$ matrix-variate normal distribution is assigned to the conditional distribution of \mathbf{B} with respect to Σ , and a inverse Wishart distribution is assigned as the prior of Σ . To understand the reasoning behind this choice, it is worth noticing that the joint prior distribution of \mathbf{B} and Σ ($p(\mathbf{B}, \Sigma)$) can be decomposed as:

$$p(\mathbf{B}, \Sigma) = p(\mathbf{B}|\Sigma)p(\Sigma) \quad (2.52)$$

where $p(\mathbf{B}|\Sigma)$ is the conditional distribution of \mathbf{B} with respect to Σ , and $p(\Sigma)$ is the prior distribution of Σ . Based on Eq. (2.52), the choice of a matrix-variate normal distribution for $p(\mathbf{B}|\Sigma)$ and an inverse Wishart distribution for $p(\Sigma)$ is related to the fact that they both represent conjugate prior distributions (i.e. a joint Normal-Wishart prior distribution for $(\mathbf{B}, \Sigma^{-1})$ will generate a Normal-Wishart posterior distribution for $p(\mathbf{B}, \Sigma^{-1})$). Therefore, $p(\mathbf{B}|\Sigma)$ can be expressed as:

$$p(\mathbf{B}|\Sigma) \sim N_{V \times M}(\mathbf{B}_0, \Sigma, \Sigma_0) \quad (2.53)$$

where \mathbf{B}_0 and Σ_0 are the initial guesses for the model parameters and residuals covariance respectively (i.e. the values that the user can set in order to make more or less informative the prior distribution). On the other hand, an inverse Wishart distribution is assigned to $p(\Sigma)$:

$$p(\Sigma) \sim W^{-1}(\mathbf{\Omega}, \nu_0) \quad (2.54)$$

where $\mathbf{\Omega}$ is called the *a priori* response scale matrix and can be interpreted as a sum of squared errors, and ν_0 (user-defined) is the number of degrees of freedom of the prior distribution. The lower ν_0 , the lower the subjectivity of the prior distribution. The analytical expression of the inverse Wishart distribution can be found in the work of Wishart (1958).

Eq (2.54) and (2.53) leads to the joint prior distribution according to Eq. (2.52). The above priors have been set in all the work presented in this Dissertation and additional details on this topic will be given in Chapter 4.

2.4.3 Posterior distribution

By combining Eq. (2.54) with (2.53) through Eq. (2.52), and combining them with Eq. (2.51), the joint posterior distribution of \mathbf{B} and $\mathbf{\Sigma}$ can be obtained from Eq. (2.44) according to the following expression:

$$\begin{aligned}
 p(\mathbf{B}, \mathbf{\Sigma} | \mathbf{X}, \mathbf{Y}) &\sim \mathcal{L}(\mathbf{B}, \mathbf{\Sigma} | \mathbf{y}) p(\mathbf{B} | \mathbf{\Sigma}) p(\mathbf{\Sigma}) \sim \\
 &\sim \det(\mathbf{\Sigma})^{-\frac{N}{2}} \exp\left(-\frac{1}{2} \text{tr}[\mathbf{\Sigma}^{-1}(\mathbf{Y} - \mathbf{XB})^T(\mathbf{Y} - \mathbf{XB})]\right) \cdot \\
 &\cdot \det(\mathbf{\Sigma})^{-\frac{V}{2}} \exp\left(-\frac{1}{2} \text{tr}[\mathbf{\Sigma}^{-1}(\mathbf{B} - \mathbf{B}_0)^T \mathbf{\Sigma}_0^{-1}(\mathbf{B} - \mathbf{B}_0)]\right) \det(\mathbf{\Sigma})^{-\frac{v_0 - 2M}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{\Omega} \mathbf{\Sigma}^{-1})\right)
 \end{aligned} \tag{2.55}$$

Sampling from the above distribution can be obtained with one of the MCMC sampling strategies discussed in section 2.3.2. The obtained PDF can then be obtained to build the PPD of a new response according to the procedure described in the next section.

2.4.4 Posterior predictive distribution of a new response

The PPD for a new response can be obtained by re-writing Eq. (2.45) adapted to the multivariate linear regression scenario, i.e.:

$$p(\hat{\mathbf{y}}_p | \mathbf{X}, \mathbf{Y}, \mathbf{x}_p) = \int_{\mathbf{B}} \int_{\mathbf{\Sigma}} p(\hat{\mathbf{y}}_p | \mathbf{B}, \mathbf{\Sigma}) p(\mathbf{B}, \mathbf{\Sigma}) d\mathbf{B} d\mathbf{\Sigma} \tag{2.56}$$

Depending to the type of priors chosen (non-informative or informative), it may be possible to give an analytical expression to Eq. (2.56). For example, if the structurally informative priors described in section 2.4.2 are used, the corresponding PPD can be obtained by plugging in Eq. (2.55) into Eq. (2.56) and Eq. (2.49) adapted to the new response vector $\hat{\mathbf{y}}_p$. In any case, the PPD is typically obtained with the sampling procedure described in section 2.3.2.

SECTION C: FEASIBILITY ANALYSIS

The concept of feasibility and flexibility analysis applied to chemical engineering problems was first introduced in the '80s by the joint studies of Morari (Lenhoff and Morari, 1982) and Grossmann (Halemane and Grossman, 1983). Although the definitions of feasibility and flexibility analysis are different from each other, they are often interchanged with a slight misuse of terminology.

Feasibility analysis describes the ability of a process to satisfy all the relevant constraints (quality constraints, production constraints, environmental constraints) in the presence of uncertainty on (i) model parameters, (ii) raw material properties (*external variability*) and (iii) critical process parameters (*external variability*). These three sources of uncertainty are typically identified as *uncertain parameters*, and the combination of these parameters that satisfy all the relevant constraints is defined as *feasible region*. The identification of the feasible region can be obtained by solving a *feasibility test problem*. Mathematically, the variables that appear in a feasibility test problem formulation can be classified as:

1. design variables \mathbf{d} [$1 \times D$]: variables related to the structure and equipment size;
2. control variables \mathbf{u} [$1 \times U$]: manipulated variables of the system;
3. uncertain parameters $\boldsymbol{\chi}$ [$1 \times C$], as defined above.

If the process is at steady-state, the set of J relevant constraints imposed on the process under investigation can be mathematically written in forms of inequalities of the type:

$$f_j(\mathbf{d}, \boldsymbol{\chi}, \mathbf{u}) \leq 0, \quad j = 1, 2, \dots, J. \quad (2.57)$$

If it is assumed that there are no control variables in the system, the maximum values of all these constraints defines the so-called feasibility function $\Psi(\mathbf{d}, \boldsymbol{\chi})$:

$$\Psi(\mathbf{d}, \boldsymbol{\chi}) = \max_j f_j(\mathbf{d}, \boldsymbol{\chi}) \quad (2.58)$$

$$\text{s.t. } \boldsymbol{\chi} \in [\boldsymbol{\chi}^{\min}; \boldsymbol{\chi}^{\max}].$$

If the process design is fixed, the feasibility function has an immediate interpretation: if $\Psi(\mathbf{d}, \boldsymbol{\chi}) < 0$, the process can be feasibly operated for the given uncertain parameters $\boldsymbol{\chi}$; if $\Psi(\mathbf{d}, \boldsymbol{\chi}) > 0$, the process is infeasible. It follows that the feasible region can be defined as the combinations of uncertain parameters that satisfy $\Psi(\mathbf{d}, \boldsymbol{\chi}) < 0$; in particular, the combinations of uncertain parameters that satisfy the condition $\Psi(\mathbf{d}, \boldsymbol{\chi}) = 0$ represent the boundary of the feasible region.

The solution of the feasibility problem (2.58) can be obtained by solving the equivalent optimization problem:

$$\begin{aligned} \Psi(\mathbf{d}, \boldsymbol{\chi}) &= \min_p p \\ \text{s. t. } f_j(\mathbf{d}, \boldsymbol{\chi}) &\leq p, \quad j = 1, 2, \dots, J. \end{aligned} \quad (2.59)$$

The solution of the optimization problem (2.59) is straightforward if the process constraints are differential closed-form constraints and their evaluation is computationally tractable. However, many applications involve black-box constraints and the evaluation of the original process model can be extremely demanding from a computational point of view. For this reason, different surrogate-based approaches have been developed to overcome these limitations (Bhosekar and Ierapetritou, 2017). These methods build a surrogate as a cheap and reliable approximation of the original process model, and compute the feasible region based on this surrogate rather than the original process model. The main concepts of surrogate-based feasibility analysis are briefly summarized below. An extension of feasibility analysis applied to dynamical systems is reported in Chapter 6.

2.5 Surrogate-based feasibility analysis

The key idea of surrogate-based feasibility analysis is to build a surrogate as an approximation of the original process model, and use this surrogate to solve the feasibility problem in order to identify the feasible region of the process. Clearly, the approximation error that is introduced with the surrogate must be kept to a minimum, in order to avoid a wrong representation of the feasible region with respect to the one that can be obtained with the original process model. Surrogate-based feasibility analysis consists of different steps, namely:

1. initial sampling of the input domain (i.e. uncertain parameters domain);
2. evaluation of the feasibility function for each of the initial samples obtained at step 1;
3. determination of a response surface for the feasibility function;
4. continuous update of the surrogate model with adaptive sampling strategies until the surrogate accuracy is deemed to be sufficient;
5. prediction of the feasible region boundary using the surrogate obtained after the adaptive sampling strategy.

Different choices can be used when applying this methodology, including the choice of the surrogate model, the type of adaptive sampling strategy adopted and the stopping criterion for the surrogate accuracy. A brief of each of these aspects is reported below.

2.5.1 Surrogate models

Different surrogate models have been proposed in the literature in order to perform a feasibility analysis exercise, including high-dimensional model representations (Banerjee and Ierapetritou, 2002; Banerjee *et al.*, 2010), Kriging surrogates (Boukouvala and Ierapetritou,

2012; Boukouvala and Ierapetritou, 2013) and Radial-Basis function (RBF) surrogates (Wang and Ierapetritou, 2017). Kriging surrogates and RBF surrogates have been proven to be more efficient than high-dimensional surrogates, while maintaining the computational burden relatively low.

2.5.1.1 Kriging surrogate

Kriging (Krige, 1951) is a well-established interpolation method that has proven to work efficiently as an approximation of many nonlinear dynamic models. The basic idea of Kriging is that the function value for an unexplored sampling point can be obtained as a weighted sum of the function values at known sampling locations. Given S design sites, the function value at the unknown sampling point \mathbf{x}^* can be obtained from the function values at the known design sites $\mathbf{x}_s, s = 1, 2, \dots, S$ according to:

$$f(\mathbf{x}^*) = \sum_{s=1}^S w_s f(\mathbf{x}_s) + \epsilon^* \quad (2.60)$$

where ϵ^* is the variance of the Kriging predictor at the sampling point \mathbf{x}^* . Eq. (2.60) can be decomposed into a regression model and a correlation model, according to the structure:

$$f(\mathbf{x}^*) = \mathbf{f}(\mathbf{x}_s)^T \boldsymbol{\gamma} + \mathbf{r}(\mathbf{x}_s)^T \boldsymbol{\delta} \quad (2.61)$$

where $\mathbf{f}(\mathbf{x}_s)^T$ is the vector collecting the function values at the S design sites and $\mathbf{r}(\mathbf{x}_s)^T$ is a vector collecting the correlations between the function values at the S design sites.

Different types of regression models and correlation models can be used in Eq. (2.61), thus leading to different types of Kriging estimators. The regression models are typically polynomial, e.g. zero-*th* order, first-order, second-order polynomials and so on. The correlation models can be of different types, including linear, exponential and Gaussian correlation models. A thorough discussion on the different correlation models of the Kriging estimator can be found in the work of Kleijnen (Kleijnen, 2009).

2.5.1.2 Radial-basis function (RBF) surrogate

RBF surrogates for feasibility analysis have recently gained attention, since it has been proven that they can outperform Kriging surrogates in many relevant situations (Wang and Ierapetritou, 2017). Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_S$ be the usual known set of sampling points and $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_S)$ the respective function values as for the Kriging surrogate. Mathematically, a RBF surrogate can be expressed with the formulation:

$$s_s(\mathbf{x}) = \sum_{s=1}^S \lambda_s \phi \|\mathbf{x} - \mathbf{x}_s\|_2 + \mathbf{b}^T \mathbf{x} + a \quad (2.62)$$

where $\|\cdot\|_2$ is the Euclidean distance, $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_S]$, \mathbf{b} and a are model coefficients that can be obtained by solving the system of equations:

$$\begin{pmatrix} \boldsymbol{\Phi} & \mathbf{S} \\ \mathbf{S}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \boldsymbol{\lambda} \\ \mathbf{c} \end{pmatrix} = \begin{pmatrix} \mathbf{F} \\ \mathbf{0} \end{pmatrix} \quad (2.63)$$

with:

$$\mathbf{S} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & 1 \\ \mathbf{x}_S^T & 1 \end{pmatrix}; \boldsymbol{\lambda} = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_S \end{pmatrix}; \mathbf{c} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_Q \\ a \end{pmatrix}; \mathbf{F} = \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_S) \end{pmatrix}. \quad (2.64)$$

The basis function ϕ can be chosen depending on the problem considered: common choices are linear, cubic and Gaussian basis functions. Additional details on this aspect can be found in the work of Forrester and Keane (2009).

2.5.2 Adaptive sampling

As explained in Section 2.5, the surrogate can be continuously updated using an adaptive sampling strategy until its accuracy reached a pre-defined threshold value. Adaptive sampling allows increasing the surrogate accuracy in contrast with the standard space-filling sampling algorithms. The basic idea of adaptive sampling for the identification of the feasible region is as follows. The new sampling points are chosen based on two criteria:

1. New points should be chosen near the boundary of the feasible region (*local search*);
2. New points should be chosen in unexplored portions of the uncertain parameters domain, i.e. where prediction uncertainty is higher (*global search*).

The balance of these local and global search criteria should allow a thorough exploration of the entire input domain *and* an accurate identification of the boundary of the feasible region. Mathematically, this can be obtained by maximizing at each iteration a modified expected improvement (EI) function (Wang and Ierapetritou, 2017):

$$\max_{\mathbf{x}} EI_{feas}(\mathbf{x}) = s \times \phi\left(\frac{-\hat{y}}{s}\right) = s \times \left(\frac{1}{\sqrt{2\pi}}\right) e^{-0.5\left(\frac{\hat{y}^2}{s^2}\right)} \quad (2.65)$$

where \hat{y} is the surrogate model prediction at \mathbf{x} and s is the standard error of prediction. If a Kriging surrogate is used, s is simply the square root of the estimated prediction variance. With

RBF surrogates, s can be estimated through an “error indicator” $1/\mu$ first proposed by Gutmann (2001) or with the modified error indicator ” $1/\mu_n$ proposed by Wang and Ierapetritou (2017). Details on the mathematical formulation of these error indicators can be found in the cited references.

2.5.3 Surrogate accuracy

Defining the accuracy of the results of a surrogate-based feasibility analysis exercise depend on the main purpose of the analysis and there are no general metrics that can be used for any situation. If the main purpose of surrogate-based feasibility analysis is that of determining with high accuracy the boundary of the feasible region of the system under investigation, there are three key requirements that can be used to assess surrogate accuracy:

1. The portion of *feasible* region that has been correctly identified by the surrogate is sufficiently high;
2. The portion of *infeasible* region that has been correctly identified by the surrogate is sufficiently high;
3. The surrogate does not considerably overestimate the actual feasible region of the system.

The third requirement is particularly useful when applied to pharmaceutical processes, as discussed in Chapter 5.

Several metrics have been proposed to assess the surrogate accuracy based on the following considerations (Rogers and Ierapetritou, 2015; Adi et al., 2016). However, all these metrics lack of meeting at least one of the above criteria. A new set of metrics that allows fulfilling all of the above requirements have been recently proposed by Wang and Ierapetritou (2017). A detailed discussion and extension of these metrics is reported in Chapter 5.

SECTION D: STATE ESTIMATION TECHNIQUES FOR ONLINE MODEL CALIBRATION

The main application of state estimators (or state *observers*) in the process industry is for process control, i.e. state estimators are used to continuously update the model-prediction of the current state of the system, given the availability of a limited number of online measurements from plant sensors. The up-to-date system state is then used to inform pre-defined control strategies (e.g. model predictive control) on the most suitable control action to be implemented at the given point in time.

Formally, the role of a state estimator can be briefly described as follows.

Let \mathfrak{N} be the system under analysis (e.g. single unit/entire manufacturing line of a pilot/full-scale plant). Let $\mathbf{f}(\mathbf{x}, \mathbf{u})$ be a mathematical model (data-driven and/or semi-empirical and/or mechanistic) that is used to describe the behavior of \mathfrak{N} . The vector $\mathbf{x}[1 \times Q]$ collects all the state variables of the system, while $\mathbf{u}[1 \times U]$ collects all the control variables. It is assumed that online measurements $\mathbf{z}[1 \times Z]$ are available from plant sensors. These measurements can be considered continuous in time $\mathbf{z}(t)$ or, in most practical situations, to be obtained at discrete points in time $\mathbf{z}_k, k = 1, 2, \dots, K$. In this Dissertation, the only scenario that will be considered is the one with discrete measurements. The following formulation is based on this underlying assumption.

A state estimator can be defined as an observer that allows updating the model-based estimate of the system state $\hat{\mathbf{x}}_k^-$ at time t_k , by accounting for the available measurements at the same time \mathbf{y}_k . The updated system state $\hat{\mathbf{x}}_k^+$ is expected to give a more accurate representation of the system at time t_k with respect to the original model prediction.

In the classical formulation, it is assumed that the system behavior can be described by a set of ordinary differential equations (ODEs) that can be expressed in a discrete-time form as:

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}, \mathbf{u}_k) + \mathbf{w}_k \quad (2.66)$$

where \mathbf{w}_k is the process error (sometimes called, in a misleading way, process *noise*) at time t_k . The measurements \mathbf{z}_k are related to the system state through a *measurement* (or *observation*) equation $\mathbf{h}(\mathbf{x}_k)$, i.e.:

$$\mathbf{z}_k = \mathbf{h}(\mathbf{x}_k) + \mathbf{v}_k \quad (2.67)$$

where \mathbf{v}_k is the measurement *noise*. Eq. (2.66) and (2.67) represent the starting point for any state estimation algorithm. However, it must be noticed that many chemical engineering-related systems are typically described by set of high-order algebraic and differential equations (i.e. high-index DAEs) whose reduction to pure ODE systems is not always possible. Therefore, the

classical formulation (2.66) does not hold true for these systems and, consequently, the available state estimation algorithms cannot be applied in a straightforward fashion. An extension of state estimation techniques to DAE systems is thoroughly described in Chapter 6. In the following, it is assumed that the system state can be described by the discrete-time ODE formulation (2.66). The three estimation algorithms that will be briefly discussed are the Extended Kalman Filter (EKF), the Unscented Kalman Filter (UKF) and the Ensemble Kalman Filter (EnKF). A discussion on other more sophisticated state estimators (e.g. particle filters (PFs), Moving Horizon Estimation (MHE) algorithms) can be found in the work of Rao *et al.*, 2001.

2.6 State estimation algorithms

The Kalman filter (Kalman, 1960) has been the precursor of most of the available state estimation algorithms and has been widely used in the process industry throughout these last decades. Although its original formulation can be applied to systems described by a linear state space model in Eq.(2.66), several extensions to nonlinear state space models have been proposed in the literature. The most famous extension is the Extended Kalman Filter (EKF), that has proved to be successful in several chemical engineering applications (Haseltine and Rawlings, 2005). In the last few years, due to the poor performance of the classical EKF formulation with highly nonlinear systems, other more sophisticated state estimation algorithms have been proposed. Among these, the Unscented Kalman Filter (UKF) and the Particle Filter (PF) have gained particular attention in the last years. Other formulations (EnKF, MHE) are out of the scope of this Dissertation.

2.6.1 Extended Kalman Filter (EKF)

Based on the discrete-time state space model (2.66) and observation model (2.67), the EKF can be formulated at each time step according to a two-step procedure (Jazwinski, 1960):

1. State prediction from time step t_{k-1} to time step t_k ;
2. Measurement update at time step t_k .

The relevant equations involved in each step can be summarized as follows.

2.6.1.1 Prediction step

The process error and the measurement noise are assumed to be white Gaussian noises with mean $\mathbf{0}$ and covariance \mathbf{Q} and \mathbf{R} respectively, i.e.:

0.

$$\mathbf{w}_k \sim N(\mathbf{0}, \mathbf{Q}_k); \mathbf{v}_k \sim N(\mathbf{0}, \mathbf{R}_k). \quad (2.68)$$

Moreover, the initial state of the system \mathbf{x}_0 is assumed to be normally distributed with initial state estimate $\hat{\mathbf{x}}_0$ and error covariance for the state estimate \mathbf{P}_0 . In all the following state

estimation algorithms, it is implicitly assumed that \mathbf{Q}_k and \mathbf{R}_k are tuning parameters to be set according to the case study considered. Additional details on the choice of \mathbf{Q}_k and \mathbf{R}_k on practical application will be thoroughly discussed in Chapter 6.

The state prediction \mathbf{x}_k^- from time t_{k-1} to time t_k is obtained through the set of model equations \mathbf{f} according to:

$$\mathbf{x}_k^- = \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{u}_k) \quad (2.69)$$

The error covariance the state prediction at time t_k \mathbf{P}_k^- is obtained through the expression:

$$\mathbf{P}_k^- = \mathbf{F}_{k-1} \mathbf{P}_{k-1}^+ \mathbf{F}_{k-1}^T + \mathbf{Q}_{k-1} \quad (2.70)$$

where \mathbf{P}_{k-1}^+ is the updated error covariance at the previous time t_{k-1} , and \mathbf{F}_{k-1} is the following Jacobian matrix:

$$\mathbf{F}_{k-1} = \left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)_{\mathbf{x} = \mathbf{x}_{k-1}^+} . \quad (2.71)$$

Eq. (2.70) is in the same form of the famous Riccati equation (Riccati, 1724) for linear systems, and it is obtained by a first order linearization of the error estimate on the systems state. It is worth noticing that this last aspect is one of the main drawbacks of the EKF, since the error covariance is propagated from one time step to another with a linearized propagation equation, thus posing a serious limitation when the behavior of the system is highly nonlinear.

2.6.1.2 Measurement update step

Given the measurement \mathbf{y}_k at time t_k , the system state prediction \mathbf{x}_k^- and error covariance prediction \mathbf{P}_k^- are updated accordingly to their new values \mathbf{x}_k^+ and \mathbf{P}_k^+ according to the following filter equations:

$$\mathbf{x}_k^+ = \mathbf{x}_k^- + \mathbf{K}_k [\mathbf{z}_k - \mathbf{h}(\mathbf{x}_k^-)] ; \quad (2.72)$$

$$\mathbf{P}_k^+ = (\mathbf{I}_{Q \times Q} - \mathbf{K}_k) \mathbf{P}_k^- \quad (2.73)$$

where \mathbf{K}_k is the *Kalman gain* and is defined as:

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k)^{-1} , \quad (2.74)$$

while $\mathbf{I}_{Q \times Q}$ is the identity matrix, and \mathbf{H}_k is the Jacobian matrix:

$$\mathbf{H}_k = \left(\frac{\partial \mathbf{h}}{\partial \mathbf{x}} \right)_{\mathbf{x} = \mathbf{x}_k^-}. \quad (2.75)$$

Eq. (2.72) suggests that the correction to the model-predicted state \mathbf{x}_k^- is proportional to the measurement residual ($\mathbf{y}_k - \mathbf{h}(\mathbf{x}_k^-)$), with the proportionality factor given by the Kalman gain of Eq. (2.74).

The combination of Eqs. (2.69)-(2.74) represents the EKF algorithm. From these equations, it is possible to notice that one big disadvantage of the EKF is that it requires the computation of (possibly) large dimensional Jacobians matrices. The numerical computation of these Jacobians can be sometimes very difficult or very expensive, thus limiting the robustness of this filter.

2.6.2 Unscented Kalman Filter (UKF)

The underlying assumption behind the EKF formulation is that, if the system state is considered as random variable, this variable is assumed to be distributed with a Gaussian PDF. The mean of this PDF is propagated from one point in time to another according to Eq. (2.69), while the covariance of this PDF is propagated according to the linearized analytical expression (2.70). As discussed in section 2.6.1.2, this approach has the notable disadvantage of introducing a linearization error and requiring the computation of expensive Jacobians.

The two key ideas behind the Unscented Kalman filter (UKF) algorithm that aims at overcoming the limitations of the EKF are:

1. It is easier to perform a nonlinear transformation on a single point rather than an entire PDF;
2. it is easy to find a set of deterministic points (called *sigma points*) whose sample PDF approximates the true PDF of the state vector.

Consistently with these statements, the UKF approximates the true PDF of the state vector with a set of deterministic (i.e. pre-allocated) sampling points, and propagates this PDF by propagating each of these sampling points through the nonlinear transformation $\mathbf{f}(\mathbf{x}, \mathbf{u})$. This procedure has the notable advantage of avoiding linearization errors and the computation of Jacobians. Mathematically, the algorithm can be decomposed in the same prediction and update step as for the EKF. The relevant equations are discussed below.

2.6.2.1 Prediction step

The key problem of the UKF is the choice of the *sigma points* whose sample PDF is assumed to be an approximation of the true PDF of the state vector. The criterion that is adopted is as follows.

Let η be the updated PDF of the state vector at time t_{k-1} , whose mean value and covariance are given by \mathbf{x}_{k-1}^+ and \mathbf{P}_{k-1}^+ respectively. In the UKF algorithm, η is approximated with the

sampling PDF of $2L$ (L is user-defined) sampling points chosen according to the following criterion:

$$\mathbf{x}_{k-1}^{-(l)} = \mathbf{x}_{k-1}^+ + \left(\sqrt{L\mathbf{P}_{k-1}^+} \right)_{(l)} \quad l = 1, 2, \dots, L \quad (2.76)$$

$$\mathbf{x}_{k-1}^{-(l)} = \mathbf{x}_{k-1}^+ - \left(\sqrt{L\mathbf{P}_{k-1}^+} \right)_{(l)} \quad l = L + 1, \dots, 2L. \quad (2.77)$$

The $2L$ sigma points are then propagated to time t_k using the model equations $\mathbf{f}(\mathbf{x}, \mathbf{u})$:

$$\mathbf{x}_k^{-(l)} = \mathbf{f}(\mathbf{x}_{k-1}^{-(l)}, \mathbf{u}_k) \quad (2.78)$$

and the predicted estimate and covariance of the state vector are computed as:

$$\mathbf{x}_k^- = \sum_{l=1}^{2L} W_s^{(l)} \mathbf{x}_k^{-(l)} ; \quad (2.79)$$

$$\mathbf{P}_k^- = \sum_{l=1}^{2L} W_c^{(l)} (\mathbf{x}_k^{-(l)} - \mathbf{x}_k^-)(\mathbf{x}_k^{-(l)} - \mathbf{x}_k^-)^T + \mathbf{Q}_{k-1}. \quad (2.80)$$

where the weighting factors are given by:

$$W_s^{(0)} = \frac{\lambda}{L + \lambda} \quad (2.81)$$

$$W_c^{(0)} = \frac{\lambda}{L + \lambda} + (1 - \alpha^2 + \beta) \quad (2.82)$$

$$W_s^{(l)} = W_c^{(l)} = \frac{1}{2(L + \lambda)} \quad (2.83)$$

$$\lambda = \alpha^2(L + \kappa) - L. \quad (2.84)$$

The two parameters α and κ are related to the spread of the sigma points of the sample PDF, while β is related to the PDF of the state vector. Although these parameters can be considered as tuning parameters (Wan *et al.*, 2000), their recommended values are typically $\alpha = 10^{-3}$, $\kappa = 0$, $\beta = 2$.

Eq. (2.79) and (2.80) represent the prediction equations for the UKF.

2.6.2.2 Measurement update step

In the update step, $2L$ sigma points are chosen using the same approach of the prediction step but considering that best state and error estimates are now \mathbf{x}_k^- and \mathbf{P}_k^- . Eqs (2.76) through (2.84) are therefore used to select the new sampling points according to these two best estimates. It is worth noticing that the covariance \mathbf{P}_k^- is augmented with the model error \mathbf{Q}_k as in Eq. (2.80). The obtained sampling points $\mathbf{x}_k^{+(l)}$, $l = 1, 2, \dots, L$ are projected through the measurement function $\mathbf{h}(\mathbf{x})$ to generate $2L$ “predicted” measurements $\mathbf{y}_k^{(l)}$ according to:

$$\mathbf{z}_k^{+(l)} = \mathbf{h}(\mathbf{x}_k^{+(l)}), \quad l = 1, 2, \dots, 2L. \quad (2.85)$$

The predicted measurements (2.85) are then combined in order to obtain the predicted measurement error \mathbf{z}_k^+ and the predicted measurement covariance \mathbf{P}_z according to:

$$\mathbf{z}_k^+ = \sum_{l=1}^{2L} W_s^{(l)} \mathbf{z}_k^{+(l)}; \quad (2.86)$$

$$\mathbf{P}_z = \sum_{l=1}^{2L} W_c^{(l)} (\mathbf{z}_k^{+(l)} - \mathbf{z}_k) (\mathbf{z}_k^{+(l)} - \mathbf{z}_k)^T. \quad (2.87)$$

The cross state-measurement covariance \mathbf{P}_{xz} is then computed:

$$\mathbf{P}_{xz} = \sum_{l=1}^{2L} W_c^{(l)} (\mathbf{x}_k^{+(l)} - \mathbf{x}_k^-) (\mathbf{z}_k^{+(l)} - \mathbf{z}_k)^T \quad (2.88)$$

And its value is used to compute the “Kalman” gain:

$$\mathbf{K}_k = \mathbf{P}_{xz} \mathbf{P}_z^{-1}. \quad (2.89)$$

The update equations for the state vector and the error covariance can therefore be written in the standard form:

$$\mathbf{x}_k^+ = \mathbf{x}_k^- + \mathbf{K}_k [\mathbf{z}_k - \mathbf{h}(\mathbf{x}_k^-)]; \quad (2.90)$$

$$\mathbf{P}_k^+ = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{P}_z \mathbf{K}_k^T. \quad (2.91)$$

Eqs (2.90) and (2.91) are the final update equations for the UKF.

2.6.3 Ensemble Kalman Filter (EnKF)

The EnKF can be considered as a Bayesian MCMC implementation of the EKF and relies on similar concepts with respect to the ones of the UKF, i.e:

1. It is easier to propagate points rather than distributions through nonlinear operators;
2. it is easy to approximate the state vector PDF with a sample PDF.

However, the most important difference between the UnKF and the EnKF is that the sampling points are not pre-allocated as in the UKF, but are drawn randomly according to a MCMC approach. The EnKF is sometimes erroneously considered as a particular type of particle filter, since it relies on the same Bayesian MCMC sampling approach. However, the EnKF is based on the assumption that the state vector follows a multivariate Gaussian PDF, while PFs do not make any assumption regarding the shape of the PDF of the state vector. For this reason, it is more convenient to consider the EnKF as a MCMC implementation of a standard EKF, rather than a specific type of particle filter.

The PDF of the state vector is approximated with S randomly-drawn sample points, with S (i.e. the ensemble size) chosen according to user's risk adversity (typical values are in the range $50 \div 1000$). Therefore, the EnKF is particular suited for systems involving a large number of state variables (i.e. $Q \gg S$), since the computational burden can be significantly reduced.

The relevant equations for the EnKF can be classified as for the previous filters in a prediction and update step.

2.6.3.1 Prediction step

A time t_k , the PDF of the state vector is approximated with an ensemble of S randomly-drawn sampling points :

$$X_k^- = (\mathbf{x}_k^{-(1)}, \mathbf{x}_k^{-(2)}, \dots, \mathbf{x}_k^{-(S)}). \quad (2.92)$$

The ensemble mean $\bar{\mathbf{x}}_k^-$ is used to approximate the state prediction \mathbf{x}_k^- at time t_k :

$$\mathbf{x}_k^- \sim \bar{\mathbf{x}}_k^- = \frac{1}{S} \sum_{s=1}^S \mathbf{x}_k^{-(s)}. \quad (2.93)$$

The ensemble error is then computed according to the following expression:

$$\mathbf{E}_k^- = \left[\left(\mathbf{x}_k^{-(1)} - \bar{\mathbf{x}}_k^- \right), \left(\mathbf{x}_k^{-(2)} - \bar{\mathbf{x}}_k^- \right), \dots, \left(\mathbf{x}_k^{-(S)} - \bar{\mathbf{x}}_k^- \right) \right]. \quad (2.94)$$

The ensemble covariance $\bar{\mathbf{P}}_k^-$ is then computed and used as an approximation of the true covariance \mathbf{P}_k^- of the state vector:

$$\mathbf{P}_k^- \sim \bar{\mathbf{P}}_k^- = \frac{1}{(S-1)} \mathbf{E}_k^- (\mathbf{E}_k^-)^T. \quad (2.95)$$

Eqs (2.93) and (2.95) are the prediction step equations for the EnKF.

2.6.3.2 Measurement update step

In this step, a set of “perturbed” measurements \mathbf{z}_k^s , $s = 1, 2, \dots, S$ is first obtained from the observed measurements at time t_k , \mathbf{z}_k , according to:

$$\mathbf{z}_k^{(s)} = \mathbf{z}_k + \mathbf{v}_k^{(s)}, \quad s = 1, 2, \dots, S \quad (2.96)$$

where $\mathbf{v}_k^{(s)}$ is randomly drawn from the PDF $N(\mathbf{0}, \mathbf{R}_k)$.

Secondly, the measurement ensemble error covariance $\mathbf{E}_{y_k}^-$ is computed according to:

$$\mathbf{E}_{y_k}^- = \left[\left(\mathbf{z}_k^{-(1)} - \bar{\mathbf{z}}_k \right), \left(\mathbf{z}_k^{-(2)} - \bar{\mathbf{z}}_k \right), \dots, \left(\mathbf{z}_k^{-(S)} - \bar{\mathbf{z}}_k \right) \right] \quad (2.97)$$

where $\bar{\mathbf{z}}_k$ is the mean of the ensemble (2.96). Given $\mathbf{E}_{y_k}^-$, the “Kalman” gain \mathbf{K}_k can then be computed according to:

$$\mathbf{K}_k = \frac{1}{S-1} \mathbf{E}_k^- (\mathbf{E}_{y_k}^-)^T \left[\frac{1}{S-1} \mathbf{E}_{y_k}^- (\mathbf{E}_{y_k}^-)^T \right]^{-1}. \quad (2.98)$$

Each member of the ensemble (2.92) can therefore be updated according to the standard expression:

$$\mathbf{x}_k^{+(s)} = \mathbf{x}_k^{-(s)} + \mathbf{K}_k \left(\mathbf{z}_k^{(s)} - \mathbf{h}(\mathbf{x}_k^{-(s)}) \right), \quad s = 1, 2, \dots, S. \quad (2.99)$$

The updated state vector \mathbf{x}_k^+ and the updated covariance \mathbf{P}_k^+ can then be computed as the mean and the covariance of this updated ensemble (2.99), using the same structural expression of Eq. (2.93) and (2.95).

The updated ensemble is then propagated to the next step in time following the set of model equations $\mathbf{f}(\mathbf{x}, \mathbf{u})$.

2.7 State estimation and online model recalibration

State estimators are typically deployed in the process industry in order to obtain an accurate prediction of the up-to-date state of the system. This information is then typically used as an input for the control system implemented in the plant (e.g. MPC). However, state estimators

can be used not only to obtain information about the current state of the system, but also to perform an online re-calibration of some of the parameters of the original process model (or to account for un-modeled disturbances that arise during the actual plant operation). In fact, during the normal operation of the plant, the representativeness of the model may change due to slow-dynamics physical phenomena that were not considered during the offline model calibration. Typical examples are fouling and catalyst deactivation. These phenomena may cause a drift in the some of the original model parameters and therefore a process-model mismatch (PMM) on the key performance indicators (KPIs) of the system. In order to account for these phenomena, state estimators can be deployed to obtain a regular re-calibration of some of the model parameters in order to remove the PMM on the KPIs. The choice of which model parameters should be re-calibrated online can be done based on previous knowledge or through sensitivity analysis on the KPIs. This aspect will be discussed in more detail in Chapter 6.

Let $\tilde{\boldsymbol{\theta}}$ be the set of model parameters that need to be re-calibrated online. The simplest way to perform a re-calibration through a state estimator is to perform a state augmentation procedure. The state-space model (2.66) is therefore re-written as:

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}, \mathbf{u}_k) + \mathbf{w}_k \quad (2.100)$$

$$\tilde{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_{k-1} + \mathbf{w}_{\tilde{\boldsymbol{\theta}}} \quad (2.101)$$

with $\mathbf{w}_{\tilde{\boldsymbol{\theta}}} \sim N(\mathbf{0}, \mathbf{R}_{\tilde{\boldsymbol{\theta}}})$ “process” error on the model parameters $\tilde{\boldsymbol{\theta}}$ (the covariance $\mathbf{R}_{\tilde{\boldsymbol{\theta}}}$ has to be considered as a tuning parameter) and initial conditions:

$$\mathbf{x}(0) = \mathbf{x}_0 \quad (2.102)$$

$$\tilde{\boldsymbol{\theta}}(0) = \bar{\boldsymbol{\theta}} \quad (2.103)$$

where $\bar{\boldsymbol{\theta}}$ is the vector collecting the nominal values of the parameters. It is worth noticing that Eq. (2.101) corresponds to the following set of equations in continuous form:

$$\frac{d\tilde{\boldsymbol{\theta}}}{dt} = \mathbf{0}. \quad (2.104)$$

The model parameters $\tilde{\boldsymbol{\theta}}$ are then used as additional states to be updated based on the measurements coming from plant sensors following one of the algorithms described above (or other state estimation algorithms that have not been considered in this Dissertation).

Chapter 3

Uncertainty back-propagation in PLS modeling for design space determination*

In this Chapter, a methodology to relate the uncertainty on the output of a PLS model to the uncertainty on the model inputs is proposed. Two uncertainty back-propagation models are formulated and critically compared and frequentist confidence regions (CRs) for the solution of the inversion problem are built. These CRs are used to target an experimental campaign for the identification of the design space (DS) of a new pharmaceutical product.

The proposed methodology is tested on three different case studies, two of which involve experimental data taken from the literature, respectively on a roller compactor and on a wet granulator. It is shown that both uncertainty back-propagation models are effective in bracketing the design space of a new pharmaceutical product, with the second model outperforming the first one in terms of shrinkage of the space within which experiments should be carried out to identify the DS.

3.1 Introduction

The quality-by design (QbD) initiative set forth by the pharmaceutical regulatory agencies (such as the U.S. Food and Drug Administration and the European Medicines Agency) encourages the adoption of systematic and science-based tools for the development and manufacturing of new pharmaceutical products (ICH, 2005). One key aspect of QbD is the determination of the design space (DS) of the process that manufactures the product under development. The DS is defined as the multidimensional combination and interaction of input variables (e.g. material attributes) and process parameters that have been demonstrated to

* Bano G., Facco, P., Meneghetti, N., Bezzo F., Barolo M. (2017) Uncertainty back-propagation in PLS model inversion for design space determination in pharmaceutical product development. *Comput. Chem. Eng.* **101**, 110-124.

provide assurance of quality (ICH, 2009). Identification of the DS of a new product is important, because “working within the DS is not considered as a change” (ICH, 2009), and as such does not require any further approval by the regulatory agencies. On the other hand, movements outside the DS would normally initiate a regulatory post-approval process.

Different model-based techniques can be used to assist a DS identification exercise. The use of static and dynamic first principles models has been reported (Pantelides *et al.*, 2009; Prpich *et al.*, 2010; Close *et al.*, 2014), as well as that of dynamic grey-box state-space models (Kishida and Braatz, 2012; Kishida and Braatz, 2014), possibly accounting for the effect of feedback control (Harinath *et al.*, 2015), and of data-driven models (MacGregor and Bruwer, 2008; Peterson, 2008). Among data-driven models, approaches based on design of experiments (Zidan *et al.*, 2007; am Ende *et al.*, 2007; Lepore and Spavins, 2008; Chatzizacharia *et al.*, 2014) and latent variable models have become particularly popular within the pharmaceutical industry community.

When a new product under development is similar to other products already developed (“old” products), historical information on the old products can be very useful for the determination of the DS for the new one. Facco *et al.* (2015) proposed a latent-variable (LV) model inversion approach (Jaeckle and MacGregor, 2000) to segment the knowledge space (KS, i.e. the space of historical information about the manufacturing of the old products) in such a way as to identify a subspace of it (called the experiment space; ES) that brackets the DS of the new product. Experiments to identify the DS can be carried out within the ES, with the advantage that this space is narrower than the KS, so that product development can be sped up. The model inverted by Facco *et al.* (2015) to determine the ES is a partial least-squares (PLS) regression one (Geladi and Kowalski, 1986; Wold *et al.*, 1983). Since this model is subject to prediction uncertainty, upon inversion the uncertainty is back-propagated to the input space, hence affecting the ES determination as shown by the authors.

Whereas the methodology proposed by Facco *et al.* (2015) is effective in segmenting the KS, it nevertheless suffers from some limitations. Firstly, the general issue of how the prediction uncertainty is allocated between the model inputs and the model parameters is not addressed. Secondly, the methodology returns an ES that is defined in the space of latent inputs rather than in the space of true inputs. While using the latent space is permitted by the regulatory agencies and can be useful to provide a compact graphical representation of a large-dimensional input space, this representation may be cumbersome for a practitioner aiming to design a set of experiments to be carried out at given values of the true inputs. A third limitation is related to the fact that the methodology only refers to products characterized by one equality constraint specification. However, in several practical situations, the product quality target may be assigned also in terms of inequality constraints.

In this study, an approach is proposed that can handle the back-propagation of uncertainty from the outputs to the inputs of a PLS model, while simultaneously allowing one to overcome the

limitations of the methodology proposed by Facco *et al.* (2015). The proposed approach makes it possible to quantify the uncertainty on the PLS inversion solution in terms of both latent variables (LVs) and real input variables. Two uncertainty back-propagation models are formulated and solved using an optimization framework. Confidence regions (CRs) for the PLS inversion solution identifying the ES are built at an assigned significance level for both back-propagation models. The resulting ES is deemed to include the DS of the product under development, thus allowing the developer to tailor his/her experimental campaign on a smaller domain of input combinations. The proposed methodology is first tested on a nonlinear mathematical example, with quality targets described by equality or inequality constraints. Then, five case studies involving experimental data taken from the literature are considered, one concerning the roll compaction of an intermediate drug load formulation (Souihi *et al.*, 2015), and the others concerning different types of granulation processes (Vemavarapu *et al.*, 2009; Oka *et al.*, 2015; Facco *et al.* (2015).

3.2 Methods

In the following, the link between the mathematical methods used in this Chapter and their application in pharmaceutical development is briefly discussed.

3.2.1 Projection to latent structures (PLS) for pharmaceutical product development

Let \mathbf{X} [$N \times V$] be a regressor (or input) matrix of N training observations and V variables, and \mathbf{Y} [$N \times L$] a response (or output) matrix of L variables. The regressor matrix collects the material attributes and process parameters and settings of the system under investigation, while the output matrix collects the product quality attributes of the pharmaceutical product to be designed. Data are assumed to be autoscaled, i.e. mean-centered and scaled to unit variance. PLS (Geladi and Kowalski, 1986; Wold *et al.*, 1983) is a multivariate regression technique that projects the regressor and response variables onto a common latent space (the model space) of A new variables, called LVs, according to the model structure:

$$\mathbf{X} = \mathbf{TP}^T + \mathbf{E}_X \quad (3.1)$$

$$\mathbf{Y} = \mathbf{TQ}^T + \mathbf{E}_Y \quad (3.2)$$

$$\mathbf{T} = \mathbf{XW}^* \quad (3.3)$$

where \mathbf{T} [$N \times A$] is the score matrix, \mathbf{P} [$V \times A$] and \mathbf{Q} [$L \times A$] are the \mathbf{X} and \mathbf{Y} loading matrices, \mathbf{E}_X and \mathbf{E}_Y the residuals; \mathbf{W}^* [$V \times A$] is the weight matrix, through which the data in \mathbf{X} are projected onto the latent space to give \mathbf{T} according to Eq.(3.3).

In this study, the algorithm used to build the PLS model is NIPALS (Geladi and Kowalski, 1986).

When a new observation $\mathbf{x}_p [1 \times V]$ is considered, its projection \mathbf{t}_p on the latent space can be obtained according to:

$$\mathbf{t}_p = \frac{\mathbf{x}_p \mathbf{W}^*}{\mathbf{p}^T \mathbf{W}^*} \quad (3.4)$$

and whether or not this new observation conforms to the calibration data set is assessed using two indices:

- the Hotelling's T^2 statistic (Hotelling, 1931):

$$T_p^2 = \sum_{a=1}^A \frac{t_{a,p}^2}{\lambda_a} \quad , \quad (3.5)$$

where λ_a is the eigenvalue of the matrix $\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}$ related to the a -th LV. T_p^2 is a measure of the Mahalanobis distance of the projection of \mathbf{x}_p from the origin of the latent space, and is usually compared to the respective 95% confidence limit:

$$T_{lim}^2 = \frac{(N-1)A}{N-A} F_{A,N-A,0.95} \quad , \quad (3.6)$$

where $F_{A,N-A,0.95}$ is the value of the 95% percentile of the F -distribution with A and $(N - A)$ degrees of freedom;

- the SPE statistic:

$$\text{SPE}_p = \mathbf{e}_{\mathbf{x},p}^T \mathbf{e}_{\mathbf{x},p} \quad , \quad (3.7)$$

which is the Euclidean distance of the projection of \mathbf{x}_p from the latent space. SPE_p is usually compared to the 95% confidence limit:

$$\text{SPE}_{lim} = \frac{\sigma_q}{2\mu} \chi_{\frac{2\mu}{\sigma_q}, 0.95}^2 \quad , \quad (3.8)$$

where μ and σ_q are the mean and standard deviation of the training set residuals, respectively, and $\chi_{\frac{2\mu}{\sigma_q}, 0.95}^2$ is the chi-square distribution with $\left(\frac{2\mu}{\sigma_q}\right)$ degrees of freedom and significance level 0.95.

In the A -dimensional latent space, T_{lim}^2 determines a 95% hyper-ellipsoidal confidence region whose semiaxis s_a along the a -th LV is given by:

$$s_a = \sqrt{\lambda_a T_{\text{lim}}^2} \quad . \quad (3.9)$$

This region is called the historical KS, and inference using the PLS model can be made only for points that lie inside this region. Therefore, in this study the input combinations, whose projections lie outside the KS, are assumed not to be described properly by the model.

This study deals with univariate responses ($L = 1$), so that \mathbf{Y} becomes $\mathbf{y} [N \times 1]$, \mathbf{Q} becomes $\mathbf{q} [1 \times A]$, and \mathbf{E}_Y becomes $\mathbf{e}_y [N \times 1]$.

Once a PLS model has been calibrated using a training set of regressors and responses according to Eqs (3.1)-(3.3), a new response $\hat{\mathbf{y}}_p [1 \times L]$ can be predicted, given a set of input data $\mathbf{x}_p [1 \times V]$:

$$\hat{\mathbf{y}}_p = \mathbf{t}_p \mathbf{Q}^T \quad (3.10)$$

where \mathbf{t}_p is the score vector of the new observation as given by Eq. (3.4).

A PLS model can also be used in its inverse form to determine the set \mathbf{x}_{new} of input combinations that yields a desired quality target y_{des} (Jaeckle and MacGregor, 2000).

From a mathematical point of view, the prediction model (10) is a linear transformation $\mathcal{L}_{PLS}(\mathbf{t})$ from an A -dimensional space Σ_T (the score or latent space) to a 1 dimensional space Σ_y (the response or y -space), whose kernel can be defined as:

$$\ker(\mathcal{L}_{PLS}) = \{\mathbf{t} \in \Sigma_T: \mathcal{L}_{PLS}(\mathbf{t}) = 0\}. \quad (3.11)$$

The rank-nullity theorem can be used to determine the dimension of $\ker(\mathcal{L}_{PLS})$:

$$\dim(\ker(\mathcal{L}_{PLS})) = \dim(\Sigma_T) - \dim(\text{Im}(\mathcal{L}_{PLS})) \quad (3.12)$$

from which, since $\dim(\Sigma_T) = A$ and (for this study) $\dim(\text{Im}(\mathcal{L}_{PLS})) = 1$, it follows that:

$$\dim(\ker(\mathcal{L}_{PLS})) = A - 1 \quad . \quad (3.13)$$

If the dimension of the latent space (i.e., A) is the same as the dimension of the y -space (i.e., 1), then $\dim(\ker(\mathcal{L}_{PLS})) = 0$ and a unique solution vector \mathbf{x}_{new} for a given target y_{des} exists:

$$\mathbf{x}_{new} = \mathbf{t}_{new} \mathbf{P}^T \quad (3.14)$$

with:

$$\mathbf{t}_{new}^T = (\mathbf{Q}^T \mathbf{Q})^{-1} \mathbf{Q}^T y_{des}^T \quad . \quad (3.15)$$

Note that \mathbf{x}_{new} (or \mathbf{t}_{new} if working on the score space) can be calculated as in Eq. (3.14) or (3.15) for any value of A . This solution will be denoted as the direct inversion solution to the inversion problem.

On the other hand, if the dimension of the latent space is greater than the dimension of the y -space (i.e. $A > 1$), then $\dim(\ker(\mathcal{L}_{PLS})) > 1$ and multiple solutions of the model inversion problem exist. The set W of possible solutions on the model space that guarantees the desired response y_{des} can be described as:

$$W = \{(\mathbf{t}_{new} + \mathbf{t}), \mathbf{t} \in \ker(\mathcal{L}_{pls})\} \quad (3.16)$$

and is called null space (Jaeckle and MacGregor, 1998). The null space can be computed analytically by determining the kernel of the y -loadings matrix \mathbf{Q} .

If the PLS model is not affected by uncertainty and the quality target is described by equality constraints only, the null space can be interpreted as the mathematical description of the DS¹⁴. However, since any model is affected by uncertainty, also the null space determined by model inversion is affected by uncertainty. By back-propagating the model prediction uncertainty onto the model space, Facco *et al.* (2015) determined a subspace of it (the ES) that most likely brackets the DS.

3.2.2 Formulation of prediction uncertainty

Consider the PLS prediction model. Given a new observation \mathbf{x}_p , the expected value of the response variable \hat{y}_p can be computed using (3.10). If the PLS model were not affected by uncertainty, \hat{y}_p would be equal to the true value of the response y_p , i.e. to the value of the response variable that would be obtained in the real process using the set of inputs defined by \mathbf{x}_p . However, since the model is affected by uncertainty, a mismatch between \hat{y}_p and y_p is observed. Three different sources of uncertainty¹⁵ can affect the prediction of a new observation:

1. measurement uncertainty in both the regressor matrix (\mathbf{X}) and the response matrix (\mathbf{Y}) used to calibrate the PLS model (Reis *et al.*, 2005);
2. uncertainty in the estimated model regression parameters (Faber and Kowalski, 1997);
3. uncertainty due to the unmodeled part of y_p . Assuming that the residual $e_p = y_p - \hat{y}_p$ belongs to a normal distribution with zero mean and variance σ^2 , the contribution to the overall prediction uncertainty given by this term is σ^2 .

¹⁴ Strictly speaking, when the PLS model is not affected by uncertainty and the quality target is described by one equality constraint, the null space represents the projection of the DS onto the PLS model space.

¹⁵ It is assumed that the variance of the residuals and the variance of the measurement errors in the input/output data are independent of the sample considered.

In this study, the prediction uncertainty is estimated according to the work of Faber and Kowalski (1997), which accounts for the second and third contributions mentioned above. An estimation of the variance $\text{var}(\hat{y}_p)$ of the prediction error is calculated using the relationship (Faber and Kowalski 1997):

$$s^2 = \text{var}(\hat{y}_p) = \sigma^2(1 + h_p) \quad , \quad (3.17)$$

where h_p is the leverage of the new observation:

$$h_p = \frac{\mathbf{t}_p \Lambda^{-1} \mathbf{t}_p^T}{N-1} \quad . \quad (3.18)$$

The variance σ^2 of the residuals is estimated as the square of standard error of calibration (SEC):

$$\sigma^2 \cong \text{SEC}^2 = \frac{\sum_{n=1}^N (y_n - \hat{y}_n)^2}{N-df} \quad , \quad (3.19)$$

where y_n and \hat{y}_n refer to the n -th calibration sample, and df is the number of degrees of freedom of the model. In this study, df has been set equal to A , although different approaches to estimate df exist (Voet, 1999).

The derivation of Eq. (3.17) involves a zeroth-order local linearization of the PLS estimator, and holds true under some assumptions (Vanlaer *et al.*, 2009). Eq. (3.17) is made up of two contributions: σ^2 is the contribution to the prediction uncertainty due to the unmodeled part of \hat{y}_p , while $\sigma^2 h_p$ accounts for the uncertainty in the model parameters. The $100(1 - \alpha)\%$ confidence interval (CI) on \hat{y}_p can be expressed as:

$$\text{CI} = \hat{y}_n \pm s t_{\frac{\alpha}{2}, N-df} \quad , \quad (3.20)$$

where α is the significance level for the confidence interval.

3.2.3 Back-propagation of uncertainty in PLS model inversion

When a PLS model is used in its direct form (i.e., to predict a response from a set of inputs) a $100(1 - \alpha)\%$ CI can be built for the prediction using the procedure outlined in section 3.2.2. The wider the CI at a given significance level, the larger the uncertainty on the output prediction. On the other hand, when a PLS model is used in its inverse form (i.e., to predict a new set \mathbf{x}_{new} of inputs that yield a desired target y_{des}), the prediction uncertainty is back-propagated to the calculated inputs. If $A > 1$, a null space exists, and the uncertainty in null-space determination must be accounted for.

A systematic methodology is now developed to back-propagate the y -space uncertainty to the latent space of input variables first, and to the real input space subsequently.

Consider the prediction model (3.10). By taking the variance/covariance of both sides it follows:

$$\text{cov}(\hat{\mathbf{y}}_p) = \text{cov}(\mathbf{t}_p \mathbf{Q}^T) \quad . \quad (3.21)$$

If model (3.17) is used to estimate the prediction uncertainty, then it follows that:

$$\text{SEC}^2(1 + h_p) = \text{cov}(\mathbf{t}_p \mathbf{Q}^T) \quad . \quad (3.22)$$

The right-hand side of Eq. (3.22) is difficult to estimate, since both the scores and the loadings are affected by uncertainty and depend on each other. However, two limiting situations can be considered (Høy *et al.*, 1998):

1. *the loadings are not affected by error.* If this is the case, Eq. (3.22) simplifies to:

$$\text{SEC}^2(1 + h_p) = \mathbf{Q}^T \text{cov}(\mathbf{t}_p) \mathbf{Q} \quad . \quad (3.23)$$

In this scenario, only the scores are assumed to contribute to the overall prediction uncertainty. In other terms, the overall prediction uncertainty on the y -space is entirely back-propagated to the score space of the inputs;

2. *the scores are not affected by error.* In this case, Eq. (3.22) simplifies as follows:

$$\text{SEC}^2(1 + h_p) = \mathbf{t}_p^T \text{cov}(\mathbf{Q}^T) \mathbf{t}_p \quad . \quad (3.24)$$

In this scenario, the overall prediction uncertainty on the y -space is fully back-propagated to the loadings of the model, i.e. it is assumed that the null space can be exactly located in the latent space when the PLS model is inverted, and the only source of uncertainty is due to the calibration step.

Following the idea implemented in the popular chemometrics software The Unscrambler™ (Camo ® Inc., Oslo, Norway), the right-hand terms of (3.23) and (3.24) can be averaged using the correction factor proposed by De Vries and Ter Braak (1995) to obtain:

$$\text{SEC}^2(1 + h_p) = \left(1 - \frac{A+1}{N}\right) \left[\underbrace{\mathbf{t}_p^T \text{cov}(\mathbf{Q}^T) \mathbf{t}_p}_{\textcircled{1}} + \underbrace{\mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q}}_{\textcircled{2}} \right] \quad . \quad (3.25)$$

Eq. (3.25) is empirical and has no theoretical foundation (Faber and Kowalski, 1996), but is nevertheless widely used (Andersson *et al.*, 1999; Cahyadi *et al.*, 2010). By using this approach, the overall prediction uncertainty is obtained as the average of the contribution due to the loadings uncertainty (① in Eq. (3.25)) and the contribution due to the scores uncertainty (② in Eq. (3.25)). In other terms, the prediction uncertainty is simultaneously determined by the uncertainty on the model “parameters” (i.e., the loadings) and the uncertainty on the location of the scores on the score space. It is reasonable to think that, when PLS model inversion is performed, only a fraction of the overall prediction uncertainty is back-propagated to the scores (i.e., to the location of the null space in the score space), whereas the remaining fraction is assigned to the loadings (i.e., to uncertainty on the model parameters).

Contribution ① can be estimated as in De Vries and Ter Braak (1995):

$$\mathbf{t}_p^T \text{cov}(\mathbf{Q}^T) \mathbf{t}_p = V_{y, \text{val}} h_p \quad , \quad (3.26)$$

where $V_{y, \text{val}}$ is the residual variance of the response variable in the validation dataset. By substituting (3.26) into (3.25), it follows:

$$SEC^2(1 + h_p) = \left(1 - \frac{A+1}{N}\right) \left[\underbrace{V_{y, \text{val}} h_p}_{\textcircled{1}} + \underbrace{\mathbf{Q}^T \text{cov}(\mathbf{t}_p) \mathbf{Q}}_{\textcircled{2}} \right]. \quad (3.27)$$

Eq. (3.27) shows that the contribution ① of the loadings to the overall prediction uncertainty is proportional to the leverage of the new observation. This is reasonable, since the uncertainty due to the non-representativeness of the model increases when the distance of the validation dataset projection from the origin on the score space increases.

The weighting factor $V_{y, \text{val}}$ in (3.27) is strongly dependent on the validation dataset available and significantly affects contribution ① to the overall prediction uncertainty. To avoid this dependency on the validation dataset available, a modification to Eq. (3.27) is proposed. Assume that the right-hand term of (3.26) can be re-written as $C h_p$, with C being a weighting factor to be tuned. Following this assumption, we will first investigate how the contribution of the scores to the overall prediction uncertainty depends on the leverage of the observation h_p . Then, we will identify the threshold value of h_p that causes the contribution of the scores to vanish. We will finally determine the value of the weighting factor C that corresponds to this limiting situation.

By substituting $V_{y, \text{val}}$ with C into (3.27) and rearranging, the following equation can be obtained:

$$\mathbf{Q}^T \text{cov}(\mathbf{t}_p) \mathbf{Q} = \frac{SEC^2}{1 - \frac{A+1}{N}} + \left[\frac{SEC^2 - C \left(1 - \frac{A+1}{N}\right)}{1 - \frac{A+1}{N}} \right] h_p \quad . \quad (3.28)$$

The left-hand term of (3.28) is the contribution to the prediction uncertainty related to the scores: if h_p increases, this contribution decreases since the dominant contribution becomes the one related to the non-representativeness of the model (i.e., related to the loadings). In particular, the contribution of the scores vanishes (i.e, all the uncertainty on the prediction is due to the non-representativeness of the model) when:

$$h_p = \frac{\text{SEC}^2}{c\left(1 - \frac{A+1}{N}\right) - \text{SEC}^2} = h^* \quad . \quad (3.29)$$

In other words, for those points on the score space that satisfy the condition $h_p < h^*$, both the scores *and* the loadings contribute to the overall prediction uncertainty; in this case, the contribution of the scores (loadings) is greater (smaller) for low (high) values of h_p . When $h_p \geq h^*$, the model is considered as non-representative, and the inversion should not be performed. To determine a value for the weighting factor C , recall that a likely region of the score space within which the model is considered representative is the KS (section 2.1). Then, if C is selected in such a way that h^* corresponds to the 95% confidence limit of the Hotelling's T^2 (i.e $h^* = h_{lim}$) then the points outside the KS will only be affected by uncertainty due to the model non-representativeness. It follows that:

$$\mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q} = \frac{\text{SEC}^2}{1 - \frac{A+1}{N}} + \left[\frac{\text{SEC}^2 - c\left(1 - \frac{A+1}{N}\right)}{1 - \frac{A+1}{N}} \right] h_p \quad \text{if } h_p < h_{lim} \quad (3.30)$$

$$\mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q} = 0 \quad \text{if } h_p \geq h_{lim} \quad (3.31)$$

The weighting factor C can then be tuned to have $h^* = h_{lim}$:

$$C = \frac{\text{SEC}^2(1 + h_{lim})}{\left(1 - \frac{A+1}{N}\right)h_{lim}} = \frac{\text{SEC}^2 \left[1 + \left(\frac{T_{lim}^2}{N-1} + \frac{1}{N} \right) \right]}{\left(1 - \frac{A+1}{N}\right) \left(\frac{T_{lim}^2}{N-1} + \frac{1}{N} \right)} \quad . \quad (3.32)$$

The resulting uncertainty back-propagation model is then given by:

$$\text{SEC}^2(1 + h_p) = \left(1 - \frac{A+1}{N}\right) [Ch_p + \mathbf{Q}^T \text{cov}(\mathbf{t}_p) \mathbf{Q}] \quad , \quad (3.33)$$

with C defined as in (3.32) .

To summarize, when a PLS model is inverted, the back-propagation of uncertainty can be done:

- model #1: according to (3.23), if only the contribution of the scores is considered. This is a conservative situation, since all the prediction uncertainty is back-propagated to the

calculated solution of the inversion problem (i.e., to the location of the null space on the score space, as in Facco *et al.*, 2015);

- model #2: according to (3.33), in the more general situation where the contributions of both the scores and the loadings to the overall prediction uncertainty are accounted for. Clearly, Eq. (3.33) is an empirical back-propagation formula whose effectiveness must be evaluated.

3.3 Uncertainty back-propagation and experiment space determination

The problem of developing a new product with a desired quality target y_{des} is addressed. Consider a PLS model built on a historical dataset of input combinations \mathbf{X} that yield products whose quality is summarized by \mathbf{y} . If the model is used for prediction, the response \hat{y}_p that corresponds to the desired input specifications \mathbf{x}_p can be computed according to the steps bracketed by the dashed black box of Fig. 3.1. On the other hand, when the model is inverted, the set of regressors that, according to the model, corresponds to the desired quality target y_{des} (that can be expressed in terms of equality or inequality constraints) can be determined.

If the number of latent variables is greater than one, the set of inputs (on the score space) that matches y_{des} is the vector space W of Eq. (3.16). As described in 3.2.3, the prediction of W is affected by different sources of uncertainty. To quantify this uncertainty, we use the systematic methodology centered on the back-propagation formulas (3.23) and (3.32)-(3.33). The rationale behind the proposed methodology is shown in Fig. 3.1, blue dashed box. The step-by-step procedure is presented in the next section.

3.3.1 Proposed methodology

In this section, the methodology will be illustrated in detail only for problems where an equality constraint is set on the product quality ($y = y_{des}$); the extension to inequality constraints is straightforward and will be discussed in section 3.5.2. Given the calibration dataset $[\mathbf{X}; \mathbf{y}]$, a PLS model with A LVs relating \mathbf{X} to \mathbf{y} is built. A can be chosen according to the eigenvalue-greater-than-one rule (Kaiser, 1960), in such a way as to explain a significant fraction of the variance of both \mathbf{y} and \mathbf{X} (to have good reliability also in model inversion). The projections of the calibration samples onto the score plot and the KS boundary (blue ellipse) are shown in Fig. 3.2a for $A = 2$.

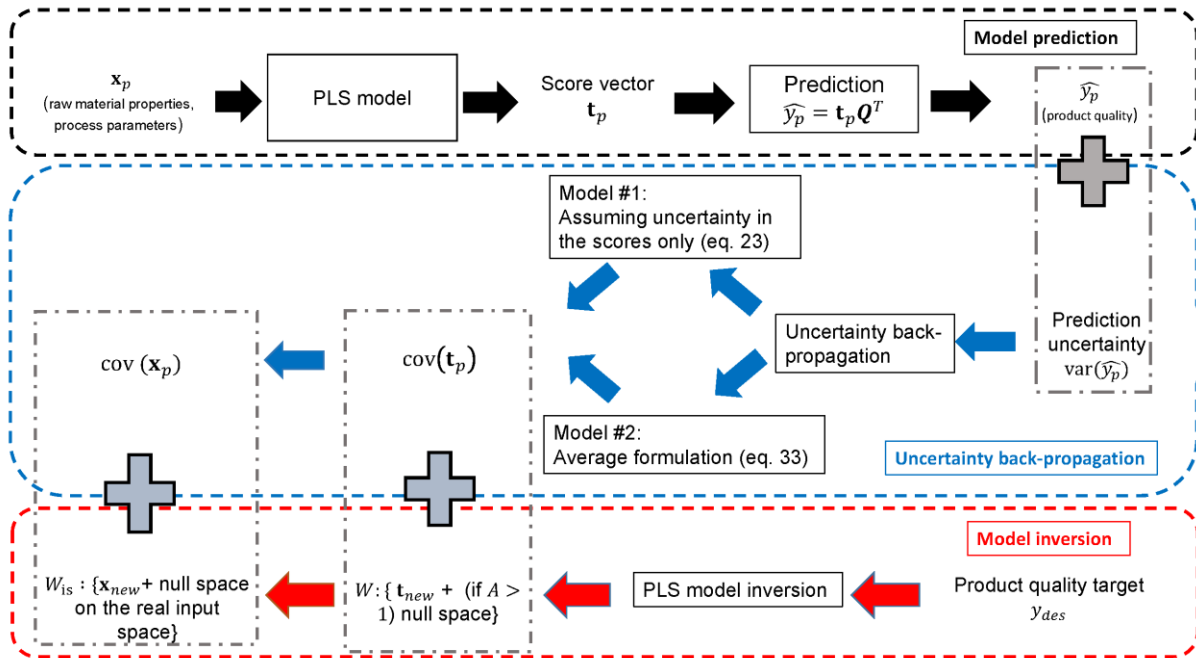


Figure 3.1. Schematic of uncertainty characterization in PLS modelling. When a PLS model is used for prediction (i.e., to predict the product quality given a set of raw material properties and process parameters), the steps are those bracketed by the dashed black box. When the model is inverted (i.e., it is used to determine the combinations of raw material properties and process parameters that give rise to a desired product quality), the steps are those bracketed by the dashed red box. The rationale behind the back-propagation of uncertainty in PLS model inversion according to the proposed methodology involves the steps bracketed by the dashed blue box.

1. A model to estimate the variance on the model prediction is chosen. The model considered in this study is given by (3.17).
2. A new product of quality y_{des} is desired. The direct inversion solution \mathbf{t}_{new} on the score space corresponding to the given quality target y_{des} is calculated according to (3.15). If $A > 1$, a null space of dimension $A - 1$ exists and can be computed as in section 3.2.2. If $A = 2$, the graphical representation of the null space in the score plot is a straight line passing through \mathbf{t}_{new} , as shown in Fig. 3.2b.
3. The null space is discretized at L score points (with L sufficiently large). The coordinates of the l -th score point along the null space are $(t_{1,l}, t_{2,l}, \dots, t_{A,l})$. For each of these points, the associated covariance matrix:

$$\text{cov}(\mathbf{t}_l) = \begin{pmatrix} \sigma_{t_{1,l}}^2 & \text{cov}(t_{1,l}, t_{2,l}) & \cdots & \text{cov}(t_{1,l}, t_{A,l}) \\ \text{cov}(t_{2,l}, t_{1,l}) & \sigma_{t_{2,l}}^2 & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(t_{A,l}, t_{1,l}) & \cdots & \cdots & \sigma_{t_{A,l}}^2 \end{pmatrix} \quad (3.34)$$

is calculated using one of the following two options:

- *model #1*: if it is assumed that the prediction uncertainty is back-propagated to the scores only, the uncertainty back-propagation formula is given by (3.23). Eq. (3.23) is solved as a constrained optimization problem:

$$\text{cov}(\mathbf{t}_l) = \min_{\text{cov}(\mathbf{t})} [\mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q} - \text{SEC}^2(1 + h_p)] [\mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q} - \text{SEC}^2(1 + h_p)]^T \quad (3.35)$$

$$\text{s.t. } \sigma_{t_{i,l}}^2 > 0, \quad i = 1, \dots, A \quad (3.36)$$

$$\text{cov}(t_{i,l}, t_{j,l}) = \text{cov}(t_{j,l}, t_{i,l}) = 0, \quad i, j = 1, \dots, A \quad (3.37)$$

Constraint (3.37) reflects the fact the scores are uncorrelated, because the NIPALS algorithm produces orthogonal \mathbf{X} -scores.

- *Model #2*: the uncertainty back-propagation formula is given by (33), meaning that the uncertainty on the prediction is propagated on both the loadings and the scores. Eq. (33) is solved as a constrained optimization:

$$\text{cov}(\mathbf{t}_l) = \min_{\text{cov}(\mathbf{t})} \left\{ \left(1 - \frac{A+1}{N}\right) [Ch_p + \mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q}] - \text{SEC}^2(1 + h_p) \right\} \left\{ \left(1 - \frac{A+1}{N}\right) [Ch_p + \mathbf{Q}^T \text{cov}(\mathbf{t}) \mathbf{Q}] - \text{SEC}^2(1 + h_p) \right\}^T \quad (3.38)$$

$$\text{s.t. } \sigma_{t_{i,l}}^2 > 0, \quad i = 1, \dots, A \quad (3.39)$$

$$\text{cov}(t_{i,l}, t_{j,l}) = \text{cov}(t_{j,l}, t_{i,l}) = 0, \quad i, j = 1, \dots, A. \quad (3.40)$$

4. The set of L points \mathbf{t}_l along the null space and their associated covariance matrices $\text{cov}(\mathbf{t}_l)$ can be used to completely identify a restricted portion of the KS that is expected to bracket the DS projection of the process. To this purpose, for each point a $(1 - \alpha)\%$ -confidence hyper-ellipsoid is built by joining the points of coordinates $\mathbf{t}[1 \times A]$ that satisfy the equation (Tomba *et al.*, 2012):

$$(\mathbf{t} - \mathbf{t}_l) \text{cov}(\mathbf{t}_l) (\mathbf{t} - \mathbf{t}_l)^T = \frac{A(N^2-1)}{N(N-A)} F_{A, N-A, \alpha}, \quad (3.41)$$

where $F_{A, N-A, \alpha}$ is the value corresponding to the α -th% percentile of the F distribution with A and $(N - A)$ degrees of freedom. An example of projection on the score space of the 95% confidence hyper-ellipsoids on the first two LVs is shown in Fig. 3.2c (model #1) and Fig. 3.2e (model #2).

5. The confidence limits of the null space can be built by calculating the envelope of the L $(1 - \alpha)\%$ -confidence hyper-ellipsoids determined in the previous step. The region bracketed by this envelope is the experiment space, namely a subspace of the

KS within which the DS projection of the process is likely to lie with a confidence of *at least* $(1 - \alpha)(\%)$. A graphical interpretation is shown in Fig. 2d for model #1 and in Fig. 3.2f for model #1 and model #2 altogether. When the scores are assumed not to be affected by uncertainty, the null space can be exactly located on the score space (black solid straight line). On the opposite, when the prediction uncertainty is entirely back-propagated to the scores, the ES is the one bracketed by the divergent red solid lines (model #1). The average situation, according to model #2, is represented by the black dashed lines: the uncertainty in the location of the null space decreases as far as the leverage of the l -th point increases, and it becomes zero outside the KS (where the model is not representative).

6. Each point \mathbf{t}_l of the null space is projected onto the multivariate V -dimensional real input space to give \mathbf{x}_l according to:

$$\mathbf{x}_l = \mathbf{t}_l \mathbf{P}^T \quad . \quad (3.42)$$

The uncertainty related to this l -th point of the null space, described by its covariance $\text{cov}(\mathbf{t}_l)$, is back-propagated to the real input space using the back-propagation formula:

$$\text{cov}(\mathbf{t}_l) = \mathbf{W}^{*T} \text{cov}(\mathbf{x}_l) \mathbf{W}^* \quad . \quad (3.43)$$

Eq. (3.43) is derived from (3.3) by assuming that the model weights are not affected by uncertainty. This is a conservative assumption since all the uncertainty related to \mathbf{t}_l is back-propagated to the multidimensional input space. Eq. (3.43) is solved for $\text{cov}(\mathbf{x}_l)$ using an optimization framework and enforcing $\text{cov}(\mathbf{x}_l)$ to be positive definite.

7. The multidimensional confidence region for the projected null space is built on the real input space following the same procedure as for steps #5 and #6, by substituting \mathbf{t}_l with \mathbf{x}_l and $\text{cov}(\mathbf{t}_l)$ with $\text{cov}(\mathbf{x}_l)$. This generates a V -dimensional space of combinations of input variables that brackets the real DS with at least $(1 - \alpha)(\%)$ confidence. To obtain a graphical interpretation of this multidimensional space, its projections onto two-dimensional plots of pairs of input variables can be plotted. It is worth remembering that this multidimensional space is *completely* defined by the set of points \mathbf{x}_l and their associated covariance matrices $\text{cov}(\mathbf{x}_l)$.

With respect to the methodology proposed by Facco *et al.* (2015), the following improvements are obtained:

- The proposed methodology is completely analytical and allows quantifying the uncertainty on the inversion solution in terms of point estimates and corresponding covariance matrices.
- As in Facco *et al.* (2015), the solution of the PLS model inversion problem and the estimation of solution uncertainty can be conveniently expressed in terms of LVs instead of true input variables, thus obtaining a strong reduction of the problem dimension. However, from a practical point of view, the approach by Facco *et al.* (2015) requires one to transform the result of the PLS inversion problem and its uncertainty in terms of real input variables. The proposed methodology allows one to project the solution back in the real input space together with the analytical quantification of its uncertainty.
- The proposed methodology is general and can also be used when *multivariate* quality targets are assigned, provided that a reliable model to estimate the uncertainty on the prediction of the multivariate response is available. The development of such models is still an open research area and this study may direct further investigation on this topic.

The methodology has been tested on three case studies. All the algorithms were implemented in MATLAB v. 2015b using a sequential quadratic programming optimization algorithm.

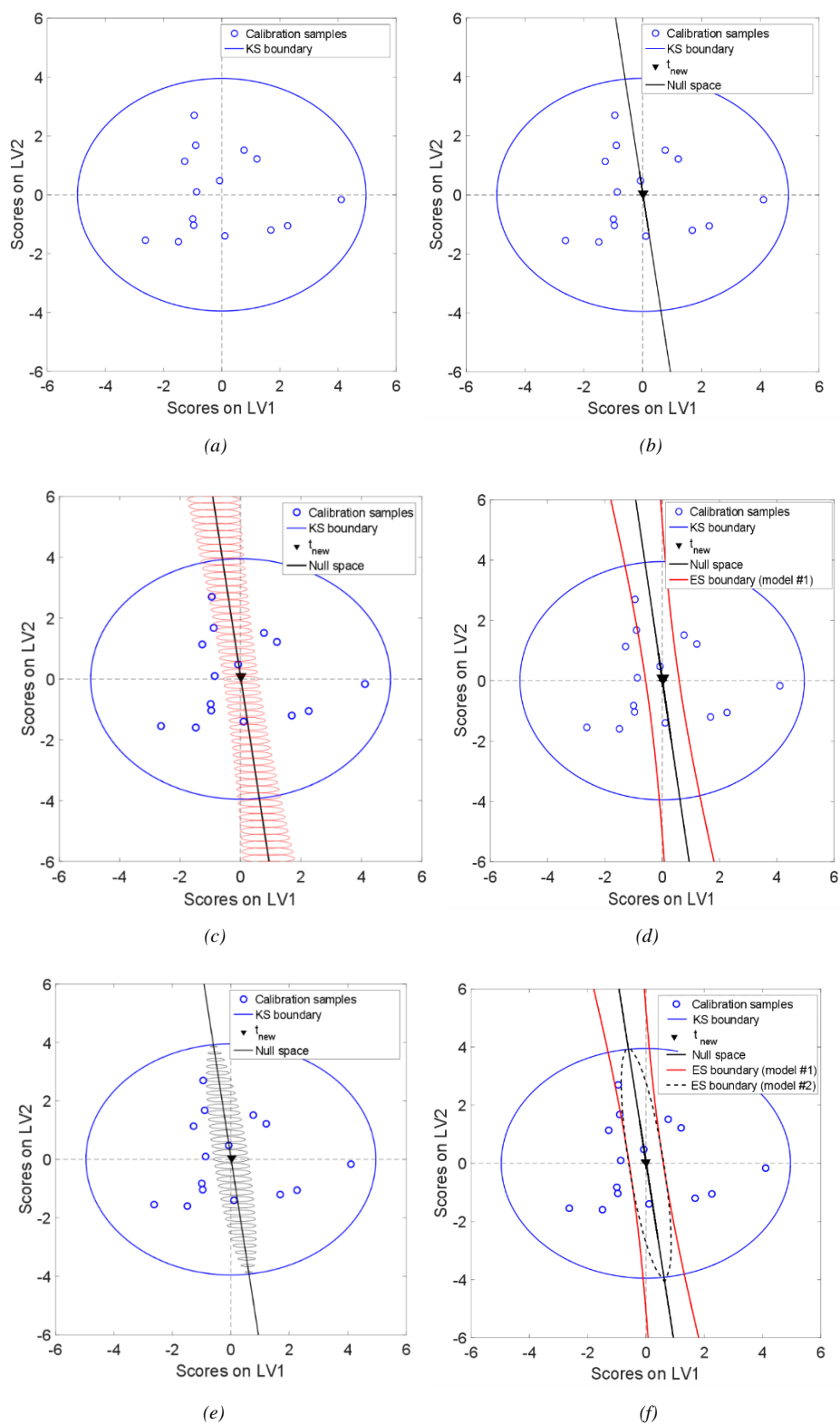


Figure 3.2. (a) Projection of the calibration samples onto the score space and KS boundary. (b) Direct inversion solution (triangle) and null space across it (solid line). (c) and (e): 95% confidence ellipses along the null space for model #1 and #2. (d) and (f): envelope of the confidence ellipses obtained with model #1 (red lines) and model #2 (dashed lines).

3.4 Case studies

3.4.1 Case study #1: mathematical example

The first illustrative example is the nonlinear mathematical example presented by Facco *et al.* (2015). The calibration dataset is made of 15 observations of five input variables \mathbf{X} [15×5] = [$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5$] and one response variable \mathbf{y} [15×1]. Two ($\mathbf{x}_1, \mathbf{x}_2$) of the five input variables are independent and are generated as normal random Gaussian distributions with mean and standard deviation as in Tab. 1. The other three ($\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5$) dependent input variables are related to \mathbf{x}_1 and \mathbf{x}_2 for any observation i according to:

$$x_{i,3} = x_{i,1}^2 \quad (3.45)$$

$$x_{i,4} = x_{i,2}^2 \quad (3.46)$$

$$x_{i,5} = x_{i,1}x_{i,2} \quad (3.47)$$

The response dataset is obtained using the model:

$$\mathbf{y} = k_0 + k_1\mathbf{x}_1 + k_2\mathbf{x}_2 + k_3\mathbf{x}_3 + k_4\mathbf{x}_4 + k_5\mathbf{x}_5 \quad (3.48)$$

with $[k_0, k_1, k_2, k_3, k_4, k_5] = [-21.0, 4.3, 0.022, -0.0064, 1.1, -0.12]$.

Table 3.1. Case study #1: characterization of the input/output data.

Variable	Mean	Standard deviation
Inputs		
\mathbf{x}_1	41.73	16.07
\mathbf{x}_2	11.13	2.97
\mathbf{x}_3	1999.15	1408.07
\mathbf{x}_4	132.63	66.93
\mathbf{x}_5	464.85	227.38
Response		
\mathbf{y}	235.99	71.35

3.4.2 Case study #2: roll compaction of an intermediate drug load formulation

This case study concerns a roll compaction of an intermediate drug load formulation with the following composition: 15% paracetamol, 55.3% mannitol, 23.7% microcrystalline cellulose, 4% croscarmellose and 2% sodium stearyl fumarate. Experimental data for this case study have been provided by Souihi *et al.* (2015) and refer to two different equipment configurations

(vertically- and horizontally- fed roll compactor). Details on the process and on the formulation can be found in the cited reference.

The historical data set includes 34 observations (11 related to the vertically-fed roll configuration, 23 to the horizontally-fed roll configuration) of 4 input variables (roll force, roll width, roll diameter, dimensionless ratio of gap to roll diameter GD) and one quality response variable (ribbon porosity)¹⁶. A summary of this dataset is reported in Table 3.2. The input variables are collected in the calibration regressor matrix \mathbf{X} [34×4], and the observations of the ribbon porosity are collected in the response vector \mathbf{y} [34×1].

Table 3.2. Case study #2: list of the input and response variables (data from Souihi *et al.*, 2015).

ID	Variable name	Units	Symbol
<i>Inputs</i>			
1	Roll force	[kN/cm]	p_{roll}
2	Roll width	[mm]	S_{roll}
3	Roll diameter	[mm]	D_{roll}
4	Ratio of the gap to roll diameter	[-]	GD
<i>Response</i>			
R1	Ribbon porosity	[%]	P

Nine external validation samples are used in the work of Souihi *et al.* (2015) to validate the model predictions. These validation samples have been used also in this study to test the effectiveness of the proposed methodology.

3.4.3 Case study #3: high-shear wet granulation of a pharmaceutical blend

Experimental data for this case study are taken from Oka *et al.* (2015) and are based on a design of experiment approach (full factorial, $3 \times 3 \times 3$). The high-shear wet granulation of a two-component (API + excipient) pharmaceutical blend is considered. The effect of three process parameters (impeller speed, liquid to solid ratio (L/S), wet massing time) on the median particle size d_{50} is considered. Twenty-seven observations of the input process parameters and of the corresponding d_{50} are available. This historical dataset is split into 20 randomly-chosen calibration samples and 7 validation samples. The input calibration matrix \mathbf{X} [20×3] collects the 20 calibration observations of the three process inputs, whereas the response calibration matrix \mathbf{y} [20×1] collects the respective 20 observations for the median particle size. A list of the input and output variables considered in this case study is given in Table 3.3.

¹⁶ In the work of Souihi *et al.* (2015), other 4 input variables (roll gap, roll speed, screw speed, ratio between the screw speed and the roll speed) are considered. However, since they have little effect on ribbon porosity (see Fig. 3.7a of the cited reference for more detail) they have not been included into the calibration dataset.

Table 3.3. Case study #3: list of the input and response variables (data from Oka *et al.*, 2015).

ID	Variable name	Units	Symbol
<i>Inputs</i>			
1	Liquid to solid ratio	[-]	L/S
2	Impeller speed	[rpm]	v_{imp}
3	Wet massing time	[min]	t_{wet}
<i>Response</i>			
R1	Median particle size	[μm]	d_{50}

3.4.4 Case study #4: dry granulation by roller compaction

Historical data for this case study were generated by Facco *et al.* (2015) using the modelling environment of gSOLIDS[®] (Process Systems Enterprise Ltd, 2014¹⁷) based on the model of Johanson (1965). Eight input variables (compressibility factor, roller diameter, roller width, roller speed, pressure force, friction angle between solid granulate and roller compactor, effective friction angle and springback factor) and one response variable (intravoid fraction of the solids out of the compactor) are considered. Details on the characterization of the input/output data for the calibration data set are reported in the cited reference.

3.4.5 Case study #5: wet granulation

This case study deals with the design of a powder product by means of high-shear wet granulation. The historical data set is taken from the experimental work of Vemavarapu *et al.* (2009) and is composed by 25 observations of seven input material properties (H_2O solubility, contact angle, H_2O holding capacity, Sauter mean diameter, distribution span, surface area, pore volume) and one response variable (fraction of granules larger than 1,4 mm). Additional details on the historical data set can be found in the cited reference.

3.5 Results for case study #1

3.5.1 Quality target described by an equality constraint

The problem of determining the set of process inputs that yields a product with $y_{des} = 300$ is addressed. First, a PLS model using the calibration datasets \mathbf{X} and \mathbf{y} as described in section 3.4.1 is built. The diagnostics of such a model is presented in Table 3.4.

¹⁷ Free software license provided by Process Systems Enterprise Ltd is gratefully acknowledged.

Table 3.4. Case study #1. Diagnostics of the PLS model.

LV #	R_X^2 (%)	$R_{X,cum}^2$ (%)	R_Y^2 (%)	$R_{Y,cum}^2$ (%)
1	60.59	60.59	96.63	96.63
2	38.47	99.06	1.48	98.11
3	0.67	99.73	0.94	99.05
4	0.19	99.92	0.71	99.77
5	0.08	100.00	0.23	100.00

Two LVs ($A = 2$) are chosen, explaining more than 99% of the total variability of \mathbf{X} and more than 98% of the total variability of \mathbf{y} . Since $A > \dim(\mathbf{y}) = 1$, a 1-dimensional null space exists. The DS corresponding to y_{des} is calculated from model (3.48) and then projected onto the score space (green line in Figs. 3.3a and 3.4a) and on the real input space (green line in Figs. 3.3b and 3.4b). Notice that, since the number of independent input variables is 2, also the dimension of the true input space is 2. The DS of the process is a curve because model (3.48) is nonlinear. The input combinations that project outside the KS (which is bracketed by the blue contours both in Figure 3.3 and in Figure 3.4) cannot be represented by the PLS model because outside the KS the model cannot be assumed to be fully representative.

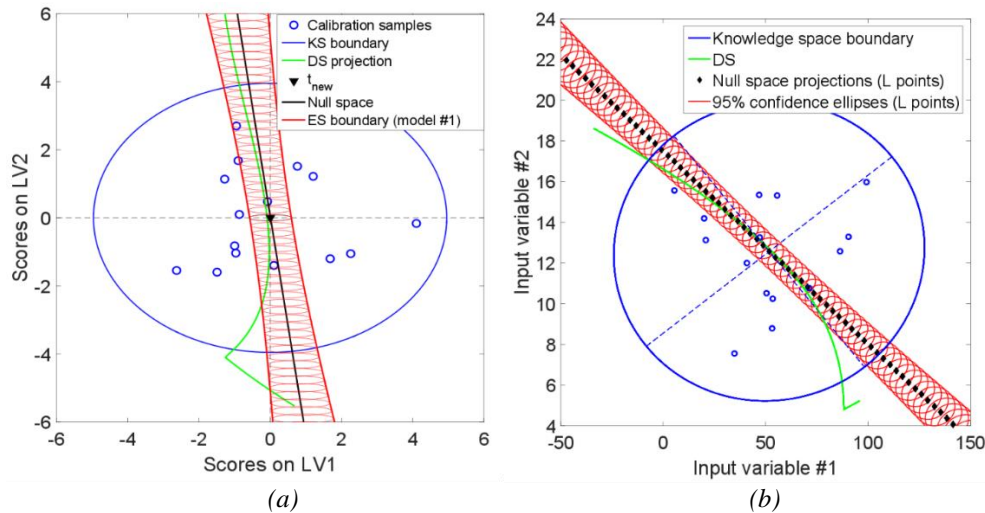


Figure 3.3. Case study #1, uncertainty back-propagation model #1: confidence limits of the null space and projection of the DS (a) onto the score space and (b) onto the real input space.

One issue deserves attention. The DS of the product is defined as the set of input combinations that corresponds exactly to a product with quality y_{des} . The green line of Figs. 3.3a and 3.4a is actually the projection of the DS onto the score space, i.e., it is the DS as it is “seen” by the PLS model. Different PLS models (i.e., different calibration datasets) would result in different DS projections. However, for a given y_{des} , the DS is unique by definition. For example, the green lines of Figs. 3.3b and 3.4b represent the DS for $y_{des}=300$. Since in these plots the DS is expressed in terms of real input variables, its representation is independent of the PLS model used to interpret the historical dataset.

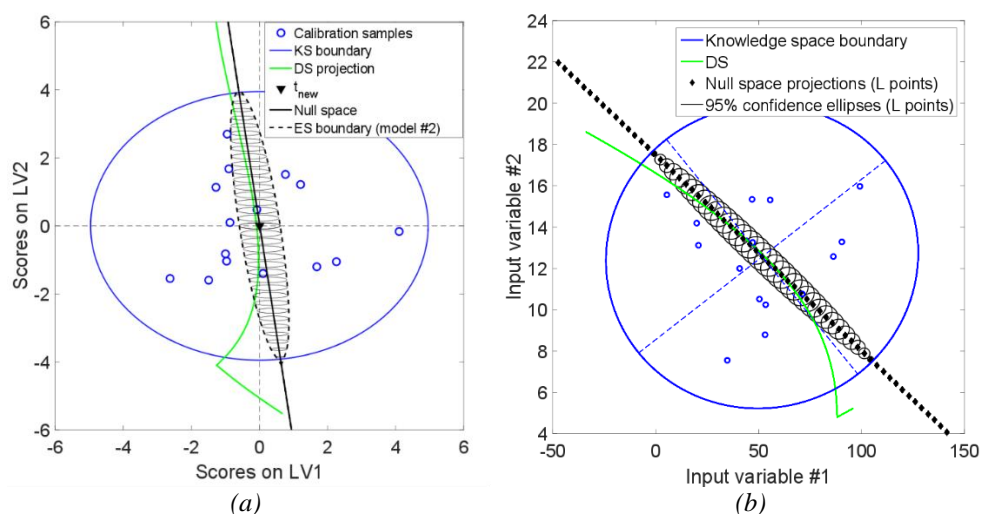


Figure 3.4. Case study #1, uncertainty back-propagation model #2: confidence limits of the null space and projection of the DS (a) onto the score space and (b) onto the real input space.

The direct inversion solution \mathbf{x}_{new} (black triangles of Figs. 3b and 4b) and its projection \mathbf{t}_{new} (black triangles of Fig. 3.3a and 3.4a) are then computed. The null space is then plotted in the score space (black solid line of Figs. 3.3a and 3.4a). A finite number of L points are chosen along the null space and their projections onto the real input space are shown as black diamonds in Figs. 3.3b and 3.4b. It is worth noticing that not all the input combinations lying along the null space may be achievable in practice because of physical or operational constraints. The same thing can be said for the DS. For each of the L points along the null space, the 95% confidence ellipse is built using the uncertainty back-propagation model #1 (red solid ellipses of Fig. 3.3a) and the modified model #2 (black solid ellipses of Fig. 3.4a). The envelope curves of this family of co-axial 95% confidence ellipses are then plotted (model #1: red solid lines of Fig. 3.3a; model #2: black dashed lines of Fig. 4a). Their projections onto the real input space are shown in Figs. 3.3b and 3.4b, respectively. The resulting ESs obtained with model #1 and model #2 are shown on the score space in Fig. 3.5a and on the real input space in Fig. 3.5b. Notice that part of the ES built using model #1 (blue region) is overlapped with the one built using model #2; namely, the red region is the portion of the KS that belongs to the CR obtained with model #1, but not with model #2. It can be noticed that the predicted ES for model #1, as well as that for model #2, effectively bracket a large fraction of the DS of the process lying within the KS. However, due to the system nonlinearity, not all the DS projection that lies inside the KS can be bracketed. The colored regions of Fig. 3.5b are the combinations of true independent input variables that should be primarily considered in the experimental campaign to determine the DS of the new product. They indeed point to a region, much narrower than the KS, within which the DS of the process is likely to lie.

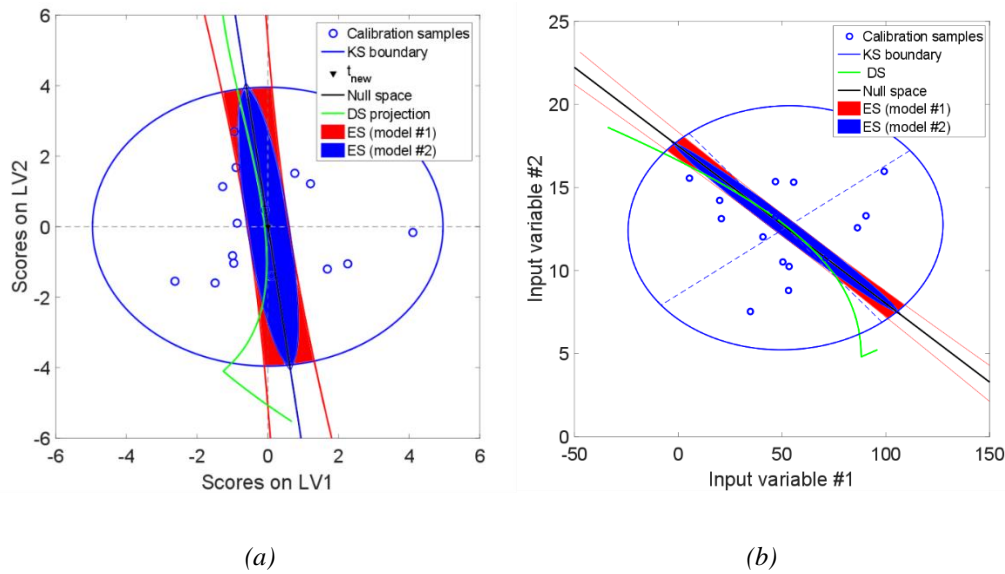


Figure 3.5. Case study #1: graphical interpretation on (a) the score space and (b) the real input space of the ES built using the uncertainty back-propagation model #1 (red + blue area) and model #2 (blue area).

From a practical viewpoint, when a new product of desired characteristics needs to be developed, the DS is not known *a priori*. The direct validation of the proposed methodology would therefore require the experimental validation of the results obtained from the model. In this mathematical example, the effectiveness of the proposed methodology has been validated using coverage probabilities (CPs) according to the following procedure. A set of $N_{val} = 1000$ validation samples $(\mathbf{X}^*, \mathbf{y}^*)$ and N_{cal} calibration samples are generated as described in section 4.1. A PLS model with $A = 2$ LVs is built with the calibration data set. For each validation sample, the null space and its confidence regions at assigned confidence levels are calculated using alternatively model #1 and model #2. It is then checked if the projection \mathbf{t}_p of the known reference value \mathbf{x}_{real} falls within the specified confidence region of the null space. The CP is calculated as the ratio between the number of validation samples that fall within the predefined CR of the null space and the total number of validation samples (CP is expressed as a percentage). If a validation sample is found to be outside the KS, it is discarded and not considered in the analysis. The CP is expected to be equal to, or greater than, the assigned confidence level, which is set to 95%.

To analyze the effect of the size of the calibration dataset, CPs have been calculated using different values for N_{cal} , but always maintaining the same 1000 validation samples. Results are reported in Table 3.5.

Table 3.5. Case study #1: observed CPs for uncertainty back-propagation models #1 and #2 with a confidence level of 95%. CPs lower than the predefined confidence level are highlighted in bold.

N_{cal}	CP model #1 (%)	CP model #2 (%)
5	92.5	86.1
10	93.6	88.5
15	94.1	91.6
20	95.1	93.2
25	95.8	95.8
30	99.0	95.4
50	96.1	95.2
100	97.8	96.1
150	98.2	96.8
200	99.3	97.2

It can be seen that both uncertainty back-propagation models perform properly even when only few calibration samples are available. Model #1 seems to perform slightly better than model #2: this is reasonable, since in this model all the prediction uncertainty is back-propagated to the scores and the CR of model #1 is larger than the one of model #2. However, an additional issue that should be considered when assessing the performance of the two uncertainty back-propagation models is quantifying the extent of the shrinkage (with respect to the original KS) of the space within which experiments should be carried out to identify the DS. We therefore propose a simple metric to quantify the extent of the shrinking of the experimental region to be spanned by experiments.

Consider the V -dimensional space obtained as the projection of the KS from the score space to the real input space. Let V_k be the hypervolume of this space. Then consider the projections of the ES boundaries located inside the KS (obtained with model #1 or #2) on the V -dimensional input space, and let $V_{e,i}$ be the hypervolume of the region bracketed by these boundaries, with $i = 1$ (model #1) or $i = 2$ (model #2). The ES shrinking (ESS) factor for model i can be quantified as:

$$ESS_i(\%) = \left[1 - \left(\frac{V_{e,i}}{V_k} \right) \right] \cdot 100 \quad (3.49)$$

The larger this factor, the more effective the segmentation of the KS obtained with a given uncertainty back-propagation model. This metric complements the CP to evaluate the model performance. An effective uncertainty back-propagation model should be reliable (i.e., with CP equal to, or greater than, the predefined confidence level) and ensure the largest possible shrinking of the ES at the same time.

In the case study under investigation, the real input space is 2-dimensional and (3.49) reduces to a ratio of areas. The ratio ESS_i as a function of the size of the calibration data set for model #1 and #2 is reported in Table 3.6.

Table 3.6. Case study #1: experimental space shrinking using the uncertainty back-propagation models #1 and #2.

N_{cal}	ESS ₁ (%)	ESS ₂ (%)
	model #1	model #2
10	79.0	93.3
15	86.0	97.2
20	81.6	96.4
25	74.3	94.3
30	77.5	92.7
50	83.1	97.0
100	86.5	97.4
200	88.9	98.0
500	91.5	98.8
1000	91.1	97.5

It can be noticed that the shrinkage of the experiment space that can be obtained with model #2 is significantly larger than the one that can be obtained with model #1. This is due to the fact that model #1 is more conservative than model #2. Moreover, ESS₁ is much more dependent on the size of the calibration data set than ESS₂.

3.5.2 Quality target described by an inequality constraint

In several practical situations, the quality target for the product to be developed is not specified as an equality constraint, but as an acceptance region (i.e., as a set of inequality constraints). In the following, we discuss the case where the set of quality constraints is assigned as $y_{lo} \leq y_{des} \leq y_{up}$, where y_{lo} and y_{up} are the lower and upper bounds for y_{des} , respectively.

Consider Fig. 3.6. The bounds for y_{des} are set to $y_{lo} = 278$ and $y_{up} = 378$. A PLS model with $A = 2$ LVs and $N = 15$ calibration samples is built as in section 3.5.1. The green area shown in Fig. 3.6a is the projection of the DS falling within the KS. This area is bracketed by the KS boundary and by the DS projections as if these projections corresponded to two distinct equality constraints, namely $y_{des} = y_{lo}$ and $y_{des} = y_{up}$. The null spaces corresponding to $y_{des} = y_{lo}$ and $y_{des} = y_{up}$ can then be computed. If no uncertainty on the location of these boundary null spaces is assumed, the DS projection that would be predicted by the PLS model is the one given by the black area of Fig. 3.6b.

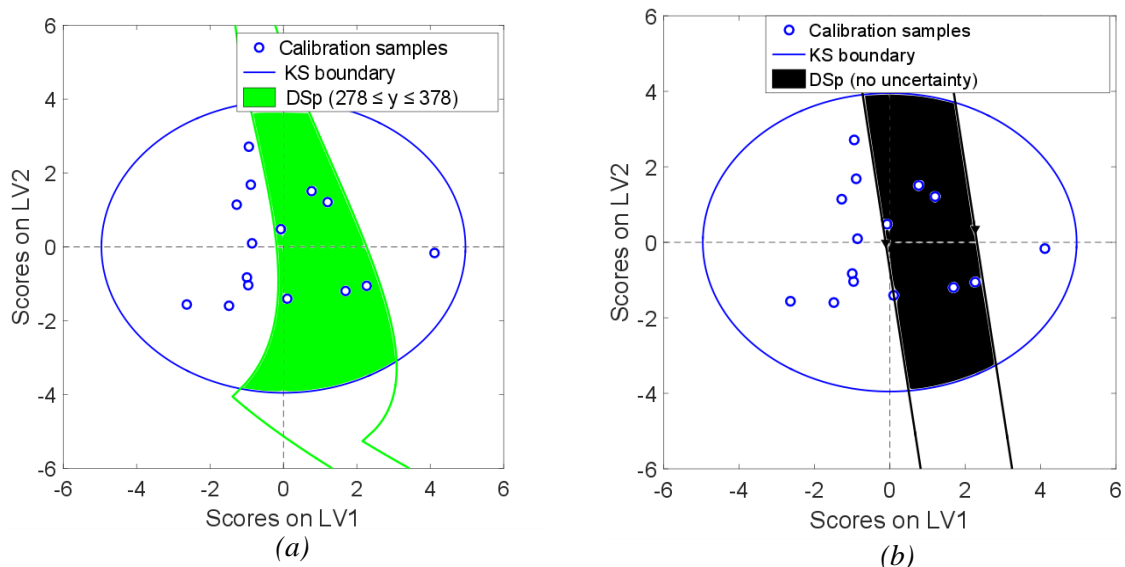


Figure 3.6. Case study #1, quality target described by $278 \leq y_{des} \leq 378$: (a) projection of the true DS on the score space (green area); (b) DS projection predicted by the PLS model when prediction uncertainty is not accounted for (black area).

If prediction uncertainty is accounted for, the envelope curves of the 95% confidence ellipses for the boundary null space locations can be built by using the uncertainty back-propagation model #1 or model#2 as shown in Figs. 3.7a and 3.7b, respectively.

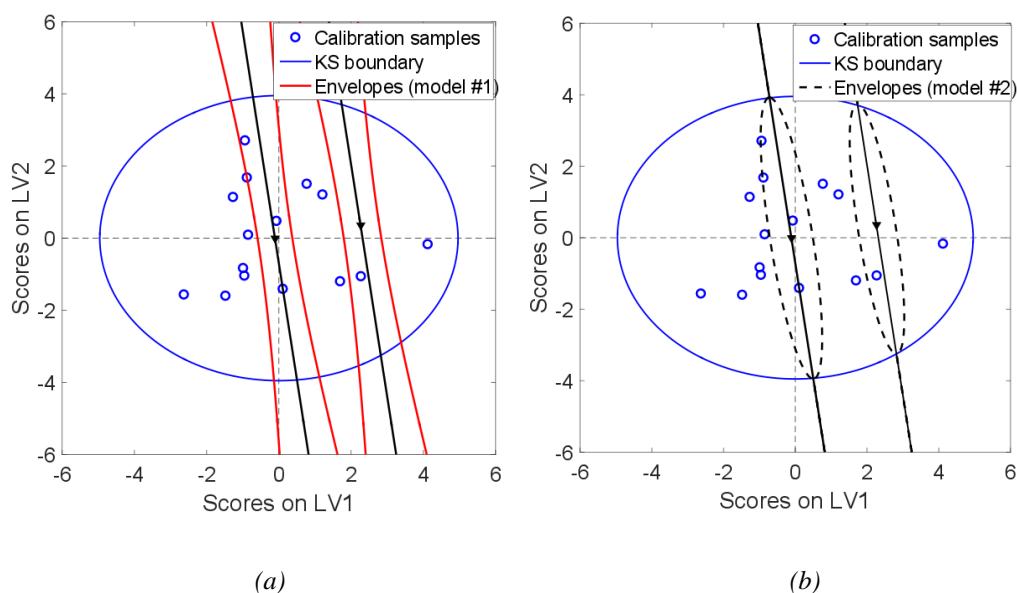


Figure 3.7. Case study #1, quality target described by $278 \leq y_{des} \leq 378$: 95% envelope curves of the 95% confidence ellipses for the boundary null spaces according to (a) model #1 and (b) model #2.

Therefore, the DS is expected to lie within the red area of Fig. 3.8a (uncertainty back-propagation model #1) or within the blue area of Fig. 3.8b (uncertainty back-propagation model #2).

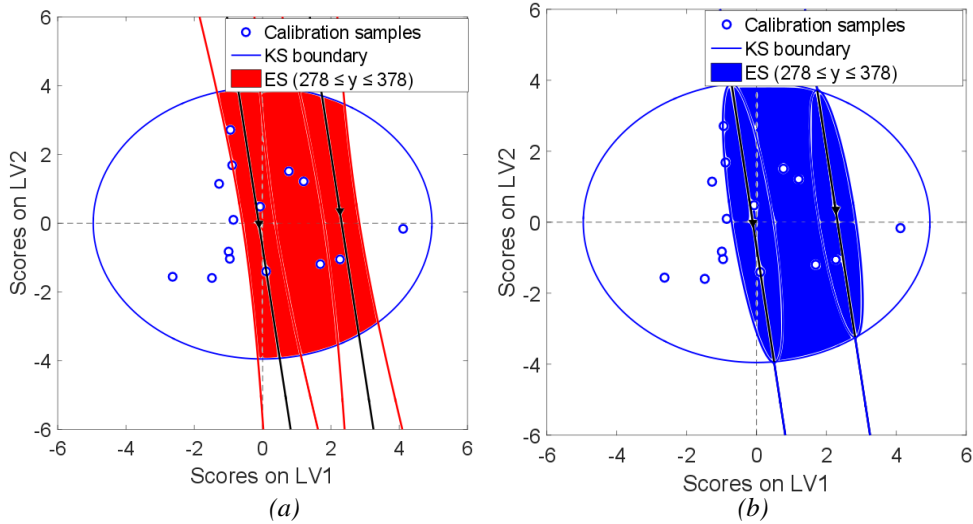


Figure 3.8. Case study #1, quality target described by $278 \leq y_{des} \leq 378$: ES according to (a) model #1 and (b) model #2.

Fig. 3.9 provides a comprehensive graphical interpretation of the overall discussion. If the segmentation of the KS is done using the PLS inversion solution and model uncertainty is not accounted for (Fig. 3.9a), a large portion of the DS projection is not bracketed and therefore will go unexplored during an experimental campaign. Note that in this case uncertainty in the determination of the DS projection mainly derives from the fact that a linear model (the PLS model) is used to predict a strongly nonlinear DS. When uncertainty in the PLS inversion solution is accounted for, a larger portion of the DS of the process is bracketed (Fig. 3.9b for model #1, Fig. 3.9c for model #2). However, a much smaller subspace of the DS is still not bracketed (green areas) by the envelope curves of the inversion solution. In fact, the contribution of the variance of the residuals to the overall prediction uncertainty in Eq. (3.17), which partially accounts for nonlinearity, is back-propagated to the calculated solution. By increasing the confidence level (i.e. by reducing the risk of not bracketing the DS), a broader ES would be obtained. In a nutshell, the greater the desired confidence level in bracketing the DS of the new product, the wider the operating window that must be investigated by experimentation.

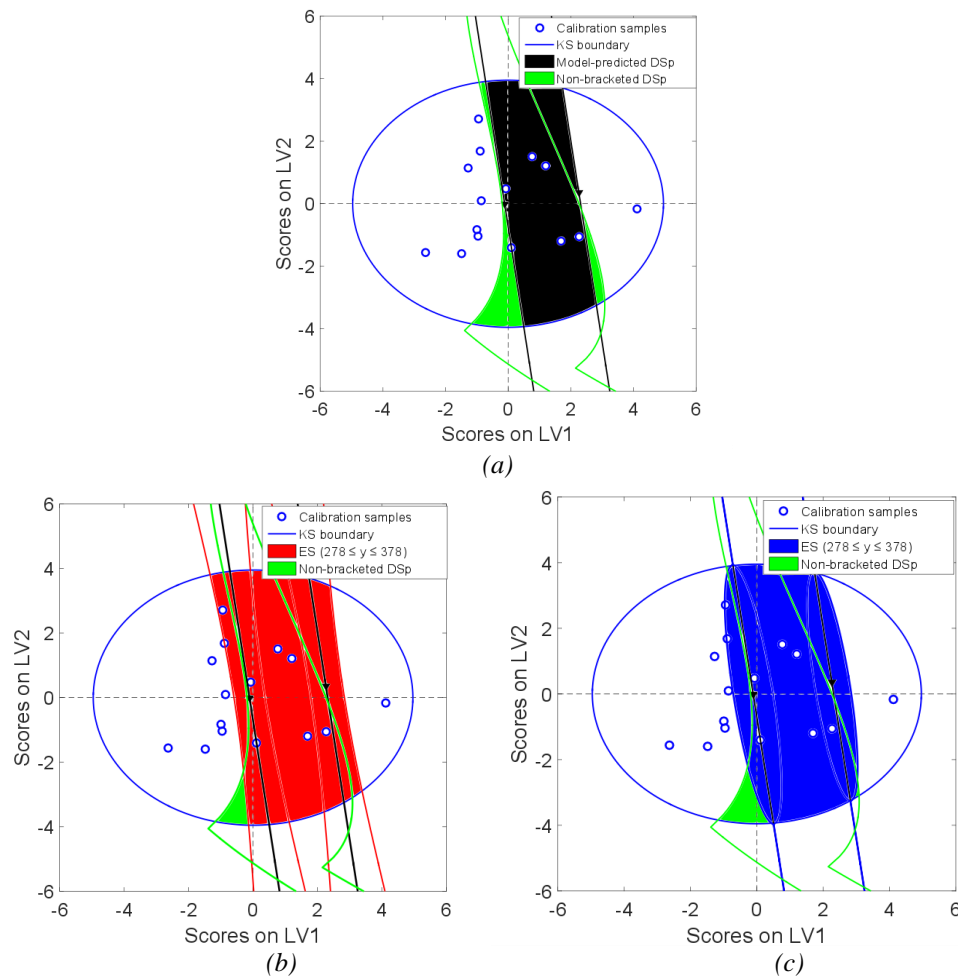


Figure 3.9. Case study #1, quality target described by $278 \leq y_{des} \leq 378$. (a) Model-predicted DS projection without accounting for uncertainty (black area), and portion of the DS projection that is not bracketed by the model predictions (green area). (b) ES (red area) according to model #1, and portion of the DS projection that is not bracketed (green area). (c) ES (blue area) according to model #2, and portion of the DS projection that is not bracketed

3.6 Results for case study #2

The same external validation data set as in Souihi *et al.* (2015) is used. First, a PLS model is built using the 34 calibration samples. The model diagnostics is shown in Table 3.9.

Table 3.9. Case study #2: diagnostics of the PLS model.

LV #	R_x^2 (%)	$R_{x,cum}^2$ (%)	R_y^2 (%)	$R_{y,cum}^2$ (%)
1	57.72	57.72	66.82	66.82
2	38.87	96.59	20.28	87.10
3	3.41	100.00	12.33	99.43
4	0.00	100.00	0.57	100.00

Using $A = 2$ LVs, 96.6% of the total variability of \mathbf{X} and 87.1% of the total variability of \mathbf{y} are explained. The proposed methodology is then tested on the 9 validation samples available. As an illustrative example, the results obtained for a product with desired average porosity $y_{des} = 22.6\%$ (validation sample #7) are shown in Fig. 3.10.

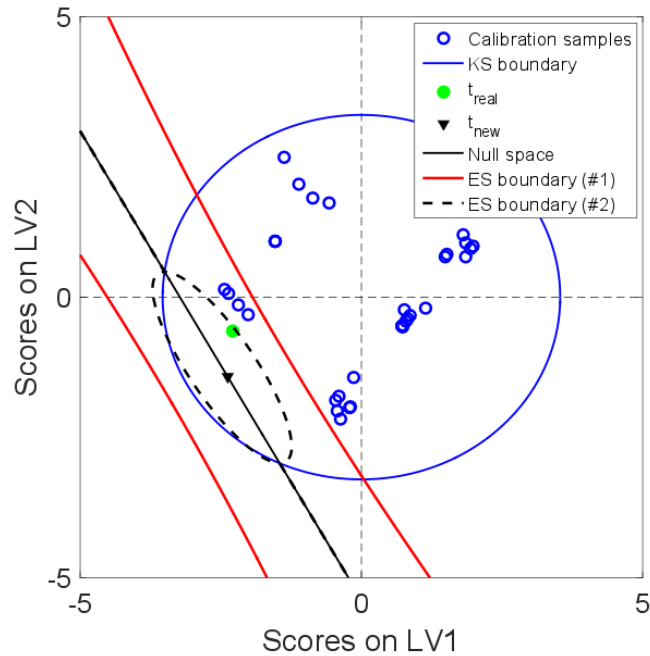


Figure 3.10. Case study #2: bracketing the DS projection of an intermediate drug load formulation with $y_{des} = 22.6\%$. The green circle is the projection of the inputs actually giving y_{des} , the black triangle is the direct inversion solution, the black solid line is the null space as calculated by the PLS model, the red solid lines are the envelope of the 95% confidence ellipses of the null space as calculated by model #1 and the black dashed lines those calculated by model #2.

The projection of the process inputs actually required to obtain the desired product quality (green circle in Fig. 3.10) lies within the ES both for model #1 and for model #2. It is worth noticing that this green circle is only one of the possible real input combinations that can lead to the desired product quality. We will refer to this point as \mathbf{t}_{real} . The null space related to the direct inversion solution \mathbf{t}_{new} (black solid line in Fig. 3.10) is quite far from \mathbf{t}_{real} . Therefore, if the PLS model were used without accounting for prediction uncertainty, a significant error on the location of the DS projection would be obtained. The ESS indices that can be obtained with the two models for this case study are $ESS_1 = 77.9\%$ and $ESS_2 = 87.7\%$ respectively. The proposed methodology has been tested on all the 9 validation samples available. The rate of success was 100% both with model#1 and with model #2.

3.7 Results for case study #3

A PLS model is first built using the randomly-chosen 20 calibration samples. $A = 2$ LVs are used, accounting for 76.3% of the total variability of \mathbf{X} and 86.5% of the total variability of \mathbf{y} . The seven validation samples are then used to test the effectiveness of the proposed methodology. An example is shown in Fig. 3.11 (validation sample #2). The green circle is the projection onto the score space of the input combination that actually yields the desired median particle size ($d_{50}^{des} = 1101 \mu\text{m}$). The direct inversion solution (black triangle) and the null space across it (black solid line) are shown in the same figure. It can be noticed that, if uncertainty is not accounted for, the model-predicted DS projection would not include the actual set of inputs. When considering the back-propagation of uncertainty on the score space, the resulting ES effectively brackets the true input combination for the desired quality target. By carrying out the experimental campaign across these restricted operating windows, an experiment space shrinkage of $ESS_1 = 69.3\%$ and $ESS_2 = 74.5\%$ can be obtained.

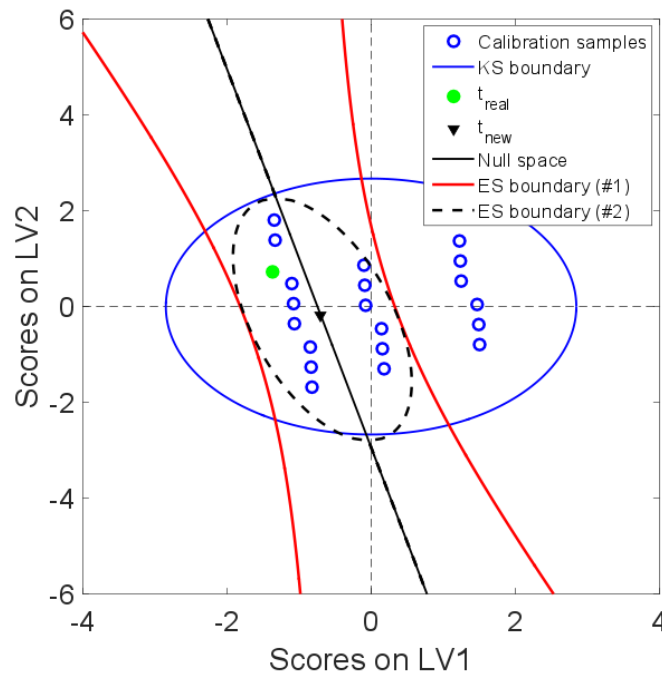


Figure 3.11. Case study #3: bracketing the true input combinations (green dot) yielding a desired median particle size granulate ($d_{50}^{des} = 1101 \mu\text{m}$) in a high-shear wet granulation process. The KS is segmented using model #1 (red convolution lines) and model #2 (black dashed convolution lines).

3.8 Results for case study #4

Differently from Facco *et al.* (2015), 30 calibration samples are selected randomly from the historical dataset (composed by 80 samples). The problem of designing a process that allows one to obtain a granulate with intravoid fraction of the solids out of the roller compactor of

$y_{des} = 0.6341 \text{ m}^3/\text{m}^3$ is considered. The results obtained with two uncertainty back-propagation models are shown in Fig. 3.12. A reduction of the experimental effort of $ESS_1 = 69\%$ can be obtained by using model #1, and of $ESS_2 = 76\%$ with model #2.

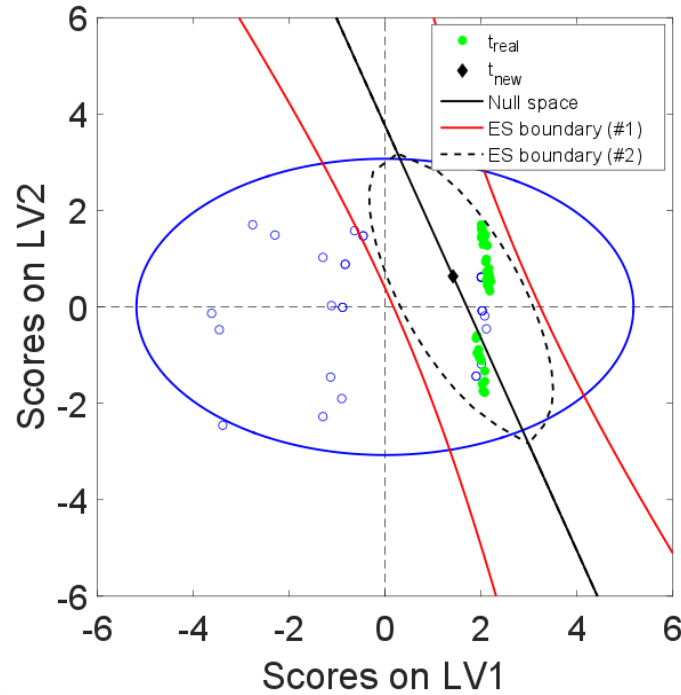


Figure 3.12. Case study #4: bracketing the true input combinations (green circle) yielding a desired intravoid fraction of solids out of the roller compactor ($y_{des} = 0.6341 \text{ m}^3/\text{m}^3$).

3.9 Results for case study #5

Fourteen calibration samples are randomly-chosen from the historical dataset, and a PLS model accounting for 67.2 % of the total variability of \mathbf{X} and 79.3 % of the total variability of \mathbf{y} is built based on this calibration data set. The problem of bracketing the knowledge space to provide guidance for the experimental campaign of granulate with a percent of oversize granules equal to 74% is considered. The results are shown in Fig. 3.13. It can be noticed that, also in this case, the segmentation of the knowledge space is effective for both uncertainty back-propagation models, being able to bracket the real input combination that yield the desired product. The reduction of the experimental effort that can be obtained for this case study are $ESS_1 = 56.1\%$ and $ESS_2 = 66.3\%$.

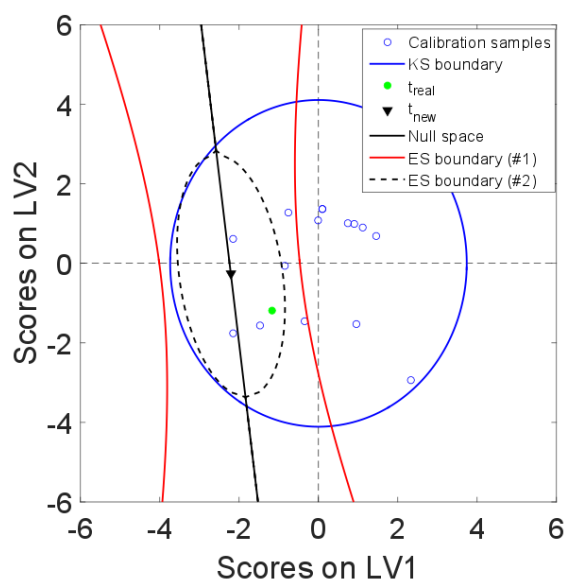


Figure 3.13. Case study #5: bracketing the true input combinations (green dot) that gives a granulate with a percent of oversize granules equal to 74% in a high-shear wet granulation process.

3.10 Conclusions

In this study, the possibility of using PLS model inversion to assist the determination of the design space of a new pharmaceutical product has been discussed. A methodology to back-propagate the uncertainty on the prediction of the PLS model to the calculated inputs has been developed, and the practical outcomes that can be obtained in pharmaceutical product development have been discussed. Two uncertainty back-propagation models have been formulated, and their performance critically compared. The first model assumes that all the prediction uncertainty is back-propagated to the PLS scores, whereas the second one attempts to allocate the overall prediction uncertainty on both the scores and the loadings. From a mathematical viewpoint, the two models allow one to determine the covariance matrix of a point belonging to the solution space of the PLS inversion problem by making some assumptions on how uncertainty back-propagates from the y -space to the latent space first, then to the real input space. The final objective of the methodology is to build confidence regions for the calculated solution of the inversion problem, thus identifying an operational space (i.e., a set of combinations of raw material properties and process parameters) that brackets the real DS of the product at the predefined confidence level. From a practitioner's perspective, the proposed methodology can be used to substantially shrink the space within which experiments must be carried out to identify the DS, since the experimental campaign can be focused on this restricted operating window instead of the entire historical decision space. Moreover, the proposed procedure returns an accurate quantification of uncertainty for quality risk management.

The effectiveness of the proposed methodology has been tested on three different case studies. The ability of the two uncertainty back-propagation models to effectively bracket the real DS of the new product and to shrink the experiment space have been assessed. The first model is more conservative than the second, in the sense that it guarantees a slightly higher performance in bracketing the true DS of the product; however, it also results in a less effective shrinkage of the experiment space.

Some areas of further investigation can be pointed out. The proposed procedure can be extended to multivariate product specifications, given the availability of a reliable model to estimate the prediction covariance matrix in PLS modelling. The development of such a model is still an open research area. Moreover, more sophisticated and rigorous relations may be derived to relate the uncertainty on the model outputs to the uncertainty on the model inputs. These formulas should account for the intrinsic reduced-rank nature of PLS modeling. Additionally, uncertainty in the input and output measurements should be accounted for.

Chapter 4

A Bayesian/latent variable approach for design space determination*

In this Chapter, a methodology that combines latent-variable modeling and multivariate Bayesian regression is presented in order to identify a subset of input combinations (process operating conditions and raw materials properties) within which the design space (DS) of a new pharmaceutical product will lie at a probability equal to, or greater than, an assigned threshold. The methodology is tested on three relevant pharmaceutical case studies, two of which involving experimental data. The ability of the proposed approach to obtain a probabilistic identification of the DS, while simultaneously reducing the computational burden for the discretization of the input domain and providing a simple graphical representation of the DS, is shown.

4.1 Introduction

In the last decade, the pharmaceutical regulatory agencies have strongly encouraged the pharmaceutical industries to gain deeper understanding of their manufacturing processes. The International Conference on Harmonization (ICH) Q8 guideline (ICH, 2009) emphasized the Quality by Design (QbD) approach, according to which quality should be built into the product since its conception, and not simply tested after the manufacturing process is completed. This requires understanding how input materials properties and process operating conditions can affect the product quality. A key concept introduced by this guideline is that of design space (DS), defined as the “multidimensional combination and interaction of input variables (e.g. material attributes) and process parameters that have been demonstrated to provide assurance of quality” (ICH, 2009). As long as raw material properties and process parameters are changed within the approved DS, no regulatory process of post-approval change is required. Product quality is expressed in terms of critical quality attributes (CQAs), which must lie within

* Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018) Probabilistic design space determination in pharmaceutical development: a Bayesian/latent variable approach. *AIChE J.* **64**, 2438- 2449.

acceptance limits defined *a priori*. “Assurance of quality” implicitly requires *quantifying* the confidence of the manufacturer at providing products with the desired quality.

In this study, we address the problem of identifying the DS of a new pharmaceutical product. From a general perspective, assisting the identification of the DS requires both experimentation and mathematical modeling. When models are used to identify the DS, different sources of uncertainty can affect the DS identification exercise, e.g. model parameter uncertainty (Faber and Kowalski, 1997), measurement uncertainty in the calibration dataset (Reis and Saravia, 2005), and structural uncertainty due to inadequacy of the model structure. However, as observed by Peterson (Peterson, 2008), most of the available modeling techniques used to assist a DS identification exercise do not account for model uncertainty, and some of them do not even fully account for the multivariate correlation between model inputs and outputs. For example, overlapping mean response surface methodologies (Anderson and Whitcomb, 1998; Zidan *et al.*, 2007) fail to account for model parameter uncertainty and for the correlation between the CQAs. The multivariate correlation between inputs and outputs is inherently considered in latent-variable modeling (LVM; Tomba *et al.*, 2013). However, accounting for model uncertainty when using this approach may be complex and at present has been shown only for products characterized by a single CQA, i.e. when quality is a univariate property (Bano *et al.*, 2017).

A comprehensive way to account for model parameter uncertainty and correlation between the CQAs, as well as to provide a metric for the assurance of quality, is to use a Bayesian posterior predictive approach. The idea of Bayesian or probabilistic DS was first introduced by Peterson (Peterson, 2004). With this approach, the posterior predictive probability to observe the CQAs within their predefined specifications is computed. DoE techniques are used to determine appropriate experimental conditions for the design of multivariate linear regression models between the input factors and the CQAs (Peterson and Yahyah, 2009; Stockdale and Cheng, 2009; Debrus *et al.*, 2011). The multivariate regression model is then used to predict the DS in a Bayesian framework. The input domain is first discretized using a sampling algorithm. Samples from the joint posterior distribution of the CQAs are then drawn using Markov-Chain Monte Carlo (MCMC) techniques for each discretization point of the input domain. The probability of the CQAs to meet their specifications for the given point is then recorded. The DS is finally identified as the subspace of the original input domain within which the probability of the CQAs to meet their specifications is higher than a user-defined acceptance level.

Several studies proved the ability of the above approach to quantify the concept of “assurance of quality” defined by ICH (Stockdale and Cheng, 2009; Debrus *et al.*, 2011; Peterson and Lief, 2010). However, two issues are still unresolved: (i) how to obtain a computationally tractable discretization of the multidimensional input domain when several input factors are involved, and (ii) how to express the results of the DS identification exercise in an intuitive and compact way (e.g., graphically). With respect to the first issue, the typical approach for the discretization

of the input domain is to use space-filling algorithms. However, these algorithms suffer from the curse of dimensionality (Cioppa and Lucas, 2007), i.e. the computational burden increases very significantly with the number of input factors. With reference to the second issue, the results of a multidimensional DS identification exercise are usually expressed in a tabular way, or as a matrix of two-dimensional probability maps by considering two inputs at a time and fixing the other factors to their nominal values. Both these approaches are difficult to interpret and fail in giving an intuitive representation of the multivariate DS. The aim of this study is to exploit the advantages of LVM to overcome both issues, while maintaining the Bayesian posterior predictive approach for DS identification.

We assume that a historical dataset of products “similar” (Jaeckle and McGregor, 1998) to the one under development is available. The historical dataset is first used to build a partial least-squares (PLS) regression model. PLS is used to obtain a linear transformation between the original multidimensional input domain and a low-dimensional latent space, and the Bayesian posterior predictive approach proposed by Peterson (Peterson, 2004) is applied to identify the probabilistic DS. The reduction of the input space dimensionality obtained with PLS is exploited to obtain a computationally efficient discretization of the input domain and to express the DS in an intuitive 2-dimensional graphical way, thus overcoming two of the main limitations of the approach proposed by Peterson. The practical outcome of the proposed methodology is the identification of the DS in a probabilistic fashion, which can be used to demonstrate the confidence the manufacturer has on the product meeting its quality targets when the manufacturing process is run within the proposed DS.

The methodology is tested on three case studies. The first one involves a model-generated historical dataset for the dry granulation of a pharmaceutical blend (taken from the work of Facco *et al.* (2015), based on the model of Johanson (1965)). In the second case study, a high-shear wet granulation process of a pharmaceutical blend is addressed, using experimental data from Oka *et al.* (2015). The third case study involves the roller compaction an intermediate drug load formulation, with experimental data taken from Souihi *et al.* (2015) .

4.2 Review of Partial least-squares (PLS) regression

Let $[\mathbf{X}; \mathbf{Y}]$ be a historical dataset of “old” products similar¹⁸ to the one under development. Matrix $\mathbf{X} [I \times Q]$ collects I samples concerned with a set of Q process inputs (raw material properties and process parameters) that affect the M quality attributes $\mathbf{Y} [I \times M]$ of the historical products.

Partial least-squares (PLS) regression (Wold *et al.*, 1983; Geladi and Kowalski, 1986) is a multivariate regression technique that projects a dataset of input and response variables onto a

¹⁸ Product similarity can be assessed using the methodology proposed by Jaeckle and MacGregor (1998).

common latent space of A new variables, called latent variables (LVs). This set of new variables is determined by finding the multidimensional direction on the \mathbf{X} -space that explains the maximum multidimensional variance direction in the \mathbf{Y} -space, according to the model structure:

$$\mathbf{X} = \mathbf{TP}^T + \mathbf{E}_X \quad (4.1)$$

$$\mathbf{Y} = \mathbf{TQ}^T + \mathbf{E}_Y \quad (4.2)$$

$$\mathbf{T} = \mathbf{XW}^* \quad (4.3)$$

where \mathbf{T} [$I \times A$] is the score matrix, \mathbf{P} [$Q \times A$] and \mathbf{Q} [$M \times A$] are the \mathbf{X} and \mathbf{Y} loading matrices, \mathbf{E}_X and \mathbf{E}_Y the residuals; \mathbf{W}^* [$Q \times A$] is the weight matrix.

The number $A \ll Q$ of LVs can be chosen in such a way as to explain a meaningful fraction of the variance of both the input and output data. In the presence of a strong input collinearity, a significant reduction of the problem dimensionality can be obtained by applying this data compression technique. This property will be exploited as later illustrated in Section 4.4.

The subspace of the input domain enclosed by the historical dataset is called the knowledge space (KS; MacGregor and Bruwer, 2008). Mathematically, the KS can be identified in the latent space as the hyper-ellipsoidal confidence region determined by the 95% confidence limit for the the Hotelling's T^2 statistic (Tracy *et al.*, 1992):

$$T_{lim}^2 = \frac{A(I^2-1)}{I(I-A)} F_{A,(I-A),0.95} \quad , \quad (4.4)$$

where $F_{A,(I-A),0.95}$ is the value of the 95% percentile of the F -distribution with A and $(I - A)$ degrees of freedom. The hyper-ellipsoid semiaxes are given by:

$$s_a = \sqrt{\lambda_a T_{lim}^2} \quad a = 1, \dots, A \quad (4.5)$$

where $\lambda_a, a = 1, \dots, A$ is the variance of the scores related to the a -th LV.

4.3 Bayesian design space: mathematical formulation

Let $\mathbf{x} = (x_1, x_2, \dots, x_Q)$ be the vector that collects the set of inputs and $\mathbf{y} = (y_1, y_2, \dots, y_M)$ the one that collects all the response variables. Let $\mathbf{h}(\mathbf{x})$ be a model relating the inputs to the outputs. From a modeling perspective, the DS of the new product can be interpreted as the region of the input domain where the corresponding predicted model response $\mathbf{y} = \mathbf{h}(\mathbf{x})$ has an acceptable probability to satisfy assigned specifications. In this study, only the subspace of the DS lying within the projection of the KS onto the latent space is explored. To simplify the notation, we will refer to this subspace as to the DS itself.

The desired quality characteristics for the product to be developed can be expressed in terms of acceptance criteria. The region where all product quality attributes meet their acceptance criteria is defined as the acceptance region (AR).

A risk-based definition of the DS is the one proposed by Peterson (Peterson, 2008):

$$DS = \{\mathbf{x} \in KS: \Pr(\mathbf{y} \in AR|\mathbf{x}, \mathbf{X}, \mathbf{Y}) \geq \theta_{th}\} \quad (4.6)$$

where θ_{th} is an assigned probability threshold for the product quality to be acceptable. Eq. (4.6) clarifies that this study considers the DS as the subspace of the historical KS where the posterior probability that the product attributes will all lie within their acceptance criteria is greater than an assigned threshold, conditionally on the historical dataset $[\mathbf{X}; \mathbf{Y}]$. The DS as defined by (4.6) will be denoted as probabilistic or Bayesian DS.

The definition of DS as in (4.6) allows for a quantification of the concept of “assurance of quality” as advocated by ICH. The benefits of this Bayesian approach have been extensively discussed (Lebrun et al., 2013). Eq. (4.6) is general and can be used whenever a model relating the product quality \mathbf{Y} to the raw material properties and process parameters \mathbf{X} is available. In this study, the model used to relate the product quality to the input factors is a multivariate linear regression model:

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E} \quad (4.7)$$

where $\mathbf{B} [Q \times M]$ is the matrix of model parameters and $\mathbf{E} [I \times M] = [\mathbf{e}_1, \dots, \mathbf{e}_I]^T$ is the matrix of the residuals. The residuals are assumed to be independent and normally distributed with mean $\mathbf{0} [I \times M]$ and covariance $\mathbf{\Sigma} [M \times M]$, i.e. $\mathbf{e}_i \propto N(\mathbf{0}, \mathbf{\Sigma}), i = 1, \dots, I$.

The posterior predictive probability included in (4.6) can be determined using a Bayesian approach with the multivariate regression model (4.7). To this purpose, a posterior predictive distribution for future responses must be obtained.

Let $f(\mathbf{y}|\mathbf{x}, \mathbf{B}, \mathbf{\Sigma})$ be the probability density function for the new response \mathbf{y} under the set of inputs \mathbf{x} , given the set of uncertain model parameters \mathbf{B} and the covariance of the residuals $\mathbf{\Sigma}$. This pdf depends on the model adopted. The posterior predictive distribution for \mathbf{y} , conditional on the historical data $[\mathbf{X}; \mathbf{Y}]$, is given by:

$$g(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int \int f(\mathbf{y}|\mathbf{x}, \mathbf{B}, \mathbf{\Sigma})p(\mathbf{B}, \mathbf{\Sigma}|\mathbf{X}, \mathbf{Y})d\mathbf{B}d\mathbf{\Sigma} \quad (4.8)$$

where $p(\mathbf{B}, \mathbf{\Sigma}|\mathbf{X}, \mathbf{Y})$ is the joint posterior distribution of the model parameters \mathbf{B} and the covariance of the residuals $\mathbf{\Sigma}$. The probability $\Pr(\mathbf{y} \in AR|\mathbf{x}, \mathbf{X}, \mathbf{Y})$ of (4.6) can then be obtained from $g(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y})$ by simple integration within the entire AR:

$$\Pr(\mathbf{y} \in AR|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int_{AR} g(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) d\mathbf{y}. \quad (4.9)$$

Eq. (4.9), combined with (4.8), indicates that the probability $\Pr(\mathbf{y} \in AR|\mathbf{x}, \mathbf{X}, \mathbf{Y})$ can be determined once the joint posterior distribution of the model parameters $p(\mathbf{B}, \mathbf{\Sigma}|\mathbf{X}, \mathbf{Y})$ is derived. This distribution can be obtained using Bayes' theorem:

$$p(\mathbf{B}, \mathbf{\Sigma}|\mathbf{X}, \mathbf{Y}) \propto \mathcal{L}(\mathbf{B}, \mathbf{\Sigma}|\mathbf{Y}) p(\mathbf{B}, \mathbf{\Sigma}) \quad (4.10)$$

where $\mathcal{L}(\mathbf{B}, \mathbf{\Sigma}|\mathbf{Y})$ is the likelihood function, given by (Tiao and Zellner, 1964):

$$\mathcal{L}(\mathbf{B}, \mathbf{\Sigma}|\mathbf{Y}) \propto |\mathbf{\Sigma}|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr} [\mathbf{\Sigma}^{-1}(\mathbf{Y} - \mathbf{XB})^T(\mathbf{Y} - \mathbf{XB})]\right) \quad (4.11)$$

and $p(\mathbf{B}, \mathbf{\Sigma})$ is the joint prior distribution of the model parameters \mathbf{B} and $\mathbf{\Sigma}$.

MCMC sampling techniques can be used to draw samples from the joint posterior distribution of the model parameters $p(\mathbf{B}, \mathbf{\Sigma}|\mathbf{X}, \mathbf{Y})$. The response posterior predictive distribution $g(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y})$ can then be obtained from this distribution and finally $\Pr(\mathbf{y} \in AR|\mathbf{x}, \mathbf{X}, \mathbf{Y})$ can be computed. The sampling algorithm that has been used in this study is the Metropolis-Hastings algorithm (Hastings, 1960).

The choice of the prior distributions for the model parameters and the covariance of the residuals is critical. Non-informative or informative priors can be used to this purpose. When informative priors are used, the amount of prior information available determines the choice of the parameters of these distributions. Strictly speaking, if no or very limited prior information is available, prior distributions that are considered informative from a structural point of view can be made uninformative by choosing uninformative prior parameters. Stated differently, when informative priors are used, the amount of prior information can be tuned by choosing the values of the parameters of these prior distributions.

In this study, informative (from a structural point of view) prior distributions for the model parameters and the covariance of the residuals have been chosen. Namely, $p(\mathbf{B}, \mathbf{\Sigma}) = p(\mathbf{B}|\mathbf{\Sigma})p(\mathbf{\Sigma})$ with $p(\mathbf{B}|\mathbf{\Sigma})$ matrix-variate normal distribution with mean \mathbf{B}_0 and covariances $\mathbf{\Sigma}$ and $\mathbf{\Sigma}_0$, i.e. $p(\mathbf{B}|\mathbf{\Sigma}) \propto N(\mathbf{B}_0, \mathbf{\Sigma}, \mathbf{\Sigma}_0)$, and $p(\mathbf{\Sigma})$ inverse Wishart distribution with scale matrix $\mathbf{\Omega}$ and degrees of freedom ν_0 , i.e. $p(\mathbf{\Sigma}) \propto W_{ish}^{-1}(\mathbf{\Omega}, \nu_0)$. The choice of the prior parameters $(\mathbf{B}_0, \mathbf{\Sigma}_0, \nu_0)$ determines the amount of prior information available. The values of these parameters have been set as uninformative as possible, since no prior information for the case studies considered was available.

A graphical representation of the discussed Bayesian approach for the determination of the posterior predictive probability of a new response is shown in Fig. 4.1.

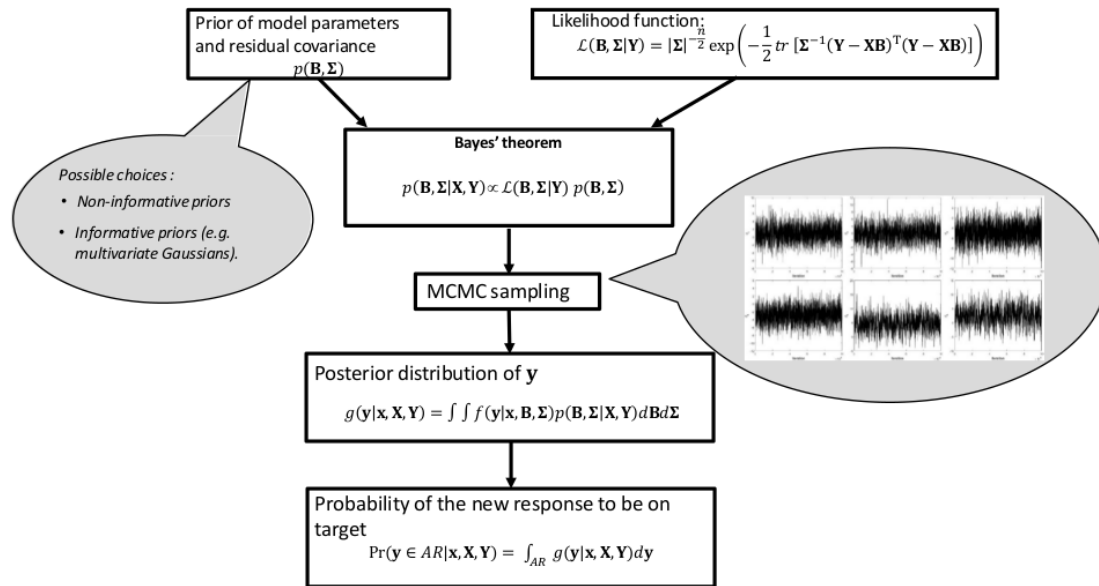


Figure 4.1. Diagram of the Bayesian approach to obtain the posterior predictive probability of a new response.

4.4 Bayesian identification of the design space

Identification of the DS according to the Bayesian definition (4.6) involves the following steps.

- 1) Identification of the critical process parameters (CPPs) and critical quality attributes (CQAs) for the new product.
- 2) Acquisition of an experimental dataset, by performing experiments according to a DoE-based plan.
- 3) Calibration of a model to relate the CPPs and raw material properties to the CQAs. In most situations, a statistical multivariate linear regression model is built at this stage.
- 4) Identification of the multidimensional knowledge space (i.e., input domain).
- 5) Discretization of the multidimensional KS. Different sampling algorithms can be used to this purpose. Space-filling algorithms are mostly used, which become computationally expensive as the number of input factors increases.
- 6) Determination (by MCMC simulations) of the joint posterior probability of the CQAs to meet their specifications for every discretization point of the KS.
- 7) Probabilistic reconstruction of the DS according to (4.6).
- 8) Representation of the multidimensional DS with a tabular approach or a matrix of 2-dimensional probability maps.

Notice that the DoE inputs are the CPPs and the raw materials properties. Hence, they may be large in number and correlated, which has some notable drawbacks. Firstly, this complicates the DoE exercise, because independent factors should be singled out (step 2). Secondly, investigating the entire KS by experiments (step 2) may be hard or even impossible. Thirdly,

the discretization of the KS with space-filling algorithms (step 5) may be computationally very intensive. Finally, visualizing the DS (step 8) may be complex. Two-dimensional posterior probability maps are often derived for each pair of model inputs, with the remaining inputs fixed at their nominal values. However, this approach does not account for the multivariate correlation structure of the model inputs, thus leading to results that are partial at best, but possibly even misleading.

As discussed in Section 4.2, PLS can be used to project the input space onto a low-dimensional latent space, while simultaneously accounting for its correlation with the output space. Coupling PLS with the multivariate linear regression Bayesian approach discussed in Section 4.3 can be a way to enhance the advantages of the latter (quantification of assurance of quality; handling of model parameter uncertainty and of the multivariate nature of the responses) with the ability of PLS to reduce the dimensionality of the multivariate input space. In this study, a systematic methodology is proposed to take full advantage of both techniques. In more detail, the ability of PLS to reduce the input space dimensionality is exploited to improve Step 5 and Step 8 of the above procedure. The methodology is discussed below.

4.4.1 Proposed methodology

Figure 4.2 provides an illustrative graphical representation of the methodology for a 2-LV space.

The step-by-step procedure is as follows.

- a. Given the historical dataset $[\mathbf{X}; \mathbf{Y}]$, a PLS model with A LVs relating \mathbf{X} to \mathbf{Y} is built. In this study, A has been chosen according to the eigenvalue-greater-than-one rule (Kaiser, 1970) in such a way as to explain a significant fraction of the variance of both \mathbf{Y} and \mathbf{X} . Once the PLS model has been calibrated, the KS is identified as the region defined by the 95% confidence limit of the Hotelling's T^2 statistic (Section 4.2). The boundary of the KS is identified by the ellipse of Fig. 4.2a for a two-dimensional latent space.
- b. The KS is discretized, e.g. by using the following approach:
 - (i) a very large (e.g., $N_a > 1000$) number of samples in the latent space is chosen inside a unit hyper-sphere centered in the origin;
 - (ii) the samples are then scaled to the confidence hyper-ellipsoid size by multiplying their coordinates by the lengths of the ellipsoid axes.

An example of KS discretization using this approach is shown in Fig. 4.2b for $N_a = 2000$.

- c. A representative number $N_s < N_a$ of samples \mathbf{t}_l , $l = 1, 2, \dots, N_s$ is chosen. The corresponding values in the original input space \mathbf{x}_l , $l = 1, 2, \dots, N_s$ are obtained according to Eq. (4.1) using the PLS loadings found in step #1. The choice of N_s must be done so as to compromise between having a good coverage of the entire KS and

reducing the computational effort. The next subsection discusses how this can be achieved. Once N_s has been determined, a space-filling algorithm is used to span the KS using the defined number of samples. In this study, the Kennard-Stone algorithm (Kennard and Stone, 1969) has been used. The Kennard-Stones algorithm is a sequential sampling algorithm that allows selecting a subset of N_s samples given a set of N_a candidates. The rationale behind this algorithm is to select each new sample by choosing the farthest sample (in terms of Euclidean distance) from the previously selected ones, starting from the two farthest points in the KS. As other space filling algorithms, the Kennard-Stone algorithm suffers from dimensionality issues. Therefore, its implementation in the latent space is much less computationally expensive than in the true input space.

An example of the selection of $N_s = 100$ samples using the Kennard-Stone algorithm is shown in Fig. 4.2c.

- d. For each of the N_s selected samples, the joint posterior predictive distribution (PPD) of the quality responses is obtained by performing a Bayesian simulation in the true input space (Section 4.3).

An example of posterior predictive distribution for a case study involving a single quality specification is shown in Fig. 4.2d.

- e. Given the desired quality target \mathbf{y}_{des} or the desired quality interval $\mathbf{y}_p < \mathbf{y}_{des}$ (or $\mathbf{y}_{lo} < \mathbf{y}_p < \mathbf{y}_{up}$, with \mathbf{y}_{lo} and \mathbf{y}_{up} lower and upper bounds for the quality variables respectively), if the condition expressed by (4.6) is satisfied, then the given sample is considered to belong to the probabilistic DS (green diamonds in Fig. 4.2e). If (4.6) is not satisfied, the given sample is rejected as it does not belong to the probabilistic DS (red diamonds of Fig. 4.2e).

A schematic flowchart involving all the different steps of the proposed methodology is shown in Fig. 4.3.

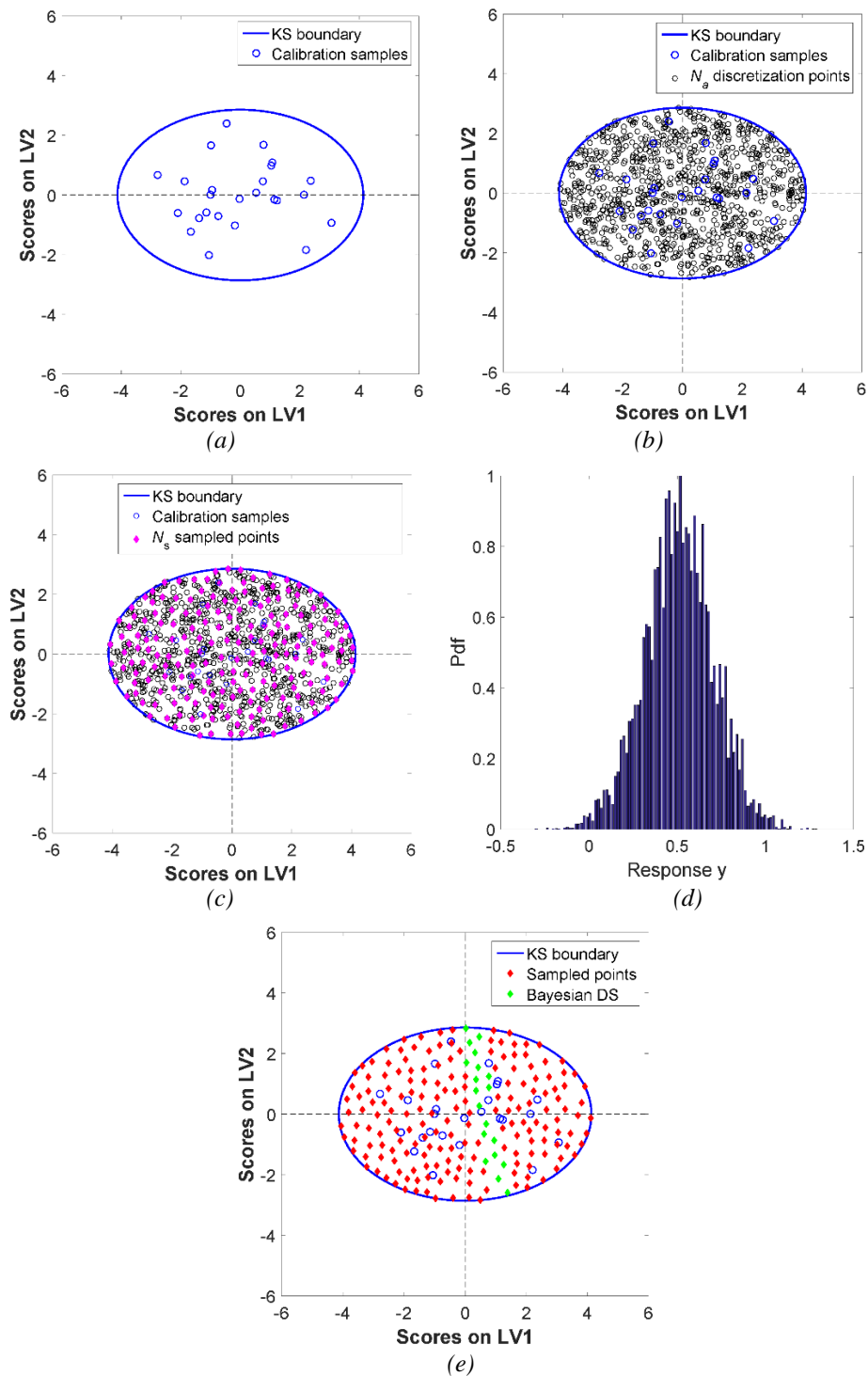


Figure 4.2. Graphical representation of the proposed methodology for probabilistic DS determination in a two-dimensional latent space: (a) KS identification in the latent space; (b) discretization of the KS; (c) sampling of the selected point with KSA; (d) posterior predictive distribution for a single sample; (e) Bayesian DS.

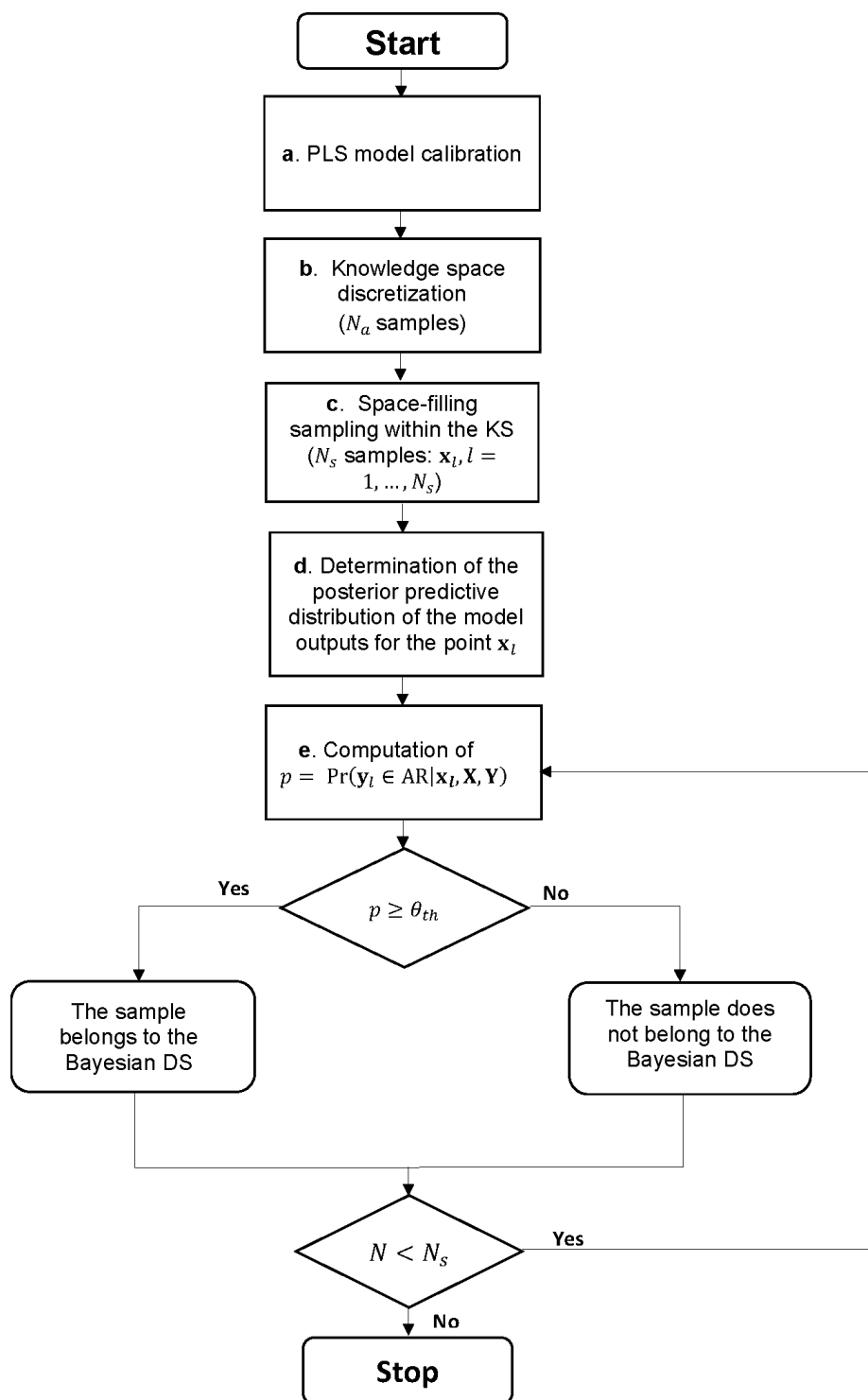


Figure 4.3. Flowchart of the proposed methodology for the determination of the Bayesian design space of a new pharmaceutical product.

4.4.2 Selection of the optimal number of samples

This section presents an automatic methodology to determine the optimal number N_s^{opt} of discretization samples of step 3 discussed in the previous section.

Consider the PLS model (4.1)-(4.3). After calibration, the PLS model can be used to predict a new response \mathbf{y}_p given a new observation \mathbf{x}_p according to :

$$\mathbf{y}_p = \mathbf{t}_p \mathbf{Q}^T \quad (4.12)$$

where \mathbf{t}_p is the score vector of the new observation \mathbf{x}_p .

The score vector \mathbf{t}_p can be related to the input vector \mathbf{x}_p through the model weights \mathbf{W} according to:

$$\mathbf{t}_p = \mathbf{x}_p \mathbf{W}^*. \quad (4.13)$$

As shown elsewhere (Zhang and Garcia Munoz, 2009) both the score vector \mathbf{t}_p and the loadings matrix \mathbf{Q}^T of (4.12) are affected by uncertainty, and accordingly weights, too. Eq. (4.13) can be seen as linear prediction model between the model “response” \mathbf{t}_p and the model “regressors” \mathbf{x}_p , through the model “parameters” \mathbf{W}^* . Therefore, the PPD of the score vector \mathbf{t}_p can be obtained by performing a Bayesian calibration of model (4.13) by feeding a new (fixed) set of model regressors \mathbf{x}_p^* .

The choice of the optimal number of discretization samples N_s can therefore be done as follows. Let $[N_1, N_2, \dots, N_M]$ be a set of discrete values for N_s , with $N_s = N_1 < N_2 < \dots < N_M$. The PPD of the score vector \mathbf{t}_p is obtained for each of these discrete values and for the given fixed value \mathbf{x}_p^* adopting the Bayesian procedure presented in Section 4.3. The conditional posterior predictive distributions (CPPDs) of each component of \mathbf{t}_p can then be obtained by picking out the contributions of every single component to the PPD.

An illustrative example of the two CPPDs obtained with a PLS model built on $A = 2$ LVs (i.e., $\mathbf{t}_p = [t_{p,1}, t_{p,2}]$) is shown in Fig. 4.4, where five discretization samples ($N_s = 5$) have been used.

The 95% credible intervals for $t_{p,a}$ ($a = 1, \dots, A$) can then be computed for each value of N_s . Mathematically, for $N_s = N_1, N_2, \dots, N_M$, the 95 % credible intervals of $t_{p,a}, a = 1, \dots, A$ can be expressed as:

$$t_{p,1}^{lo}(N_s) < t_{p,1} < t_{p,1}^{up}(N_s) \quad (4.14)$$

...

$$t_{p,a}^{lo}(N_s) < t_{p,a} < t_{p,a}^{up}(N_s) \quad (4.15)$$

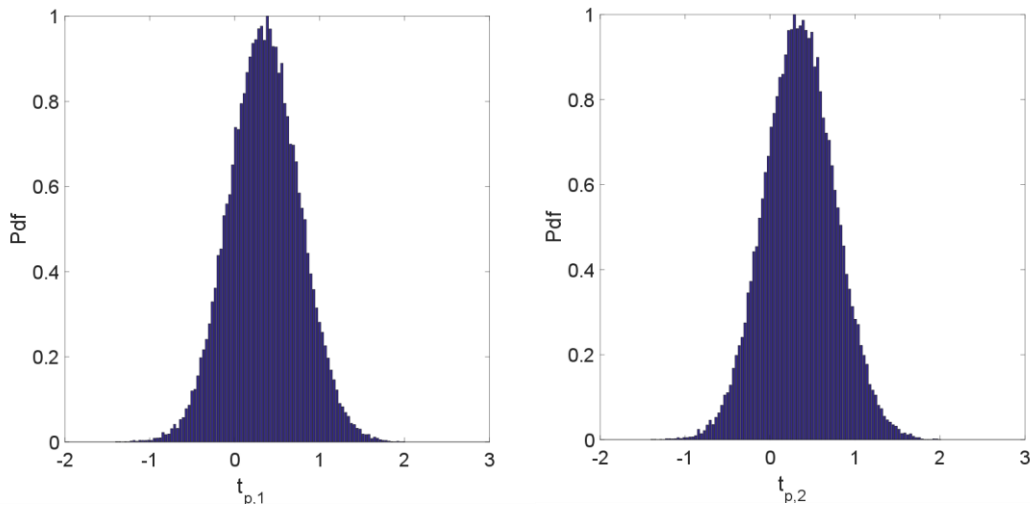


Figure 4.4. Illustrative example of the CPPDs of the two components of the score vector \mathbf{t}_p with 5 discretization samples.

where $t_{p,a}^{lo}$ and $t_{p,a}^{up}$ are the lower and upper bound of the credible interval for the a -th component of \mathbf{t}_p .

Let $w_a(N_s) = t_{p,a}^{up}(N_s) - t_{p,a}^{lo}(N_s)$, $a = 1, \dots, A$ be the width of the credible intervals for the a -th component of the score vector, and s_a the semiaxis of the KS on the a -th LV as defined by Eq. (4.5). The ratio:

$$\alpha_a(N_s) = \frac{w_a(N_s)}{s_a} \quad a = 1, \dots, A \quad (4.16)$$

can be used as a measure of the amount of variability captured by the N_s discretization samples for the a -th LV of the PLS model. For each value of N_s , the respective values of α_a , $a = 1, \dots, A$ can then be computed. Simple regression models f_a , $a = 1, 2, \dots, A$ can then be built to relate the behavior of α_a , $a = 1, \dots, A$ with respect to the number of discretization points N_s . Mathematically, the following regression models can be obtained:

$$\alpha_a(N_s) = f_a(N_s) \quad a = 1, \dots, A \quad (4.17)$$

Notice that regression models (4.17) can be built with a small set of values of N_s . Moreover, the maximum value assumed by N_s (i.e. N_M) is relatively small (typically, less than 25). Therefore, the derivation of the regression models (4.17) is computationally very fast and does not significantly contribute to the overall computational burden.

We define the total amount of variability $\Psi(N_s)$ captured by the N_s discretization samples on the latent space as a weighted sum of the metrics $\alpha_a(N_s)$ for all the a -th components, the weights being the amount of \mathbf{X} -variability explained by the given LV. Mathematically:

$$\Psi(N_s) = \sum_{a=1}^A R_{X,a}^2 \alpha_a(N_s). \quad (4.18)$$

The optimal number N_s^{opt} of discretization points can be derived as the number of samples that maximizes $\Psi(N_s)$ (i.e. that maximizes the total amount of variability captured by the N_s discretization samples) while keeping the total computational time to a minimum. For simplicity, the total computational time is considered as the product of the time τ_s required to perform a simulation for a single sample and the number of samples N_s :

$$T_s = N_s \cdot \tau_s. \quad (4.19)$$

N_s^{opt} can then be computed as:

$$N_s^{opt} = \min_{N_s} [-\Psi(N_s) + N_s \cdot \tau_s]. \quad (4.20)$$

The above strategy for the selection of the optimal number of discretization points has been implemented in MATLAB[®] and has been used for all the Bayesian simulations involved in this study. All the Bayesian Monte Carlo simulations for the multivariate linear regression model (4.7) have been implemented in MATLAB v. 2015b. Parallel computing (12 parallel threads) has been used to speed up the simulations on an Intel[®] Core[™] I7-6700 CPU@3.40GHz processor with 16.0 GB RAM.

4.5 Case studies

4.5.1 Case study #1: dry granulation of a pharmaceutical blend by roller compaction

The first case study concerns the dry granulation of a pharmaceutical blend by roller compaction. Historical data for this case study were generated by Facco *et al.* (2015) based on the model of Johanson (1965). The modeling environment of gSOLIDS[®] (Process Systems Enterprise, 2014) was used to this purpose.

The historical dataset is composed by eight input variables (compressibility factor, roller diameter, roller width, roller speed, pressure force, friction angle between solid granulate and roller compactor, effective friction angle and springback factor) and one response variable (intravoid fraction of the solids out of the compactor). A summary of the input/output variables and the characterization of the input dataset is reported in Table 4.1. Full details on the derivation of the historical dataset can be found in the cited reference.

Table 4.1. Case study #1: list of the input and response variables (data from Facco *et al.* (2015), based on the model of Johanson (1965)) and characterization of the input dataset (columns 5 and 6).

ID	Variable name	Units	Symbol	Mean	Std.
<i>Inputs</i>					
1	Compressibility factor	[-]	K	9.85	2.53
2	Roller width	[m]	S_{roll}	0.13	0.02
3	Roller diameter	[m]	D_{roll}	0.40	0.07
4	Roller speed	[Hz]	v_{roll}	0.1707	0.1072
5	Pressure force	[kN]	F_{roll}	13866.67	6951.19
6	Friction angle	[rad]	γ_{FR}	27.51	8.78
7	Effective friction angle	[rad]	γ_{EFF}	48.17	31.86
8	Springback factor	[-]	F_{sb}	0.11	0.03
<i>Response</i>					
R1	Intravoid fraction of solids	$[m^3/m^3]$	ϕ_s	[-]	[-]

The historical dataset presented in Facco *et al.* (2015) includes 80 calibration samples. In this study, 20 calibration samples have been randomly selected from the historical dataset. These data have been collected in the input calibration matrix \mathbf{X} $[20 \times 8]$, and in the response calibration vector \mathbf{y} $[20 \times 1]$.

4.5.2 Case study #2: high-shear wet granulation of a pharmaceutical blend

This case study concerns the high-shear wet granulation of a pharmaceutical blend (API+excipient). Experimental data can be found in the study of Oka *et al.* (2015) and are based on a full factorial DoE ($3 \times 3 \times 3$). Three input variables are considered (impeller speed, liquid to solid ratio L/S , wet massing time) and their effect on the particle size distribution (PSD) of the granulate has been studied. Twenty-seven samples are available.

Three different scenarios of increasing complexity are considered:

1. *scenario 1 (univariate)*: only the median particle size d_{50} is considered to characterize the quality of the granulate;
2. *scenario 2 (bivariate)*: the quality of the granulate is characterized by the median particle size d_{50} and the 90-th percentile of the PSD (d_{90});
3. *scenario 3 (trivariate)*: the quality of the granulate is characterized by d_{50} , d_{90} and the 10-th percentile of the PSD (d_{10}).

The respective calibration datasets have been collected in the input calibration matrix \mathbf{X} $[27 \times 3]$ and the response matrices \mathbf{y}_1 $[27 \times 1]$, \mathbf{y}_2 $[27 \times 2]$ and \mathbf{y}_3 $[27 \times 3]$ for the three scenarios considered.

A summary of the input-output variables involved for the three scenarios is reported in Table 4.2. In the same table, the target values of the CQAs that have been considered in the three scenarios are reported.

Table 4.2. Case study #2: list of input and response variables for the three scenarios considered (data from Oka *et al.*, 2015).

ID	Variable name	Units	Symbol	Target value
<i>Inputs</i>				
1	Liquid to solid ratio	[-]	L/S	
2	Impeller speed	[Hz]	v_{imp}	[-]
3	Wet massing time	[s]	t_{wet}	[-]
<i>Response</i>				
Scenario 1				
R1	Median particle size	$[\mu m]$	d_{50}	1129
Scenario 2				
R1	Median particle size	$[\mu m]$	d_{50}	1129
R2	90-th percentile	$[\mu m]$	d_{90}	1849
Scenario 3				
R1	Median particle size	$[\mu m]$	d_{50}	1129
R2	90-th percentile	$[\mu m]$	d_{90}	1849
R3	10-th percentile	$[\mu m]$	d_{10}	732

The purpose of analyzing three different scenarios were:

- 1) to test the ability of the proposed methodology to identify the probabilistic DS when the number of quality specifications increases;
- 2) to study the robustness of the methodology, namely to analyze how numerical issues (e.g., convergence of the Markov chain for every Bayesian simulation) may be affected when the problem dimensionality increases.

4.5.3 Case study #3: roll compaction of an intermediate drug load formulation

Experimental data for this case study are taken from the study of Souihi *et al.* (2015). Four input variables (roller force, roller width, roller diameter, dimensionless ratio of gap to roller diameter GD) and 3 product quality variables (ribbon envelope density, ribbon relative density and ribbon porosity) are considered. The historical dataset is composed by 34 calibration samples. A summary of the input-output variables for this case study is reported in Table 4.3.

Table 4.3. Case study #3: list of the input and response variables (data from Souihi et al.(2015))

ID	Variable name	Units	Symbol	Target values
<i>Inputs</i>				
1	Roller force	[kN/cm]	p_{roll}	[-]
2	Roller width	[mm]	s_{roll}	[-]
3	Roller diameter	[mm]	D_{roll}	[-]
4	Ratio of the gap to roll diameter	[-]	GD	[-]
<i>Response</i>				
R1	Ribbon envelope density	[g/cm ³]	ρ_e	1065
R2	Ribbon relative density	[-]	ρ_r	0.72
R3	Ribbon porosity	[%]	P	28

The target values for the CQAs that have been investigated are reported in the same table.

4.6 Results

4.6.1 Results for case study #1

The problem addressed in this case study is the development of a granulate with an intravoids fraction of solids leaving the roller compactor equal to $0.6341 \text{ m}^3/\text{m}^3$. The ability of the proposed methodology is tested to identify a subspace of the historical knowledge space within which the probability that the product specification will be on target is equal to, or greater than, 90%. First, a PLS model is calibrated using the historical dataset. Two LVs are used, explaining 95.1 % of the variability of \mathbf{X} and 89.8 % of the variability of \mathbf{y} . The PLS model allows identifying the KS onto the latent space (Fig. 4.5a). The KS is then discretized according to step 2 of the proposed methodology, and the discretization is shown in Fig. 4.5b. A total of $N_s^{opt} = 242$ representative samples are then chosen according to the procedure discussed in Section 4.4.2. Fig. 4.5c shows the optimal sampling points (red diamonds). Once the sampling point have been chosen, an MCMC simulation is performed for each optimal sample and the posterior predictive probability is calculated. A summary of the computational specifications for this case study is shown in Table 4.4. A multivariate Gaussian distribution is chosen for the model parameters, and an inverse Wishart distribution for the covariance of the residuals. The values of the prior parameters have been set as uninformative, since no prior information on the system was available. $1e5$ MCMC iterations are used for each selected input combination

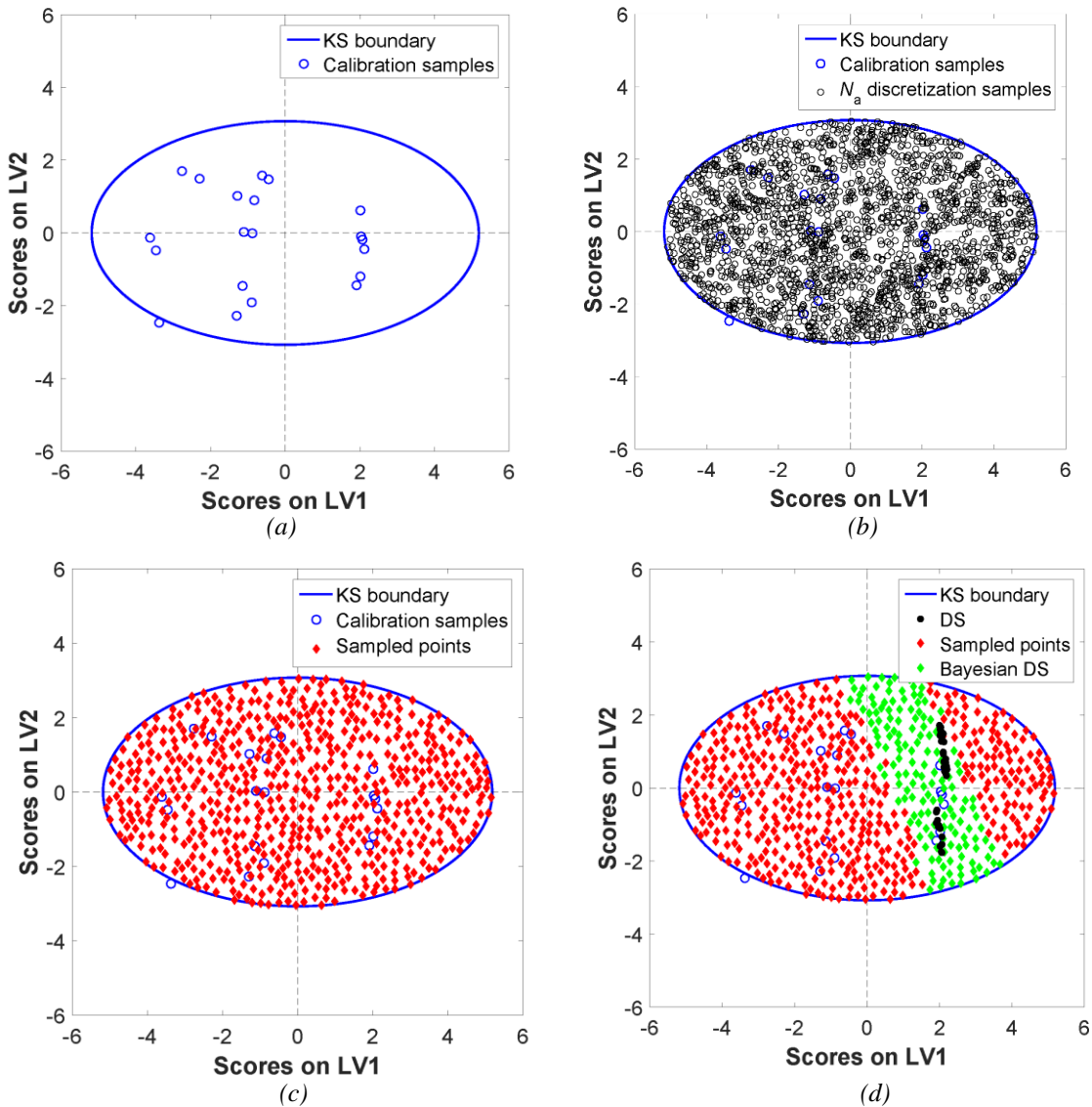


Figure 4.5. Case study #1. (a) KS boundary and calibration samples projected onto the latent space. (b) KS discretization. (c) Sampling within the KS. (d) Bayesian DS and rejected points.

Of them, $1e3$ are discarded (“burn-in” iterations). The optimal number of sampling points for the Kennard-Stone algorithm is 242.

Table 4.3. Case study #1: problem specifications.

Specification	Value
No. of calibration samples	20
No. of model parameters	8
N_s^{opt}	242
Priors of model parameters	Non informative Gaussians
No. of iterations per sample	$1e5$
No. of burn-in iterations	$1e3$
Total simulation time (parallel computing)	3h 35min

The probability of the given sampling points to give a granulate with the desired intravoid fraction of solids is then computed. Finally, the samples (i.e., combinations of raw material properties and process parameters) that allow obtaining the desired granulate with 90% probability are determined. These samples are represented by green diamonds in Fig. 4.5d.

Being this a simulated case study, the first-principles model can be used to determine the set of input combinations that belong to the DS. A trial-and-error approach was used to determine a subset of them, and the result shown by the black circles in Fig. 4.5d was obtained. Note that we denote this subset as the actual DS, even though it more properly represents a subspace of it. It can be seen that the Bayesian DS returns a reliable estimation of the actual DS. The role of the threshold value for the probability can be understood from Fig. 4.6: the larger the probability required for the final product to be on target (i.e., the smaller the risk to manufacture a product *not* meeting its quality specifications), the narrower the subspace of input combinations that can be used.

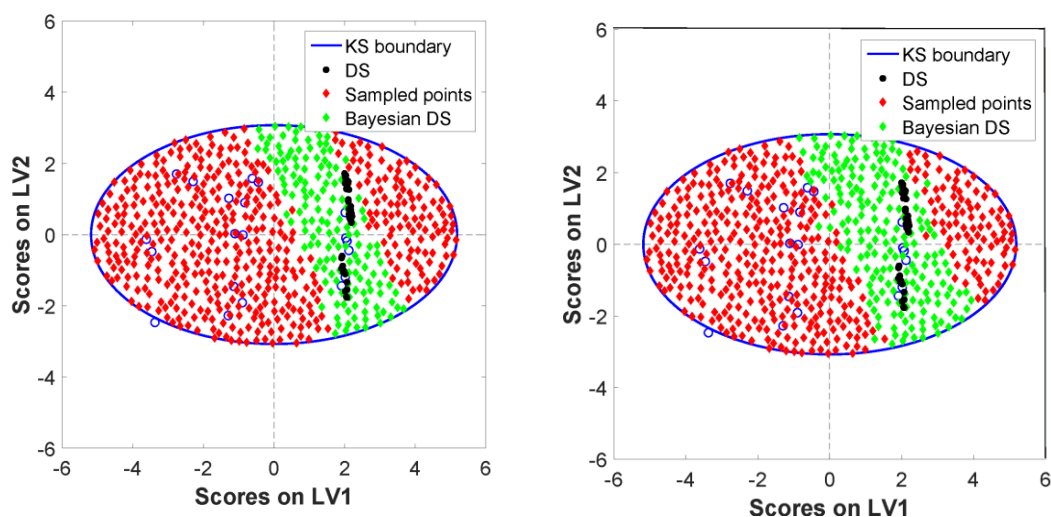


Figure 4.6. Case study #1: effect of the threshold value on the Bayesian DS. (a) Threshold = 90%. (b) Threshold = 80%.

The threshold value should be determined by engineering judgement. Too large values for the threshold may result in rejecting all the selected samples within the knowledge space, since the desired probability cannot be achieved at the given model accuracy (which is limited by definition, since different sources of uncertainty affect the model predictions). On the other hand, too small values for the threshold mean that the DS is identified with a larger uncertainty. In all the case studies considered in the following, a threshold value of 90% is used.

4.6.2 Results for case study #2

One of the 27 historical samples (namely, sample #17) is selected to validate the proposed methodology. The PLS model is then calibrated using the remaining 26 samples. Two LVs are

chosen, accounting for 97.2 % of the total variability of \mathbf{X} and 91.8 % of the total variability of \mathbf{Y} . As discussed in Section 5.2, three scenarios are considered. The target values assigned to the CQAs for each scenario are reported in Table 4.2.

Table 4.5 reports the problem specifications for the three scenarios. Non-informative Gaussian distributions are used for the model parameters, and inverse Wishart distribution for the covariance of the residuals.

Table 4.5. Case study #2: problem specifications for the three scenarios considered.

Specification	Scenario #1	Scenario#2	Scenario #3
No. of calibration samples	26	26	26
No. of model parameters	3	6	9
N_s^{opt}	390	485	491
Priors of model parameters	Non informative Gaussians	Non informative Gaussians	Non informative Gaussians
No. of iterations per sample	1e4	1e5	1e6
No. of burn-in iterations	1e3	1e4	2e5
Total simulation time (parallel computing)	2h 05min	2h 15 min	2h 20min.

The number of required iterations for each sampling point to obtain convergence of the Markov chain increases (from 1e4 to 1e6) with the number of quality specifications. Additionally, a larger number of iterations is required to build the joint PPD of the product quality attributes for the multivariate scenarios with respect to the univariate one. This is reasonable since the number of model parameters for the linear regression model increases with the number of model outputs, also resulting in a slower convergence of the Markov chain for each iteration.

The probabilistic DS projections that can be obtained with the proposed methodology for the three scenarios are shown in Fig. 4.7. The green diamonds represent the sampling points (i.e., input combinations) that have a probability equal to, or higher than, 90% to give the desired granulate. For the univariate scenario (Fig. 4.7a), a narrow portion of the historical KS is determined as the probabilistic DS for the product. This subset of input combinations does include the actual input combination that allows obtaining the desired median particle size (circle in Fig. 4.7). It is worth noting that this circle is only one of the possible real input combinations that can lead to the desired product quality. When two quality specifications are assigned (Fig. 4.7b), a smaller portion of the historical KS is identified as a probabilistic DS. The subspace further shrinks in the trivariate scenario (Fig. 4.7c).

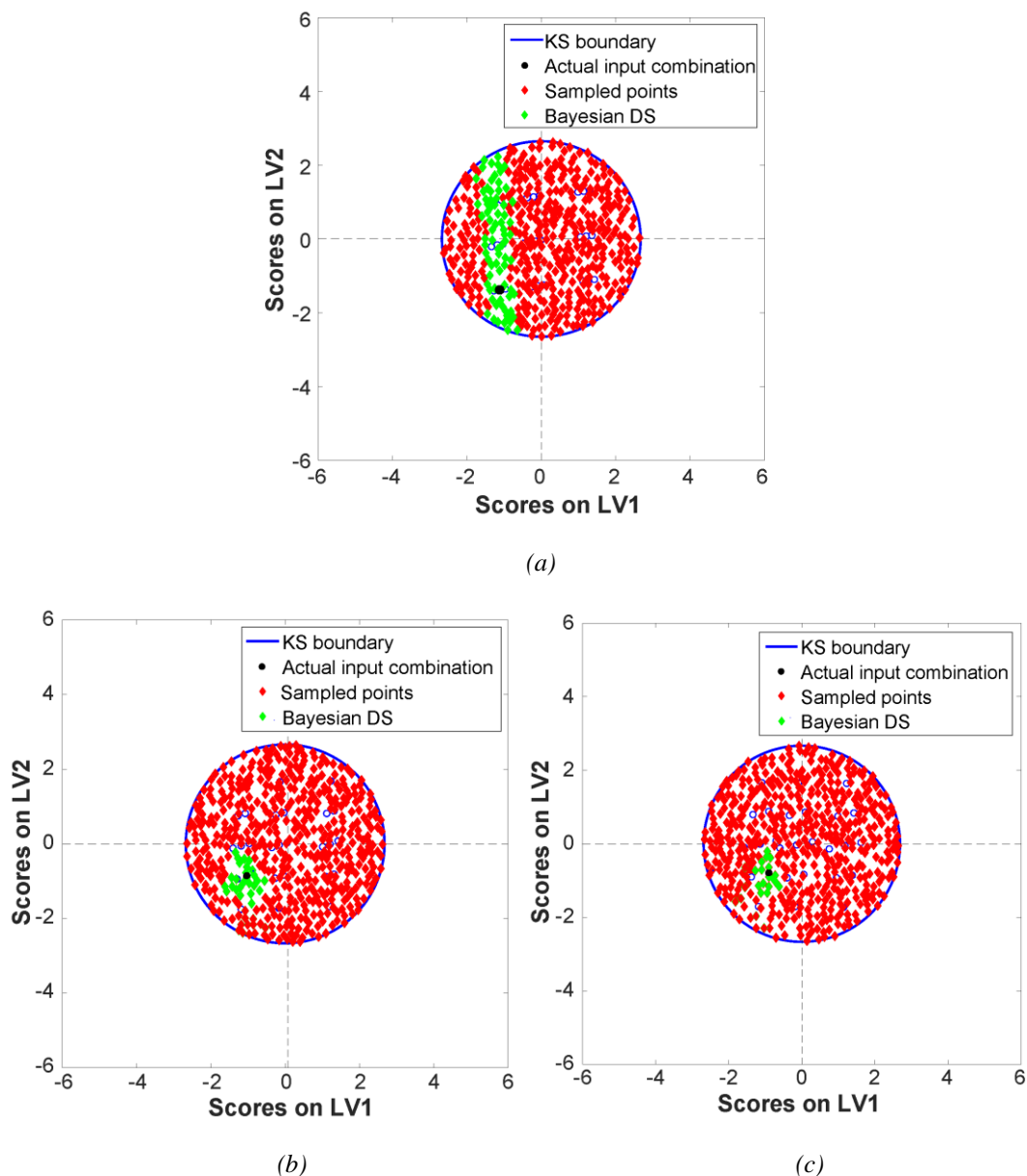


Figure 4.7. Case study #2: Bayesian DS at 90% probability (green diamonds) for: (a) a granulate with assigned $d_{50} = 1129 \mu\text{m}$; (b) a granulate with assigned $d_{50} = 1129 \mu\text{m}$ and $d_{90} = 1849 \mu\text{m}$; (c) a granulate with assigned $d_{50} = 1129 \mu\text{m}$, $d_{90} = 1849 \mu\text{m}$ and $d_{10} = 732 \mu\text{m}$.

The different shapes of the Bayesian DS obtained for the three scenarios can be explained by considering the reduction of dimensionality obtained with the PLS model. Since the number of LVs used in the PLS model is 2, when a single quality specification is assigned a 1-dimensional null space (Jaeckle and MacGregor, 1998) exists in the space of the model inputs, and the projection of the model inputs onto the latent space is therefore a straight line. This explains why, in the latent space, the probabilistic DS appears to be bounded by straight lines. When the number of quality specifications is greater than 1, a null space does not exist in the latent space. In this case, if no uncertainty is accounted for, a PLS model would predict (after inversion) a

single point in the latent space for a given set of quality specifications. This is why no shape of the probabilistic DS is apparent in the latent space. Also note that even if a shape of the DS may be apparent in the latent space, this does not imply that the shape is maintained in the real input space.

The results obtained for this case study confirm the ability of the proposed methodology to identify a set of input variables combinations for which the probability to meet the quality specifications is greater than the assigned thresholds, thus providing a quantification of the “assurance of quality” for the investigated product also for products characterized by multivariate quality targets. The computational effort increases as the number of quality specifications increases, but remains fully tractable

4.6.3 Results for case study #3

A historical dataset of 34 samples is available. One sample is chosen as a validation sample, and the other 33 are used as calibration samples for the PLS model. As for the other two case studies, 2 LVs are chosen, accounting for 96.1 % of the total variability of \mathbf{X} , and 87.7 % of the total variability of \mathbf{y} . The target values assigned to the CQAs are reported in Table 4.3.

The problem specifications are summarized in Table 4.6. Note that the number of optimal samples for this case study is 234. $1e5$ iterations are used for each iteration, and non informative Gaussians are set as prior distributions for the model parameters.

Table 4.6. Case study #3: problem specifications.

Specification	Value
No. of calibration samples	33
No. of model parameters	12
N_s^{opt}	234
Priors of model parameters	Non informative Gaussians
No. of iterations per sample	$1e5$
No. of burn-in iterations	$1e4$
Total simulation time (parallel computing)	1h 15min

The projection onto the latent space of the input combinations that allow obtaining the desired quality with a probability equal to, or higher than, 90% is shown in Fig. 4.8 (green diamonds). The probabilistic DS captures the available true input combination leading to the desired quality target, and offers additional combinations of input materials properties and process operating conditions that can be used to manufacture the desired product, with a given probability to meet the desired quality targets. The probabilistic DS identification exercise is carried out in a reasonably short time (1.25 h).

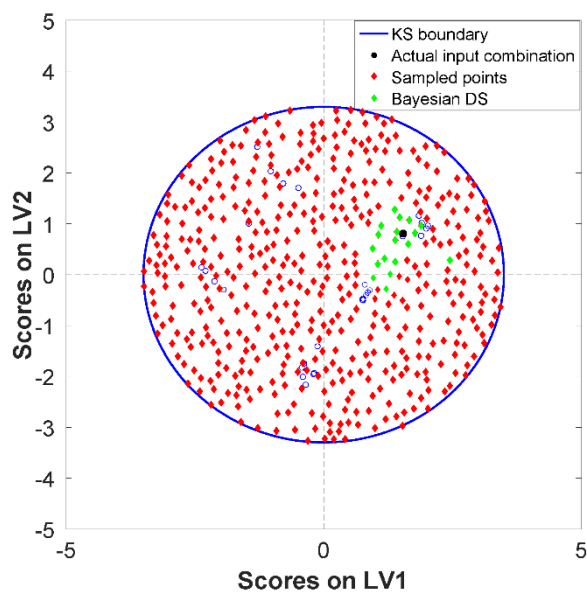


Figure 4.8. Case study #3: Bayesian DS at 90% probability (green diamonds) and actual input combination that gives the desired multivariate quality target (black circle).

4.7 Discussion

Two issues deserve further discussion at this point, namely *i*) the relation between Bayesian DS (as discussed in this study) and confidence region for the DS (as assumed by frequentist approaches), and *ii*) some limitations on the use of the proposed methodology.

With respect to the first issue, it should be pointed out that the meaning of the Bayesian DS is intrinsically different from the concept of confidence region for the DS predicted with traditional frequentist approaches. For example, Facco *et al.* (2015) build a frequentist confidence region for the model prediction of the DS and name this restricted portion of the KS as the “experiment space”. This experiment space represents a subspace of the KS, which is deemed to include the DS at the given confidence level: the greater the confidence level, the wider the experiment space that must be investigated to determine the DS. On the other side, the Bayesian DS discussed in this study can be interpreted as the portion of the KS within which the probability for the product to be on target is greater than (or equal to) an assigned threshold: the greater the probability threshold, the narrower the Bayesian DS.

The difference between the two approaches relies on the difference between the concepts of confidence region (in frequentist statistics) and credible region (in Bayesian statistics). The confidence region (i.e., the experiment space) determined by Facco and coworkers considers the model prediction of the DS as fixed, and the bounds of the experiment space as random variables. On the other side, the credible region (i.e., the Bayesian DS) as discussed in this study considers the bounds of the probabilistic DS as fixed, and the model prediction of the DS as a random variable. The practical outcome is that, whereas the experiment space proposed by

Facco and coworkers cannot be considered as a model-based representation of the DS, the Bayesian DS proposed in this study can be deemed as such, given that the intrinsic probabilistic nature of model predictions is recognized (Pantelides *et al.*, 2012).

With regard to the limitations of the proposed methodology, it must be noted that the amount of \mathbf{X} -variability explained by the PLS model for the chosen number of LVs plays a key role for its successful implementation. Indeed, the smaller the amount of cumulative \mathbf{X} -variability explained by the LVs, the greater the error on the projection of the KS onto the latent space, and the worse the representation of the original input space by the latent space. In practice, when the amount of explained \mathbf{X} -variability is small, the risk of not considering portions of the KS that may belong to the DS of the process increases. As a rule of thumb, a limiting value for the amount of cumulative \mathbf{X} -variability explained by the PLS model may be set as 90%, but this value can possibly be decreased according to the user's risk adversity. The cumulative \mathbf{X} -variability may be increased by increasing the number of LVs of the PLS model (e.g., for $A > 2$), at the expense of a larger computational burden for the KS discretization as well as of a higher complexity of the graphical interpretation of the DS (possibly even losing the possibility to provide a graphical interpretation if $A > 3$). On the other side, a small amount of cumulative \mathbf{y} -variability explained by the PLS model does not preclude the use of the methodology, even though in certain situations it may be a warning that significant nonlinearity exists in the original dataset that is not captured adequately by the linear PLS model.

4.8 Conclusions

The aim of this Chapter was to propose a methodology that allows identifying the design space of a new pharmaceutical product with a risk-based Bayesian posterior predictive approach, while reducing the problem dimensionality by using PLS modeling.

A methodology, based on PLS modeling, that allows automatic determination the optimal number of sampling points needed to cover the entire knowledge space has been developed. This methodology compromises between reduction of the computational effort and achievement of a satisfying coverage of the historical knowledge space. The joint posterior predictive probability of each sampling point was obtained with a multivariate Bayesian linear regression model, and the probability that product quality will meet its specifications for the given point was computed. By performing this analysis for each of the sampling points of the historical knowledge space, the probabilistic (or Bayesian) DS of the product under investigation was determined. The results were then expressed with a simple and intuitive graphical representation on the low-dimensional latent space.

Three case studies were considered to illustrate the effectiveness of the proposed methodology. The obtained probabilistic DS represents a way to demonstrate (i.e., quantify) the level of

“assurance of quality” the manufacturer can guarantee for a given product, as advocated by the pharmaceutical regulatory agencies.

Possible areas of further investigation can be pointed out. First, the methodology is completely general and can be applied not only to linear multivariate regression model, but also to nonlinear data-driven models or nonlinear mechanistic models. Therefore, the combined effect of a reduction of problem dimensionality given by PLS and the identification of the Bayesian DS using a detailed mechanistic model for the process can be investigated. Secondly, the effect of measurement uncertainty (e.g., noise) on the Bayesian DS can be studied (incorporation of measurement uncertainty in the Bayesian framework is straightforward). Finally, the effect of uncertainty propagation in manufacturing processes involving multiple units (thus with significant dimensionality issues) can be assessed with the proposed methodology.

Chapter 5

Feasibility analysis and latent variable modeling for design space determination*

In this Chapter, a methodology that exploits first-principles (or semi-empirical) and latent variable modeling for the identification of the design space of a new pharmaceutical product is presented. Specifically, feasibility analysis is coupled with projection to latent structures (PLS) to obtain a computationally tractable identification of the DS of complex pharmaceutical processes. PLS is used to obtain a linear transformation between the original multidimensional input space and a lower dimensional latent space. Radial Basis Function (RBF) adaptive sampling feasibility analysis is then used on this lower dimensional space to identify the feasible region of the process. The accuracy and robustness of the results is assessed with three metrics, and the criteria that should be adopted for the choice of the number of latent variables are discussed. The performance of the methodology is tested on three simulated case studies, one of which involving the continuous direct compaction of a pharmaceutical powder.

5.1 Introduction

The design space (DS) of a pharmaceutical process is defined as the multidimensional combination and interaction of input variables (e.g. material attributes) and process parameters that have been demonstrated to provide assurance of quality (ICH, 2009). A modeling technique that can be used to assist a DS identification exercise is feasibility analysis (Boukouvala *et al.*, 2010; Wang and Ierapetritou, 2017).

The aim of feasibility analysis (Halemane and Grossman, 1983) is to quantify the capability of a process design to be feasibly operated over the whole domain of input factors, thus including raw material properties and process parameters. The final objective is to determine the

* Bano G., Wang, Z., Facco, P., Bezzo, F., Barolo, M., Ierapetritou, M. (2018) A novel and systematic approach to identify the design space of pharmaceutical processes. *Comput. Chem. Eng.* **115**, 309-322

multivariate region of the input domain within which the process is considered to be feasible. Hence, the concept of feasible region is strictly related to the one of DS for pharmaceutical processes.

Different mathematical approaches can be used to identify the feasible region. When the model of the process is computationally inexpensive, the feasible region can be directly determined from the model itself (Prpich *et al.*, 2010; Close *et al.*, 2014). When the original model is computationally expensive or the computation of its derivatives is cumbersome, surrogate-based approaches can be used (Boukouvala *et al.*, 2010; Grossmann *et al.*, 2014; Wang and Ierapetritou, 2017; Wang *et al.*, 2017). The underlying idea behind these methods is to build a surrogate as a computationally cheap and reliable approximation of the original model, and use this surrogate to identify the feasible region. In this regard, different types of surrogate models have been used (Goyal and Ierapetritou, 2002; Banerjee *et al.*, 2010; Boukouvala *et al.*, 2010; Rogers and Ierapetritou, 2015). Recently, Wang and Ierapetritou (2017) proposed a radial basis function (RBF) adaptive sampling approach that outperforms all the other surrogate-based approaches for low dimensional test problems. The adaptive sampling is a technique that was first proposed in the optimization literature (Jones *et al.*, 1998). In surrogate-based feasibility analysis, it is used to maximize the potential of the sampling budget and to simultaneously explore regions of the original input domain that are close to the boundary of the feasible region and less explored regions.

Although feasibility analysis can be a valuable tool to identify the feasible region of a manufacturing process, it suffers from one main limitation, namely the curse of dimensionality (Shan and Wang, 2010). When a large number of input factors is involved, as it is often the case with large and complex integrated flowsheet models, the computational cost of feasibility analysis has a potential to increase significantly. The solution of a feasibility analysis problem in this scenario can be extremely complicated if not impossible to solve. When the number of input factors is large, the computational burden could be so high or the results so inaccurate as to effectively preclude the application of this methodology. Moreover, a visualization problem of the feasible region in high dimensions arises, given the impossibility of graphically representing a N -dimensional space ($N > 3$).

A recent work of Wang *et al.* (2017) tried to solve this problem by transforming a multidimensional feasibility analysis problem into a series of disjoint 2-dimensional problems, and presenting the results as a matrix of 2-dimensional feasibility contour plots. However, this approach does not account for the multivariate correlation between the original input factors and is thus incomplete.

The input factors of a pharmaceutical process are often correlated to each other (Tomba *et al.*, 2013) and have a different impact on the process output (e.g. product quality; Saltelli *et al.*, 2010). In most common situations, not all the combinations of the original input factors have a strong effect on the output, i.e. some “driving forces” can be identified that predominantly affect

the responses (Jaeckle and McGregor, 2000). A class of statistical models that can identify these underlying driving forces are latent variable models (LVMs). In particular, partial least-squares (PLS) regression (Geladi and Kowalski, 1986; Wold *et al.*, 1983) is a multivariate latent variable technique that can be used to capture the variability of the input and output spaces by means of few meaningful variables, called latent variables, thus reducing the problem dimensionality. The latent variables are chosen by performing a simultaneous decomposition of the input and output space, such that these variables explain as much as possible of the covariance between the two spaces. Stated differently, PLS identifies linear combinations of the original input factors that best describe the correlation between the input factors *and* their effect on the model responses. PLS can be used by itself as a black-box data-driven modelling technique to assist a DS identification exercise when a historical dataset of the process under investigation is available (Facco *et al.*, 2015; Bano *et al.*, 2017). However, when a historical dataset is not available or the process behavior cannot be captured by a linear regression model, the ability of PLS to “compress” the original multidimensional input and output spaces (Mevik, *et al.*, 2004) can be coupled with more sophisticated modelling techniques to assist a DS identification exercise.

The aim of this work is to overcome the curse of dimensionality of feasibility analysis identification by exploiting the ability of PLS to reduce the input space dimensionality. A PLS model is used to perform a linear transformation from the original multidimensional input space to a lower dimensional latent space. Based on the PLS diagnostics, different scenarios are identified in which the proposed methodology is deemed to be profitable.

RBF-based adaptive sampling feasibility analysis is then performed on the latent space and the accuracy and robustness of the results are assessed with three appropriate metrics. The performance of the methodology is tested on three simulated case studies, two of them involving two- and three-unit multidimensional test problems, and one involving a continuous direct compaction of a pharmaceutical powder. In all case studies, the ability of the methodology to give accurate and robust results by simultaneously reducing the computational burden for the feasibility analysis problem is shown.

The rest of this Chapter is organized as follows. PLS and RBF-based adaptive sampling feasibility analysis are briefly reviewed in Section 5.2. The proposed methodology to couple PLS and feasibility analysis for the identification of the feasible region of a manufacturing process is then discussed in Section 5.3. Section 5.4 collects the case studies considered in this work and Section 5.5 shows the results obtained with the proposed methodology.

5.2 Mathematical background

In the following, we briefly review the mathematical techniques that we will use throughout this Chapter.

5.2.1 Partial Least Squares Regression (PLS)

PLS (Geladi and Kowalski, 1986; Wold *et al.*, 1983) is a multivariate regression technique that projects a historical dataset of input \mathbf{X} [$I \times Q$] and response variables \mathbf{Y} [$I \times M$] onto a common latent space of A latent variables (LVs), according to the model structure:

$$\mathbf{X} = \mathbf{T}\mathbf{P}^T + \mathbf{E}_X \quad (5.1)$$

$$\mathbf{Y} = \mathbf{T}\mathbf{Q}^T + \mathbf{E}_Y \quad (5.2)$$

$$\mathbf{T} = \mathbf{X}\mathbf{W}^* \quad (5.3)$$

where \mathbf{T} [$I \times A$] is the score matrix, \mathbf{P} [$Q \times A$] and \mathbf{Q} [$M \times A$] are the \mathbf{X} and \mathbf{Y} loading matrices, \mathbf{E}_X and \mathbf{E}_Y the residuals, \mathbf{W}^* [$Q \times A$] is the weight matrix.

The new set of latent variables is determined by finding the multidimensional directions on the \mathbf{X} -space that explain the maximum multidimensional variance direction in the \mathbf{Y} -space. A can be chosen in such a way as to explain a significant fraction of the variance of both the input and output data.

The PLS model calibration described by Eqs. (5.1)-(5.3) can be interpreted as a linear transformation L_{PLS} between the original set of input factors (x_1, x_2, \dots, x_Q) and the new set of A latent variables, identified by the model scores (t_1, t_2, \dots, t_A), according to Eq. (5.1). Therefore, PLS can be used as a projection technique to express the original set of input factors (x_1, x_2, \dots, x_Q) in terms of a reduced set of variables (t_1, t_2, \dots, t_A).

The latent representation of the original input space obtained with the PLS model captures most of the variability of the original input space that is *correlated* with the model outputs. This is a key advantage with respect to other decomposition techniques such as principal component analysis (PCA; Hotelling, 1933) or Singular Value Decomposition (SVD), in which the latent representation of the original input space only accounts for the variability of the original input dataset. In the presence of a strong input collinearity and/or in the presence of a strong sensitivity of the model outputs to a reduce set of input combinations, a significant reduction of the problem dimensionality can be obtained by exploiting the PLS model (5.1)-(5.3) (i.e. $A \ll Q$). The linearity of the transformation (5.1)-(5.3) allows expressing the original input factors as linear combinations of the latent variables, as opposed to other more sophisticated nonlinear methods such as Kernel PCA (Schölkopf *et al.*, 1998). For the purpose of this study, this has the advantage of keeping the dimensionality reduction step as simple as possible, since process nonlinearity is tackled at a subsequent stage of the proposed methodology.

5.2.2 Feasibility analysis

Feasibility analysis is a mathematical tool that allows investigating the portion of the input space within which the process under investigation is considered to be feasible. For pharmaceutical processes, this is equivalent to determining the design space of the process (ICH, 2009). A rigorous mathematical description of the feasibility of a process can be found in Wang and Ierapetritou (2017). In a nutshell, given a process with J constraints $f_j(\mathbf{d}, \mathbf{x})$, $j = 1, \dots, J$, where \mathbf{d} is the vector of the design variables (which are constant since the process is determined) and \mathbf{x} is the vector collecting the input variables (raw material properties and process parameters), the process is considered to be feasible when all its constraints are met, i.e. when $f_j(\mathbf{d}, \mathbf{x}) \leq 0 \forall j \in [1, \dots, J]$. In order to verify if any of the process constraints is violated, it is enough to check the maximum value of all the constraint function values. The maximum value of the constraint function values is defined as a “feasibility function” $\psi(\mathbf{d}, \mathbf{x})$ (Grossmann, 2014) according to

$$\psi(\mathbf{d}, \mathbf{x}) = \max_{j \in J} f_j(\mathbf{d}, \mathbf{x}). \quad (5.4)$$

If $\psi(\mathbf{d}, \mathbf{x}) \leq 0$, it means that, for the given set of design variables and input variables, the process is feasible. If $\psi(\mathbf{d}, \mathbf{x}) > 0$, it means that one or more of the constraints are violated, i.e. the process is not feasible. If $\psi(\mathbf{d}, \mathbf{x}) = 0$, it means that we are at the boundary of the feasible region. The objective of feasibility analysis is to identify the subspace of input combinations of the original input domain within which the process is feasible, i.e. to identify the DS of the process under investigation.

5.2.3 Surrogate-based feasibility analysis

The feasibility analysis problem described in Section 5.2.2 is based upon the availability of a model to describe the single unit and/or the overall process under investigation. When the model complexity is high (e.g. in the case of an integrated flowsheet model) feasibility analysis can be rather difficult due to the computational burden required by the simulation. In this scenario, a surrogate-based method can be used to efficiently identify the feasible region. With this approach, a surrogate model is built as a computationally cheap approximation of the original model. The feasible region is then determined based on this surrogate. Adaptive sampling (Rogers and Ierapetritou, 2015) can then be used in order to improve the accuracy of the surrogate by making the best possible use of the sampling budget. New points are sampled near the boundary of the feasible region as well as in less explored regions of the historical input domain. In this work, the feasibility analysis problem is solved using a radial basis function

(RBF)-based adaptive sampling method proposed by Wang and Ierapetritou (2017). The main features of this approach are briefly discussed below; additional details can be found in the original reference.

5.2.4 Radial basis function (RBF) surrogate model

RBFs (Björkman and Holmström, 2000) can be used to predict the value of the original function at an unsampled point according to:

$$y_n(\mathbf{x}) = \sum_{i=1}^{nN} \lambda_i \xi \left(\|\mathbf{x} - \mathbf{x}_i\|_2 \right) + \mathbf{b}^T \mathbf{x} + a \quad (5.5)$$

where $\mathbf{x}_i \in \mathbf{R}^Q, i = 1, \dots, N$ are N distinct sample points of known function values $f(\mathbf{x}_i)$; $\|\cdot\|_2$ is the Euclidean norm and ξ is the basis function that in this study is chosen as the cubic basis function:

$$\xi(r) = r^3. \quad (5.6)$$

The RBF surrogate model with cubic basis function has been proven to outperform other surrogate models by Wang and Ierapetritou (2017).

The coefficients λ_i , \mathbf{b} and a of Eq. (5.5) can be determined by solving the following equation:

$$\begin{pmatrix} \Phi & \mathbf{S} \\ \mathbf{S}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \boldsymbol{\lambda} \\ \mathbf{c} \end{pmatrix} = \begin{pmatrix} \mathbf{F} \\ \mathbf{0} \end{pmatrix}. \quad (5.7)$$

where Φ is the $N \times N$ matrix whose elements are defined by $\Phi_{ij} = \xi \left(\|\mathbf{x} - \mathbf{x}_i\|_2 \right)$ and:

$$\mathbf{S} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & 1 \\ \mathbf{x}_N^T & 1 \end{pmatrix}; \boldsymbol{\lambda} = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_N \end{pmatrix}; \mathbf{c} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \vdots \\ \mathbf{b}_Q \\ a \end{pmatrix}; \mathbf{F} = \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_N) \end{pmatrix}. \quad (5.8)$$

In order to describe how well the region near an unsampled point \mathbf{x}^* has been explored, an error indicator $1/\mu_n(\mathbf{x}^*)$ first introduced by Gutmann (2001) is used. The larger the value of $1/\mu_n(\mathbf{x}^*)$, the less the nearby region of the unsampled point has been explored. Full details on the derivation of this estimator can be obtained in the cited reference and in Wang and Ierapetritou (2017).

5.2.5 Adaptive sampling

Adaptive sampling is a technique that allows improving the accuracy of the surrogate without spanning the whole historical input domain of the original model (Boukouvala and Ierapetritou, 2012; Rogers and Ierapetritou, 2015). Various adaptive sampling techniques have been developed for different purposes, such as enhancing the accuracy of a global surrogate model (Garud, *et al.* 2017), improving the approximation of distribution functions (Martino *et al.* 2015), and facilitating surrogate-based global optimization (Huang *et al.* 2006; Picheny *et al.* 2013; Boukouvala and Ierapetritou 2014).

In this Chapter, the sampling strategy proposed by Boukouvala and Ierapetritou (2014) was used. With this technique, the search for new sample points is directed towards the boundary of the feasible region as well as less explored regions of the historical input domain. New sampling points are chosen by maximizing a modified expected improvement (EI) function:

$$\max_{\mathbf{x}} EI_{feas}(\mathbf{x}) = s \times \phi\left(\frac{-\hat{y}}{s}\right) = s \times \left(\frac{1}{\sqrt{2\pi}}\right) e^{-0.5\left(\frac{\hat{y}^2}{s^2}\right)} \quad (5.9)$$

where $EI_{feas}(\mathbf{x})$ is the modified expected improvement function at \mathbf{x} ; s is the standard error of the predictor; \hat{y} is the surrogate model prediction; $\phi(*)$ is the standard normal density function. For RBF-based methods, the prediction error can be estimated with the error indicator $1/\mu$ discussed in Section 5.2.4, by introducing a scale factor to balance the magnitude of $1/\mu$ with that of the surrogate value \hat{y} . As shown by Wang and Ierapetritou (2017), s can be estimated as:

$$s = \sqrt{\frac{1/\mu}{scale}} \quad (5.10)$$

where:

$$\sqrt{\frac{1/\mu}{scale}} = \sqrt{\frac{\left(\frac{1}{\mu}\right)}{\left(\frac{1}{\mu_0}\right)_{\max}}} \times \frac{RBF_0^{\max}}{V_0}. \quad (5.11)$$

with $\left(\frac{1}{\mu_0}\right)_{\max} = \max\left(\frac{1}{\mu_0}\right)$ maximum value of $\frac{1}{\mu}$ with the initial surrogate model; $RBF_0^{\max} = \max(RBF_0)$ is the maximum value of the initial cubic RBF prediction; V_0 is the number of initial sample points.

By taking the derivative of $EI_{feas}(\mathbf{x})$ with respect to \hat{y} , it can be shown (Wang and Ierapetritou, 2017) that its value is large when the surrogate prediction \hat{y} is close to zero (i.e. near the boundary of the feasible region) and/or when prediction uncertainty s is high. Therefore, the

adaptive sampling allows sampling new points near the boundary of the feasible region for the surrogate as well as in those regions of the historical input domain where prediction uncertainty is high (i.e. less explored regions).

A schematic representation of the RBF-based adaptive sampling algorithm is shown in Figure 5.1. This algorithm was first proposed by Wang and Ierapetritou (2017). Three steps are involved: in the first step, a space-filling sampling algorithm is used to span the historical input domain and build the initial surrogate according to a Design of Experiments (DoE) approach. The scale factor of Eq. (5.11) is computed at this step. In the second step, the surrogate accuracy is improved by sequentially adding new sampling points obtained with the adaptive sampling strategy discussed above. The adaptive sampling stops when the number of iterations exceeds a maximum defined by user. A typical value for the maximum number of iterations is 100. In the third and last step, the surrogate is used to predict the feasible region (i.e. the DS) of the process under investigation.

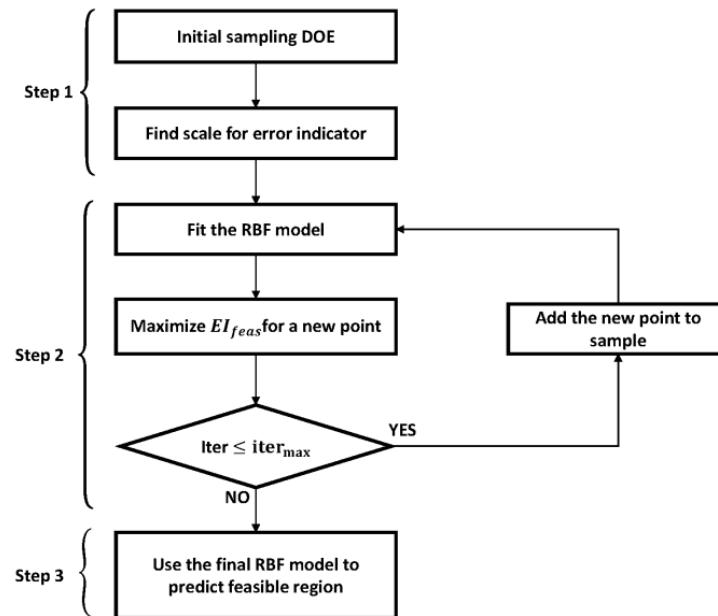


Figure 5.1. RBF-based adaptive sampling for feasibility analysis. Originally presented by Wang and Ierapetritou (2017).

5.2.6 Surrogate accuracy

The results of the surrogate-based feasibility analysis can be deemed to be accurate if three conditions are simultaneously satisfied:

1. The feasible region is correctly identified by the surrogate;
2. The infeasible region is correctly identified;

3. The surrogate does not considerably overestimate the actual feasible region of the process.

Different metrics have been proposed to quantify the surrogate accuracy, based on the computation of simple coverage probabilities (Rogers and Ierapetritou, 2015) or comparison between the predicted and theoretical hyper-volumes of the feasible region (Adi *et al.*, 2016). However, while these techniques aim at quantifying the first two requirements listed above, they fail at addressing the third one. A set of simple metrics accounting for all three aspects is the one that has been recently proposed by Wang and Ierapetritou (2017). A graphical and analytical interpretation of these metrics can be found in the cited reference.

The metrics are:

1. CF (%): percentage of the feasible region for the original function which has been correctly discovered by the surrogate model;
2. CIF (%): percentage of infeasible region for the original function which has been correctly discovered by the surrogate model;
3. NC (%): percentage of feasible region that has been overestimated by the surrogate model.

The first two metrics can be considered as “quality” metrics: they describe how well the input space has been explored and classified with respect to feasibility. The third metric can be considered as a measure of the “conservativeness” of the results, i.e. it describes how large is the portion of the input domain that has been wrongly classified as feasible by the surrogate with respect to the overall feasible region predicted by the surrogate itself. In summary, it can be said that the feasible region has been classified thoroughly and conservatively if CF(%) and CIF(%) are close to 100, and NC(%) is close to 0.

5.3 Feasibility analysis and reduction of input space dimensionality

In the following, a systematic methodology to couple PLS input space projection and feasibility analysis for the identification of the feasible region of a process involving a (possibly) large number of input factors is proposed.

5.3.1 Proposed methodology

The proposed methodology involves multiple steps. Although the discussion will be general, an illustrative case study shown in Figure 5.2 will be used to give a simple interpretation of the methodology. The example involves two units. Let x_1 and x_2 being the input factors for the first unit. Let ψ_1 be the feasibility function for the first unit, i.e. the maximum values of all the constraints on the outputs of the first unit. Let u_1 being an additional input for the second unit and ψ_2 be the feasibility function for the second unit. The objective is to determine the set of

input combinations $\mathbf{x} = (x_1, x_2, u_1)$ that satisfy $\psi_2(\mathbf{x}) \leq 0$. In the test problems, it is assumed that ψ_2 monotonically increases with ψ_1 . It is worth noticing that this condition does not imply that $\psi_1 \leq 0$ if $\psi_2 \leq 0$; in practical terms, the outputs of the second unit may satisfy all their constraints, but the outputs of the first unit may violate some of their own constraints. This reflects a practical situation where a process may satisfy the quality constraints on the final product (output of unit #2), but not all the constraints on the outputs of all the units of the manufacturing line.

In practical situations, the methodology assumes the availability of an integrated flowsheet model of the process (we will refer to this model as the “original” model), with quality constraints imposed on the final product. The maximum value of all these constraints will be denoted by the feasibility function $\psi(\mathbf{x})$.

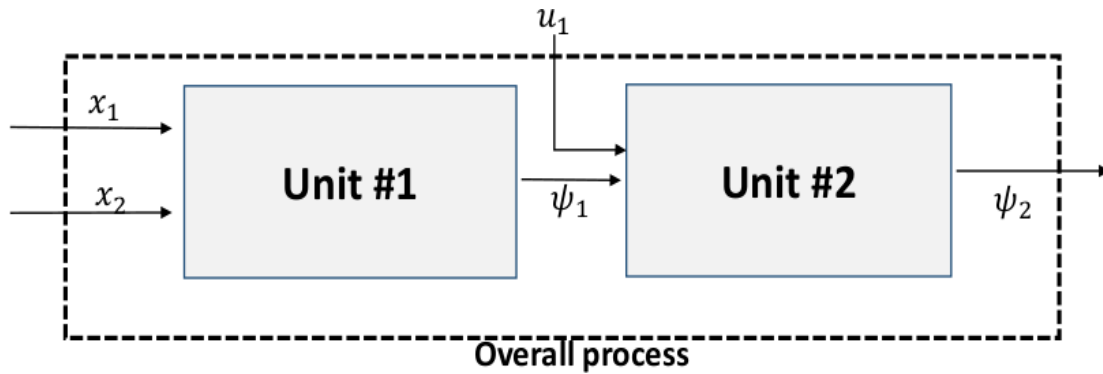


Figure 5.2. Illustrative example for the proposed methodology: block flow diagram.

The flowchart of the methodology is shown in Figure 5.3. The steps are as follows.

1. **Step #1:** collection of simulated data from the model. A dataset is generated using the original model (i.e., simulated data). First, the input space is spanned using the Kennard-Stones’s sampling algorithm (KSA; Kennard and Stone, 1969). KSA is a space-filling sampling strategy that selects each new sample by choosing the farthest sample (in terms of Euclidean distance) from the previously selected ones, starting from the two farthest points on the historical input domain. Each selected sample corresponds to a selected input combination. Let $\mathbf{x}_i, i = 1, \dots, Q$ be the vector that collects all the selected values for the i -th input factor. The selected samples can then be collected in the input matrix $= [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_Q]$. The number of samples can be selected by the user and can be possibly increased to improve the PLS model accuracy. In all case studies, we generated a number of sample points equal to the number of sample points that we used to build the initial surrogate.

Then, the values of the feasibility function $\psi(\mathbf{x})$ are collected according to the original process model. The vector that collects all these values is denoted by Ψ . For the

illustrative example of Fig. 5.2, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{u}_1]$, and the feasibility function for the overall process corresponds to the feasibility function for the second unit (i.e. $\Psi = \Psi_2$).

2. **Step #2:** PLS model calibration. A PLS model (Eqs.(5.1)-(5.3)) is calibrated with the dataset generated in the previous step. A total number of A latent variables are chosen according to the desired dimensionality reduction of the input space. In most situations, two or three LVs are able to capture most of the variability of the original input space. As explained in section 5.2.1, the PLS model is then used in the form of Eq. (5.1) to obtain a linear transformation L_{PLS} between the original set of input factors (x_1, x_2, \dots, x_Q) and a new set of A latent variables, identified by the PLS model scores (t_1, t_2, \dots, t_A) , according to Eq. (5.1). In the illustrative example considered, assuming $A = 2$ latent variables, the linear transformation $L_{PLS} : (x_1, x_2, u_1) \leftrightarrow (t_1, t_2)$ can be written as:

$$\begin{aligned} x_1 &= t_1 p_{11} + t_2 p_{12} + e_1 \\ x_2 &= t_1 p_{21} + t_2 p_{22} + e_2 \\ u_1 &= t_1 p_{31} + t_2 p_{32} + e_3 \end{aligned} \quad (5.12)$$

where e_i is the residual (i.e. projection error) for the i -th input factor.

3. **Step #3:** PLS model diagnostics. Some diagnostics on the representativeness of the PLS model are evaluated. To this purpose, the amount of \mathbf{X} - and \mathbf{y} -variability explained by the PLS model is used, which is quantified by the determination coefficients R_X^2 and R_Y^2 (together with their cumulative values $R_{X,cum}^2$ and $R_{Y,cum}^2$):

$$R_X^2 = 1 - \frac{\sum_{i=1}^I \sum_{q=1}^Q (x_{i,q} - \hat{x}_{i,q})^2}{\sum_{i=1}^I \sum_{q=1}^Q (x_{i,q})^2} \quad (5.13)$$

$$R_Y^2 = 1 - \frac{\sum_{i=1}^I \sum_{m=1}^M (y_{i,m} - \hat{y}_{i,m})^2}{\sum_{i=1}^I \sum_{m=1}^M (y_{i,m})^2} \quad (5.14)$$

where $\hat{x}_{i,q}$ is the element of the i -th row and q -th column of the matrix $\hat{\mathbf{X}} = \mathbf{TP}^T$ reconstructed through the PLS model; $\hat{y}_{i,m}$ is the element of the i -th row and m -th column of the matrix $\hat{\mathbf{Y}} = \mathbf{TQ}^T$. The two metrics (5.13)-(5.14) can be computed for every latent variable included in the model and their cumulative values can be reported. If the chosen number of latent variables is A , $R_{X,cum}^2$ and $R_{Y,cum}^2$ represent the cumulative amount of input and output variability respectively captured by the PLS model with A LVs.

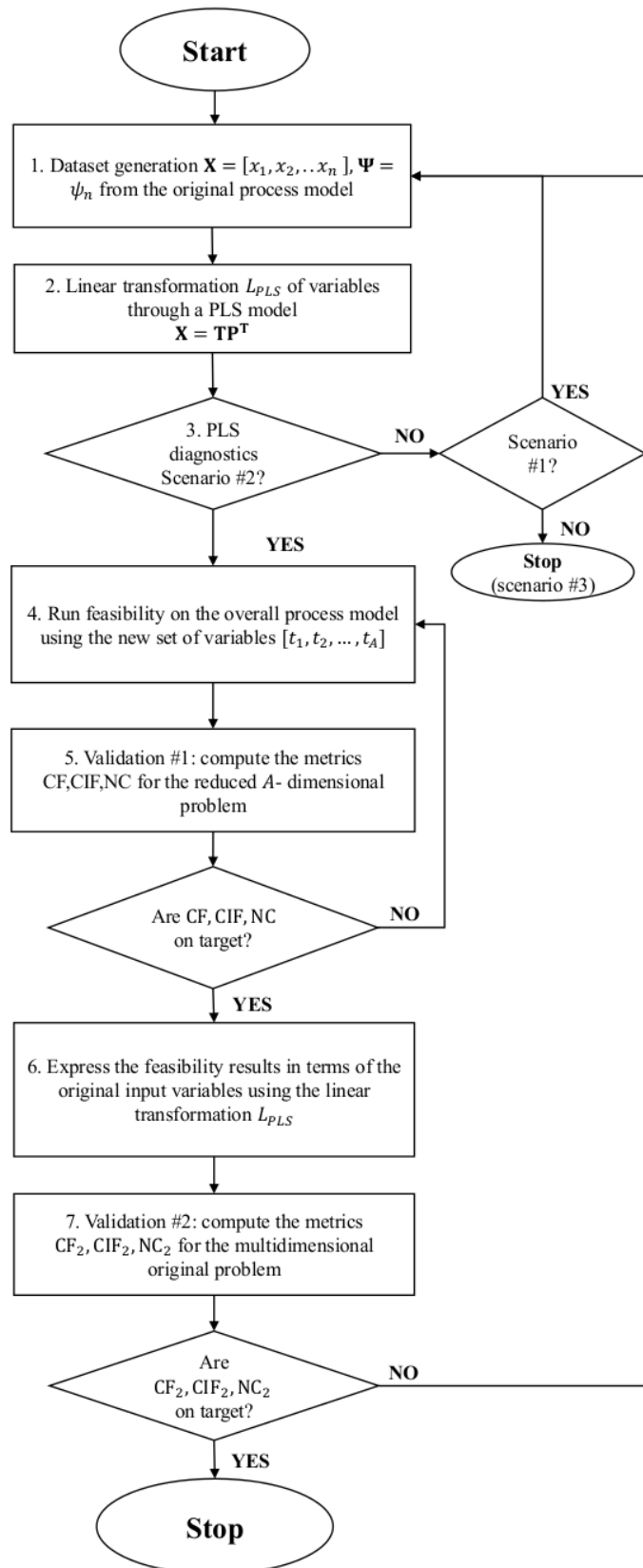


Figure 3. Flow chart of the proposed methodology.

Three possible levels for the metrics $R_{X,cum}^2$ and $R_{y,cum}^2$ are identified, namely “low” when they are smaller than 70%, “medium” when they are between 70% and 90%, and “high” when they are greater than 90%. These are indicative values dictated by experience and can be adjusted according to user’s risk adversity.

According to the values of $R_{X,cum}^2$ and $R_{y,cum}^2$, three different possible scenarios can be identified, namely Scenario #1, Scenario #2 and Scenario #3 respectively.

Scenario #1: the PLS model cannot explain a large portion of both the input and the output variability for the chosen number of LVs. In other words, the predictive ability of the PLS model is limited (low $R_{y,cum}^2$) and a large portion of the original input space is not captured by the PLS model (low $R_{X,cum}^2$). In order for the methodology to be effective, the number of latent variables A must be increased to increase the cumulative \mathbf{X} -variability explained by the PLS model. The higher $R_{X,cum}^2$, the better the description of the original input space by the latent space, the lower the projection error of the linear transformation of variables L_{PLS} . The value of A should be increased to obtain a $R_{X,cum}^2$ close to or higher than 90% (in order to switch to scenario #2). The greater the number of LVs chosen, the smaller the reduction of the dimensionality of the original input space, the higher the computational burden for the feasibility analysis. If feasibility analysis is computationally infeasible for the number of LVs required to satisfy the conditions above, the proposed methodology is not useful and should not be used.

Scenario #2: for the chosen number of LVs, the amount of \mathbf{X} -variability explained by the PLS model is high, while the amount of \mathbf{y} -variability explained is low. In this scenario, the predictive ability of the PLS model is low (low $R_{y,cum}^2$), but the original input space is well described by the new latent space, i.e. the projection error of the linear transformation L_{PLS} is low. This is an ideal scenario for the application of the methodology. In fact, recall from section 5.2.1 that, in this study, the PLS model is only used as a projection method and not for prediction. Therefore, a small value of $R_{y,cum}^2$ indicates that PLS (i.e. a simple linear regression model) is not a suitable model for the description of the process under investigation. On the other hand, a high value of $R_{X,cum}^2$ indicates that the original input space is well described by the latent space identified by the new set of latent variables (t_1, t_2, \dots, t_A) . Since feasibility analysis is performed on the *original* process model by applying the linear transformation of variables L_{PLS} , this scenario represents the ideal situation for the methodology to be effective. The feasibility analysis problem is transformed from a Q -dimensional problem to an A -dimensional problem, with $A \ll Q$.

Scenario #3: the amount of \mathbf{y} -variability explained by the PLS model is high, independently of the amount of \mathbf{X} -variability captured by the A LVs. In this situation, the predictive ability of the PLS model is high, no matter how large is the projection error of the linear transformation L_{PLS} . This suggests that PLS (which is a linear

regression model) would be enough to describe the input-output correlation for the original process, and could be directly used to predict the feasible region. In this scenario, the use of a more sophisticated integrated model for the overall process is not necessary to determine the feasible region. Therefore, feasibility analysis is not required and the PLS model can be used by itself without going on with the next steps of the methodology. Details on the direct use of PLS modelling for the determination of the feasible region are reported elsewhere (e.g. Tomba *et al.*, 2012; Facco *et al.*, 2015; Bano *et al.*, 2017).

Figure 5.4 shows the possible combinations of the values of $R_{X,cum}^2$ and $R_{Y,cum}^2$ and their corresponding scenarios. The cells with two indicated scenarios (e.g. Scenario 1/ Scenario 3) represent transitional situations where it is up to the user deciding which scenario to be chosen between the two, and therefore which corrective actions to be implemented as indicated in Figure 5.5. This gives the user a degree of freedom on the choice of the most suitable techniques/corrective actions to be implemented for the case study under investigation.

4. **Step #4:** feasibility analysis. Once verified that the methodology can be effective (scenario #2), the surrogate-based feasibility analysis previously discussed is performed on the original process model, by applying the linear transformation of variables L_{PLS} . Going back to the illustrative example, the surrogate-based feasibility analysis is performed on the original model structure for units 1 and 2, but applying the linear transformation of the original input variables described by Eq. (5.14). The surrogate-predicted feasible region is then identified on the new A -dimensional latent space, i.e. it is identified in terms of the new set of variables $[t_1, t_2, \dots, t_A]$.
5. **Step #5:** results validation on the latent space. The surrogate accuracy is validated on the latent space of the A new variables. The three metrics CF(%), CIF(%), NC(%) are computed to test the accuracy and robustness of the results. It is worth noticing that these metrics, as computed at this step, simply quantify the agreement between the feasible region predicted by the surrogate and the projection of the actual feasible region of the process onto the latent space. If these metrics satisfy the threshold values specified by the user, the analysis of the surrogate accuracy with respect to the original input space can be performed (step #6). If the metrics do not satisfy the requirements, the accuracy of the results can be improved by increasing the maximum number of iterations of the feasibility analysis or by increasing the number of initial sample points to build the surrogate. As a rule of thumb, the following stopping criterion can be used for the accuracy metrics: $CF(\%) > 97, CIF(\%) > 97, NC(\%) < 3$.

Cumulative X – variance explained \ Cumulative y – variance explained	Low ($R_{y,cum}^2 < 70\%$)	Medium ($70\% \leq R_{y,cum}^2 \leq 90\%$)	High ($R_{y,cum}^2 > 90\%$)
Low ($R_{x,cum}^2 < 70\%$)	Scenario 1	Scenario 1 / Scenario 3	Scenario 3
Medium ($70\% \leq R_{x,cum}^2 \leq 90\%$)	Scenario 1 / Scenario 2	Scenario 2 / Scenario 3	Scenario 3
High ($R_{x,cum}^2 > 90\%$)	Scenario 2	Scenario 2 / Scenario 3	Scenario 3

Figure 5.4. Application of the proposed methodology: possible PLS modeling scenarios.

6. Step #6: PLS model inversion. The results of the feasibility analysis are expressed in terms of the original input variables by applying the inverse linear transformation L_{PLS}^{-1} . In the illustrative example considered, the results of the feasibility of step #5 (expressed in terms of (t_1, t_2)) are expressed in terms of (x_1, x_2, u_1) .

7. Step #7: results validation on the original input space. The accuracy and robustness of the results are evaluated by computing the metrics CF(%), CIF(%), NC(%) in terms of the original input factors (x_1, x_2, \dots, x_Q) . To emphasize that these metrics are different from the ones of step #5, they have been defined as $CF_2(\%), CIF_2(\%), NC_2(\%)$. It must be noticed that $CF_2(\%) \leq CF(\%), CIF_2(\%) \leq CIF(\%)$ and $NC_2(\%) \geq NC(\%)$ always result, due to the projection error of the linear transformation L_{PLS} . In the limiting case where there is no projection error with the PLS model, then $CF_2(\%) = CF(\%), CIF_2(\%) = CIF(\%)$ and $NC_2(\%) = NC(\%)$. If the metrics satisfy the requirements defined by the user, the results can be considered accurate and the algorithm stops. If these metrics do not satisfy the requirements, the projection error should be decreased by increasing the number of LVs of the PLS model (thus increasing the computational burden). A possible set of threshold values that can be imposed for the accuracy metrics at this step is $CF_2(\%) > 90, CIF_2(\%) > 95, NC_2(\%) < 5$.

Scenario #	Conditions	Actions
1	<ul style="list-style-type: none"> Low PLS predictive ability A large portion of the historical input space cannot be investigated for the given number of LVs 	<ul style="list-style-type: none"> Increase the number of LVs to increase $R_{X,cum}^2$ if computationally feasible If computationally infeasible, a different modelling approach is needed
2	<ul style="list-style-type: none"> Low PLS predictive ability A large portion of the historical input space can be investigated for the given number of LVs 	<ul style="list-style-type: none"> The proposed methodology is effective Increasing $R_{X,cum}^2$ (i.e. by increasing the number of LVs) is effective, but increases the computational burden Select A as a compromise between projection accuracy and computational burden
3	<ul style="list-style-type: none"> High PLS predictive ability A large portion of the historical input space can be investigated for the given number of LVs 	<ul style="list-style-type: none"> Use the PLS model to predict the feasible region Feasibility analysis is not required The original model structure can possibly be simplified

Figure 5.5. Description and corrective actions for the different PLS modeling scenarios.

The methodology presented above was implemented in MATLAB[®] 2015, and the simulations were performed on an Intel[®] Core[™] I7-5600U CPU@2.60GHz processor with 16.0 GB RAM. A schematic of the different MATLAB[®] routines and their interactions are shown in Figure 6. Two main routines (PLS and feasibility routines) interact with each other to obtain the linear transformation of variables L_{PLS} . In the PLS routines, the dataset generation and the PLS diagnostics subroutines are sequentially called. In the main feasibility routine, the unit models and the results validation subroutines are involved. The feasibility routine interacts with the PLS routine to transform the reduce set of latent variables to the original input variables through the PLS loadings, thus computing the results validation on the original multidimensional input space.

5.4 Case studies

The proposed methodology was tested on three case studies of increasing complexity. In the first case study, a low dimensional two-unit mathematical test problem involving 3 input factors is considered. The second case study is a three-unit high dimensional test problem with 7 input factors. The third case study involves a direct compaction manufacturing process of a pharmaceutical powder. Six units and six relevant input factors are considered with 40 constraints on the model output.

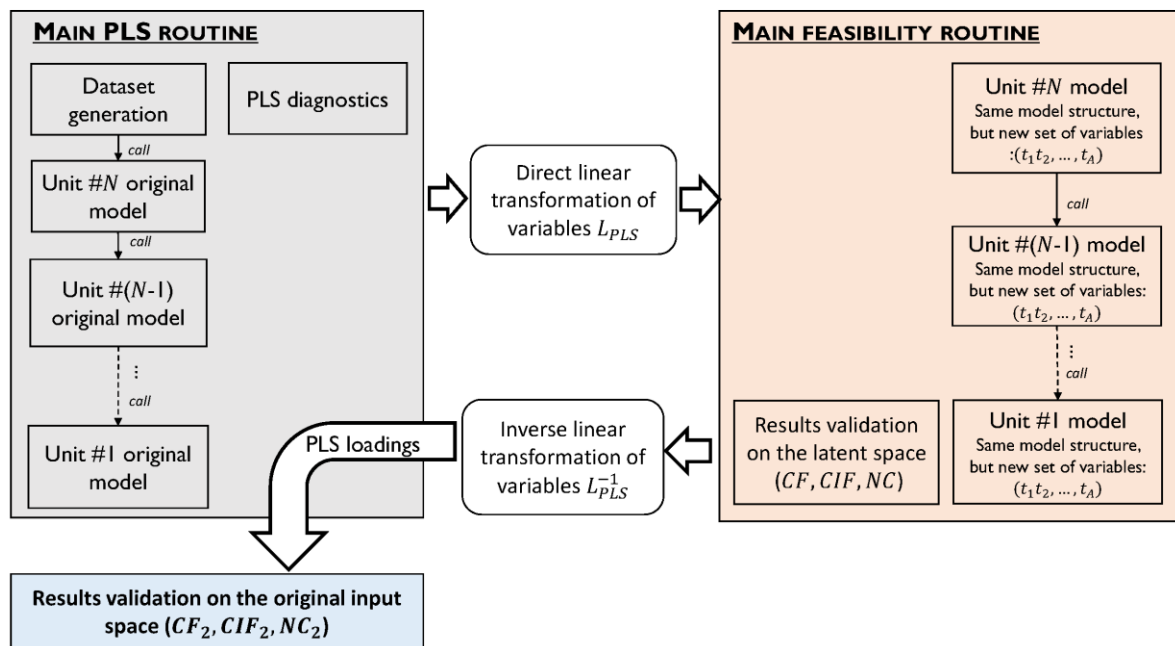


Figure 6. Set of MATLAB[®] routines used to implement the methodology and their interaction

5.4.1 Case study #1: low dimensional test problem

The case study is a simplified 2-unit process, as schematically shown in Figure 5.7. The first unit is characterized by two input factors (x_1, x_2), and the maximum constraint value for its output is denoted by ψ_1 . An additional input u_1 enters the second unit, whose feasibility function is ψ_2 .

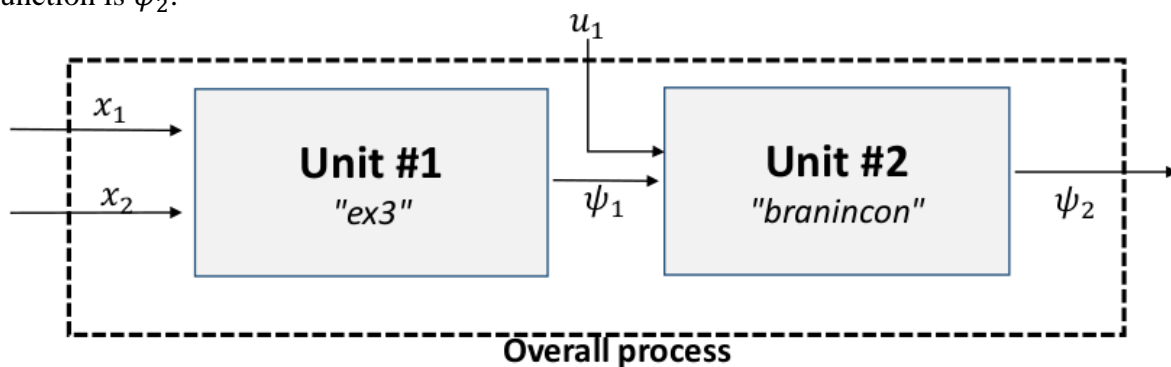


Figure 5.7. Case study #1: block flow diagram.

The models for the two units have been taken from Wang and Ierapetritou (2017), based on the previous work of Rogers and Ierapetritou (2015). Mathematically, the feasibility problem for the first unit (“ex3”) can be expressed as:

$$-2x_1 + x_2 - 15 \leq 0 \tag{5.15}$$

$$\frac{x_1^2}{2} + 4x_1 - x_2 - 5 \leq 0 \quad (5.16)$$

$$-\frac{(x_1-4)^2}{5} - \frac{x_2^2}{0.5} + 10 \leq 0. \quad (5.17)$$

$$10 \leq x_1 \leq 5 \quad (5.18)$$

$$-15 \leq x_2 \leq 15 \quad (5.19)$$

while the feasibility problem for the second unit (“*branincon*”) can be expressed as:

$$a(\psi_1 - bu_1^2 - cu_1 - d)^2 + h(1 - ff)\cos(u_1) - 5h \leq 0 \quad (5.20)$$

$$-5 \leq \psi_1 \leq 10 \quad (5.21)$$

$$0 \leq u_1 \leq 15 \quad (5.22)$$

where

$$a = 1; b = \frac{5.1}{4\pi^2}; c = \frac{5}{\pi}; d = 6; h = 10; ff = \frac{1}{8\pi}. \quad (5.23)$$

The overall process consists of three input factors (x_1, x_2, u_1), thus being a 3-dimensional problem.

5.4.2 Case study #2: high dimensional test problem

The case study is a high dimensional test problem involving three units and a total of 7 input factors for the overall process (Figure 5.8).

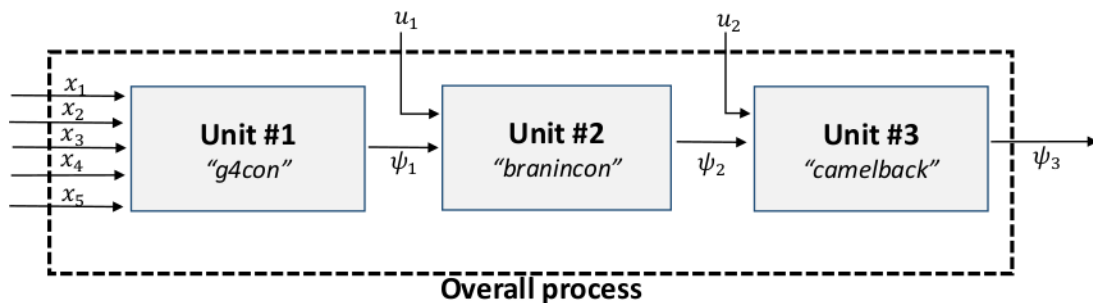


Figure 5.8. Case study #2: block flow diagram.

The model for the first unit (“*g4con*”) is taken from Koziel and Michalewicz (1999) and can be mathematically described as:

$$0 \leq 85.334407 + 0.0056858x_2x_5 + 0.0006262x_1x_4 - 0.0022053x_3x_5 \leq 92 \quad (5.24)$$

$$90 \leq 80.51249 + 0.0071317x_2x_5 + 0.002955x_1x_2 + 0.0021813x_3^2 \leq 110 \quad (5.25)$$

$$20 \leq 9.300961 + 0.0047026x_3x_5 + 0.0012547x_1x_3 + 0.0019085x_3 \leq 25 \quad (5.26)$$

$$78 \leq x_1 \leq 102, 33 \leq x_2 \leq 45, 27 \leq x_i \leq 45, i = 3,4,5. \quad (5.27)$$

The second unit is the “branincon” test model described for case study #1. The model for the third unit (“camelback”) is taken from Goel *et al.* (2007) and is given by:

$$\left(4 - 2.1\psi_2^2 + \frac{\psi_2^4}{3}\right)\psi_2^2 + \psi_2u_2 + (-4 + 4u_2^2)u_2^2 \leq 0 \quad (5.28)$$

$$-3 \leq \psi_2 \leq 3 \quad (5.29)$$

$$-2 \leq u_2 \leq 2. \quad (5.30)$$

The feasibility function for the third unit is denoted as ψ_3 . The input factors for the overall process are $(x_1, x_2, x_3, x_4, x_5, u_1, u_2)$. The case study is thus a 7-dimensional test problem.

5.4.3 Case study #3: simulation of continuous direct compaction for pharmaceutical production

This case study involves the continuous direct compaction (DC) for the production of powder based drugs and is based on the work recently presented by Wang *et al.* (2017). The flowsheet model of the manufacturing line is based on a continuous DC pilot plant situated at ERC-C-SOPS Rutgers University and is shown in Figure 5.9.

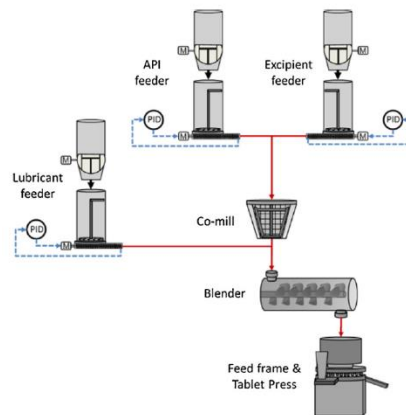


Figure 5.9: Case study #3: process flow diagram of the continuous direct compaction manufacturing line. Originally presented by Wang *et al.* (2017).

The integrated model was developed in gPROMS[®] (Process Systems Enterprise, 2014) and consists of 6 units (API, excipient and lubricant feeders, co-mill, blender and tablet press). The behavior of each unit is described by a semi-empirical model. Full details on the models of each unit can be found in the work of Wang *et al.* (2017).

In this DC process, the API and the excipient are first fed to the co-mill. The lubricant is added after the co-mill to improve the flowability of the powder mixture. The API and the excipient, together with the lubricant, are then mixed in a continuous convective blender and then transported to a rotary feed frame, which feeds the mixture into dies where the press compacts the powder into tablets.

The integrated flowsheet model consists of 22 input factors, including raw material properties, process parameters and tablet geometry, and 20 output variables, which include blends material properties, operation safety, mixing characterization and tablet product qualities. A detailed list of the input and output variables can be found in the cited reference.

As shown by Wang *et al.* (2017), sensitivity analysis can be performed to determine the most influential input factors with respect to the desired output variables. Wang *et al.* (2017) identified 10 influential input factors, 4 of them being raw material properties that are assumed to be fixed at their nominal values throughout the entire process operation. Therefore, only 6 input factors are investigated by feasibility analysis, as reported in Table 5.1.

Table 5.1. Case study #3: list of the influential input factors and the fixed raw material properties. Originally presented by Wang *et al.* (2017).

Input #	Variable name	LB	UB	Nominal value	Measurement unit
<i>Significant input factors</i>					
1	API flow rate	2.85	3.15	[-]	[kg/h]
2	Excipient flow rate	25.365	28.035	[-]	[kg/h]
3	Co-mill blade speed	1064	1176	[-]	[RPM]
4	Blender blade speed	237.5	262.5	[-]	[RPM]
5	Tablet press fill depth	0.0095	0.0105	[-]	[m]
6	Tablet press thickness	0.002375	0.002625	[-]	[m]
<i>Fixed material properties</i>					
MP1	Excipient bulk density	[-]	[-]	400	[kg/m ³]
MP2	Excipient true density	[-]	[-]	2500	[kg/m ³]
MP3	Excipient d_{50}	[-]	[-]	120	[μ m]
MP4	Excipient d_{90}	[-]	[-]	250	[μ m]

A total of 20 double-sided inequality constraints (i.e. 40 single-sided inequality constraints) are assigned to the 20 output variables of the flowsheet model, thus making the feasibility problem

as a 6-dimensional problem with 40 inequality constraints. The maximum value of all these constraints is the feasibility function $\psi(\mathbf{x})$ of this case study.

5.5 Results

In the following, the results of the proposed methodology for the three case studies discussed above are presented and critically discussed.

5.5.1 Results for case study #1

The original feasibility problem for this case study was a 3-dimensional problem with quality constraints set for the output of the second unit. In other words, it has been imposed that the quality constraints must not be violated for the output of the second unit, while some of the constraints described by Eqs. (5.15)-(5.17) for the output of the first unit may possibly be violated.

Mathematically, the feasible region on the 3-dimensional input domain identified by (x_1, x_2, u_1) was considered by accounting only for the feasibility function ψ_2 . The objective was to understand if the proposed methodology can reduce the original 3-dimensional problem to a lower-dimensional problem, without substantially affecting the accuracy of the prediction of the original feasible region.

Following the procedure of Fig. 5.3, a simulated dataset was first built using the single unit models described in section 5.4.1. 49 = 7² KSA-selected calibration samples were used to build the PLS model. The obtained PLS diagnostics is reported in Table 5.2.

Table 5.2. Case study #1. Diagnostics of the PLS model

LV #	R_X^2 (%)	$R_{X,cum}^2$ (%)	R_Y^2 (%)	$R_{Y,cum}^2$ (%)
1	53.12	53.12	23.12	23.12
2	42.01	95.13	28.08	51.20
3	4.87	100.00	7.76	58.96

From Table 5.2, it can be seen that with $A = 2$ latent variables the PLS model explains 95.13 % of the input variability and 51.20 % of the output variability. Recalling Figure 5.4, the case study falls into scenario #2 and the methodology can be effective to reduce the original 3-dimensional problem to a 2-dimensional problem. The linear transformation $L_{PLS}: (x_1, x_2, u_1) \leftrightarrow (t_1, t_2)$ was therefore applied and feasibility analysis was executed with this new set of latent variables. A rectangular-grid sampling with 49 initial samples (as we did to build the PLS model) was used to build the initial surrogate onto the 2-dimensional latent space. The contour plot of the original feasibility function ψ_2 projected onto the 2-dimensional

latent space is shown in Figure 5.10a. The grey-shaded areas represent the projection of the true feasible region of the process onto the latent space. It is important to notice that these areas, even if computed with the original feasibility function, do not correspond exactly to the 3-dimensional feasible region of the process, due to the projection error of the linear transformation L_{PLS} .

Adaptive sampling (section 5.2.5) was used to iteratively sample new points and improve the surrogate accuracy. A scale factor for the adaptive sampling 10 times larger than the default value identified by Eq. 5.11 was used, thus enhancing local search with respect to global search. In other words, it was preferred optimizing the sample budget without facing the risk to over explore the latent space. In all the case studies considered, it was verified *a posteriori* that all the disjoint feasible regions were correctly identified by the surrogate using the increased scale factor, thus confirming the validity of our choice. 100 iterations were fixed as the maximum number of iterations for the adaptive sampling. The surrogate accuracy after 100 iterations can be seen in Figure 5.10b, where the blue dots represent the initial sampling points and the red circles represent the adaptive sampling points after 100 iterations. From the same figure, it can be seen that the adaptive sampling points are correctly located near the boundaries of the disjoint feasible regions onto the latent space.

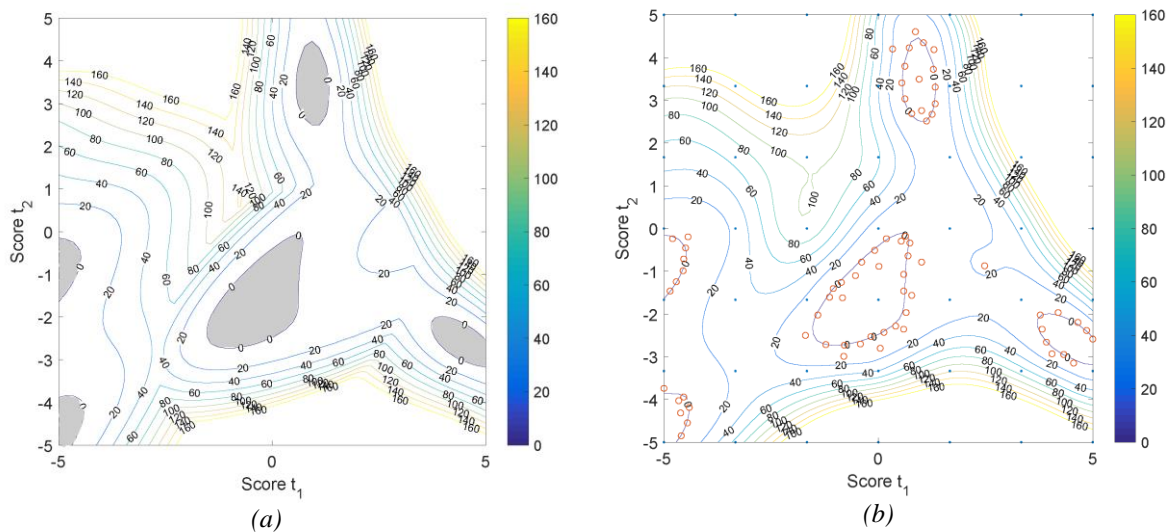


Figure 5.10. Case study #1. (a) Projection of the original function onto the latent space. The grey-shaded area represents the projection of the feasible region. (b) RBF model after 100 iterations. The blue dots represent the initial samples. The red circles represent the adaptive sampling points after 100 iterations.

The three metrics $CF(\%)$, $CIF(\%)$, $NC(\%)$, that quantify the accuracy and robustness of the results with respect to the projection of the actual feasible region onto the latent space, were then computed. After 100 iterations, the following values were obtained: $CF(\%) = 99.1$, $CIF(\%) = 99.9$, $NC(\%) = 1.20$. These values confirmed that the surrogate can identify with high accuracy the projections of both the feasible and infeasible regions (over 99% of the

feasible regions were correctly discovered). They also confirmed that the percent of overestimated feasible regions is very low (1.20 %).

As discussed in section 5.2.6, these metrics quantify the accuracy of the results on the 2-dimensional latent space and are computed just as an intermediate step for the final assessment of the results.

Since the values of these metrics resulted to be very satisfactory, the results in terms on the original input factors were computed, by applying the inverse linear transformation L_{PLS}^{-1} . The values of the new metrics $CF_2(\%)$, $CIF_2(\%)$, $NC_2(\%)$, that describe the accuracy and robustness of the results with respect to the original input space, were: $CF_2(\%) = 98.2$, $CIF_2(\%) = 99.90$, $NC_2(\%) = 1.32$. As expected, the values of these metrics are worse than the previous ones, due to the projection error of L_{PLS} . However, these values confirm that over 98% of the actual 3-dimensional feasible regions are correctly discovered, and that the methodology overestimates only 1.32% of them.

In conclusion, the results are very satisfactory: the methodology was able to transform a 3-dimensional problem into a 2-dimensional problem, reducing the computational burden for feasibility analysis, but at the same time guaranteeing the prediction of the actual feasible region with an accuracy higher than 98 %. Table 5.3 shows a summary of the results obtained for this case study.

Table 5.3. Case study #1. Surrogate accuracy onto the latent and original input space after 100 iterations.

Metric	Latent space (2-D)	Original input space (3D)
CF(%)	99.10	98.20
CIF(%)	99.90	99.90
NC(%)	1.20	1.32

For this case study, the total simulation time for feasibility analysis was 11 min and 32 sec.

5.5.2 Results for case study #2

In the second case study, the original problem is 7-dimensional and it would make feasibility analysis computationally difficult if not impossible. Moreover, giving a compact representation (i.e. graphic) of a 7-dimensional feasible region is impossible and makes the interpretation of the results cumbersome.

The PLS model was built with 49 KSA-selected calibration samples and the diagnostics of Table 5.4 were obtained.

As can be seen in Table 5.4, even though the original number of input factors is 7, the “driving forces” that affect the process output are much less. In fact, with 2 latent variables, 88.14 % of the variability of the original input space can be captured, and 41.15 % of the variability of

output space is described. It is worth noticing that this represents an intermediate situation between scenario #1 and scenario #2, since the amount of \mathbf{X} - variability explained is smaller than 90 %.

Table 5.4. Case study #2. Diagnostics of the PLS model

LV #	$R_X^2(\%)$	$R_{X,cum}^2(\%)$	$R_Y^2(\%)$	$R_{Y,cum}^2(\%)$
1	51.12	51.12	26.89	26.89
2	37.02	88.14	14.26	41.15
3	5.99	94.13	6.11	47.26
4	1.09	95.22	1.12	48.38
5	1.78	97.00	3.29	51.67
6	0.15	97.15	1.01	52.68
7	2.85	100.00	2.15	54.83

Therefore, if the case study is treated as belonging to scenario #2, feasibility analysis can be applied to a reduced 2-dimensional problem and the accuracy of the results must be assessed. On the other hand, if we consider scenario #1, the number of latent variables should be increased to increase $R_{X,cum}^2$. From Table 5.4, it can be noticed that if the number of latent variables is increased to 3, $R_{X,cum}^2$ becomes 94.13%, while $R_{Y,cum}^2$ does not increase significantly. In this case, feasibility analysis must be performed on a 3-dimensional latent space thus increasing the computational burden. To compare the balance between the reduction of the projection error and the increase of the computational burden, the methodology was tested with both 2 and 3 latent variables and compared the results.

In the 2-dimensional case, the initial surrogate was built with 49 calibration samples selected according to a rectangular-grid sampling. The projection of the original function onto the latent space is shown in Figure 5.11a. The surrogate accuracy after 100 iterations is shown in Figure 5.11b.

It can be noticed that the adaptive sampling points correctly identify the boundary of the feasible regions for the process. The metrics CF(%), CIF(%), NC(%) resulted to be 92.15, 98.87 and 2.97 respectively. When computed on the original 7-dimensional input space, they become $CF_2(\%) = 90.11$, $CIF_2(\%) = 98.15$, $NC_2(\%) = 3.26$. These results show that, with 2 latent variables, over 90% of the feasible regions are correctly identified by the surrogate, and the percent of overestimated feasible regions is smaller than 4%.

In the 3-dimensional case, a Latin hypercube sampling strategy with 64 samples was used to build the initial surrogate. We noticed that the accuracy and robustness of the results increases as expected. In fact, the surrogate accuracy after 100 iterations with respect to the original input space is $CF_2(\%) = 94.13$, $CIF_2(\%) = 99.15$, $NC_2(\%) = 1.8$. However, the total computational time for the feasibility analysis significantly increases from the 2-dimensional to the 3-dimensional problem (from 36 min. to 1h 26 min), as shown in Table 5.5.

This case study shows how the choice of the number of latent variables should be done as a compromise between the desired accuracy of the results and the reduction of the computational burden for the feasibility analysis. For most practical applications, an accuracy of around 90%

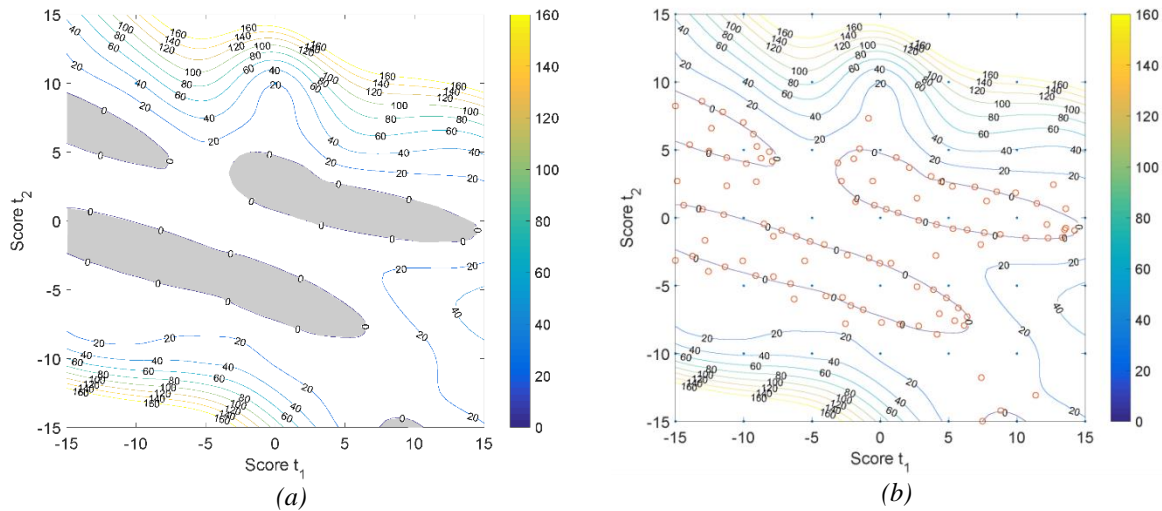


Figure 5.11. Case study #2. (a) Projection of the original function onto the latent space. The grey-shaded area represents the projection of the feasible region. (b) RBF model after 100 iterations. The dots represent the initial samples. The red circles represent the adaptive sampling points after 100 iterations.

Table 5. Case study #2. Comparison of the surrogate accuracy and simulation time after 100 iterations with 2 and 3 latent variables respectively.

Metric	2 latent variables (2-D)	3 latent variables (3-D)
CF ₂ (%)	90.11	94.13
CIF ₂ (%)	98.15	99.15
NC ₂ (%)	3.26	1.89
Total computational time	32 min	1h 26 min

in the identification of the feasible region is acceptable and therefore the original 7-dimensional problem can be reduced to a 2-dimensional problem. If more accuracy is required, this has to be done by increasing the computational burden for the feasibility analysis.

5.5.3 Results for case study #3

This case study is a comprehensive practical exercise of design space identification for a continuous manufacturing pharmaceutical process. The step-by-step methodology was applied to this 6-dimensional problem, involving 40 inequality constraints for the model output. The “historical” dataset was generated using the original integrated flowsheet model by selecting

25 KSA calibration samples. A PLS model was then built for the overall process. The PLS diagnostics are presented in Table 5.6.

Table 5.6. Case study #3. Diagnostics of the PLS model.

LV #	R_X^2 (%)	$R_{X,cum}^2$ (%)	R_Y^2 (%)	$R_{Y,cum}^2$ (%)
1	62.12	62.12	36.45	36.45
2	29.60	91.72	12.13	48.58
3	4.12	95.84	11.16	59.74
4	3.16	99.00	4.21	63.95
5	0.17	99.17	4.01	67.96
6	0.83	100.00	1.04	69.00

Table 5.6 results illustrate that the PLS model is able to capture 91.72% of the \mathbf{X} -variability and 48.58 % of the \mathbf{y} -variability with only two latent variables, thus giving scenario #2. The explained \mathbf{X} -variability would increase to 95.84% if 3 latent variables were chosen. Considering $R_{X,cum}^2 = 91.72\%$ a good projection accuracy, only 2 latent variables were used for this case study. Feasibility analysis was then performed on the reduced 2-dimensional problem with the new set of latent variables. A total of $25 = 5^2$ samples were used to build the initial surrogate according to a rectangular grid sampling scheme. The increased scale factor discussed for case study #1 was also used. The adaptive sampling routine was applied and the surrogate RBF model onto the latent space after 100 iterations is shown in Figure 5.12b. The projection of the original function onto the latent space is shown in Figure 5.12a. It can be noticed that the adaptive sampling points clearly identifies the boundary of the feasible region (i.e. the DS) of the process. The values of the three metrics CF(%), CIF(%), NC(%) after 100 iterations are 95.22, 97.95 and 1.13 respectively.

According to the ICH Q8(R2) guideline, the DS of the process can be described as a “combination of variables such as components of a multivariate model” (ICH, 2009). Moreover, the DS “can be explained mathematically through equations describing relationships between parameters for successful operation” (ICH, 2009). Being the original process parameters linear combinations of the PLS scores $[t_1, t_2]$ of Figures 5.12a and 5.12b according to Eq. (5.1), the representation of the DS of Figure 5.12b fully complies with the regulatory requirements. However, from a practitioner’s point of view, describing the DS in terms of latent variables may be difficult to translate into input materials properties and process operating conditions.

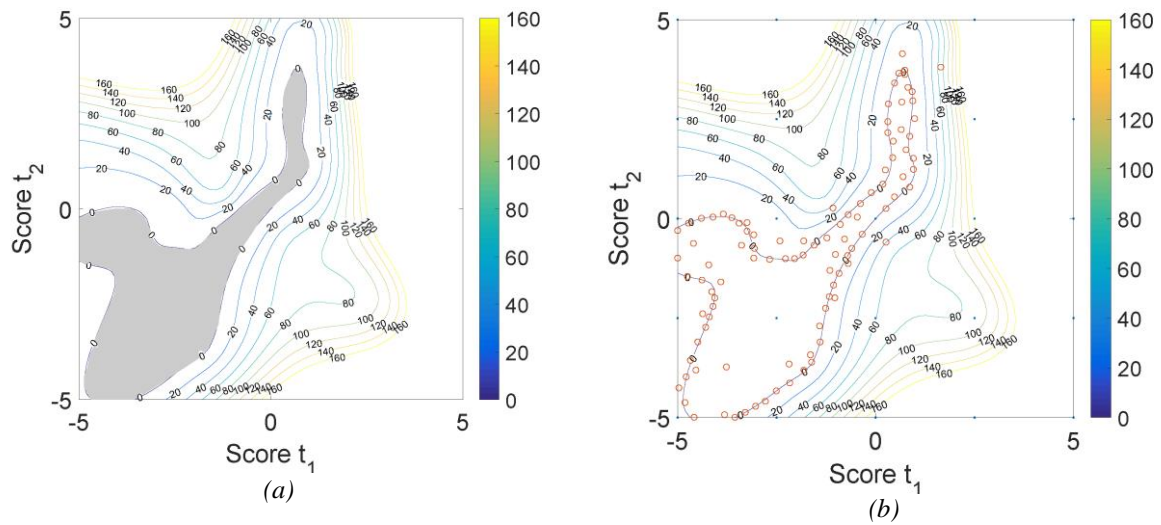


Figure 12. Case study #3. (a) Projection of the original function onto the latent space. The grey-shaded area represents the projection of the feasible region. (b) RBF model after 100 iterations. The blue dots represent the initial samples. The red circles represent the adaptive sampling points after 100 iterations.

It must be recalled that each point of Figure 5.12 corresponds to a combination of the six input factors of Table 5.1, and the resulting DS can be easily translated into practical information for the user by applying the inverse of the linear transformation L_{PLS} (step 6 of Figure 6). For example, let $\mathbf{t}^f = [t_1^f, t_2^f]$ be a point inside the feasible region identified the RBF model of Figure 12b. The corresponding 6-dimensional feasible operating condition \mathbf{x}^f can be obtained by applying the linear transformation (5.1) using the loadings of the PLS model (step 2 of Figure 5.3). The procedure can be repeated for any desired number of points falling within, or at the boundary of, the feasible region identified by the RBF model (Figure 5.12b), thus translating the latent space results into practical operating conditions. As an example, let $\mathbf{t}^f = [-2.12, -3.18]$ be one of the points inside the feasible region of Figure 5.12b. The corresponding input combination that can be obtained by applying L_{PLS}^{-1} is reported in Table 5.7. The metrics $CF_2(\%)$, $CIF_2(\%)$ $NC_2(\%)$ on the original 6-dimensional space were 93.50, 97.20 and 2.30 respectively. These results show that more than 93% of the feasible region is correctly identified by the RBF model in the original input space, and the percent of overestimated feasible region is a little bit higher than 2%. A summary of the results is presented in Table 5.8.

Table 5.7. Case study #3. Example of feasible point expressed in terms of original input factors, given its coordinates on the latent space $t^f = [-2.12, -3.18]$.

Original input factor combination	Numerical value	Measurement unit
API flow rate	3.05	[kg/h]
Excipient flow rate	26.925	[kg/h]
Co-mill blade speed	1122	[RPM]
Blender blade speed	253.2	[RPM]
Tablet press fill depth	0.00998	[m]
Tablet press thickness	0.002502	[m]

Table 5.8. Case study #3. Surrogate accuracy onto the latent and original input space after 100 iterations.

Metric	Latent space (2-D)	Original input space (6-D)
CF(%)	95.22	93.50
CIF(%)	97.95	97.20
NC(%)	1.13	2.30

The results are very satisfactory considering the dimensionality of the original problem. In a previous work, Wang *et al.* (2017) used a matrix of 2-dimensional feasibility plots to assess the accuracy of the results, without accounting for the multivariate correlation between the input factors. With the proposed methodology, a high accuracy in the prediction of the feasible region is obtained while fully taking into account the multivariate nature of the original problem.

5.6 Conclusions

In this work, a novel and systematic methodology to identify the feasible region of an integrated process by applying a linear transformation of the original input factors into a new set of PLS-determined latent variables, and then applying feasibility analysis on this low dimensional problem was proposed. Given the original integrated flowsheet model, a PLS model was first built based on a simulated dataset obtained from the original flowsheet model. According to the amount of input and output variability explained by the PLS model for the chosen number of latent variables, three possible scenarios were identified. In the first scenario, the methodology can be applied only if the number of latent variables can be increased while simultaneously maintaining the feasibility analysis computationally applicable. In the second scenario, the methodology can be directly used and can show its potential; in the latter scenario, feasibility analysis is not needed and the feasible region can be directly identified with the PLS model. The PLS model was used as a linear transformation between the original set of input variables

and the new set of latent variables, and RBF-based adaptive sampling feasibility analysis on the reduced latent space was performed. The quality and robustness of the results were assessed with three metrics, that describe how well the feasible and infeasible regions are identified by the RBF model, and how large is the portion of overestimated feasible region. The methodology was tested on three relevant case studies, the first two involving multiple unit test problems with 3 and 7 input factors respectively, and the last one involving a continuous manufacturing direct compaction process of a pharmaceutical powder. For all the case studies, the ability of the methodology to reduce the problem dimensionality while maintaining a good accuracy and robustness of the results was shown. It was also shown the ability of the methodology to give a simple and compact representation of the feasible regions for the process under investigation. For the second case study, it was proved how the choice of the number of latent variables must be done by finding a compromise between the desired accuracy and the reduction of the computational burden.

Chapter 6

Design space maintenance by online model adaptation in pharmaceutical manufacturing*

In this Chapter, a model adaptation strategy that allows obtaining an up-to-date representation of the design space in the presence of process/model mismatch is proposed. First, the motivation and the problem statement are presented. Then, the state-of-the art for design space determination using classical feasibility analysis is briefly discussed, together with a critical review of its main limitations. The proposed methodology is then presented and tested on two simulated case studies. Finally, the results of the proposed design space adaptation strategy are shown and critically analyzed.

6.1 Introduction

Although not compulsory, a detailed description of the design space (DS) of a pharmaceutical product is a key step of the Quality by Design (QbD) paradigm advocated by the pharmaceutical regulatory agencies. The design space is defined as the multidimensional combination and interaction of raw materials properties and critical process parameters that have been demonstrated to meet assigned specifications on the critical quality attributes of the final product (ICH, 2009). If the submitted design space is approved, the manufacturing process can be run anywhere within the design space without initiating a regulatory post approval change procedure. This increases process flexibility, since manufacturers can decide to change the process parameters and/or the raw material characteristics within the approved design space without the need to obtain a new regulatory endorsement.

The use of mathematical models has been extensively exploited to assist design space description during the last decade, both in academia and in industry. In the following, we will

* Bano G., Facco, P., Ierapetritou, M., Bezzo, F., Barolo, M. (2018) Design space maintenance by online model adaptation in pharmaceutical manufacturing. Submitted to: *Comput. Chem. Eng.*

refer to the prediction of the design space that can be obtained using a process model as to the *model-based* DS, to distinguish it from the *actual* (and unknown) DS. In the context of data-driven models, the use of response surface methodologies (Chatzizacharia and Hatzivramidis, 2014; Kumar *et al.*, 2014), multivariate experimental design (Charoo *et al.*, 2012), high-dimensional model representations (Li *et al.*, 2001; Boukouvala *et al.*, 2010), kriging methodologies (Jia *et al.*, 2009; Boukouvala *et al.*, 2010) and latent-variable model inversion (Tomba *et al.*, 2012) have been reported to support DS description, with application to the manufacture of drug substances, intermediates and drug products. With reference to first-principles models (or semi-empirical ones), techniques such as Monte-Carlo simulations (Burt *et al.*, 2011; Brueggemeier *et al.*, 2012; Garcia-Munoz *et al.*, 2015; Garcia-Munoz *et al.*, 2018), surrogate-based feasibility analysis (Banerjee *et al.*, 2010; Rogers and Ierapetritou, 2015; Wang *et al.*, 2017a; Wang *et al.*, 2017b; Bano *et al.*, 2018) and operability analysis (Vinson and Georgakis, 2000; Uztürk and Georgakis, 2002) have been used to assist DS determination. Steady-state systems (Lima *et al.*, 2010; Boukouvala and Ierapetritou, 2012) as well as dynamic systems (Garcia-Munoz *et al.*, 2015; Rogers and Ierapetritou, 2015) have been considered.

Whether data-driven or first-principles, DS description using a mathematical model requires prior knowledge of the uncertainty affecting the model and of the disturbances that may affect plant operation. However, neither all sources of uncertainty nor all disturbances can be modeled *a priori*. Additionally, it should be reminded that the behavior of an industrial production plant is typically the outcome of a complex interaction of physical phenomena, and is affected by downstream and upstream units, as well by the external environment (Pantelides, 2009). Laboratory or pilot-scale measurements often do not embed enough information about these interactions, and this may result in deficiencies in the plant model. Finally, drifts may occur during actual operation of the plant due to (typically slow) physical phenomena, such as equipment fouling, leaks, catalyst deactivation and the like (Pantelides and Renfro, 2013). All of the above occurrences can result in process-model mismatch, which impacts on the accuracy of the critical quality attributes predicted by the model, thus in turn undermining the appropriateness of the design space described by means of the model as originally developed. Process-model mismatch can sometimes be mitigated or even removed by adjusting some of the original model parameters in real time. This continuous model adaptation methodology is long known in process engineering (Lainiotis, 1971; Cheng *et al.*, 1993; Ogunnaike, 1994; Pantelides, 2009) and, following the regulatory parlance, may also be defined as online model maintenance. Maintenance of a process model also affects the DS described using that model, an issue that is also addressed by the regulatory agencies in the ICH Q8(R2) guideline (ICH, 2009), where it is noted that "...for certain design spaces using mathematical models, periodic maintenance could be useful to ensure the model's performance..." and that "...expansion, reduction or redefinition of the design space could be desired upon gaining additional process knowledge...". The fact that a DS should be adapted during the product lifecycle has started

being acknowledged by the academic community (Herwig *et al.*, 2017). However, no indications have been provided so far on *how* DS maintenance should be carried out and reported. This study is intended to offer a contribution to the development of a DS maintenance strategy. A discussion on the regulatory implications of DS maintenance is beyond of the scope of the study, and therefore this issue will not be addressed in this paper.

In this Chapter, a methodology that allows continuous maintenance of the model-based representation of the DS using online model adaptation is proposed. The main focus is on pharmaceutical processes described by first-principles models or semi-empirical models, and the availability of plant measurements from online sensors is assumed. Recent advances in online applications of first-principles models (commonly referred to as tracking simulators; Patwardhan *et al.*, 2012; Pantelides *et al.*, 2016) are coupled with advances in the area of feasibility analysis and surrogate-based feasibility analysis for DS determination (Grossmann *et al.*, 2014; Bhosekar and Ierapetritou, 2017). Specifically, a dynamic state estimator suitable for differential-algebraic equation (DAE) systems (Cheng *et al.*, 1997) is exploited for online model adaptation, and an online update of the model-based DS representation using feasibility analysis techniques is obtained. Situations where the impact of process-model mismatch can be removed online by adjusting some of the original model parameters are analyzed.

The proposed methodology is tested on two simulated case studies: the former involves the roller compaction of microcrystalline cellulose based on the model of Hsu *et al.* (2010), and the latter the fermentation of penicillin in a fed-batch bioreactor based on the model of Riascos and Pinto (2004).

6.2 Methods

In the following, the two modelling techniques that will be jointly exploited in this study are briefly reviewed. It is assumed that a first-principles model or semi-empirical model of the pharmaceutical unit or process under investigation is available.

6.2.1 State estimator for DAE systems

Most of the models describing a pharmaceutical system are governed by index-1 or higher-index system of nonlinear DAEs of the form:

$$\begin{bmatrix} \dot{\mathbf{x}}_1(t) \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{f}_1(\mathbf{x}_1(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \\ \mathbf{f}_2(\mathbf{x}_1(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \end{bmatrix} + \begin{bmatrix} \mathbf{w}_1(t) \\ \mathbf{w}_2(t) \end{bmatrix} \quad (6.1)$$

subject to the initial conditions:

$$\mathbf{x}_1(0) = \mathbf{x}_{1,0} + \mathbf{w}_1(0) \quad (6.2)$$

where the N -dimensional system state $\mathbf{x}(t)$ is split into an N_1 -dimensional differential state $\mathbf{x}_1(t)$ and an N_2 -dimensional algebraic state $\mathbf{x}_2(t)$. The model error $\mathbf{w}(t)$ is a zero-mean random process with unknown statistics that, in accordance with $\mathbf{x}(t)$, can be split into the N_1 -vector $\mathbf{w}_1(t)$ and the N_2 -vector $\mathbf{w}_2(t)$. The I -dimensional vector $\mathbf{u}(t)$ is the control vector. Online measurements are assumed to be available at discrete sampling instants t_k and are collected in the M -dimensional vector $\mathbf{y}(t_k)$. The online observations are related to the state vector at time t_k according to:

$$\mathbf{y}(t_k) = \mathbf{h}(\mathbf{x}_1(t_k), \mathbf{x}_2(t_k), t_k) + \mathbf{v}(t_k); k = 1, 2, \dots, K \quad (6.3)$$

where $\mathbf{v}(t_k)$ is the M -dimensional vector of unknown measurement noise. Eq. (6.1) is the state-space model of the system, while Eq. (3) is the observation model.

In principle, index-1 DAE systems can always be reformulated in standard ODE form, and higher index DAEs can be reduced to 1-DAE systems e.g. using Pantelides' algorithm (Pantelides, 1988) or dummy derivative index reduction (Mattsson and Söderlind, 1993). Therefore, standard state estimators based on pure ODE state-space models may be used for these systems. However, in practical applications, obtaining consistent initial conditions for the system state and the initial model error covariance for index-1 DAE systems can be very cumbersome. In this case, derivation of state estimators based on the original DAE form can be beneficial to overcome this issue. In this study we use the sequential state estimator proposed by Cheng *et al.* (1997), which is directly obtained from the original DAE formulation of the system.

Given the discrete time measurements $\mathbf{y}(t_k)$ in the time interval $0 \leq t_k \leq T$, it is required to estimate $\mathbf{x}_1(t)$ and $\mathbf{x}_2(t)$ at a given time t . The estimation criterion is the minimization of the least-squares error functional (Cheng *et al.*, 1997):

$$\begin{aligned} I = & \frac{1}{2} [\mathbf{x}_1(0) - \mathbf{x}_{1,0}]^T \mathbf{P}_{11}^{-1}(0) [\mathbf{x}_1(0) - \mathbf{x}_{1,0}] + \\ & + \frac{1}{2} \int_0^T \left\{ \begin{bmatrix} \dot{\mathbf{x}}_1(t) \\ 0 \end{bmatrix} - \begin{bmatrix} \mathbf{f}_1(\mathbf{x}_1(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \\ \mathbf{f}_2(\mathbf{x}_1(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \end{bmatrix} \right\}^T \mathbf{Q}^{-1}(t) \left\{ \begin{bmatrix} \dot{\mathbf{x}}_1(t) \\ 0 \end{bmatrix} \right. \\ & \left. - \begin{bmatrix} \mathbf{f}_1(\mathbf{x}_1(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \\ \mathbf{f}_2(\mathbf{x}_1(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \end{bmatrix} \right\} dt + \\ & + \frac{1}{2} \int_0^T [\mathbf{y}(t) - \mathbf{h}(\mathbf{x}_1(t), \mathbf{x}_2(t), t)]^T \mathbf{R}_d^{-1}(t) [\mathbf{y}(t) - \mathbf{h}(\mathbf{x}_1(t), \mathbf{x}_2(t), t)] dt \end{aligned} \quad (6.4)$$

where the first term minimizes the square error of the initial estimate of $\mathbf{x}_1(0)$, the second term minimizes the integral square model error and the third term minimizes the integral square measurement error. Although the three matrices $\mathbf{P}_{11}^{-1}(0)$, $\mathbf{Q}^{-1}(t)$, $\mathbf{R}^{-1}(t)$ do not have statistical

meaning, they can be tuned by the user to reflect the errors on the initial estimate of the differential state, the process model and the measurement device, respectively. Their meaning is analogous to the covariances of the initial state errors, the process noise and the measurement noise respectively of the linear Kalman filter (Kalman, 1960). Matrix $\mathbf{R}_d^{-1}(t)$ can be computed according to:

$$\mathbf{R}_d^{-1}(t) = \sum_{k=1}^K \mathbf{R}^{-1}(t_k) \delta(t - t_k) \Delta_k \quad (6.5)$$

where Δ_k is the sampling interval at sample k and $\delta(t - t_k)$ is the Dirac delta function.

Application of the calculus of variations to (6.4) allows obtaining a sequential solution to the estimation problem. The relevant equations that can be obtained are reported in Appendix A. The sequential estimator that can be obtained simplifies to the standard extended Kalman filter (EKF) (Ray, 1981) equations when the system is described by ordinary differential equations (ODEs). Therefore, the approach considered in this study can be considered as an extension of the standard EKF estimation to DAE systems. The robustness of this estimator has been recently proved to be satisfactory even with large (227 differential variables plus 14268 algebraic variables) high-order nonlinear DAE systems (Pantelides *et al.*, 2016).

6.2.2 Feasibility analysis for dynamical systems

Feasibility analysis (Halemane and Grossman, 1983) is a mathematical technique that can be used to identify the subset of combinations of uncertain parameters that guarantee to satisfy all the relevant process constraints (e.g. quality constraints, production constraints, environmental constraints). These combinations represent the feasible region of the process. According to the standard terminology of feasibility analysis, the uncertain parameters collect three different sources of uncertainty (Dimidiatris and Pistikopoulos, 1995):

- a) variability on external variables such as raw material properties (e.g., composition of available feedstock);
- b) variability on internal variables such as critical process parameters (e.g., fluctuations of the flowrate of the inlet streams);
- c) uncertainty on model parameters (e.g., heat/mass transfer coefficients, kinetic constants).

These sources of uncertainty can potentially affect the key performance indicators (KPIs) of the process.

From a general perspective, the physical behavior of a dynamic system can be described by the systems of DAEs of Eq. (1). This set of equations can be written in the implicit form:

$$\mathbf{F}(\mathbf{d}, \dot{\mathbf{x}}(t), \mathbf{x}(t), \mathbf{u}(t), \boldsymbol{\theta}(t), t) = \mathbf{0}; \quad \mathbf{x}(0) = \mathbf{x}_0 \quad (6.6)$$

where the vector collecting the design variables \mathbf{d} (which is fixed, since we assume that the process design is fixed) and the vector collecting the uncertain parameters $\boldsymbol{\theta}(t)$ are shown explicitly.

The key performance indicators of the process are subject to a set of constraints (e.g. quality constraints). These constraints can be path constraints, i.e. of the form:

$$\mathbf{g}(\mathbf{d}, \dot{\mathbf{x}}(t), \mathbf{x}(t), \mathbf{u}(t), \boldsymbol{\theta}(t), t) \leq \mathbf{0} \quad (6.7)$$

or a set of $p = 1, 2, \dots, P$ point constraints, i.e.:

$$\mathbf{g}_p(\mathbf{d}, \mathbf{x}(t_p), \mathbf{u}(t_p), \boldsymbol{\theta}(t_p), t_p) \leq \mathbf{0}; \quad p = 1, \dots, P. \quad (6.8)$$

A common type of point constraints is the end-time point constraint:

$$\mathbf{g}_{\text{end}}(\mathbf{d}, \mathbf{x}(t_{\text{end}}), \mathbf{u}(t_{\text{end}}), \boldsymbol{\theta}(t_{\text{end}}), t_{\text{end}}) \leq \mathbf{0}. \quad (6.9)$$

where t_{end} may be, for example, the total duration of a batch or the total duration of a production campaign. The characterization the feasible region of the process can be obtained by solving a flexibility test problem (Halemane and Grossmann, 1983). The conventional formulation of the flexibility test problem, under the assumption that no control variables are manipulated (open-loop scenario), is as follows:

$$\begin{aligned} \chi(\mathbf{d}) &= \max_{\boldsymbol{\theta}(t) \in \mathbf{T}(t)} \Psi(\mathbf{d}, \boldsymbol{\theta}(t)) \\ \text{subject to: } \Psi(\mathbf{d}, \boldsymbol{\theta}(t)) &= \max_{i \in I} (g_i(\mathbf{d}, \mathbf{x}(t), \boldsymbol{\theta}(t), t)), \quad i = 1, 2, \dots, I. \\ \mathbf{F}(\mathbf{d}, \dot{\mathbf{x}}(t), \mathbf{x}(t), \boldsymbol{\theta}(t), t) &= \mathbf{0}; \\ \mathbf{x}(0) &= \mathbf{x}_0 \\ \mathbf{T}(t) &= \{\boldsymbol{\theta}(t) | \boldsymbol{\theta}^L(t) \leq \boldsymbol{\theta}(t) \leq \boldsymbol{\theta}^U(t)\} \end{aligned} \quad (6.10)$$

where $\Psi(\mathbf{d}, \boldsymbol{\theta}(t))$ is the feasibility function, i.e. the maximum value of all the I constraints on the process KPIs. The practical meaning of the feasibility function is straightforward: if $\Psi(\mathbf{d}, \boldsymbol{\theta}(t)) > 0$, there is at least one constraint that is being violated for the given design and the process is not feasible at the point in time considered. If $\Psi(\mathbf{d}, \boldsymbol{\theta}(t)) < 0$, the process is feasible at time t . If $\Psi(\mathbf{d}, \boldsymbol{\theta}(t)) = 0$, vector $\boldsymbol{\theta}(t)$ identifies the boundary of the feasible region at the given point in time. We will refer to the feasible region of the process as to the set of uncertain parameter values that satisfies:

$$R(t) = \{\boldsymbol{\theta}(t) \mid \Psi(\mathbf{d}, \boldsymbol{\theta}(t)) \leq 0\} \quad , \quad (6.11)$$

namely, the combinations of uncertain parameters $\boldsymbol{\theta}(t)$ that satisfy *all* the constraints on the key performance indicators.

The solution of the bi-level optimization problem (6.10) can be deterministic, i.e. by directly using the model equations (6.6), or can be obtained exploiting surrogate-based methods (Boukouvala and Ierapetritou, 2012).

In the former case, the solution can be obtained by embedding a differential solver into the optimization algorithm or by converting the differential equations into algebraic equations, e.g. by means of the direct quadrature (DC) method (Adi and Chang, 2013). In both cases, the solution can be computationally very demanding if the process model is expensive to evaluate or in the presence of black-box constraints. In the latter case, a surrogate is used as a cheap approximation of the original process model, and the optimization problem (6.10) is solved using this surrogate. This clearly results in a lower computational cost at the expense of introducing an approximation error. A detailed surrogate-based approach for dynamical systems based on kriging regression was developed by Rogers and Ierapetritou (Rogers and Ierapetritou, 2015).

6.3 Feasibility analysis for design space determination

In the following, a systematic description of the application of feasibility analysis for design space determination is given. First, the terminology used in classical feasibility analysis formulation is adapted to the one used by the regulatory documents for design space determination. Secondly, the application of feasibility analysis for offline model-based DS identification is described. Thirdly, the typical procedure adopted for offline model-based DS identification is presented and the motivation that leads to an adaptive online model-based DS determination is discussed.

6.3.1 Feasible region and design space

The definition of feasible region as described by Eq. (6.11) is closely related to the definition of design space of a pharmaceutical product. To make the terminology adopted by the pharmaceutical regulatory documents consistent with the one adopted in the classical feasibility analysis formulation, it must be observed that:

- i) when applied to a design space identification, the constraints on the key performance indicators that need to be satisfied are the constraints on the critical quality attributes of the pharmaceutical product considered (e.g., mean particle size of a pharmaceutical granulate, ribbon density etc.). These constraints reflect the quality specifications that are imposed on the final product.

- ii) The regulatory definition of design space strictly refers to the multidimensional combinations of critical process parameters and raw material properties (uncertainty sources of the types a) and b) of the previous section) that have been demonstrated to meet the specifications on the critical quality attributes of the product.. On the other hand, in the classical formulation of the feasibility region (6.11), the uncertainty vector $\boldsymbol{\theta}(t)$ also collects the uncertain parameters of the model that is used to describe the process behavior. This source of uncertainty is unrelated to the concept of design space and should therefore be considered separately from the other two sources.

In view of the above, the uncertainty vector $\boldsymbol{\theta}(t)$ in two parts can be split according to:

$$\boldsymbol{\theta}(t) = [\mathbf{q}(t); \tilde{\mathbf{p}}] \quad (6.12)$$

where $\mathbf{q}(t)$ collects the raw material properties and critical process parameters that directly affect the product critical quality attribute (i.e. the input variables), whereas $\tilde{\mathbf{p}}$ collects the subset of uncertain model parameters of the first-principles model describing the pharmaceutical process considered. The model-based DS (\widehat{DS}) can then be defined from a feasibility analysis perspective as:

$$\widehat{DS} = \{\mathbf{z}(t) \mid \Psi(\mathbf{d}, \mathbf{q}(t), \tilde{\mathbf{p}}) \leq 0\} \quad (6.13)$$

where the feasibility function $\Psi(\mathbf{d}, \mathbf{q}(t))$ is, in this context, the maximum value of all the constraints imposed on the productcritical quality attributes.. The definition of design space as described by the feasibility analysis formulation (6.13) is now fully consistent with the regulatory definition.

6.3.2 Limitations of existing model-based design space description methodologies

Design space description using a mathematical model is typically performed as an offline activity following four sequential steps, namely:

- 1) Parameter estimation and model validation.
- 2) Design space determination.
- 3) Uncertainty quantification on the estimated model parameters and uncertainty propagation to the predicted product critical quality attributes.
- 4) Uncertainty back-propagation from the product critical quality attributes to the design space obtained at step #2.

During the parameter estimation and model validation step, the information available from designed laboratory or pilot-scale experiments is exploited to obtain an accurate estimate of the model parameters (Pantelides *et al.*, 2012). Well-established techniques are available to

perform this activity for data-driven models as well as for first-principles models (Asprey and Macchietto, 2000). In the design space determination stage, a representation of the design space is obtained using one selected modeling strategy. In the uncertainty quantification step, a determination of the different sources of uncertainty (i.e., uncertainty on the model parameters, measurement error on the calibration dataset, external disturbances) that may affect the model-based design space representation obtained at the previous step is obtained. Uncertainty can be quantified with frequentist (Zhang and Garcia-Munoz, 2009; Garcia-Munoz *et al.*, 2010), or Bayesian (Peterson and Yahyah, 2009; Stockdale and Cheng, 2009; Bano *et al.*, 2018a) methodologies. During the uncertainty back-propagation step, uncertainty is propagated from the output space (i.e. the space spanned by the product critical quality attributes) to the design space using suitable uncertainty back-propagation strategies (Facco *et al.*, 2015; Bano *et al.*, 2017).

If feasibility analysis is used, steps 2-4 are merged together and the uncertainty on the model parameters and modeled disturbances is expressed in terms of expected deviations from their nominal values (Pistikopoulos and Mazzucchi, 1990; Straub and Grossmann, 1990; Dimitriadis and Pistikopoulos, 1995; Pistikopoulos and Ierapetritou, 1995). The model-based DS is then computed accounting for these expected uncertainty ranges (Wang and Ierapetritou, 2017b).

Despite uncertainty quantification, process-model mismatch may still exist due to the following reasons:

1. Not all of the disturbances that may affect plant operation can be known and modeled *a priori*;
2. Model calibration and validation (step #1) is typically performed on laboratory or pilot-scale data. During a production campaign, the plant behavior is affected by the interaction of downstream and upstream units as well as of external environmental factors. Laboratory-scale (or pilot-scale) data typically do not include enough information to capture these phenomena and therefore affect the model-based representation of the design space obtained with the calibrated model of step#1. This source of uncertainty may, in principle, be included in the uncertainty quantification analysis (step #3) and therefore back-propagated to the model-based DS (step #4). However, this would require at least a qualitative understanding of these phenomena (Pantelides and Renfro, 2013) and how they can affect the model prediction of the product critical quality attributes. This is usually a very difficult task to implement in practice and requires a strong empirical reasoning.
3. Some of the model parameters may drift during plant operation due to (typically slow) physical phenomena (e.g., catalyst deactivation, fouling, etc.).

Since process-model mismatch can potentially affect the appropriateness of the design space described using the original model, maintenance of the model-based DS would be useful

whenever process-model mismatch is encountered. In the next section, we propose a methodology to tackle this problem.

6.4 Design space maintenance by online model adaptation

An extension to the offline model-based DS identification described in section 6.3.2 is here proposed. A 4-step systematic methodology is derived that performs sequentially: *i*) an offline model-based DS identification using feasibility analysis; *ii*) an online maintenance of the model-based DS by jointly exploiting the state estimator described in section 6.2.1 and feasibility analysis.

Let us consider the production of a pharmaceutical product with a process with a fixed design and assume that a first-principles model (or semi-empirical model) describing the dynamic process is available. Mathematically, the model can be described by the DAE system (6.15) in its fully-implicit form, or by Eq. (6.1) (without the error terms) in its semi-explicit form. It is assumed that the raw material properties and critical process parameters (collected in the vector $\mathbf{q}(t)$) and the product critical quality attributes have already been identified at a previous stage of the product development. Quality specifications are imposed on the product critical quality attributes in the form of single- and/or double-sided path inequality constraints or single- and/or double-sided point inequality constraints. In the former case, the quality constraints can be expressed in the implicit form of Eq. (6.7), whereas in the latter case they can be expressed in the implicit form of Eq. (6.8). Notice that double-sided inequality constraints can be decoupled into a system of two single-sided inequality constraints. The analysis will be focused only to end-time specifications imposed on the product critical quality attributes, but the methodology is general and can be applied to any type of quality specifications. In this scenario, quality specifications are expressed by Eq. (6.9).

It is assumed that the first-principles model has been validated with laboratory experimental data and the nominal values of the model parameters have been estimated and/or taken from the literature accordingly. It is also assumed that a set of measurements $\mathbf{y}(t_k), k = 1, \dots, K$ will be available online from the plant unit at each time instant t_k . In most situations, the product critical quality attributes are not measured online and therefore do not belong to $\mathbf{y}(t_k)$. The methodology can be summarized in 4 steps, as schematically shown in Fig. 6.1.

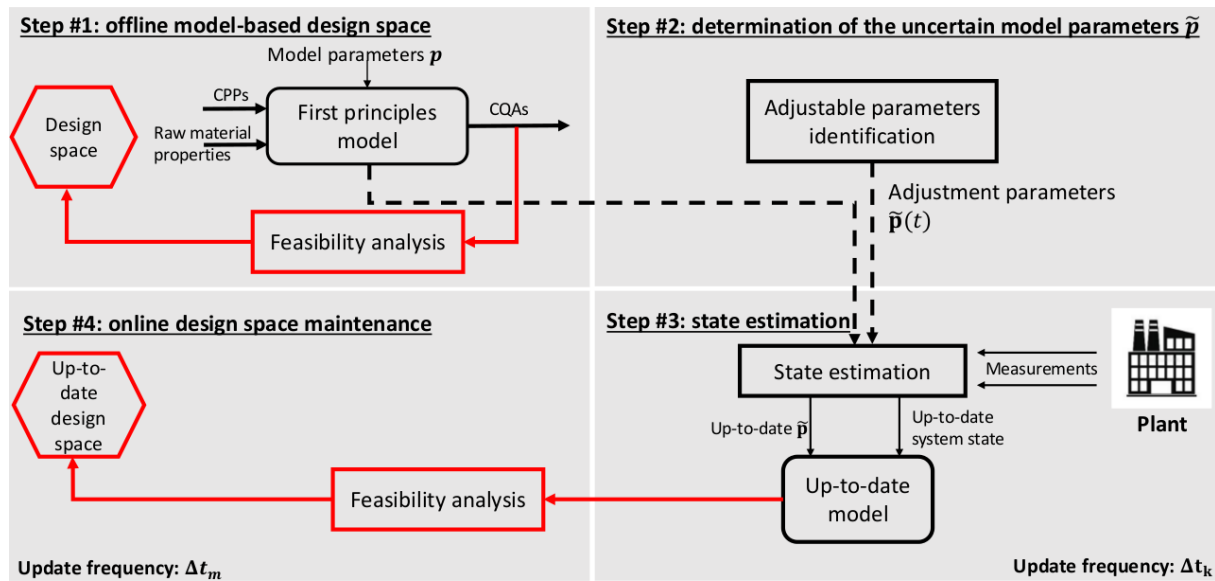


Figure 6.1. Schematic representation of the online design space maintenance methodology.

6.4.1. Proposed methodology

The step-by-step methodology is summarized as follows.

1. Step #1: offline model-based DS identification. Feasibility analysis is used to identify the design space of the process with the model whose parameters are set at their nominal values. The design space is determined as per the regulatory definition (6.13) by solving the optimization problem (6.10) and replacing $\theta(t)$ with $\mathbf{z}(t)$. The model-based DS is predicted in a deterministic way directly from the first-principles model and does not account for the uncertainty on model parameters.
2. Step #2: determination of the uncertain model parameters $\tilde{\mathbf{p}}$ to be calibrated online. As discussed in section 3, the underlying assumption that we make is that the process-model mismatch can be removed online by simultaneously updating the system state and recalibrating some of the model parameters. If severe structural deficiencies are detected in the model, the methodology is not useful: the model should be corrected and the methodology re-started from step#1.

The choice of the parameters to be calibrated online can be made following two criteria:

- a. if prior knowledge of the un-modeled physical phenomena or disturbances that may cause the process-model mismatch in the actual plant is available, the model parameters that are mostly affected by these disturbances can be determined by a qualitative preliminary diagnosis of the model;
- b. if case a) does not hold true, a preliminary sensitivity analysis can be performed offline in order to determine the model parameters that most affect the prediction of the product critical quality attributes. Since the main target is to remove the process-model mismatch affecting the prediction of the critical quality

attributes, only the most influential parameters will be calibrated during plant operation.

3. Step #3: state estimation and online parameter recalibration. Given the online measurements $\mathbf{y}(t_k)$, the state estimator described in section 6.2.1 is used to obtain, at each time t_k , both the current state vector of the system, an up-to-date and an updated (recalibrated) model providing an accurate prediction of the product critical quality attributes.

The online model recalibration is obtained through a state augmentation procedure (Ricker and Lee, 1995; Smith *et al.*, 2013). The differential subset $\mathbf{x}_1(t)$ of the state vector is augmented with a new set of time-dependent parameters $\tilde{\mathbf{p}}(t)$, leading to the overall state vector:

$$\mathbf{x}(t) = \begin{bmatrix} \mathbf{x}_1(t) \\ \tilde{\mathbf{p}}(t) \\ \mathbf{x}_2(t) \end{bmatrix}. \quad (6.14)$$

The system of equations (6.1) can then be re-written as:

$$\begin{bmatrix} \dot{\mathbf{x}}_1(t) \\ \tilde{\mathbf{p}}(t) \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{f}_1(\mathbf{x}_1(t), \tilde{\mathbf{p}}(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \\ \mathbf{0} \\ \mathbf{f}_2(\mathbf{x}_1(t), \tilde{\mathbf{p}}(t), \mathbf{x}_2(t), \mathbf{u}(t), t) \end{bmatrix} + \begin{bmatrix} \mathbf{w}_1(t) \\ \mathbf{w}_{\tilde{\mathbf{p}}}(t) \\ \mathbf{w}_2(t) \end{bmatrix} \quad (6.15)$$

with initial conditions:

$$\mathbf{x}_1(0) = \mathbf{x}_{1,0} + \mathbf{w}_1(0) \quad (6.16)$$

$$\tilde{\mathbf{p}}(0) = \bar{\mathbf{p}} + \mathbf{w}_{\tilde{\mathbf{p}}}(0) \quad (6.17)$$

where $\bar{\mathbf{p}}$ is the vector collecting the nominal values of the parameters to be recalibrated (i.e., the ones obtained during the offline model calibration procedure); $\mathbf{w}_{\tilde{\mathbf{p}}}(t)$ is a random process with mean $\mathbf{0}$ representing the uncertainty on the parameters $\tilde{\mathbf{p}}$. Based on the augmented state (6.14), the estimation criterion and the sequential estimator equations are the same as in section 6.2.1 and Appendix A. It is important to notice that, with the above state augmentation procedure, the number of parameters that can be calibrated online cannot exceed the number of online measurements. The model can be updated with the same frequency of online measurements or at a user-defined time step Δt_k .

4. Step #4: online model-based DS maintenance. Feasibility analysis is used to periodically update the design space prediction with the up-to-date model returned by the state estimator. The design space maintenance frequency Δt_m must be compatible

with the computational time required to solve the optimization problem (6.10). In this regard, surrogate-based feasibility analysis provides a competitive advantage in terms of the total computational time (typically few seconds) required when complex models are used. The condition $\Delta t_m \geq \Delta t_k$ must hold true in any event.

In order to quantify the accuracy of the model-based DS identification with respect to the actual DS (which is however unknown in practical applications), the three metrics proposed by Wang and Ierapetritou (2017b) are exploited, namely:

- (i) $CF(\%)$ = percentage of the design space of the process that has been correctly estimated by the model;
- (ii) $CIF(\%)$ = percentage of the infeasible region (i.e. portion of the input domain that does not belong to the design space) that has been correctly estimated by the model;
- (iii) $NC(\%)$ = percentage of design space predicted by the model that does not belong to the actual DS of the process.

The closer $CF(\%)$, $CIF(\%)$ are to 100 and $NC(\%)$ to 0, the better the accuracy in the design space identification. Details on the derivation of these metrics can be found in the cited reference.

The feasibility analysis step can be solved using the classical formulation (i.e., by using deterministic NLP techniques) or surrogate-based approaches (e.g., by using kriging, or radial-basis function, or other surrogates), based on the computational requirements for the online implementation. Notice that the methodology can be first tested offline once a set of measurements are collected in the plant and then implemented online once its robustness is assessed.

The above methodology was implemented in MATLAB[®] and all simulations discussed in this study were performed on an Intel[®] Core™ I7-5600U CPU@2.60GHz processor with 16.0 GB RAM. A sequential quadratic programming (SQP) optimizer was used in all simulations.

6.5 Case studies

In the following, two simulated case studies are discussed to test the proposed methodology. The following notation will be used:

1. the “process” denotes the set of equations used to generate synthetic measurements to be fed to the state estimator (i.e., the set of equations representing the actual plant behavior);
2. the “model” denotes the set of equations used to simulate the plant behavior.

Process-model mismatch was enforced as a parametric mismatch in Case study #1 and as a structural mismatch in Case study #2.

6.5.1 Case study #1: roller compaction of microcrystalline cellulose

The first case study involves the roll compaction of microcrystalline cellulose grade Avicel PH102 and is based on the model proposed by Hsu *et al.* (2010a), which combines Johanson's rolling theory (Johanson, 1965) with a dynamic material balance. The set of model equations is as follows:

$$\frac{d}{dt} \left(\frac{h_0}{R} \right) = \frac{\omega [\rho_{in} \cos(\theta_{in}) (1 + \frac{h_0}{R} - \cos(\theta_{in})) (\frac{u_{in}}{\omega R}) - \rho_{exit} (\frac{h_0}{R})]}{\int_0^{\theta_{in}} \rho(\theta) \cos(\theta) d\theta} \quad (6.18)$$

$$P_h = \frac{W}{A} \frac{\sigma_{exit} R}{(1 + \sin(\delta))} \int_0^\alpha \left[\frac{\frac{h_0}{R}}{(1 + \frac{h_0}{R} - \cos(\theta)) \cos(\theta)} \right]^K \cos(\theta) d\theta \quad (6.19)$$

$$\sigma_{exit} = C_1 \rho_{exit}^K \quad (6.20)$$

The product critical quality attributes are the ribbon density at the outlet ρ_{exit} and the ribbon thickness h_0 . The roll pressure P_h , the roll speed ω and inlet feed speed u_{in} are process parameters that affect the critical quality attributes. However, as proved by Hsu *et al.* (2010a), the roll pressure has a significantly stronger influence on the product quality than the other two factors. Therefore, the roll pressure was considered as a critical process parameter and the other two variables as control variables that can be manipulated in order to guarantee the product specifications. In this regard, Hsu *et al.* (2010a) proposed to simulate the actuator dynamics for the two control variables with a first-order dynamics:

$$\tau_\omega \frac{d\omega}{dt} + \omega = \omega_d \quad (6.21)$$

$$\tau_u \frac{du_{in}}{dt} + u_{in} = u_d \quad (6.22)$$

where τ_ω and τ_u are the time constants, and ω_d and u_d are the set-points for the roll speed and feed speed respectively. The inlet density of the bulk material ρ_{in} is a raw material property that directly affects the critical quality attributes. The meaning and nominal values of all the variables and parameters of the model are reported in Table 6.1. An approximated analytical expression for the denominator of Eq. (19) can be found in the cited reference.

Table 6.1. Case study #1: list of model parameters for roll compaction process (taken from Hsu et al., 2010a).

Parameter	Definition	Units	Value
δ	Effective angle of friction	[rad]	0.7068
R	Roll radius	[m]	0.125
A	Compact surface area	[m ²]	0.01
W	Roll width	[m]	0.05
C_1	Compression parameter	[Pa/(kg m ³) ^{4.97}]	7.5×10^{-8}
K	Compression parameter #2	[-]	4.97
α	Nip angle	[rad]	0.173
τ_ω	Time constant for ω	[s]	6
τ_u	Time constant for u_{in}	[s]	6
ω_d	Set-point for ω	[rad/s]	0.5236
u_d	Set-point for u_{in}	[m/s]	3.27×10^{-2}

The design space for this system can be defined as the multidimensional combinations of roll pressure (critical process parameter) and inlet bulk density (raw material property) that guarantee to meet the desired constraints on the product critical quality attributes (ribbon density and ribbon thickness). With reference to the dynamic model (18)-(22), the following issues deserve attention.

- It has been pointed out (Hsu et al., 2010a; Hsu et al., 2010b) that the estimate of the inlet angle θ_{in} is affected by a high degree of uncertainty. Rogers and Ierapetritou (2015) showed that θ_{in} has a strong effect on the critical quality attributes and consequently on the prediction of the design space. For instance, a $\pm 5^\circ$ (0.09 rad) error on the inlet angle can result in a $\pm 50\%$ error in the inlet mass flow rate, thus causing a significant process-model mismatch on the critical quality attributes.
- The ribbon density ρ_{exit} can be measured online with near-infrared (NIR) spectroscopy at a minimum time frequency of 0.1s (Hsu et al., 2010b).
- Since the roll compactor is meant to operate at steady state (continuous manufacturing), the dynamic model (6.18)-(6.22) describes transitions (e.g. start-up, set-point changes) between different steady-state operations of the unit. Fast process dynamics coupled with model computational complexity act as severe challenges towards online model applications (Hsu et al., 2010b).

In view of the above, the “process” is denoted as the set of equations (6.18)-(6.22) where nominal values for θ_{in} and initial state vector are used. On the other hand, the “model” is given by the set of equations (6.18)-(6.22) with a wrong initial estimate of the system state and a wrong initial estimate of θ_{in} (parametric mismatch). Namely, a scenario is considered where the value of θ_{in} obtained during the offline validation procedure is wrongly estimated by $\pm 30\%$, and the initial steady-state values of state variables ρ_{exit} and h_0 are wrongly estimated (by +5% with respect to the nominal values). The error in the initial estimate of θ_{in} is consistent with the uncertainty on the value of θ_{in} that is reported in the literature (Rogers and Ierapetritou,

2015). Note that this scenario represents a conservative worst-case one for the estimation of θ_{in} .

6.5.2 Case study #2: penicillin fermentation

In this case study, the fermentation of penicillin in a fed-batch bioreactor is considered. The simplified mechanistic model used by Riascos and Pinto (2004) complemented with the CO₂ production equation proposed by Birol *et al.*, (2002) was used as the reference “model”. The model consists of 5 differential equations on the following differential states: volume $V(t)$, biomass concentration $B(t)$, substrate concentration $S(t)$, penicillin concentration $P(t)$, moles of CO₂ produced per liter of broth [CO₂]. The model equations are:

$$\frac{dB(t)}{dt} = \mu(t)B(t) - \left(\frac{B(t)}{S_F V(t)}\right)U \quad (6.23)$$

$$\frac{dP(t)}{dt} = \rho(t)B(t) - K_{deg}P(t) - \left(\frac{P(t)}{S_F V(t)}\right)U \quad (6.24)$$

$$\frac{dS(t)}{dt} = -\mu(t)\left(\frac{B(t)}{Y_{B/S}}\right) - \rho(t)\left(\frac{B(t)}{Y_{P/S}}\right) - \left(\frac{m_S S(t)}{K_m + S(t)}\right)B(t) + \frac{\left(1 - \frac{S(t)}{S_F}\right)U}{V(t)} \quad (6.25)$$

$$\frac{dV(t)}{dt} = \frac{U}{S_F} \quad (6.26)$$

$$\frac{d[CO_2](t)}{dt} = \alpha_1 \frac{dB(t)}{dt} + \alpha_2 B(t) + \alpha_3 \quad (6.27)$$

with:

$$\mu(t) = \mu_{max} \left(\frac{S(t)}{K_X B(t) + S(t)} \right) \quad (6.28)$$

$$\rho(t) = \rho_{max} \left(\frac{S(t)}{K_P + S(t) \left(1 + \frac{S(t)}{K_{in}} \right)} \right). \quad (6.29)$$

The nominal values of the model parameters and the initial conditions for the state variables are reported in Table 6.2.

Table 6.2. Case study #2: list of model parameters for the penicillin fermentation model (data from Riascos and Pinto (2004) and Birol *et al.* (2002)).

Parameter	Definition	Units	Value
μ_{\max}	Maximum specific biomass growth rate	$[\text{h}^{-1}]$	0.11
ρ_{\max}	Maximum specific production rate	$[\text{g}_P/(\text{g}_B \text{h})]$	0.0055
K_X	Saturation parameter for biomass growth	$[\text{g}_S/\text{g}_B]$	0.006
K_P	Saturation parameter for production	$[\text{g}_S/\text{L}]$	0.0001
K_{in}	Inhibition parameter for production	$[\text{g}_S/\text{L}]$	0.1
K_{deg}	Product degradation rate	$[\text{h}^{-1}]$	0.01
K_m	Saturation parameter for maintenance consumption	$[\text{g}_S/\text{L}]$	0.0001
m_S	Maintenance consumption rate	$[\text{g}_S/(\text{g}_B \text{h})]$	0.029
$Y_{B/S}$	Yield factor for substrate to biomass	$[\text{g}_B/\text{g}_S]$	0.47
$Y_{P/S}$	Yield factor for substrate to product	$[\text{g}_P/\text{g}_S]$	1.2
α_1	Constant relating CO_2 to growth	$[\text{mmol CO}_2/\text{g}_B]$	0.143
α_2	Constant relating CO_2 to maintenance energy	$[\text{mmol CO}_2/(\text{g}_B \text{h})]$	4×10^{-7}
α_3	Constant relating CO_2 to penicillin production	$[\text{mmol CO}_2/(\text{Lh})]$	10^{-4}

By plugging Eqs. (6.28)-(6.29) into Eqs (6.23)-(6.27), the original DAE system can be converted into a pure ODE system whose state vector is given by $\mathbf{x}_1(t) = [B(t); P(t); S(t); V(t); [\text{CO}_2](t)]$. Note that in this scenario, the estimator of section 6.2.1 simplifies to the standard EKF observer.

The model described by Eqs. (6.23)-(6.27) is based on some assumptions and simplifications that do not hold true in industrial penicillin fermenters (structural mismatch). These assumptions and simplifications can be summarized as follows.

- The effects of environmental factors, such as pH and temperature in the bioreactor, are not modeled. These variables actually play a key role in penicillin fermentation, having a direct impact primarily on the biomass growth rate and penicillin production rate.
- The model assumes saturation conditions for the oxygen dissolved in the broth. However, in industrial applications, oxygen limitation is a key issue that directly affects the biomass growth rate and the penicillin production rate.
- In industrial fermenters, a significant amount of the broth inside the reactor (up to ~10–20 % in one week of operation, as observed by Birol *et al.*, 2002) is lost due to evaporation. In addition, the volume inside the reactor is affected by addition of acid/base solutions to control the broth pH. The model does not account for any of these phenomena.

A detailed set of equations that take into account all the phenomena described above in penicillin fermentation is the one proposed by Birol *et al* (2002), based on a previous work of Bajpai and Reuss (1980). A description of these equations is reported in the supplementary material. This set of equations was used as the reference “process”, from which synthetic measurements to be fed to the state estimator were generated. Summarizing:

1. the “process” is the set of equations proposed by Birol *et al* (2002) that can be found in the supplementary material. These equations account for all the relevant physical phenomena described above;
2. the “model” is the set of equations (6.23)-(6.29), in which the phenomena described above are not explicitly considered in the model structure.

6.6 Results and discussion

6.6.1 Results and discussion for Case study #1

6.6.1.1 Simulation set-up

By plugging Eq. (6.20) into Eq. (6.19), and solving Eq. (6.19) with respect to the ribbon density ρ_{exit} , the dynamic model of the granulator is in the general DAE form of Eq. (6.1), where the only differential state is given by the ribbon thickness $h_0(t)$, the only algebraic state is $\rho_{exit}(t)$, and the control vector is given by the inlet feed speed $\mathbf{u}(t) = [\omega(t); u_{in}(t)]$. As remarked in section 6.5.1, the inlet angle $\theta_{in}(t)$ is the uncertain parameter that needs online recalibration. During normal operation of the roller compactor, measurements of the ribbon density can be obtained online through NIR spectroscopy. To the purpose of testing the proposed methodology, synthetic process measurements of $\rho_{exit}(t_k)$ ($k = 1, \dots, K$) were generated at a frequency of 1 Hz and a nominal value $\bar{\theta}_{in} = 0.4$ rad. White Gaussian noise $(t_k) \sim (0, 0.9^2)$, corresponding to a standard deviation of 0.1% of the nominal value $\bar{\rho}_{exit} = 900$ kg/m³, was added to the ρ_{exit} measurements.

6.6.1.2. Offline model-based DS identification (step #1)

Wrong identification of θ_{in} and of the initial state of the system causes a process-model mismatch that strongly impacts on the prediction of the product critical quality attributes, as shown in Fig. 6.2.

In order to identify offline the model-based DS, the following upper (superscript up) and lower (superscript lo) specifications were assigned to the product critical quality attributes (ρ_{exit} and h_0) for a new steady-state reached at time t_{SS} :

$$\rho_{exit}(t_{SS}) - \rho_{exit}^{up} \leq 0 \quad (6.30)$$

$$\rho_{exit}^{lo} - \rho_{exit}(t_{SS}) \leq 0 \quad (6.31)$$

$$h_0^{lo} - h_0(t_{SS}) \leq 0 \quad (6.32)$$

$$h_0(t_{SS}) - h_0^{up} \leq 0. \quad (6.33)$$

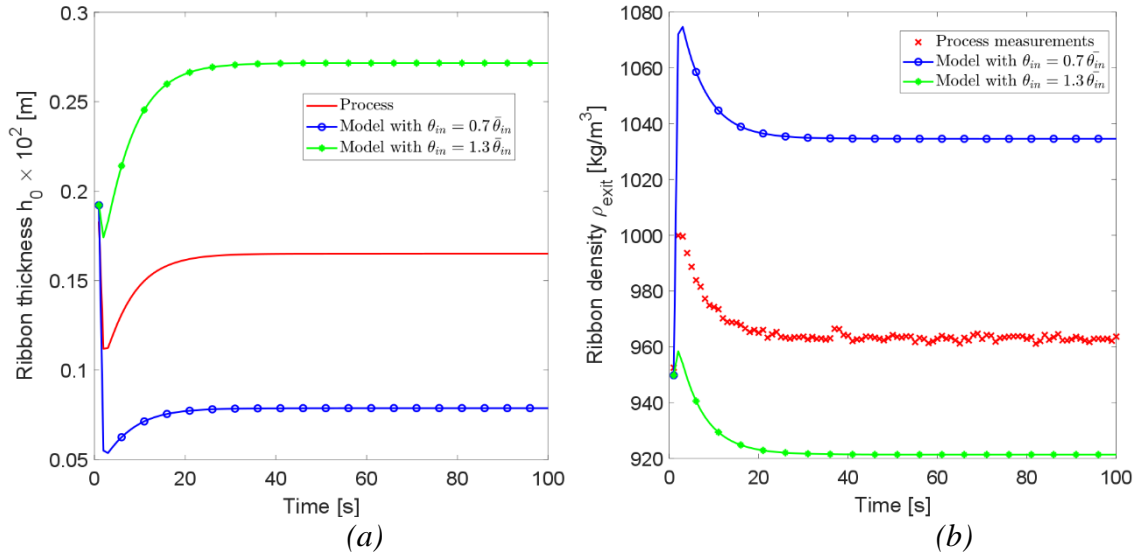


Figure 6.2. Case study #1: effect of process-model mismatch on (a) prediction of the ribbon thickness and (b) prediction of the ribbon density.

According to the feasibility analysis formulation (6.10), the model-based DS can be formulated as follows:

$$\chi(\mathbf{d}) = \max_{P_h, \rho_{in}} \Psi(\mathbf{d}, P_h, \rho_{in}) \quad (6.34)$$

$$\text{s.t.} \quad \Psi(\mathbf{d}, P_h, \rho_{in}) = \max \{ \rho_{exit}(t_{SS}) - \rho_{exit}^{up} \leq 0;$$

$$\rho_{exit}^{lo} - \rho_{exit}(t_{SS}) \leq 0;$$

$$h_0^{lo} - h_0(t_{SS}) \leq 0;$$

$$h_0(t_{SS}) - h_0^{up} \leq 0 \}; \quad (6.35)$$

$$\frac{d}{dt} \left(\frac{h_0}{R} \right) - \frac{\omega [\rho_{in} \cos(\theta_{in}) (1 + \frac{h_0}{R} \cos(\theta_{in})) (\frac{u_{in}}{\omega R}) - \rho_{exit}(\frac{h_0}{R})]}{\int_0^{\theta_{in}} \rho(\theta) \cos(\theta) d\theta} = 0 \quad (6.36)$$

$$P_h - \frac{W}{A} \frac{\sigma_{exit} R}{(1 + \sin(\delta))} \int_0^\alpha \left[\frac{\frac{h_0}{R}}{(1 + \frac{h_0}{R} \cos(\theta)) \cos(\theta)} \right]^K \cos(\theta) d\theta = 0 \quad (6.37)$$

$$\sigma_{exit} - C_1 \rho_{exit}^K = 0 \quad (6.38)$$

$$h_0(0) = h_0^{in} \quad (6.39)$$

$$\rho_{in} \in [\rho_{in}^{\min}; \rho_{in}^{\max}] \quad (6.40)$$

$$P_h \in [P_h^{min}; P_h^{max}]. \quad (6.41)$$

All the numerical values (initial conditions; product quality specifications; input variables domain) are reported in Table 6.3. In problem formulation (6.34)-(6.41), we assumed time-invariant input variables P_h and ρ_{in} ; the control variables $\omega(t)$ and $u_{in}(t)$ do not appear because their trajectories are explicit functions of (P_h, ρ_{in}) . Notice that, although we steady-state quality constraints were assigned and time-invariant input variables were considered, formulation (34)-(41) is inherently dynamic since it requires integration until the new steady-state is reached.

Table 6.3. Case study #1: values of the parameters appearing in problem formulation (34)-(41).

Parameter	Units	Value
h_0^{in}	[cm]	0.183
h_0^{up}	[cm]	0.190
h_0^{lo}	[cm]	0.170
ρ_{exit}^{up}	[kg/m ³]	950
ρ_{exit}^{lo}	[kg/m ³]	850
ρ_{in}^{max}	[kg/m ³]	550
ρ_{in}^{min}	[kg/m ³]	150
P_h^{max}	[MPa]	1.1
P_h^{min}	[MPa]	0.9

The model was updated every $\Delta t_k = 1$ s, whereas the design space was updated every $\Delta t_m = 5$ s. Accordingly, a solution to the optimization problem (6.34)-(6.41) must be obtained within 5 s. However, the optimization problem in its classical formulation can only be solved in a time frame 10 times larger than required, and classical feasibility analysis would therefore make the proposed methodology ineffective. As discussed in section 6.2.2, the approach proposed by Rogers and Ierapetritou (2015) was implemented, which builds a kriging surrogate as a cheap approximation of model (6.18)-(6.22).

The ability of the surrogate approach to give a reliable representation of the actual DS with respect to the classical formulation approach was first tested offline. The boundary of the model-based DS obtained with the classical and kriging-based feasibility analysis are shown in Fig. 6.3. It can be noticed that the boundary of the model-based DS obtained with the kriging approach (obtained after 10 iterations) are in perfect agreement with the ones obtained with the classical formulation. However, the computational time with the kriging surrogate is 30 times faster than with the original model, thus making online implementation of the proposed methodology possible.

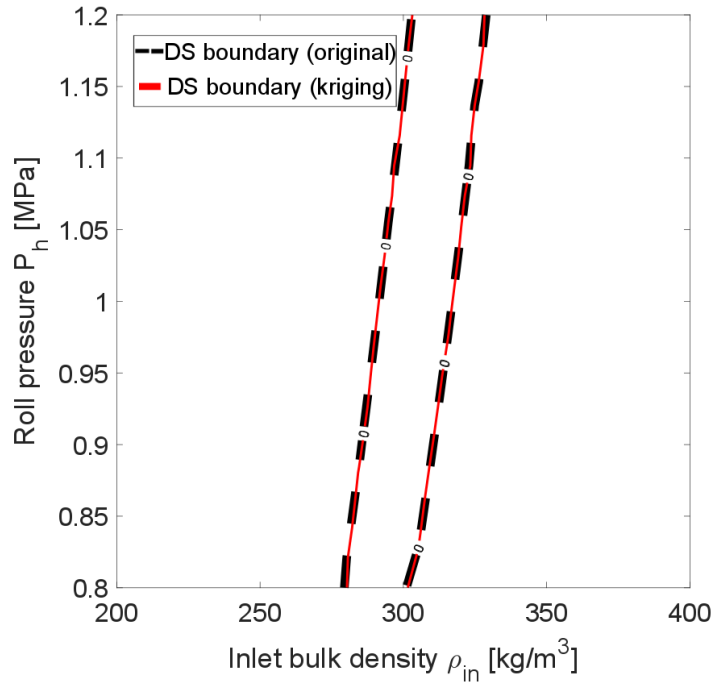


Figure 6.3. Case study #1: boundary of the model-based DS of the process obtained with the classical formulation (42) and with the kriging surrogate approach.

Fig. 6.3 also shows that there is a direct relationship between the effect of the roll pressure and inlet bulk density with respect to the product critical quality attributes. The higher the roll pressure, the greater the inlet bulk density the process can tolerate to guarantee the final product specifications. This is in agreement with the experimental work of Hsu *et al.* (2010a).

6.6.1.3 State estimation for online model adaptation (steps #2 and #3)

The case study represents a situation where the model parameter that needs online re-calibration (i.e. θ_{in}) is known *a priori*. Implementation of step #2 of the proposed methodology is therefore straightforward.

With respect to step #3, the following conditions were set for the state estimator:

$$\mathbf{P}_{11}(0) = \begin{pmatrix} 0.01^2 & 0 \\ 0 & 1^2 \end{pmatrix}; \quad \mathbf{Q} = \begin{pmatrix} 0.01^2 & 0 \\ 0 & 1^2 \end{pmatrix}; \quad \mathbf{R} = 0.9^2; \quad \Delta t_k = 1 \text{ s}. \quad (6.42)$$

The value of the initial error covariance, model error covariance and measurement covariance were treated as tuning parameters. An alternative systematic procedure that can be used to set the values of these matrices is the one proposed by Schneider and Geoargakis (2012). Note that we used a model error for ρ_{exit} ($= 1^2$) slightly greater than its measurement error ($= 0.9^2$), in order to reflect the greater confidence that we have on the measured values of ρ_{exit} rather than their model prediction.

The results obtained with the state estimator for the two scenarios considered ($\pm 30\%$ estimation error on θ_{in} , $+5\%$ error on initial state) are shown in Fig. 6.4. Note that this figure refers to a situation where P_h and ρ_{in} are set at their nominal values and simply illustrates the estimator performance.

It can be noticed that the state estimator quickly removes the effect of process-model mismatch from both ρ_{exit} (measured variable) and h_0 (inferred state variable, not available to measure). Additionally, the alignment of θ_{in} to its actual value occurs quickly.

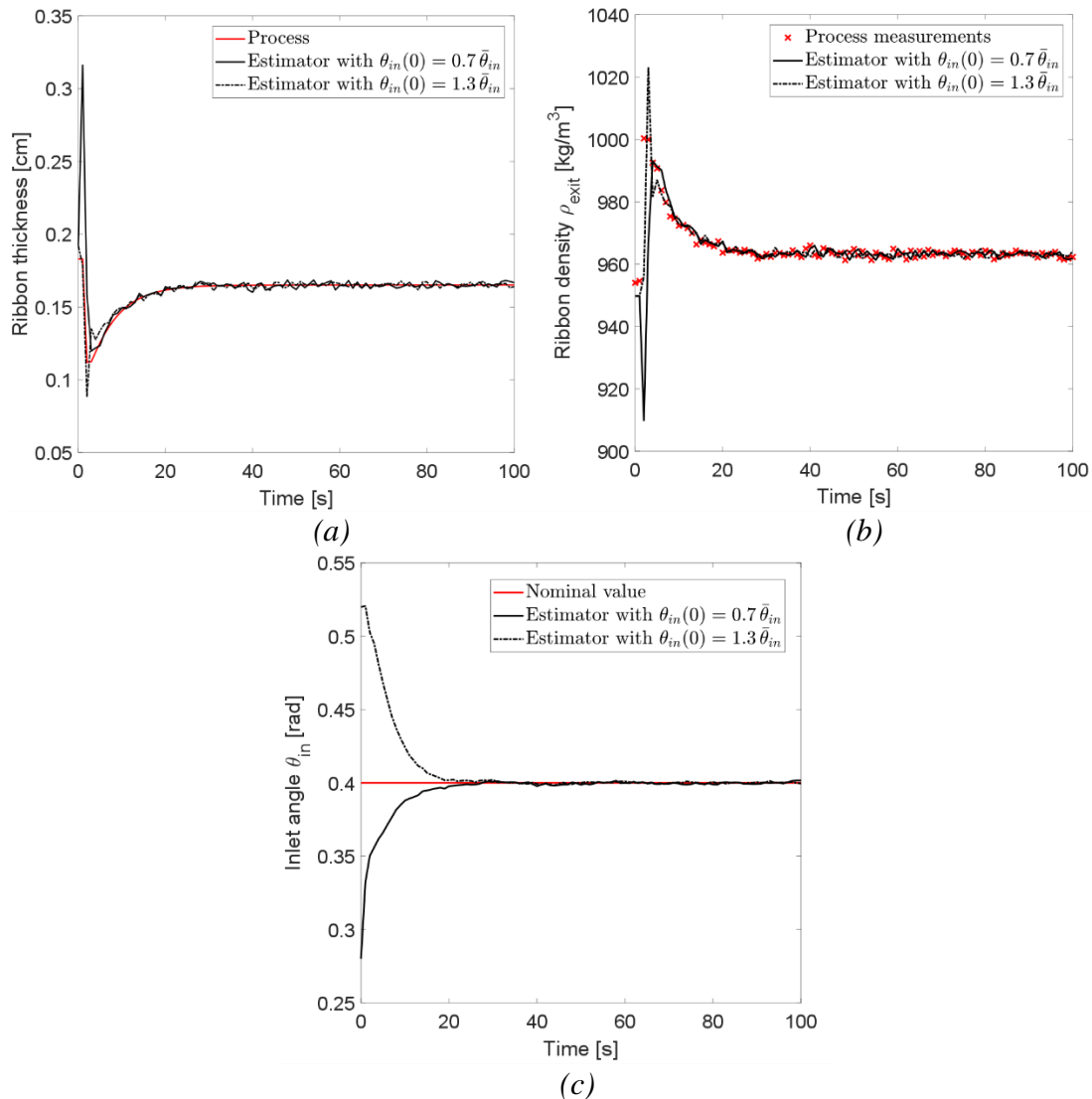


Figure 6.4. Case study #1: state estimator predictions, process states and measurements. (a) Ribbon thickness. (b) Ribbon density. (c) Inlet angle.

6.6.1.4 Online model-based DS maintenance (step #4)

The online model-based DS maintenance can be formulated at each time-step t_m using the kriging-based feasibility analysis as:

$$\chi(\mathbf{d}) = \min_{P_h, \rho_{in}} U_{p,m}(P_h, \rho_{in}, \theta_{in}(t_m)) \quad (6.51)$$

$$\text{s.t.} \quad \rho_{in} \in [\rho_{in}^{\min}; \rho_{in}^{\max}] \quad (6.52)$$

$$P_h \in [P_h^{\min}; P_h^{\max}]. \quad (6.53)$$

where $U_{pm}(\cdot)$ is the surrogate (kriging) feasibility function with initial conditions for the state variables given by the up-to date state vector $\mathbf{x}(t_m)$. The design space identification with formulation (6.51)-(6.53) is accurate (as proved in Fig. 6.3) and requires a short computational time (1.67 s in less than 10 iterations).

The offline model-based DS that would be obtained with the -30% and $+30\%$ error on the estimate of θ_{in} are compared with the actual DS of the process in Fig. 6.5. The figure shows that a wrong estimate of θ_{in} strongly affects the model-based DS identification. Namely, θ_{in} strongly influences the prediction of the variability of the inlet bulk density that the process can tolerate to guarantee the final product quality.

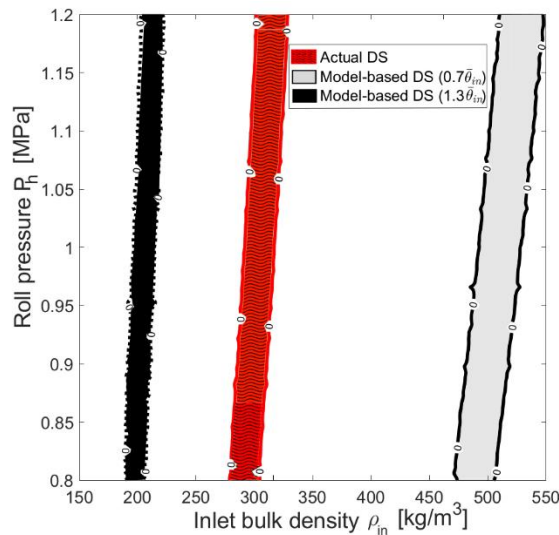
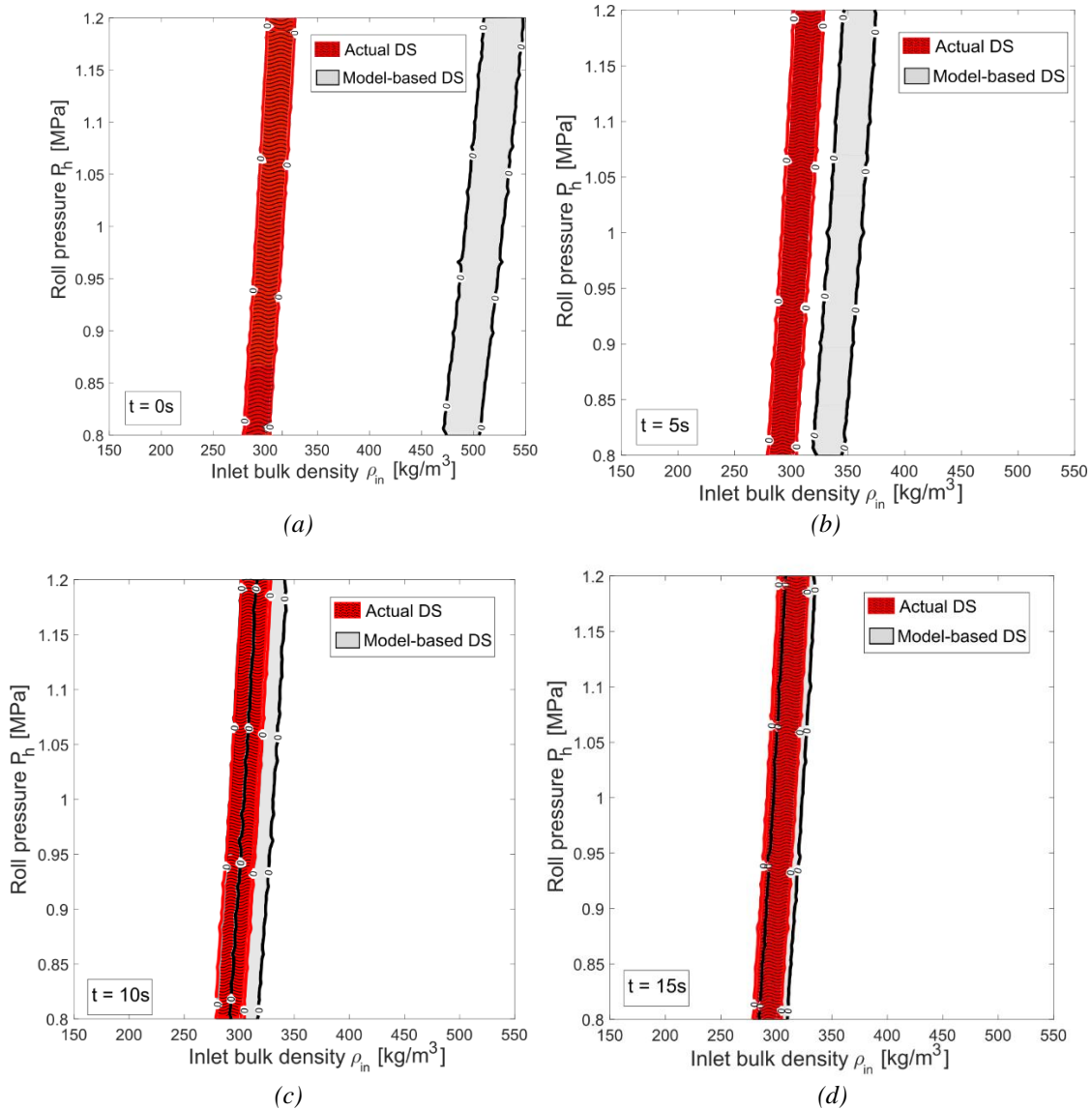
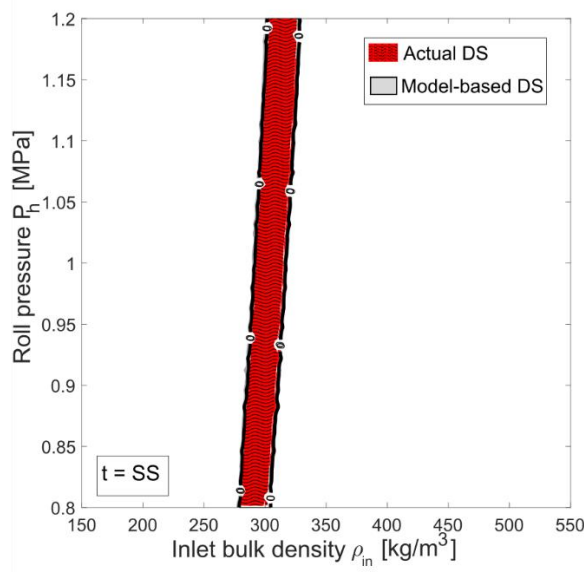


Figure 6.5. Case study #1: comparison between the actual DS and the offline model-based DS with $+30\%$ and -30% estimate on the bulk inlet angle θ_{in} and 5% error on the initial state.

The evolution of the model-based DS that can be obtained with the proposed online DS maintenance procedure is shown in Fig. 6.6. The figure refers to the scenario with a -30% wrong estimate on θ_{in} and $+5\%$ error on the initial state of the system. The results for the case with $+30\%$ error on θ_{in} are analogous and can be found in the supplementary material. The online model-based DS is shown after 0 (i.e., initial offline model-based DS), 5, 10, 15 and 35 s (i.e., the time the process takes to reach its new steady-state) seconds. Consistently with Fig. 6.6, it can be noticed that it takes less than 20 s for the model-based DS to be in complete

agreement with the actual DS. For this lab-scale roller compactor, the methodology is able to restore the correct model-based DS identification in less than half of the time the process needs to reach its new steady state. The values of the metrics $CF(\%)$, $CIF(\%)$ and $NC(\%)$ are 99.95, 99.97 and 0.03 respectively, thus confirming the perfect agreement between the model-based DS and the actual DS. The total computational time required for the sequential state estimation step is 0.44 s, whereas the kriging-based feasibility step requires 1.67 s.





(e)

Figure 6. Case study #1: up-to date model-based DS obtained after (a) 0 s, (b) 5 s, (c) 10 s, (d) 15 s and (e) 35 s. ($= t_{SS}$) for a -30% deviation in the inlet angle estimate and $+5\%$ error in the initial state

6.6.2 Results and discussion for Case study #2

6.6.2.1 Simulation set-up

A scenario was considered where a first principles model developed and calibrated at the laboratory scale (Eqs (6.23- 29)) fails to account for some physical phenomena that occur at a different scale and/or that were not observed during the offline model validation procedure (type 2 uncertainty as described in section 6.4). In industrial applications, two variables that can be measured online with little effort are the moles of CO_2 in the outlet gas per liter of broth the volume (y_1), and the volume $V(t)$ (y_2). Therefore, synthetic process measurements were generated for these two variables with a measurement interval of $\Delta t_k = 10$ min. These measurements were corrupted with white Gaussian noise $\mathbf{v}_k \sim N(\mathbf{0}, \mathbf{R})$ with:

$$\mathbf{R} = \begin{pmatrix} 0.1^2 & 0 \\ 0 & 0.5^2 \end{pmatrix}. \quad (6.43)$$

The concentration of penicillin $P(t)$ was considered as a product critical quality attribute, and the concentration of substrate s_F in the feed as a raw material property. The critical process parameter that was taken into account was the inlet feed flowrate U . As in the previous case study, both variables were considered as time-invariant, but the analysis can be easily extended to time-varying inputs.

A single end-point specification was set on the penicillin concentration for a batch duration of $t_{end} = 200$ h:

$$P^{lo} - P(t_{end}) \leq 0 . \quad (6.44)$$

6.6.2.2 Offline model-based DS identification (step #1)

The model-based DS of this process can be described in terms of a feasibility analysis problem as:

$$\chi(\mathbf{d}) = \max_{U, s_F} \Psi(\mathbf{d}, U, s_F) \quad (6.45)$$

$$\text{s.t.} \quad \Psi(\mathbf{d}, U, s_F) = P^{lo} - P(t_{end}) \leq 0 ;$$

$$\text{eqs. (6.23) - (6.29)}$$

$$U \in [U^{lo}; U^{up}] \quad (6.46)$$

$$s_F \in [s_F^{lo}; s_F^{up}] . \quad (6.47)$$

The values of the parameters that appear in the formulation (6.45) are reported in Table 6.4.

Table 6.4. Case study #2: values of the parameters appearing in the formulation (6.45).

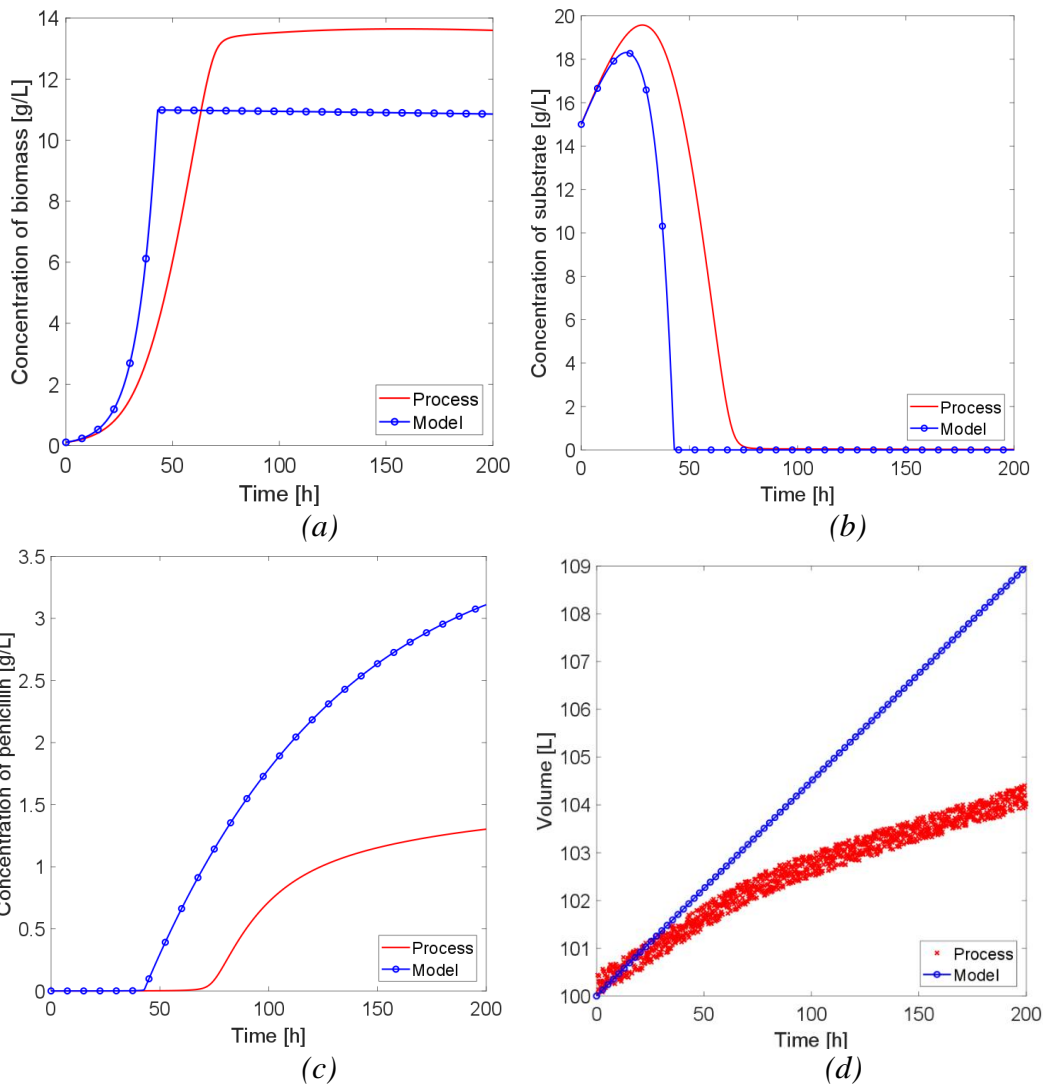
Parameter	Units	Value
P^{lo}	[g _S /L]	1.25
U^{lo}	[L/h]	5
U^{up}	[L/h]	80
s_F^{lo}	[g _S /L]	300
s_F^{up}	[g _S /L]	900

Fig. 6.7 shows that the un-modeled phenomena cause a strong process-model mismatch that impacts on the penicillin concentration (Fig 6.7c) as well as all other state variables. The plots of Fig. 6.7 refer to the initial conditions (assumed to be estimated correctly in the first principles model) and nominal values U and s_F reported in Table 6.4. Note that the model significantly overestimates the final concentration of penicillin of the actual process (Fig. 6.7c).

The solution of the feasibility problem (6.45)-(6.47) requires the integration of the model equations (6.23)-(6.27) embedded in the optimization algorithm, and it takes 2.6 min to get a solution. However, since the process dynamics is very slow and the design space is updated every 24 h, the computational time for the solution of the feasibility problem (6.45)-(6.47) is consistent with online implementation. Therefore, surrogate-based approaches are not needed for this case study.

Table 6.5. Case study #2: initial conditions for the state variables and nominal values of the substrate feed concentration and feed flowrate.).

Variable	Symbol	Units	Initial value	Nominal value
Biomass concentration	B	$[g_B/L]$	0.1	[-]
Substrate concentration	S	$[g_S/L]$	15	[-]
Penicillin concentration	P	$[g_P/L]$	0	[-]
Volume	V	[L]	100	[-]
CO ₂ concentration	$[CO_2]$	[mmol/L]	0.5	[-]
Feed flowrate	U	$[g_S/h]$	[-]	27
Substrate feed concentration	s_F	$[g_S/L]$	[-]	600



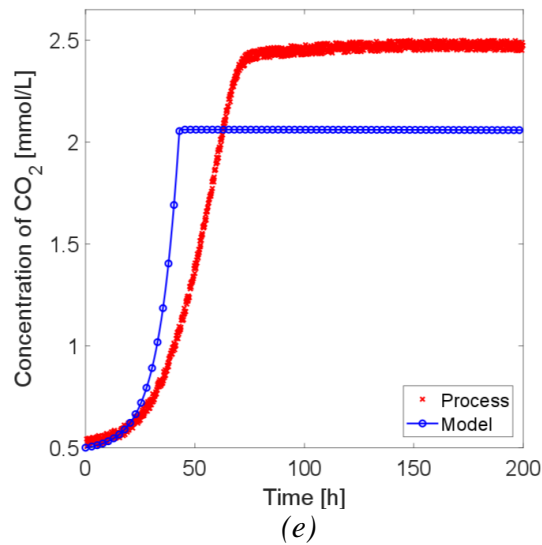


Figure 6.7. Case study #2: impact of process-model mismatch for the penicillin fermenter on the (a) concentration of biomass, (b) concentration of substrate, (c) concentration of penicillin, (d) volume and (e) concentration of CO₂.

6.6.2.3 Identification of uncertain model parameters (step #2)

According to step #2 of the proposed methodology, the parameters that need online recalibration should be identified in order to remove the process-model mismatch. Note that no more than two parameters can be recalibrated, since we assumed that only two online measurements are collected in the plant.

Qualitatively, it is expected that the environmental factors (pH, temperature) and the dissolved oxygen concentration will have a remarkable effect on the biomass growth rate as well as on the penicillin production rate (Birol *et al.*, 2002). Since these are not included explicitly in the model formulation, an adjustment parameter $\phi_p(t)$ was introduced in the expression of the penicillin production rate (i.e., the relevant critical quality attribute in this case) in order to account for this lack of the “model”:

$$\rho(t) = \rho_{max} \phi_p(t) \left(\frac{S(t)}{K_P + S(t) \left(1 + \frac{S(t)}{K_{in}} \right)} \right). \quad (6.48)$$

$\phi_p(t)$ is treated as an additional state that is used with the purpose of increasing the adaptive capability of the state estimator. Note that $\phi_p(t)$ was introduced without directly adjusting ρ_{max} in order to preserve its physical meaning. No claim is made on the physical meaning of $\phi_p(t)$. To account for broth evaporation, a second adjustment parameter $\phi_V(t)$ was introduced in the volume balance equation:

$$\frac{dV(t)}{dt} = \left(\frac{U}{S_F}\right) (1 - \phi_V(t)). \quad (6.49)$$

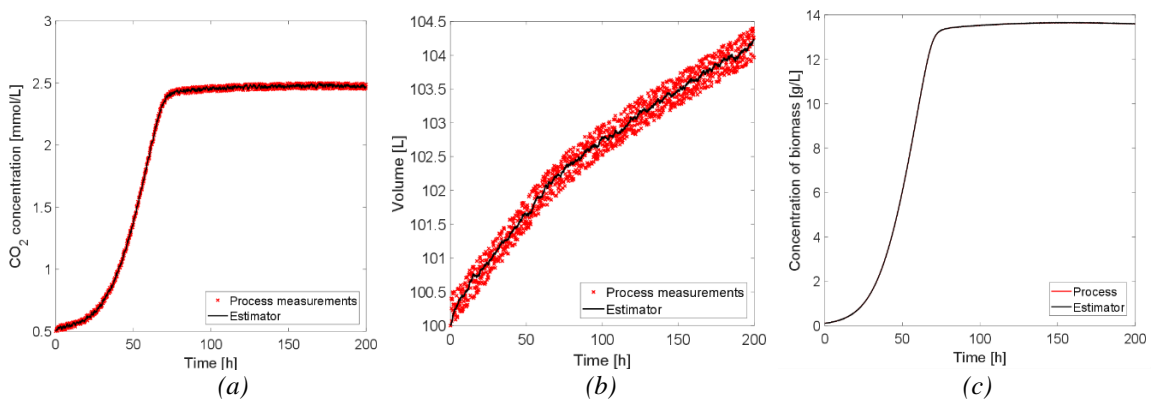
The same considerations made for $\phi_p(t)$ hold true for $\phi_V(t)$. The initial (i.e., nominal) values of ϕ_p and ϕ_V are 1 and 0, respectively.

6.6.2.4 State estimation for online model adaptation (step #3)

State estimation techniques and control applications have been extensively studied for this process (Birol *et al.*, 2002; Georgakis, 2013; Kager *et al.*, 2018). The initial model error covariance and the model error covariance were tuned with a trial-and-error approach in order to obtain the best performance of the estimator. The following values were chosen (1=B; 2=P; 3=S; 4=V, 5=[CO₂], 6= ϕ_p , 7= ϕ_V):

$$\mathbf{P}_{11}(0) = \begin{pmatrix} 0.1^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.1^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.1^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.5^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.1^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1^2 \end{pmatrix}; \mathbf{Q} = \begin{pmatrix} 0.1^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.1^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.1^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.1^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.5^2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.5^2 \end{pmatrix}. \quad (6.50)$$

The ability of the state estimator to remove the process-model mismatch is shown in Fig. 6.8 for all the state variables considered (U and S_F are set at their nominal values in these plots). Very good agreement between the estimator predictions and the simulated process profiles for the state variables is obtained.



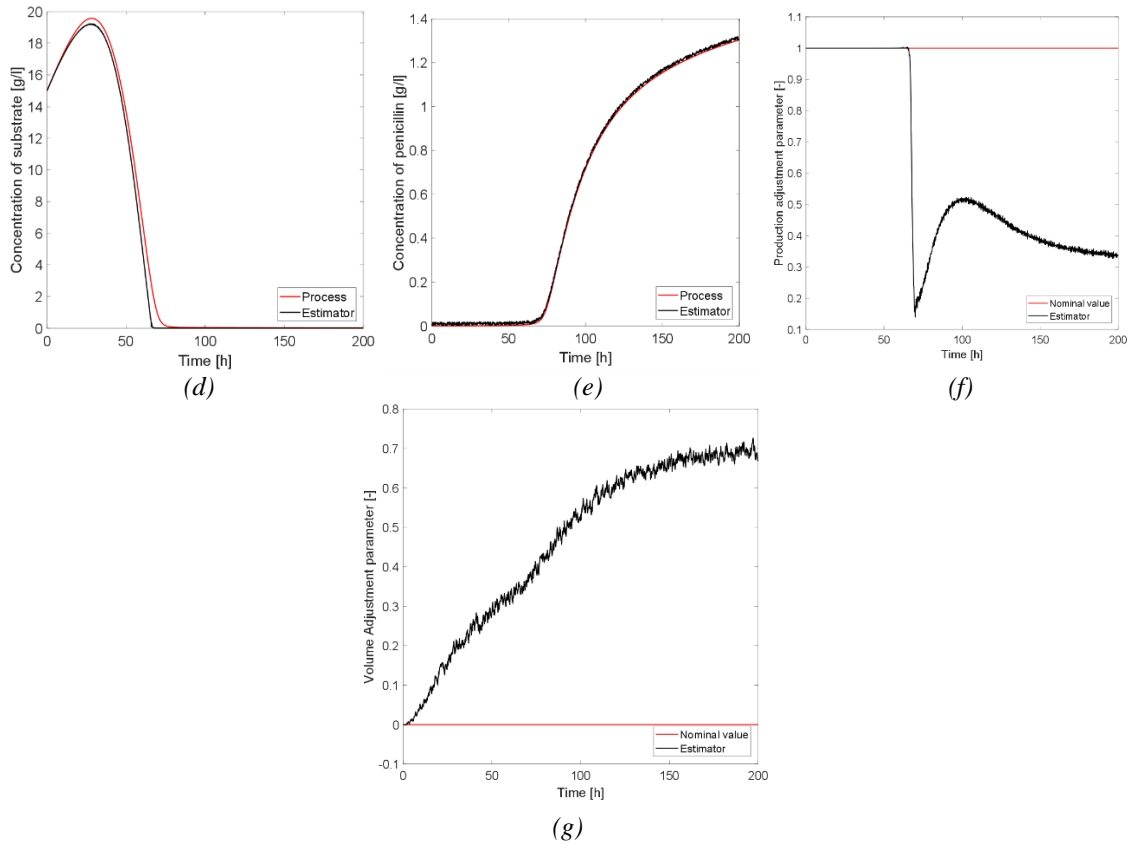


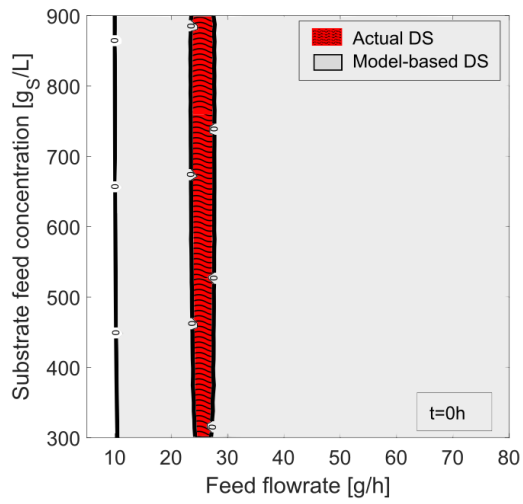
Figure 6.8. Case study #2: state estimator predictions and process states and measurements. (a) CO_2 concentration. (b) Volume. (c) Concentration of biomass. (d) Concentration of substrate. (e) Concentration of penicillin. (f) Growth adjustment parameter. (g) Volume adjustment parameter.

6.6.2.5 Online model-based DS maintenance (step #4)

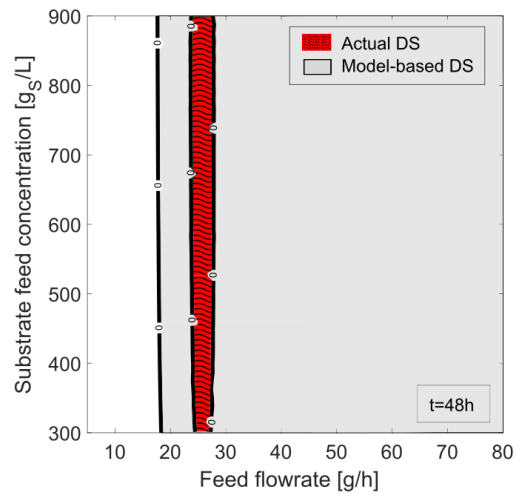
As discussed earlier, the update interval for the design space was set to $\Delta t_m = 24$ h. The online model-based DS can be obtained by solving the feasibility problem (6.45)-(6.47) with the initial conditions given by the up-to-date state vector $\mathbf{x}_1(t_m)$ returned by the state estimator at time t_m . The solution of the feasibility problem requires integration from time t_m till the end of the batch (200 h).

Fig. 6.9 shows the actual DS of the process and the evolution of the model-based DS after 0, 48, 72, 96 and 120 h.

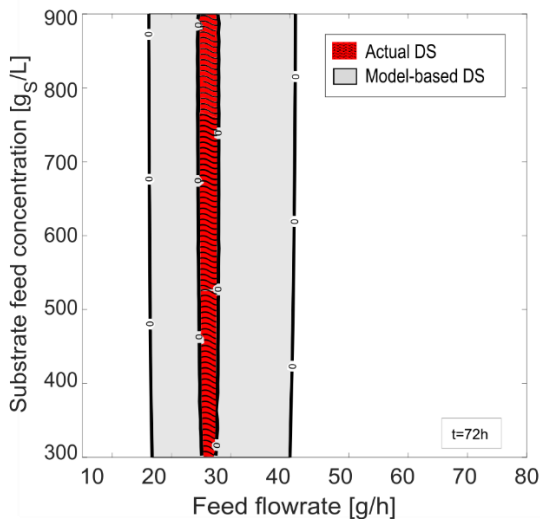
The values of the three metrics $CF(\%)$, $CIF(\%)$, $NC(\%)$ for all the time steps involved during the fermenter operation is reported in Table 6.6. Fig. 6.9a shows that the offline model-based DS (i.e., $t = 0$ h) is significantly wider than the actual DS of the process. This is in agreement with Fig. 6.7c, which shows that the original model significantly overestimates the concentration of penicillin at the end of the batch, and it is confirmed by the value of the metric $NC(\%)$, which is significantly greater than 0. The model-based DS keeps being significantly wider than the actual DS till the time step $t_m = 48$ h: this can be easily understood (Fig. 6.8c) with the fact that there is not a process-model mismatch impact on the concentration of penicillin during the first 48 h (penicillin is still not being produced).



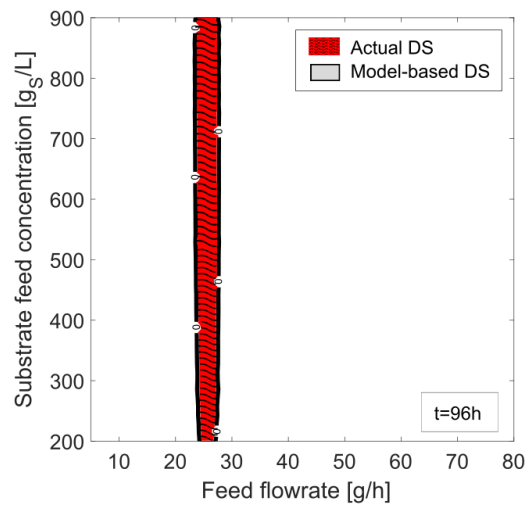
(a)



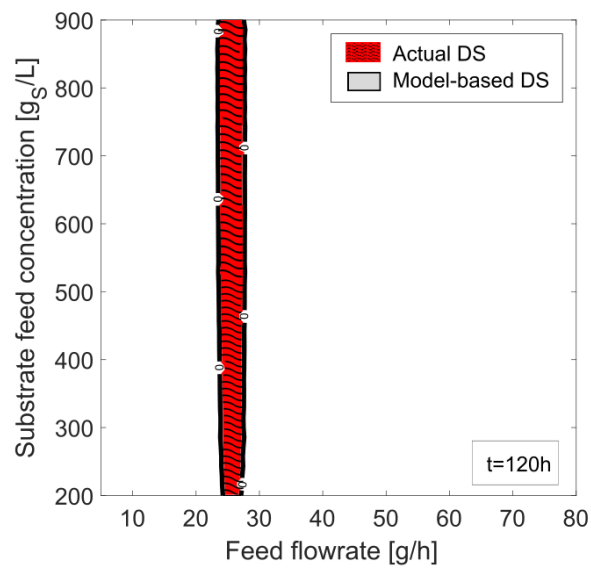
(b)



(c)



(d)



(e)

Consistently, the adjustment parameter $\phi_p(t)$ is not adjusted by the state estimator during this time frame. The accuracy of the model-based DS identification significantly increases after 72 h, since a substantial adjustment of the system state and $\phi_p(t)$ is performed (Fig. 6.8f).

Figure 6.9. Case study #2: actual DS of the process (red wavy region) and model-based DS (light gray regions) obtained after (a) 0 h., (b) 48 h., (c) 72 h., (d) 96 h and (e) 120 h.

Table 6.6. Case study #2: accuracy metrics for the DS maintenance tool.

Time[h]	CF (%)	CIF (%)	NC(%)
0	100.0	15.3	81.25
24	100.0	17.4	77.23
48	100.0	21.2	72.5
72	100.0	71.2	25.0
96	100.0	99.1	0.22
120	100.0	99.0	0.23
144	100.0	99.1	0.21
168	100.0	99.2	0.22
192	100.0	99.0	0.21

The DS identification accuracy stabilizes at satisfactory values ($CF(\%) \cong 100.0$, $CIF(\%) \cong 99.0$, $NC(\%) \cong 0.2$) after 96 h, and this accuracy is maintained throughout the entire fermenter operation. It is worth noticing the model-based DS obtained with the online maintenance tool is slightly conservative ($NC(\%) > 0$), which is desirable in terms of operational flexibility. Summarizing, after the first 50 h (where penicillin is still not being produced and therefore the adaptive ability of the state estimator is still not exploited in terms of critical quality attributes prediction), the design space maintenance methodology quickly recovers the correct DS. Note that the DS updating frequency can be increased after the first 50 h of operation to obtain an even prompter DS adaptation.

6.7. Conclusions

In this study, a methodology was proposed for the maintenance of the design space of a pharmaceutical product using online model adaptation. By combining the predictions of a first-principles model with measurements made available by plant sensors, the methodology exploits dynamic state estimation for online model adaptation and feasibility analysis to obtain an up-to-date representation of the design space during plant operation.

The effectiveness of the proposed methodology was proved on two simulated case studies, involving the roller compaction of microcrystalline cellulose and the fermentation of penicillin in a bioreactor. For both the case studies considered, the ability of the proposed approach to timely track the design space was proved.

Applications of this up-to date design space include, but are not limited to:

- (i) management of real-time raw materials variability;
- (ii) search of the optimum operation point within the design space by real-time dynamic optimization based on the up-to-date first-principles model;
- (iii) open-loop decision support within the boundary of the design space (e.g., manual intervention on critical process parameters' set-points);
- (iv) closed-loop decision support within the boundary of the design space (e.g., definition of optimal set-points or set-point trajectories).

Chapter 7

Model-based design of experiments in pharmaceutical freeze drying

DISCLOSURE STATEMENT

This work was done under a cooperative research and development agreement between GlaxoSmithKline Biologicals SA and Dipartimento di Ingegneria Industriale (DII) at University of Padova, in the context of the project entitled “Verso la digitalizzazione dell’industria farmaceutica: generazione di dati ad alto contenuto d’informazione per l’ottimizzazione di processi industriali di liofilizzazione. DIGI-LIO”. The study is co-funded by GSK and University of Padova.

All the numerical results reported in this Chapter have been normalized with non-disclosed numerical values due to confidentiality reasons. The normalization does not affect the discussion and conclusions that can be drawn from the reported results.

In this Chapter, model-based design of experiments (MBDoE) techniques are exploited to design optimal experiments for the correct identification of the most critical parameters in a pharmaceutical freeze drying first-principles model. A preliminary analysis is carried out in order to identify the structural consistency and the most influential parameters of the model, and MBDoE is subsequently used to design an optimal experiment based on the results of the preliminary analysis. The designed experiment is first simulated *in silico*, then performed on a real small-scale equipment, and finally the model parameters are identified based on the experimentation.

The Chapter is organized into six sections. In the first section, a brief description of the freeze drying process is given, with a specific focus on pharmaceutical applications. Section 7.1 collects a detailed description of the primary drying step of the process, which is the focus of this Chapter. In section 7.2, a brief overview of the mathematical formulation of MBDoE techniques is presented. Section 7.3 thoroughly describes the mathematical model used for the description of the primary drying stage. In section 7.4, a preliminary assessment of the performance of this model is presented. Section 7.5 collects the results obtained with the MBDoE activity. Finally, in section 7.6 the final remarks and the future areas of research are discussed.

7.1 Freeze drying: process description

Freeze drying (or *lyophilization*) is a process that is used to remove a solvent (typically water) from a frozen product, in order to (Pikal, 1980):

1. guarantee product stability both in terms of its physicochemical properties and microbiological quality attributes;
2. allow an easy reconstitution of the dried product by water addition (the dried product has a high surface area and can be easily re-hydrated).

It is extensively used in pharmaceutical manufacturing to recover the active pharmaceutical ingredient (API) and the excipients from an aqueous solution, in all those situations where other simpler drying technologies cannot be applied. In pharmaceutical applications, the frozen product is typically processed in vials, placed over shelves in a drying chamber. A schematic representation of a typical freeze-dryer, with a short description of all its components, is shown in Fig. 7.1.

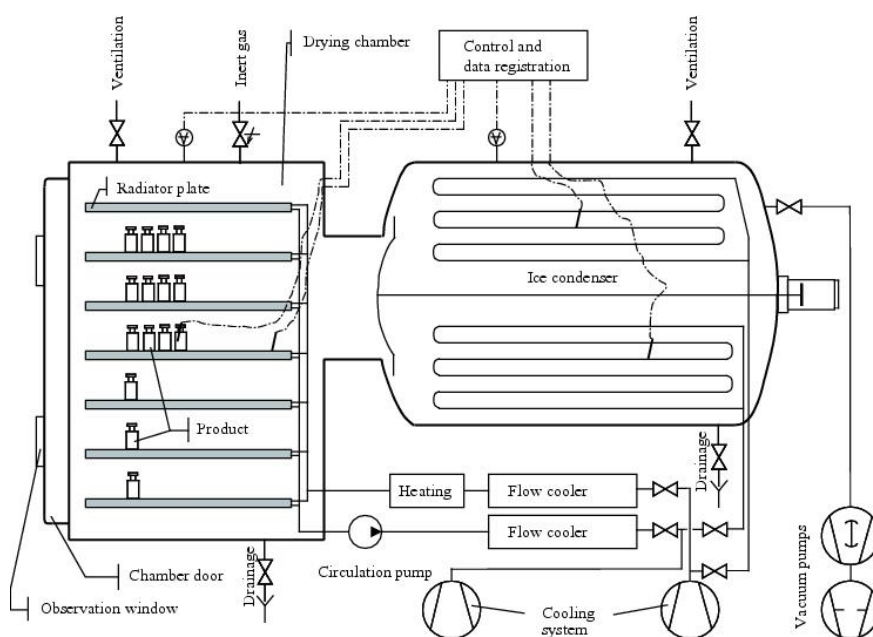


Figure 7.1. Schematics of a freeze dryer and its components (Adapted from: <http://www.gmpua.com>).

The process involves low temperatures and high vacuum, and can thus be considered an energy-intensive one. A full drying cycle typically consists of seven steps, as schematically shown in Fig. 7.2:

1. **Loading.** The vials are loaded on the shelves of the drying chamber. The loading can be full (i.e. all the shelves are loaded with vials) or partial (only one or more of the shelves are loaded). The loading level (i.e. partial or full) can affect the subsequent stages of the drying cycle.

2. Freezing. The product temperature is lowered below the solution freezing point at atmospheric pressure and the solvent freezes, forming ice crystals. The main physical mechanisms of this step are crystal nucleation and growth: in most situations, stochastic nucleation (i.e. non-controlled nucleation) is obtained. In this scenario, the freezing protocol strongly affects the crystal size (i.e., the product morphology), hence all the subsequent drying stages. Some freeze dryers allow keeping under control the nucleation of crystal sites (i.e., nucleation is not stochastic anymore). In this scenario, product morphology can be controlled and therefore the subsequent drying stages can be sped up. In both scenarios, part of the solvent can remain bounded to the product and must therefore be desorbed. The freezing stage typically involves 5-10% of the total duration of the drying cycle.
3. Condenser cooling and evacuation. The condenser temperature is cooled down and vacuum is introduced in the drying chamber. This is an intermediate step that is performed to prepare the system to the subsequent stage.
4. Primary drying. In this step, ice sublimation occurs at low temperatures and high vacuum. The energy supply for ice sublimation is given by the heating fluid (typically a silicone oil) flowing through the coils inserted in the shelves of the drying chamber. Typical temperatures are in the range (-40) - (-10) °C, and typical operating pressures in the drying chamber are between 50-150 mT (milli Torr, where 1 Torr corresponds to 133.322 Pa or 1/760 atm). This step is the most energy-intensive stage and involves 60-80% of the total duration of the drying cycle.
5. Secondary drying. After ice sublimation, the residual particles of solvent that are bounded to the dried product are desorbed in order to lower the residual solvent content (defined as residual moisture content if the solvent is water) of the final product. This is achieved by increasing the temperature of the shelves (typical temperatures range between 20-40 °C) and (possibly) decreasing the chamber pressure. The desorption time strongly depends on the product morphology obtained after freezing, and typically involves 10-20% of the total duration of the drying cycle.
6. Backfilling. The vacuum pump is stopped, and the drying chamber and the condenser are filled with an inert gas at an assigned pressure. A typical backfill operation is performed between 500-900 mbar of dry nitrogen. The advantage of backfilling is that oxidation of the dried product during storage is avoided.
7. Stoppering. The shelves are closed by hydraulic means and the partially inserted stoppers are fully inserted in the vials.

The most critical steps of a drying cycle are step # 2 (freezing), step #4 (primary drying) and step #5 (secondary drying).

Freezing is crucial since the freezing rate directly affects the morphology of the ice crystals and therefore all the subsequent drying steps. Primary drying involves the largest portion of the entire drying duration, and its optimization is therefore pivotal to reduce the overall cycle time. As regard secondary drying, a satisfactory removal of the residual moisture in the product is critical to guarantee its stability over its entire life cycle duration, thus making this step crucial for the entire drying operation.

The focus of this Chapter is to parametrically identify a mathematical model to be subsequently used to optimize the behavior of the system during the primary drying operation, given its strong impact on the entire drying duration. A thorough description of the other two steps and the mathematical models that can be used to describe them can be found in Pikal, 1985 and Fissore *et al.*, 2012.

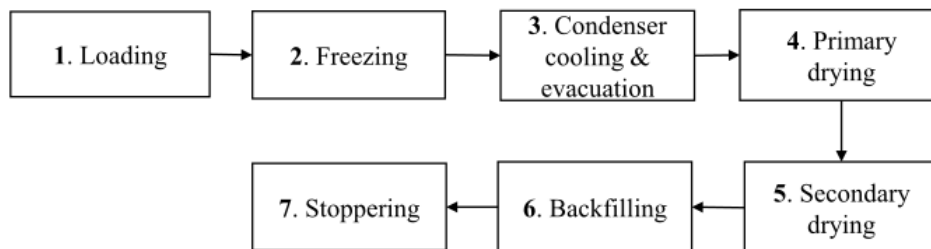


Figure 7.2. Schematics of the different steps of a freeze-drying cycle.

7.1.1 Primary drying

As described in Section 7.1, during primary drying the frozen product is subject to low temperature and high vacuum, and ice sublimation occurs at the interface between the frozen layer and the dried layer. The sublimation front progresses inwards as long as primary drying proceeds, leaving a layer of dried material behind. A schematic representation of the evolution of the dried layer with time is shown in Fig.7.3.

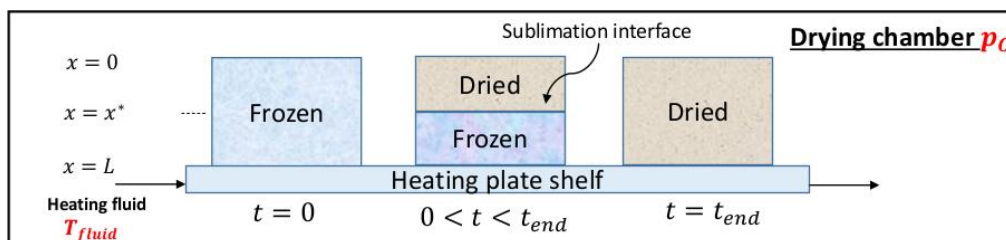


Figure 7.3. Evolution of ice sublimation during primary drying.

In order to avoid product melting (i.e., solid/liquid transition), the partial pressure of water in the drying chamber (but not the total pressure in the chamber) is kept below the triple point of the solvent (6.63 mbar in case of water). Energy supply for ice sublimation at the sublimation interface is provided by the silicone oil flowing inside the shelves and involves different mechanisms of heat transfer, namely (Scutellà *et al.*, 2017a):

1. heat conduction through the gas in the drying chamber;
2. heat conduction between “solids”;
3. radiation.

The first mechanism involves heat conduction through low-pressure water vapor in the gap between the vials and the metallic rails that are used to load the vials in the chamber. Recent studies (Pikal *et al.*, 2016; Scutellà *et al.*, 2018) have shown that this mechanism plays a key role in the overall heat transfer mechanism to the sublimation interface.

The second mechanism involves heat transfer by conduction through solids, namely (i) heat transfer between the silicone oil and the shelf surface, (ii) heat transfer between the shelf surface and the vial bottom through direct contact points, and (iii) heat transfer between the vial bottom and the sublimation interface through the frozen layer.

The third mechanism involves all the radiant effects that occur in the drying chamber. These effects are responsible for the drying heterogeneity between vials of the same shelf and between different shelves, and involve different phenomena such as:

- (a) radiation from top and bottom shelves with respect to the vial considered;
- (b) radiation from the chamber walls to the vials facing the walls or to the rails;
- (c) radiation from the top shelf to the rails;
- (d) radiation from the internal walls of the vials and the sublimation interface.

The above considerations suggest that the description of the heat transfer mechanisms inside the drying chamber can be very complex and some simplifying assumptions are often considered to capture the overall behavior of the freeze dryer. These mechanisms are not the same for all the vials within a single shelf and within the vials of different shelves, and this results in a strong heterogeneity of the drying behavior between vials. As an example, vials that are close to the chamber walls are subject to higher sublimation rates than vials that are placed in internal positions within the shelf, due to heat radiation through the chamber walls and conduction through the gas in the drying chamber.

Besides the issue of drying heterogeneity, the dryer operation during primary drying is affected by two constraints (related both to the equipment and to the product) that must be fulfilled.

1. The product temperature must be kept below the glass transition temperature (for amorphous materials) or eutectic temperature (for crystalline materials) in order to avoid the “collapse”/“melting” of the product (i.e., irreversible alteration of the product morphology). Within a vial, the fraction of the product that is subject to a higher temperature is located at the vial bottom, since it is in contact with the heated shelf.

Monitoring of the vial bottom temperature is therefore critical to avoid the collapse of the cake. Moreover, vials that are close to the chamber walls or to the window (if present) of the freeze dryer, i.e. the so-called “edge” vials, are subject to higher temperatures and slower freezing rates than the internal vials, thus being potentially critical for product collapse.

2. The sublimation rate must be compatible with the condenser capacity and choked flow must be avoided in the duct that connects the drying chamber to the condenser (Searles, 2004). This prevents loss of pressure control inside the drying chamber, with related effects on both the product and the safety of the equipment operation.

The presence of operational and quality constraints on the dried product, coupled with the complexity of heat transfer mechanisms occurring inside the drying chamber, make the mathematical description of primary drying very complex. Moreover, other factors such as equipment geometry, vials geometry, as well as duct, valve and condenser design can affect the drying operation and the final product quality (Pikal, 1985). Several modeling approaches have been proposed to describe the behavior of the system during primary drying. A short classification of these models, with focus on mono-dimensional lumped parameter models, is reported in the next section.

7.2 Mathematical modeling of primary drying

Several mathematical models have been proposed to describe the primary drying stage of freeze-drying. These models aim at predicting the key performance indicators of the system (e.g., vial bottom temperature, amount of water sublimated, etc...) and at describing how primary drying progresses with time. Model classification can be performed using different criteria, e.g. based on the modeling assumptions adopted (pseudo steady-state / dynamic models), based on the ways the system geometry is described (mono- and bi-dimensional models), or based on the type of model parameters involved (lumped/distributed parameter models).

A rough classification of freeze-drying models involves the following categories:

1. Mathematical models based on a detailed description of the multi-physics phenomena of the system, typically (but not necessarily) implemented on computationally fluid dynamics (CFD) software;
2. Mathematical models involving distributed parameters, i.e. parameters that depend on the spatial position within the drying chamber;
3. Mathematical models involving lumped parameters, i.e. parameters that do not depend on the spatial position within the drying chamber.

The first class of models (Liapis and Bruttini, 1995; Mascarenhas *et al.*, 1997; Sheehan and Liapis, 1998; Pisano *et al.*, 2011; Ramšak *et al.*, 2017) involves detailed CFD simulations and

high model complexity (i.e. large number of model equations and model parameters). This model complexity results in two notable drawbacks:

- large computational time, that partially or totally prevents online use of these models for real-time parameter estimation or real-time process optimization;
- large number of model parameters to be estimated. In most cases, the values of these parameters are taken from the literature and may have been obtained under different equipment configuration and/or different product formulation.

The second class of models collects mono- and bi-dimensional models that aim at capturing the heterogeneous behavior within the drying chamber by introducing a spatial dependence on some of the model parameters. The model complexity is kept relatively low and the models can be used, under certain circumstances, for process simulation and control. However, their online use is still limited by identifiability issues (Brülls and Rasmuson, 2002; Tsinontides *et al.*, 2004; Zhai *et al.*, 2005 ; Hottot *et al.*, 2006; Gieseler *et al.*, 2007; Trelea *et al.*, 2007; Duralliu *et al.*, 2018).

The third class of models collects several mono-dimensional lumped-parameter models (e.g., Pikal, 1985; Millman *et al.*, 1985; Sadikoglu and Liapis, 1997; Velardi and Barresi, 2008). involving few parameters and requiring small computational times. This is obtained at the expense of neglecting such phenomena as radial gradients of temperature and composition and heat transfer between the product and the side wall of the vial, as well as other radiation phenomena. Despite these limiting assumptions, it has been shown (Fissore *et al.*, 2010; Fissore *et al.*, 2011a; Fissore *et al.*, 2011b; Bosca *et al.*, 2013; Bosca *et al.*, 2015) that these models can capture the key performance indicators (KPIs) of the system during primary drying with reasonable accuracy, and therefore they represent a good starting point towards a macroscopic mechanistic description of the process. Moreover, the complexity of these models can be modulated by introducing and/or neglecting some (heat-transfer) phenomena that are considered to be relevant/irrelevant for the end purpose, thus providing some flexibility to the modeling framework.

A simplified mono-dimensional lumped parameter model that has gained increasing attention is the one proposed by Velardi and Barresi (2008), based on a previous study (Pikal, 1985). This model is one-dimensional and collects the mass balance for the frozen layer and the macroscopic heat- and mass-transfer equations for the frozen product. The relevant model equations are as follows.

The frozen product is heated by the silicone oil flowing through the shelf, and the driving force for heat transfer is given by the difference between the temperature $T_{fluid}(t)$ of the heating fluid and the temperature $T_B(t)$ at the bottom of the frozen layer. The heat flux $J_q(t)$ [W/m²] between the silicone oil and the vial bottom can therefore be expressed as:

$$J_q = K_v(T_{fluid} - T_B) \quad , \quad (7.1)$$

where K_v [W/(m²K)] is the heat transfer coefficient. It is worth noticing that Eq. (7.1) only accounts for the heat transfer between the silicone oil and the vial bottom. However, as explained in section 7.1.1., other heat transfer mechanisms affect the overall energy supply to the frozen layer, but they are not taken into account explicitly in the heat flux formulation (7.1). The heat transfer coefficient of Eq. (7.1) strongly depends on the pressure p_c [Pa] in the chamber. The dependence of K_v with respect to p_c is expressed as:

$$K_v = C_1 + \frac{C_2 p_c}{1 + C_3 p_c} \quad , \quad (7.2)$$

where C_1, C_2, C_3 are parameters to be estimated.

Given the heat supply expressed by Eq. (7.1), the frozen product sublimates at the sublimation front and the sublimation flux $J_w(t)$ [kg/(m²s)] can be expressed as:

$$J_w = \frac{1}{R_p} (p_{w,i} - p_{w,c}) \quad , \quad (7.3)$$

where $p_{w,i}$ [Pa] is the partial pressure of water at the sublimation front, $p_{w,c}$ [Pa] is the partial pressure of water in the chamber (which is typically assumed equal to the total pressure p_c of the chamber), and R_p [m/s] is the resistance to the mass transfer.

The interpretation of Eq. (7.3) is straightforward: the driving force for the sublimation flux is the difference between the vapor partial pressure at the sublimation interface and the vapor partial pressure in the chamber, and the proportionality factor is the inverse of the mass transfer resistance to the vapor flow. The mass transfer resistance R_p is a function of the thickness L_{frozen} of the frozen layer according to the expression:

$$R_p = R_0 + \frac{R_1 L_{frozen}}{1 + R_2 L_{frozen}} \quad , \quad (7.4)$$

where R_0, R_1 and R_2 are parameters to be estimated.

The partial pressure $p_{w,i}$ of water at the sublimation interface is a function of the temperature T_i [K] at the interface. This dependence can be expressed with the Goff-Gratch equation (Goff, 1946); a simplified expression for the values of pressure and temperature that are typically involved during primary drying has been proposed by Fissore and Barresi (2011):

$$p_{w,i} = \exp \left(-\frac{6150.6}{T_i} + 28.932 \right) \quad . \quad (7.5)$$

The heat that is exchanged between the oil flowing in the shelves and the product is assumed to be entirely used for ice sublimation. The energy balance at the sublimation interface can be written as:

$$J_w = \frac{1}{\Delta H_s} J_q = \frac{1}{\Delta H_s} K_v (T_{fluid} - T_B) \quad , \quad (7.6)$$

where ΔH_s [J/kg] is the heat of sublimation, which is dependent on the partial pressure of water in the chamber (the pressure dependency can often be neglected). Since no heat accumulation is assumed in the frozen layer, the energy balance for the frozen product can be written as:

$$K_v (T_{fluid} - T_B) = \frac{T_{fluid} - T_i}{\left(\frac{1}{K_v} + \frac{L_{frozen}}{\lambda_{frozen}}\right)} \quad , \quad (7.7)$$

where λ_{frozen} [W/(mK)] is the thermal conductivity of the frozen layer. Eq. (7.7) can be rewritten as:

$$T_B = T_{fluid} - \frac{1}{K_v} \left(\frac{1}{K_v} + \frac{L_{frozen}}{\lambda_{frozen}}\right)^{-1} (T_{fluid} - T_i) \quad , \quad (7.8)$$

which allows relating the temperature T_B at the bottom of the vial (measurable variable) to the temperature T_i at the sublimation interface (unmeasurable variable).

The time evolution of the frozen layer thickness can be obtained by solving the material balance:

$$\frac{dL_{frozen}}{dt} = - \frac{1}{\rho_{frozen} - \rho_{dried}} J_w \quad , \quad (7.9)$$

where ρ_{frozen} [kg/m³] is the density of the frozen product, and ρ_{dried} [kg/m³] is the density of the dried product.

The total amount M [kg] of water sublimated during primary drying can be computed from the sublimation flux J_w according to the equation:

$$M = A_v \int_0^{t_d} J_w(t) dt \quad , \quad (7.10)$$

where t_d is the total duration of primary drying, and A_v is the internal cross-sectional area of the vial.

The mono-dimensional model consists of Eq. (7.1)-(7.10). From a structural point of view, the model is a system of one differential equation (Eq. (7.9)), one integral equation and seven algebraic equations. The model inputs, outputs and parameters are listed in Table 1.

Table 7.1. *Input variables, output variables and model parameters for the mono-dimensional model (7.1)-(7.10).*

ID	Variable name	Units	Symbol	Type
<i>Inputs</i>				
1	Chamber pressure	[Pa]	p_c	Algebraic
2	Heating fluid temperature	[K]	T_{fluid}	Algebraic
<i>Outputs</i>				
1	Length of the frozen layer	[m]	L_{frozen}	Differential
2	Cumulative amount of water sublimated	[kg]	M	Algebraic
3	Heat flux	[W/m ²]	J_q	Algebraic
4	Sublimation flux	[kg/(m ² s)]	J_w	Algebraic
5	Vial bottom temperature	[K]	T_B	Algebraic
6	Sublimation interface temperature	[K]	T_i	Algebraic
7	Water partial pressure at the sublimation interface	[Pa]	$p_{w,i}$	Algebraic
<i>Parameters</i>				
1	Heat transfer coefficient - parameter #1	[W/(m ² K)]	C_1	[-]
2	Heat transfer coefficient - parameter #2	[W/(m ² KPa)]	C_2	[-]
3	Heat transfer coefficient - parameter #3	[W/(m ² KPa)]	C_3	[-]
4	Mass transfer resistance - parameter #1	[m/s]	R_0	[-]
5	Mass transfer resistance - parameter #2	[1/s]	R_1	[-]
6	Mass transfer resistance - parameter #3	[1/m]	R_2	[-]

With reference to Table 7.1, some comments are in order.

1. The two variables that can be manipulated during the dryer operation are the pressure in the chamber and the heating fluid temperature. These variables can be manipulated according to a time-invariant, piecewise constant or (only the heating fluid temperature) piecewise linear behavior. Additional details will be given in section 7.5.
2. Among the output variables, one of the variables that can typically be measured during a drying cycle is the vial bottom temperature. This is done by placing a thermocouple inside the vial, and measurements can be obtained online with relatively low frequency. Thermocouples are typically placed inside different vials located at selected positions within the dryer in order to account for their different thermal behaviors within the same shelf.
3. The dryer operation is subject to both product constraints and equipment constraints. From a modeling perspective, the product constraints are imposed on the vial bottom temperature T_B , which must be kept below the critical temperature (glass transition or eutectic temperature) of the product for the entire drying operation. The equipment constraints are imposed on the sublimation flux J_w , that must be kept below a limiting

value to avoid choked flow through the duct linking the drying chamber to the condenser.

In view of the above, obtaining good predictions for the vial bottom temperature and sublimation flux from the model is crucial to guarantee product quality and safe operation. Therefore, T_B and J_w are considered as KPIs for the model. Moreover, the amount of water sublimated during primary drying can be considered as an additional KPI, since its value is directly related to the overall duration of the primary drying stage: obtaining a good prediction of M is therefore essential for process optimization.

4. The model involves 6 parameters, that describe the dependency of the global heat transfer coefficient K_v and the mass transfer resistance R_p with respect to the chamber pressure and length of the frozen layer respectively. A good estimation of the values of these parameters is required in order to obtain reliable predictions for the KPIs described above.

The step-by-step methodology that has been used to calibrate and validate the above model is reported in the next section.

7.3 Model-based design of experiments (MBDoe): brief overview

In the following, a short description of the mathematical background of model-based design of experiments (MBDoe) techniques is reported. A thorough mathematical formulation can be found in the work of Franceschini and Macchietto (2008) and Galvanin (2013).

7.3.1 Design of an optimal experiment

The goal of MBDoe is to find the optimal initial conditions, the optimal profile of manipulated inputs, the optimal sampling instants for the measured variables, and the optimal duration of a new experiment in order to maximize the information that can be extracted from the experiment for the purpose of model identification.

Mathematically, given a dynamic system described by a system of differential and algebraic equations (DAEs) of the form:

$$\begin{aligned} \mathbf{f}(\dot{\mathbf{x}}(t), \mathbf{x}(t), \mathbf{u}(t), \mathbf{w}, \boldsymbol{\theta}, t) &= 0 \\ \mathbf{y}(t) &= \mathbf{h}(\mathbf{x}(t)) \end{aligned} \quad (7.11)$$

where $\boldsymbol{\theta} [1 \times N_\theta]$ is the set of parameters to be estimated, $\mathbf{x}(t)$ is the state vector, $\mathbf{u}(t)$ and \mathbf{w} are the time-dependent and time-invariant control variables respectively, $\mathbf{y}(t)$ is the output vector, the target is to find the optimal value of the design vector:

$$\boldsymbol{\varphi} = [\mathbf{y}_0, \mathbf{u}(t), \mathbf{w}, \mathbf{t}^{sp}, \tau] \quad (7.12)$$

that maximizes the information for the identification of $\boldsymbol{\theta}$. Eq. (7.12) uses the following notation:

- \mathbf{y}_0 is the set of initial conditions for the measured variables;
- $\mathbf{t}^{sp} = [t_1, t_2, \dots, t_{sp}]$ is the vector collecting the different sampling times in which measurements are collected;
- τ is the total duration of the experiment.

The time-dependent input variables $\mathbf{u}(t)$ are typically approximated as piecewise constant or piecewise linear variables for a given set of intervals splitting the entire experiment duration. The optimal value of the design vector is found as the value that minimizes a given metric α of the variance-covariance matrix of the model parameters $\mathbf{V}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi})$ according to:

$$\boldsymbol{\varphi}_{opt} = \arg \min \{ \alpha [\mathbf{V}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi})] \} \quad . \quad (7.13)$$

The variance-covariance matrix $\mathbf{V}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi})$ is defined as:

$$\mathbf{V}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi}) = [\mathbf{H}_\theta^0 + \mathbf{H}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi})]^{-1} \quad , \quad (7.14)$$

where $\mathbf{H}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi})$, defined as the dynamic information matrix, is a measure of the information that can be exploited from the experiment for the identification of $\boldsymbol{\theta}$, while \mathbf{H}_θ^0 is the initial (i.e., preliminary) value of this matrix. For a set of N_{exp} experiments, $\mathbf{H}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi})$ is expressed as:

$$\mathbf{H}_\theta(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \sum_{k=1}^{N_{exp}} \sum_{i=1}^{N_y} \sum_{j=1}^{N_y} s_{ij|k} \mathbf{Q}_{i|k}^T \mathbf{Q}_{j|k} + \mathbf{H}_\theta^0 \quad , \quad (7.15)$$

where $\mathbf{Q}_{i|k}$ is the dynamic sensitivity matrix of the i -th measured response in the k -th experiment, whose elements for the k -th experiment are given by:

$$\mathbf{Q}_i = \left[\frac{\partial y_{il}}{\partial \theta_m} \right] \quad l = 1, 2, \dots, n_{sp} ; m = 1, 2, \dots, N_\theta \quad , \quad (7.16)$$

and $s_{ij|k}$ is the ij -th element of the inverse matrix of the measurement errors in the k -th experiment. According to the metric α that is optimized in Eq. (7.13), different design criteria can be identified:

- A-optimal: α is defined as the trace of the variance-covariance matrix;
- E-optimal: α is defined as the largest eigenvalue of the variance-covariance matrix;
- D-optimal: α is defined as the determinant of the variance-covariance matrix.

Based on the design criterion adopted, the solution of optimization problem (7.13) allows obtaining all the information needed to perform a new optimal experiment for the system under investigation.

7.3.2 Additional MBD_{oE} activities

The design of a new experiment according to the formulation described in the previous section is followed by two sequential activities:

1. Execution of the designed experiment on the experimental equipment;
2. Estimation of parameters based on the data collected from the experiment.

If the parameter estimation is statistically satisfactory, the MBD_{oE} activity can be considered as concluded. On the other hand, if the parameter estimation is not satisfactory, a new experiment must be designed based on the updated estimates of the model parameters, and the overall procedure must be iterated. A schematic representation of the iterative procedure is shown in Fig. 7.4.



Figure 7.4. Simplified flowsheet for a standard MBD_{oE} activity.

The stopping criterion for the MBD_{oE} activity can be dictated by different (and often case-dependent) occurrences, including:

- attainment of a desired level of *precision* on the estimated values of the model parameters;
- attainment of a desired level of model *accuracy*, to be intended as goodness of fit for the desired responses;
- attainment of the maximum experimental budget, typically dictated by both operational and cost-dependent constraints.

Once one or more of the above conditions is reached, the MBD_{oE} activity is stopped. The final target is the attainment of a well-calibrated model to be used for decision-support scenarios.

7.4 Preliminary assessment of freeze-drying model performance

In order to facilitate a successful application of MBD_{oE} techniques with the mono-dimensional model described in Section 7.2, two preliminary activities are required, namely:

1. assessment of the adequacy of the model structure with respect to the experimental observations on the real equipment. In other words, it must be assessed if the (potential) mismatch between the model predictions and the experimental observations can be removed by adjusting (i.e., calibrating) the model parameters. Alternatively, structural

deficiencies (e.g., due to a lack of fundamental knowledge or to simplifying assumptions introduced during model development) may be identified in the original model. In the latter case, the model needs structural adjustments before using it for MBDoE purposes (yet, MBDoE can potentially be used also to discriminate among rival model structures).

2. identification of reasonable preliminary values for the model parameters, in order to speed up the MBDoE activity (Fig. 7.4) and reduce the total number of experiments to be performed for model calibration. These values may be taken (if available) from the literature or estimated from historical data.

For the case study reported in this Chapter, historical data of primary drying were available. A detailed description of each step, together with the organization of the available historical data, is reported in the next sections.

7.4.1 Description of the historical datasets

Five historical experimental datasets have been used to preliminarily calibrate and validate the model parameters. The datasets differ from each other for the product considered, the type of vials used during drying, the loading of the shelves in the freeze dryer, the number of measurements available, and the experimental protocol adopted. All experiments the datasets refer to were carried out in the same equipment, and this equipment is the same to be used within the current MBDoE activity.

Three historical datasets (namely, datasets A, B and C) were used for model calibration. The remaining two datasets (D, E) were used to validate the model. A detailed description of each dataset is reported in the following.

7.4.1.1 Calibration datasets

The following three datasets were used to calibrate the model.

- Dataset A: data refer to three incomplete primary drying experiments (i.e., the primary drying was stopped before reaching its completion) at three different values of chamber pressure, namely at 67%, 100% and 133% of its nominal value. Siliconized vials were used in each experiment. The vials were filled with 0.6 mL of microfiltered water and only one shelf out of five (i.e., 20% partial loading) was loaded. Two different vial frames were loaded for each shelf, as schematically shown in Fig. 7.6. From the same figure, it can be seen that the vial arrangement in each shelf has been partitioned in 12 different zones, each of which with expected similar thermal behavior. Measurements of vial bottom temperature were collected by placing one thermocouple in one single vial for each zone (light blue-shaded vials in Fig. 7.6). In other words, the behavior of the vial inside which a thermocouple was placed is considered as representative of the behavior of all the vials within the same zone. This is a simplification dictated by

operational limits, and some studies (Pikal, 2000; Pikal *et al.*, 2016) have also shown that the very presence of a thermocouple alters the thermal behavior of the vial with respect to the adjacent ones.

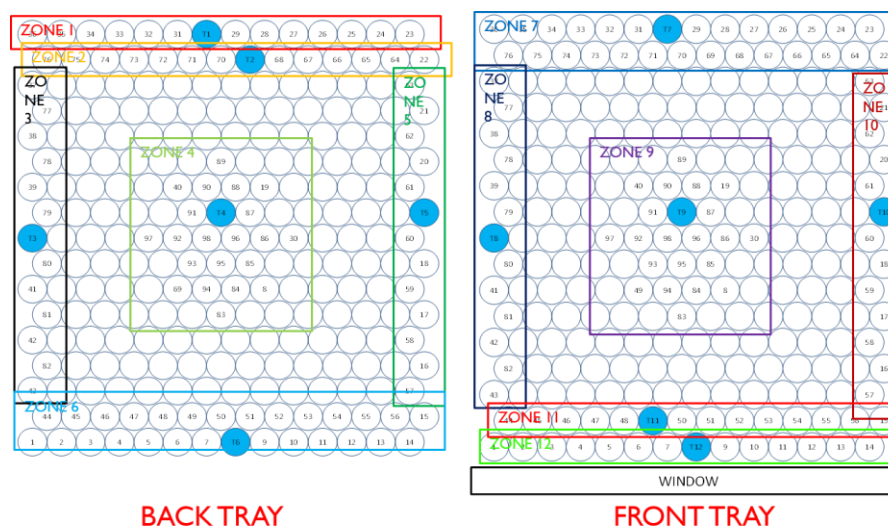


Figure 7.6. Arrangement of the vials on the shelves. Vials tagged with a number were weighted before and after drying. Thermocouples were placed inside the vials shaded in light blue. The rectangular boxes indicate the logic partitioning of vials considered in each shelf.

Measurements of chamber pressure (obtained with both a Pirani gauge and a capacitance manometer), heating fluid temperature at the inlet of the shelf and vial bottom temperature at the 12 locations described above were collected with a measurement interval of 30 s for the first two experiments, and of 20 s for the third experiment. Moreover, the vials tagged with a number in Fig. 7.6 were weighted at the beginning and at the end of the experiment, thus giving an additional measurement of the total amount of water sublimated during primary drying. The time profiles for the input variables (chamber pressure and heating fluid temperature) that were set during the experiments are reported in Table 7.2

Table 7.2. Time profiles of the input variables for dataset A.

Step #	Interval duration [min]	p_c/\bar{p}_c [-]	$T_{fluid}/\bar{T}_{fluid}$ [-]
<i>Experiment 1</i>			
1	15	0.667	0.898 (constant)
2	30	0.667	0.898-1.000 (linear)
3	327	0.667	1.000 (constant)
<i>Experiment 2</i>			
1	15	1	0.898 (constant)
2	30	1	0.898-1.000 (linear)
3	186	1	1.000 (constant)
<i>Experiment 3</i>			
1	15	1.333	0.898 (constant)
2	30	1.333	0.898-1.000 (linear)
3	100	1.333	1.000 (constant)

- **Dataset B:** the structure of this dataset is the same as for dataset A (with three experiments at three different pressures). The formulate is a 5% sucrose-based placebo and the vials used were non-siliconized. The same measurements collected for dataset A were collected also for this dataset. The time profiles of the input variables set during the three experiments are reported in Table 7.3.

Table 7.3. Time profile of the input variables for dataset B.

Step #	Interval duration [min]	p_c/\bar{p}_c [-]	$T_{fluid}/\bar{T}_{fluid}$ [-]
<i>Experiment 1</i>			
1	15	0.667	0.898 (constant)
2	30	0.667	0.898-1.000 (linear)
3	331	0.667	1.000 (constant)
<i>Experiment 2</i>			
1	15	1	0.898 (constant)
2	30	1	0.898-1.000 (linear)
3	180	1	1.000 (constant)
<i>Experiment 3</i>			
1	15	1.333	0.898 (constant)
2	30	1.333	0.898-1.000 (linear)
3	120	1.333	1.000 (constant)

- **Dataset C:** as for the previous datasets, data refer to an incomplete drying cycle, but were obtained with a 5% sucrose-based formulation placebo at only one value for the chamber pressure (equal to its nominal value). The same vials format and the same thermocouple arrangement as in the previous datasets experiments were used, but vials

were in this case siliconized. As for the previous case, measurements of chamber pressure (Pirani and capacitance manometers), heating fluid temperature and vial bottom temperature were collected every 30 s. The tagged vials were weighted before and after drying in order to measure the total amount of water sublimated during primary drying. The recipe used for this experiment is reported in Table 7.4.

Table 7.4. Time profiles of the input variables for dataset C.

Step #	Interval duration [min]	$p_c/\bar{p}_c [-]$	$T_{fluid}/\bar{T}_{fluid} [-]$
1	15	1 (constant)	0.898 (constant)
2	30	1 (constant)	0.898-1.000 (linear)
3	175	1 (constant)	1.000 (constant)

7.4.1.2 Validation datasets

The following two datasets have been used to validate the model.

- **Dataset D:** data have been extracted from a dataset of a complete drying cycle (involving all the seven steps described in section 7.1.2) and refer to the primary drying step of this cycle. The dataset was obtained with a 5% sucrose-based formulation placebo using non-siliconized vials. Only one shelf (the central one) out of five was loaded with 238 vials. Each vial was filled with 0.6 mL of placebo. Two thermocouples were placed to monitor the bottom temperature of two vials placed in the middle of the first row and last row of the first half of the central shelf, according to the schematic representation of Fig. 7.7. The time profile for the input variables (chamber pressure and heating fluid temperature) during primary drying is the one reported in Table 7.5. The dataset collects measurements of chamber pressure (obtained with a Pirani gauge and a capacitance manometer), heating fluid temperature at the inlet of the shelf and vial bottom temperature at the two locations described above with a measurement interval of 30 s.

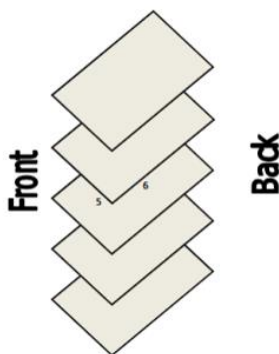
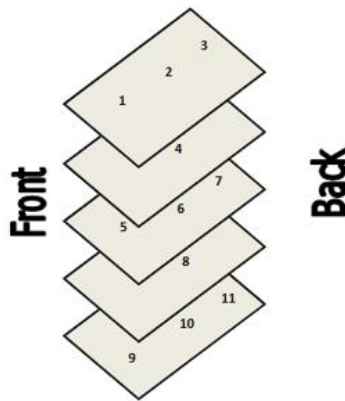


Figure 7.7. Arrangement of the thermocouples (numbered as 5 and 6) on the central shelf for dataset D.

Table 7.5. Times profile of the input variables for dataset D.

Step #	Interval duration [min]	p_c/\bar{p}_c [-]	$T_{fluid}/\bar{T}_{fluid}$ [-]
1	15	1 (constant)	0.910 (constant)
2	30	1 (constant)	0.910-1.000 (linear)
3	990	1 (constant)	1.000 (constant)
4	180	1 (constant)	1.237 (constant)

- **Dataset E:** as for the previous dataset, data of primary drying have been extracted from a complete freeze-drying cycle of a 5% sucrose-based placebo. The dataset was obtained by loading 2380 siliconized vials (full loading) on all the five shelves placed in the drying chamber. Eleven thermocouples were inserted in 11 vials placed according to the schematic representation of Fig. 7.8. The filling volume was set to 0.6 mL as for the previous datasets. The recipe for the primary drying is reported in Table 7.6. The time profiles for the input variables are the same as for the previous dataset, with a significant increase in the duration of the third time interval. Measurements of chamber pressure (Pirani and capacitive), heating fluid temperature and vial bottom temperature for the 11 locations described above were collected every 30 s.

**Figure 7.8.** Arrangement of the eleven thermocouples (numbers) on the five shelves for dataset E.**Table 7.6.** Time profile of the input variables (i.e. “recipe”) for dataset E.

Step #	Interval duration [min]	p_c/\bar{p}_c [-]	$T_{fluid}/\bar{T}_{fluid}$ [-]
1	15	1 (constant)	0.910 (constant)
2	30	1 (constant)	0.910-1.000 (linear)
3	2430	1 (constant)	1.000 (constant)
4	180	1 (constant)	1.237 (constant)

The step-by-step methodology used for this preliminary analysis is shown in Fig. 7.5.

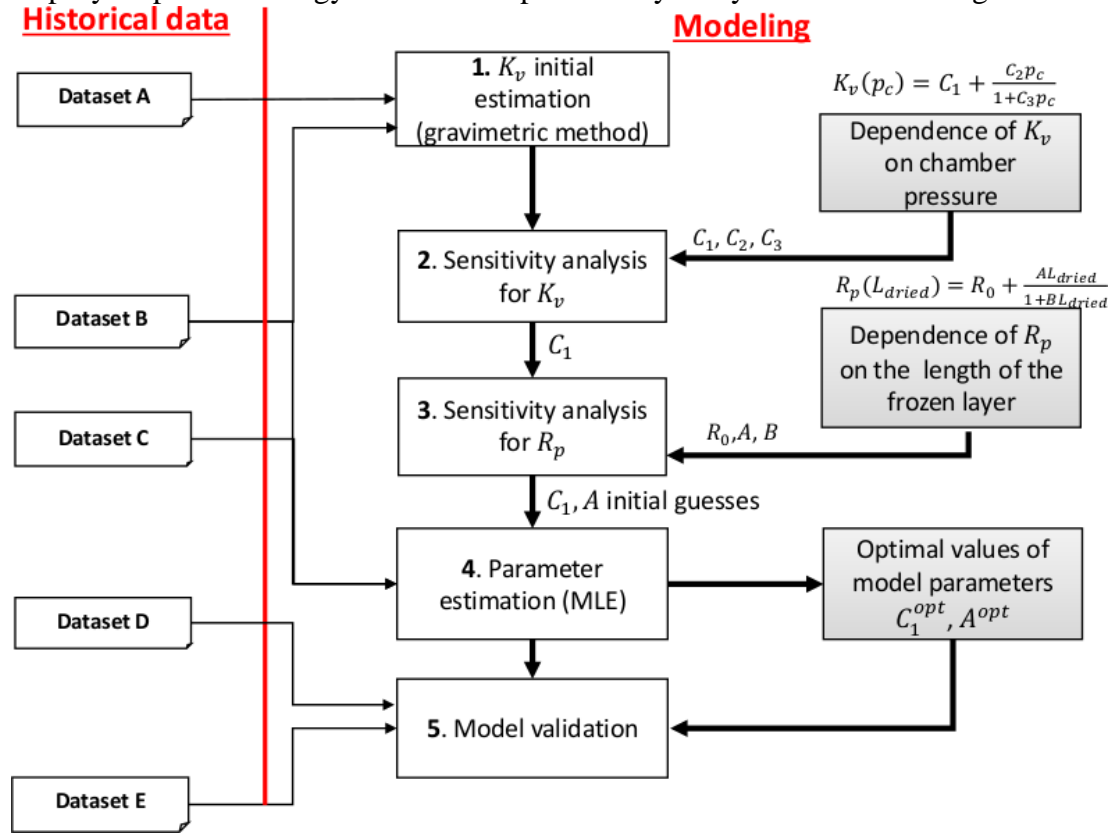


Figure 7.5. Schematic of the different steps for the preliminary assessment of model performance.

7.4.2 Step #1: gravimetric estimation of K_v

In the first step of the proposed methodology, an estimation of the parameters C_1, C_2 and C_3 for the heat transfer coefficient K_v has been obtained through a standard gravimetric method (Fissore and Barresi, 2010). The purpose of this gravimetric estimation for the scope of this work is to facilitate the numerical solution of the maximum likelihood estimation of these parameters (step # 4), by assigning reasonable initial guesses to the estimator. In fact, it has been verified that bad initial guesses of C_1, C_2 and C_3 may lead the estimator to diverge.

The gravimetric method is a standard procedure adopted in the characterization of freeze drying equipment to estimate the heat transfer coefficient K_v . This coefficient is estimated from incomplete drying experiments (such as the ones collected for datasets A and B), given the availability of measurements of vial bottom temperature $T_B(t)$ and total amount of water sublimated M for a given drying time t_d , according to:

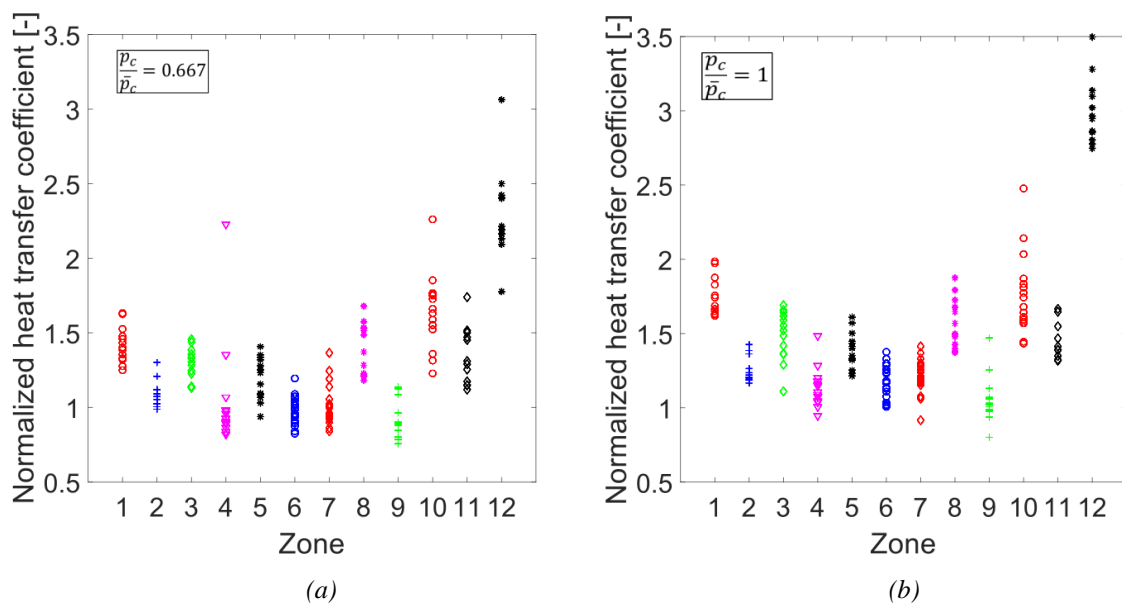
$$K_v = \frac{M \Delta H_s}{A_v \int_0^{t_d} (T_{fluid}(t) - T_B(t)) dt} \quad (7.17)$$

As discussed in section 7.2, the $T_B(t)$ profile, as measured by a thermocouple located inside a given zone, is typically assumed to be the same for all the vials within that zone. Therefore, Eq. (7.17) allows obtaining a different value of K_v for each vial characterized by a measurement of M . This allows one to obtain a distribution of values of K_v both within the same zone and between different zones of a given shelf. This in turn allows a quantitative characterization of the thermal behavior of each zone.

If the experiment is repeated for at least three different values of chamber pressure, the values of C_1 , C_2 and C_3 of Eq. (7.2) can be obtained by simple fitting. It is worth noticing that the values of K_v do not depend (in principle) on the product considered, but are affected by the type of vials (i.e. siliconized or non-siliconized in this study). In fact, the presence of a thin layer of silicone on the internal surface of a vial affects the overall heat transfer to the frozen product. Based on the previous considerations, a preliminary estimation of the heat transfer coefficient K_v , as well as of the parameters C_1 , C_2 and C_3 , have been obtained from the historical datasets A (siliconized vials) and C (non-siliconized vials). The results are presented in the next sections.

7.4.2.1 Results for dataset A

The distribution of the values of K_v for each zone and between different zones is shown in Fig. 7.9 for the three different values of chamber pressures described in section 7.4.1. The values of K_v have been normalized with the value obtained for a vial placed at the same location (zone 7) at $p_c/\bar{p}_c = 0.667$.



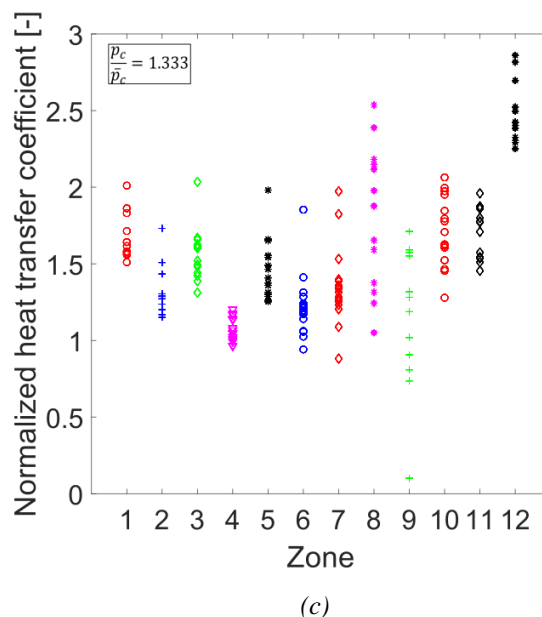


Figure 7.9. Values of the heat transfer coefficient for the different vials on a shelf at three different chamber pressures.

The mean values of (normalized) K_v are reported in Tables 7.7-7.9 together with their standard deviations, for each zone and at the three pressures at which the experiments were carried out. Note that standard deviations are computed with respect to the normalized values of K_v , and that the mean values have been normalized with the same nominal values used for Fig 7.9. The vials have been divided into two categories, namely “edge” (the ones belonging to zones 1, 2, 3, 5, 8, 10, 11) and “central” (the ones belonging to zones 4, 6, 7, 9), in order to reflect the geometric arrangement shown in Fig. 7.6.

Table 7.7. Mean values of the heat transfer coefficient and standard deviations for dataset A at $p_c/\bar{p}_c = 0.667$.

Zone #	Normalized mean value of K_v [-]	Standard deviation	Vial type
1	1.420	0.121	Edge
2	1.089	0.090	Edge
3	1.279	0.113	Edge
4	1.036	0.154	Central
5	1.192	0.138	Edge
6	0.977	0.128	Central
7	1.015	0.171	Central
8	1.396	0.131	Edge
9	0.888	0.142	Central
10	1.627	0.242	Edge
11	1.343	0.192	Edge
12	2.286	0.299	Edge (window)

Table 7.8. Mean values of the heat transfer coefficient and standard deviations for dataset A at $p_c/\bar{p}_c = 1$.

Zone #	Normalized mean value of K_v [-]	Standard deviation	Vial type
1	1.747	0.098	Edge
2	1.255	0.066	Edge
3	1.495	0.125	Edge
4	1.134	0.097	Central
5	1.390	0.087	Edge
6	1.168	0.082	Central
7	1.202	0.085	Central
8	1.572	0.121	Edge
9	1.052	0.113	Central
10	1.765	0.204	Edge
11	1.430	0.095	Edge
12	2.988	0.164	Edge (window)

Table 7.9. Mean values of the heat transfer coefficient and standard deviations for dataset A at $p_c/\bar{p}_c = 1.333$.

Zone #	Normalized mean value of K_v [-]	Standard deviation	Vial type
1	1.690	0.122	Edge
2	1.321	0.131	Edge
3	1.563	0.139	Edge
4	1.069	0.062	Central
5	1.460	0.157	Edge
6	1.211	0.141	Central
7	1.338	0.185	Central
8	1.810	0.359	Edge
9	1.149	0.369	Central
10	1.700	0.180	Edge
11	1.697	0.133	Edge
12	2.476	0.158	Edge (window)

The numerical values of K_v obtained with the gravimetric method suggest some considerations.

1. As expected, the vial thermal behavior across different zones changes remarkably. Namely, vials at the edge of a shelf are characterized by greater values of K_v with respect to the central ones. This results in larger sublimation rates (hence smaller drying times), but also in greater product temperatures. Product temperature control for these vials is therefore crucial to avoid irreversible alteration of the product morphology.

2. Vials directly exposed to the equipment window (zone 12) show considerable greater values of K_v with respect to other vials. This is clearly related to the heat exchange through the window, which speeds up the sublimation process. Following the reasoning of the previous remark, this means that these vials are significantly exposed to the risk of product collapse/melting.
3. The variability of K_v across vials belonging to the same zone is not negligible, as the standard deviations of Tables 7.7-7.9 suggest. However, this variability is significantly smaller than the one that would be obtained by considering a single value of K_v for *all* the vials placed on the shelf, as can be easily inferred from Table 7.10. In other words, whereas assuming a uniform thermal behavior for the vials belonging to a given zone is reasonable, the same cannot be said for *all* the vials located on the same shelf. This means that the mono-dimensional model (7.1)-(7.10) must be calibrated *for each thermal zone*, since significantly different values of K_v are expected across zones.

Table 7.10. Mean values and standard deviations of the “global” (i.e. computed for all the 12 zones) heat transfer coefficient at the three different values of chamber pressure considered.

Normalized chamber pressure [-]	Normalized mean value of K_v [-]	Standard deviation
0.667	1.347	0.765
1	1.577	0.854
1.333	1.614	0.895

In view of the above considerations, the values of parameters C_1 , C_2 and C_3 of Eq. (7.2) can be computed for each zone. Although the numerical values of these parameters cannot be reported for confidentiality reasons, it was noted that:

- the value of C_1 significantly differs between the different zones. This is justified by the fact that, as proved by Pikal and coworkers (Pikal, 2000), this parameter primarily accounts for the radiation contributions from the top and bottom shelves, as well as from the chamber walls. These contributions change with the location in the freeze dryer (Tables 7.7-7.9).
- The values of C_2 and C_3 do not significantly vary between zones. This is in agreement with the work of Pikal (Pikal, 2000), who justifies this behavior by noting that C_3 mainly depends on the distance between the bottom of the vial and the shelf (that is constant in the equipment considered), whereas C_2 mainly accounts for conductive effects through the sidewall of the vial (which are constant within the shelf) and through the gas entrapped between the vial bottom and the shelves.

Based on the above observations, the twelve different values of C_1 obtained with the gravimetric method have been used as initial guesses for the twelve different parameter (preliminary)

estimation exercises required at step #4 of the proposed methodology. On the other hand, C_2 and C_3 have been used as fixed values.

The procedure described above for dataset A was repeated for dataset C, in order to account for the different type of vials considered (siliconized vs. non-siliconized). The numerical results can be found in Appendix A. The results show that, surprisingly, the presence of the silicone layer does not significantly affect the values of the heat transfer coefficient, even though they significantly affect the values of the mass transfer resistance, as will be discussed in more detail in section 7.4.5.

7.4.3. Step #2: sensitivity analysis for K_v

As discussed in the previous section, due to the variability of the heat transfer coefficient between the different zones of a shelf, twelve different model calibrations are required in order to capture the relevant behavior of each zone of the shelf. It is worth noticing that, in principle, whereas parameter C_1 requires to be estimated with the maximum likelihood estimator on each zone, parameters C_2 and C_3 need to be estimated for one zone only, since their values are expected to be the same for all zones.

Notwithstanding this aspect, there would be at least one model calibration activity that would require the joint maximum likelihood estimation of parameters C_1, C_2 and C_3 , as well as of mass transfer resistance parameters R_0, R_1, R_2 . As discussed in section 7.4, it was verified that the joint estimation of these six parameters is affected by convergence issues due to the limited number of measured response variables (vial bottom temperature and end-point cumulative amount of water sublimated). This prevented the possibility of jointly estimating all parameters together. Therefore, the possibility to reduce the total number of parameters to be estimated was considered in order to facilitate the estimator convergence.

Sensitivity analysis (Cruz and Perkins, 1964) can be used to understand the influence of model parameters with respect to the KPIs of the system. The greater the sensitivity of a given parameter with respect to a given KPI, the greater the impact of a wrong calibration of that parameter on the model prediction of the KPI. In other words, it is always desirable to obtain accurate estimations of the parameters with highest sensitivities with respect to the outputs of interest, since small variations on their values can significantly affect the model predictions. On the other hand, parameters with smaller sensitivities can be set to their nominal (e.g., literature) values without significantly impacting on the model performance.

During primary drying, the KPIs of interest are the temperature T_B at the bottom of the vials, the sublimation flux J_w , and the cumulative amount M of water sublimated during the drying stage. The first two KPIs are dynamic variables, while the third one is an end-point variable. Therefore, the sensitivities of C_1, C_2, C_3 with respect to the first two KPIs are time-dependent (i.e., they are defined as *dynamic* sensitivities) and are given by:

$$s_{T,i}(t) = \frac{\partial T_B(t)}{\partial C_i} ; s_{J_w,i}(t) = \frac{\partial J_w(t)}{\partial C_i} \quad i = 1,2,3 \quad , \quad (7.18)$$

while the sensitivities with respect to M do not depend on time:

$$s_{M,i} = \frac{\partial M}{\partial C_i} ; \quad i = 1,2,3 \quad . \quad (7.19)$$

As regard the dynamic sensitivities, the maximum absolute values of their time profiles have been considered

Table 7.11. Sensitivities of vial bottom temperature, sublimation flux and cumulative amount of water sublimated with respect to the three parameters C_1, C_2, C_3 . The values refer to a +1% deviation on the nominal value of the C_i parameters and sensitivities are expressed as % of the nominal value of the KPI considered. Values in boldface indicate the strongest sensitivities.

Parameter	$ s_{T,i}^{max} $ [%]	$ s_{J_w,i}^{max} $ [%]	$ s_{M,i} $ [%]	Parameter	$ s_{T,i}^{max} $ [%]	$ s_{J_w,i}^{max} $ [%]	$ s_{M,i} $ [%]
<i>Zone 1</i>				<i>Zone 7</i>			
C_1	0.752	5.789	26.96	C_1	0.512	4.874	24.12
C_2	0.186	1.012	0.312	C_2	0.182	0.945	0.289
C_3	0.163	0.198	0.059	C_3	0.158	0.164	0.045
<i>Zone 2</i>				<i>Zone 8</i>			
C_1	0.645	5.124	26.12	C_1	0.647	5.478	27.56
C_2	0.178	1.123	0.365	C_2	0.187	1.058	0.395
C_3	0.166	0.201	0.062	C_3	0.159	0.241	0.056
<i>Zone 3</i>				<i>Zone 9</i>			
C_1	0.654	5.212	26.54	C_1	0.403	3.215	23.11
C_2	0.183	1.056	0.251	C_2	0.184	0.665	0.221
C_3	0.169	0.195	0.069	C_3	0.163	0.121	0.039
<i>Zone 4</i>				<i>Zone 10</i>			
C_1	0.478	4.786	25.13	C_1	0.785	5.632	27.89
C_2	0.184	0.996	0.354	C_2	0.184	1.026	0.397
C_3	0.157	0.164	0.044	C_3	0.161	0.226	0.055
<i>Zone 5</i>				<i>Zone 11</i>			
C_1	0.612	5.478	27.14	C_1	0.635	5.124	27.56
C_2	0.179	1.047	0.321	C_2	0.183	1.023	0.378
C_3	0.164	0.212	0.065	C_3	0.162	0.225	0.074
<i>Zone 6</i>				<i>Zone 12</i>			
C_1	0.496	4.324	25.56	C_1	1.245	6.451	29.23
C_2	0.180	0.987	0.276	C_2	0.201	1.451	0.452
C_3	0.168	0.156	0.048	C_3	0.189	0.321	0.087

. The results obtained with a +1% deviation on nominal values (obtained by the gravimetric method) of the parameters are reported in Table 7.11 for all the twelve groups of vials¹⁹ considered. The values reported in Table 7.11 have been computed with the model inputs set at their nominal values, i.e. $p_c/\bar{p}_c = 1$ and $T_{fluid}/\bar{T}_{fluid} = 1$. The general trend of the sensitivities is the same independently of the values set for these inputs.

The results of Table 7.11 show that parameter C_1 has the strongest influence on all the three KPIs of interest. In particular, the total amount of water sublimated show a very strong sensitivity with respect to C_1 (values are highlighted in bold), whereas very limited sensitivities with respect to C_2 and C_3 .

Based on these results, it was decided to set the values of C_2 and C_3 at their nominal (gravimetric) values, and to perform a maximum likelihood estimation of C_1 for each zone. The estimation results are discussed in detail in section 7.4.6.

7.5.4 Step #3: sensitivity analysis for R_p

Once C_1 has been identified as the most relevant parameter to be estimated for an accurate description of the heat transfer coefficient, an estimation of the mass transfer resistance R_p must be obtained in order to complete the preliminary model calibration activity.

The resistance to the vapor flow during ice sublimation is the sum of three main contributions, namely:

1. the resistance to the vapor flow due to the dried layer, whose thickness increases as the ice sublimation progresses;
2. the resistance to the vapor flow exerted by the elastomeric stoppers that are partially inserted in the neck of the vials during the drying process. These stoppers can have different opening configurations to allow the vapor flow during drying operation and are fully inserted at the end of the overall drying cycle to guarantee product sterility;
3. the resistance to the vapor flow due to path from the drying chamber to the condenser.

Several experimental studies (Pikal, 2000; Rambhatla *et al.*, 2004; Pikal *et al.*, 2016) have shown that the first contribution accounts for more than 80% of the total mass transfer resistance, whereas the stoppers and chamber-to-condenser contributions account for no more than 3-6% each. Therefore, the latter contributions are typically neglected and the mass transfer resistance reported in Eq. (7.3) only accounts for the contribution of the growing dried layer. The dependence of R_p with respect to the thickness L_{frozen} of the dried layer is typically expressed by Eq. (7.4), and the three parameters appearing in this equation, which need to be estimated, are R_0 , R_1 and R_2 .

¹⁹ Note that the fact that C_2 and C_3 are the same for each zone does *not* mean that their sensitivities are the same, since their values are affected by the nominal value of C_1 (which is different for each zone).

R_0 is the mass transfer resistance at the beginning of the ice sublimation process and depends on the product formulation. The interpretation of this parameter has been attributed to the mass transfer resistance of the top layer of the frozen product, but a rigorous physical explanation is still missing (Rambhatla *et al.*, 2004). R_1 describes the “speed” at which the mass transfer resistance increases with the dried layer thickness. The contribution of R_2 is sometimes neglected (Pikal, 2000), and the dependence of R_p with respect to L_{frozen} is expressed as a linear relationship.

Some experimental studies (Pikal, 2000; Rambhatla *et al.*, 2004) have shown that the evolution of the mass transfer resistance during ice sublimation strongly depends on the product morphology, which is in turn affected by the freezing protocol adopted. In fact, the greater the nucleation temperature and degree of supercooling during the freezing step, the greater the size of ice crystals and dried pores inside the product, the smaller the mass transfer resistance to the vapor flow during ice sublimation. Different techniques (e.g., controlled ice nucleation and introduction of an annealing step during freezing) can be used to increase the dimensions of the ice crystals and therefore to reduce the value of R_p during primary drying.

The estimation of the parameters R_0 , R_1 and R_2 is typically performed with *ad-hoc* experimental techniques that all suffer from several limitations. The typical procedure consists of two sequential steps.

1. First, the mass transfer resistance R_p is determined experimentally using the relationship:

$$R_p = \frac{A_v(p_i - p_c)}{J_w} \quad (7.20)$$

The sublimation flux J_w can be obtained experimentally by inserting a weighing device inside the drying chamber (Fissore *et al.*, 2010), by performing a pressure rise test (Fissore *et al.*, 2010), or by using a non-invasive spectroscopic method called tunable diode laser absorption spectroscopy (Cassidy and Reid, 1982). However, the experimental values of J_w (hence of R_p) obtained with these methods are only “local” values (i.e., valid for a single vial in the case of the weighing device) or “global” values (i.e., a single value of the mass transfer resistance is obtained for *all* the vials inside the drying chamber, in the case of pressure rise test and tunable diode laser absorption spectroscopy). Additionally, these values do not account for the different mass transfer behavior of the vials inside the drying chamber. Note also that the application of Eq. (7.20) also requires measurements of the temperature at the bottom of the vials, since temperature T_i at the sublimation interface can be inferred from them according to Eq. (7.8); after that, also p_i can be inferred, according to Eq. (7.5). As will be shown in section 7.6, this is a serious (and still not resolved) limitation of these methods, since

the estimation of the mass transfer resistance depend on the estimation of the heat transfer coefficient. A wrong estimation of K_v (which is typically obtained with the gravimetric method described in section 7.4.1) would result in a wrong estimation of R_p , and, most importantly, the (possible) presence of a high correlation between those two parameters would completely invalidate the ability of these experimental methods (gravimetric for K_v together with experimental identification of R_p) to correctly identify the model parameters.

- Parameters R_0 , R_1 and R_2 are then estimated by looking for the best fit between a reference curve and the experimental values of R_p vs. the dried layer thickness L_{dried} that can be obtained with one of the experimental strategies described above.

Despite the aforementioned experimental techniques, which all suffer from some limitations, the parameters of the mass transfer resistance R_p equation can also be estimated directly with the maximum likelihood estimator once measurements of vial bottom temperature and total amount of water sublimated are available. In other words, R_0 , R_1 and R_2 can be estimated together with the heat transfer parameter C_1 at step #4 (Fig. 7.5) of the proposed methodology. For the same reasons discussed in the previous section, it is always desirable to study if some of these parameters can be set to their nominal values without affecting the overall model performance, so as to facilitate the estimation. To this purpose, a sensitivity analysis for the system KPIs with respect to R_0 , R_1 and R_2 was performed.

The nominal values of these parameters are reported in the literature for the same formulation (5% sucrose-based placebo) that was used in the historical experiments. In this study, the values reported by Mortier *et al.* (2016) were used as nominal values for R_0 , R_1 and R_2 . Note that these values are assumed to be the same for all the vials, and therefore sensitivity analysis can only be performed by assuming a “global” behavior of the mass transfer resistance (i.e., “global” nominal values are available for R_p , whereas “local” nominal values were available for K_v thanks to the gravimetric measurements). However, in this work, R_p was not assumed to have the same values for all zones: instead, a value of R_p was determined for each zone through the maximum likelihood estimator. As a consequence, since the results of the sensitivity analysis are only available on a “global” level, it is implicitly assumed that the sensitivities of the KPIs with respect to the three mass transfer parameters are not significantly affected by their “local” nominal values.

The results of the sensitivity analysis are shown in Table 7.12. These values refer to the same drying recipe used for the sensitivity analysis on K_v and were obtained by keeping K_v constant and equal to the gravimetric value.

Table 7.12. Sensitivities for the vial bottom temperature, sublimation flux and cumulative amount of water sublimated with respect to the three parameters R_0, R_1, R_2 . The values refer to a +1% deviation on the nominal value of the parameters and sensitivities are expressed as % of the nominal value of the KPI considered. The strongest sensitivity value is shown in boldface.

Parameter	$ s_{T,i}^{max} [\%]$	$ s_{w,i}^{max} [\%]$	$ s_{M,i} [\%]$
R_0	0.020	0.009	0.000
R_1	0.076	1.321	0.199
R_2	0.024	0.000	0.000

From the results shown in Table 7.12, it can be concluded that:

- R_2 has a negligible impact on all the KPIs of interest and can therefore be neglected (as confirmed by several experimental studies);
- R_1 has a stronger influence on the KPIs of interest than R_0 , especially on the sublimation flux and the total amount of water sublimated (as one would expect). Therefore, accurate estimation of this parameter is crucial to obtain a reliable prediction of the sublimation flux and duration of the ice sublimation.

In view of the above, parameter R_0 was set to its nominal value, and R_1 was treated as an additional model parameter to be preliminarily estimated. Note that the nominal value of R_0 depends on the product formulation; therefore, if its value cannot be found in the literature for the product under investigation, this parameter needs to be calibrated together with C_1 and R_1 at the next step of the proposed methodology.

7.4.5 Step #4: parameter estimation

Based on the analyses performed in the previous steps, the maximum likelihood estimation of the parameters of model (7.1)-(7.10) is reduced to a joint estimation of:

- the heat transfer parameter C_1 , with initial guess set as the corresponding gravimetric value determined during step #1;
- the mass transfer parameter R_1 , with initial guess set as the corresponding literature value.

The two parameters C_1, R_1 must be estimated for each one of the 12 zones, with the other parameters C_2, C_3, R_0, R_2 set at their nominal values.

The estimation activity has to be repeated for dataset B (non-siliconized vials) as well as dataset C (siliconized vials). The results of the parameter estimation activity for all the 12 zones considered is reported in Table 7.13 for dataset B. The results for dataset C can be found in appendix A. Note that the values of the parameters are reported as % of their nominal (i.e. initial guess) values.

Table 7.13. Maximum likelihood parameter estimation for the 12 zones considered. Dataset B (non-siliconized vials).

Parameter	% nominal value	95% t-value	Reference t-value	Parameter	% nominal value	95% t-value	Reference t-value
<i>Zone 1</i>				<i>Zone 7</i>			
C_1	103.65	74.41	1.65	C_1	103.12	78.12	1.65
R_1	104.13	63.58	1.65	R_1	103.89	66.32	1.65
<i>Zone 2</i>				<i>Zone 8</i>			
C_1	109.14	75.87	1.65	C_1	106.82	77.52	1.65
R_1	102.12	66.41	1.65	R_1	104.82	65.56	1.65
<i>Zone 3</i>				<i>Zone 9</i>			
C_1	102.45	77.89	1.65	C_1	107.74	74.52	1.65
R_1	105.69	69.45	1.65	R_1	106.16	63.56	1.65
<i>Zone 4</i>				<i>Zone 10</i>			
C_1	106.23	67.89	1.65	C_1	104.56	78.98	1.65
R_1	109.12	62.12	1.65	R_1	102.12	67.85	1.65
<i>Zone 5</i>				<i>Zone 11</i>			
C_1	112.25	70.25	1.65	C_1	102.31	75.65	1.65
R_1	103.64	64.12	1.65	R_1	106.14	66.32	1.65
<i>Zone 6</i>				<i>Zone 12</i>			
C_1	121.53	72.14	1.65	C_1	108.12	71.21	1.65
R_1	114.73	65.9	1.65	R_1	114.56	59.87	1.65

Table 7.13 shows that the two parameters can be identified by the maximum-likelihood estimator with high precision²⁰, since their 95% t-values are significantly greater than the reference t-value. The good agreement between the experimental data and the model prediction is shown in Fig. 7.10 for $p_c/\bar{p}_c = 0.667$. Similar results have been obtained for the other values of chamber pressure. Although these plots refer to the model calibrated for zone # 6 of the batch of vials, similar results have been obtained for all the other 11 zones.

²⁰ Note that a high precision on the parameter estimation does not necessarily mean a high accuracy on the estimated values.

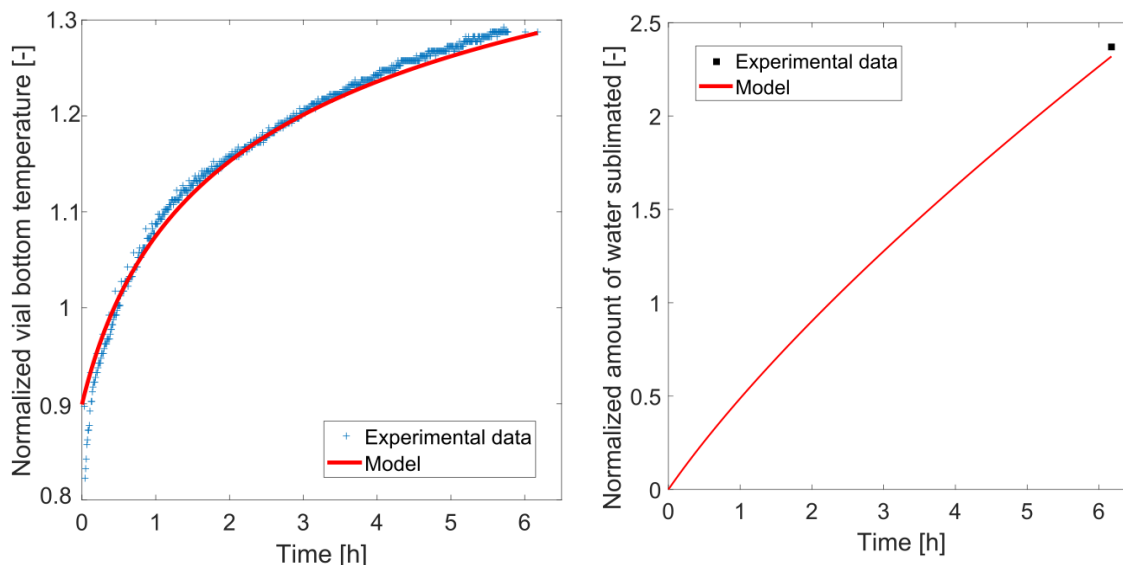


Figure 7.10. Comparison between experimental observations (light blue markers) and model predictions (red lines) for zone 6 and $\frac{p}{\bar{p}_c} = 0.667$.

A good fit of the historical experimental data is obtained with the available model. Particularly, the prediction of the end-point value of the cumulative amount of water sublimated is in good agreement with the experimental value (the model prediction is 0.87% smaller than the averaged experimental value). Moreover, the model is slightly underestimating this variable, which provides a conservative prediction of the total time needed to complete the ice sublimation. The state variables profiles predicted by the calibrated model in the zone 6 of the shelf are shown in Fig. 7.11.

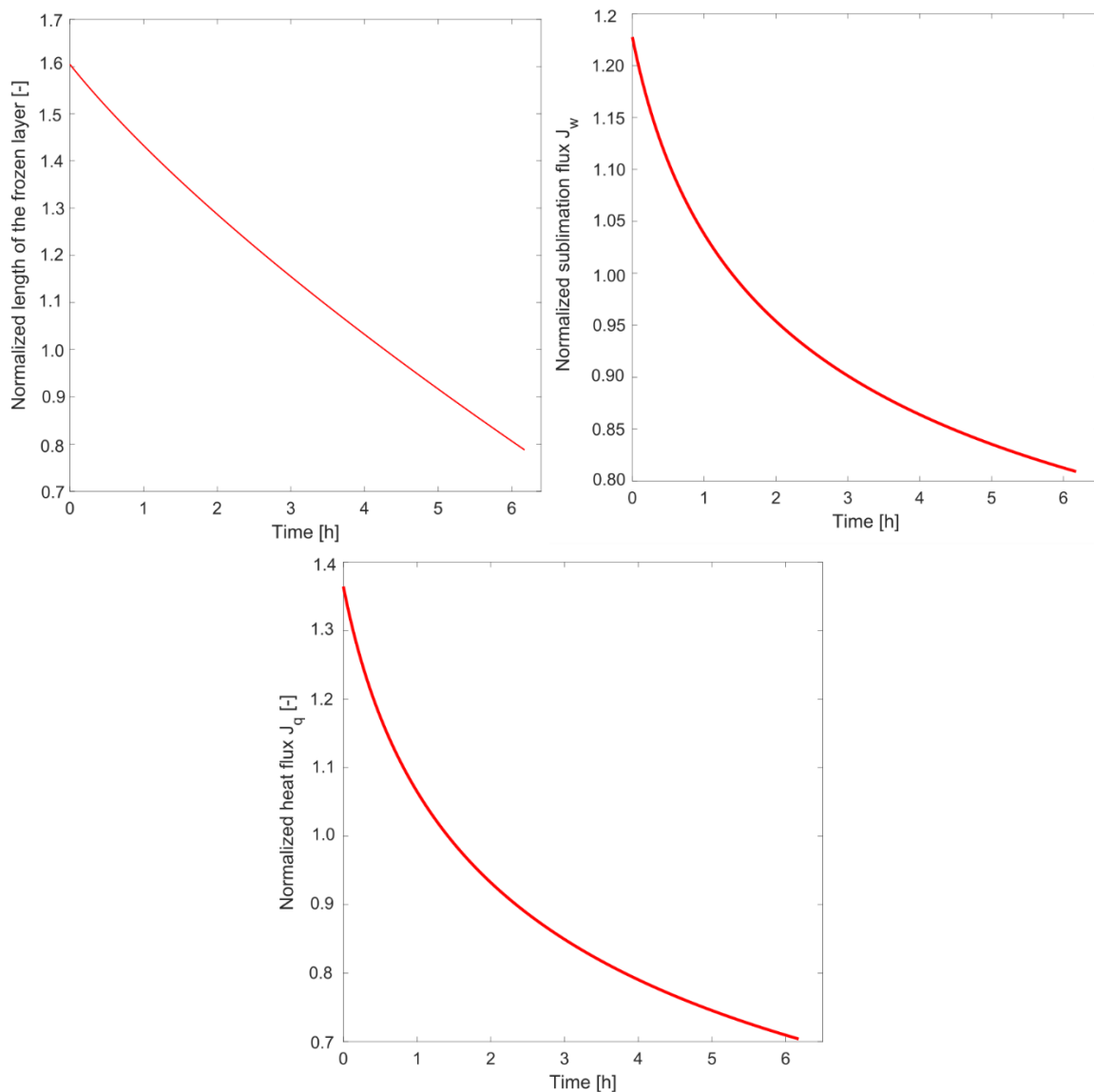


Figure 7.11. Time profiles of the (a) length of the frozen layer, (b) sublimation flux, and (c) heat flux as predicted by the calibrated model for zone #6.

7.4.6 Step #5: model validation

Once the model has been preliminarily calibrated, the model performance can be assessed using a dataset that was *not* used for model calibration. Two validation datasets have been used to this purpose, namely dataset D for non-siliconized vials, and dataset E for siliconized vials.

The ability of the model to give a reliable estimation of the KPIs of interest has been investigated. A slight complication to this model validation activity is due to a difference in the location of the thermocouples between the calibration and validation datasets, as well as the different loading between the calibration and validation datasets that could impact on the temperature dynamics.

For example, for non-siliconized vials, the validation dataset D collects measurements of product temperature for two vials placed in the middle of the first and last row of the first half of the shelf. According to Fig. 7.7, these locations should correspond, in principle, to zone 12 and zone 7 (respectively) of the calibration dataset. Therefore, the model calibrated for these two zones should be used to challenge the model predictions against the experimental observations. However, since the position of the thermocouple itself is different between the two datasets, the models calibrated for zones adjacent to zones 12 and 7 have also been used to test the performance. Note that the only measurement available in the validation datasets is the vial bottom temperature, and the accuracy of the model prediction can therefore be assessed only for this variable.

The comparison between the experimental values for the two temperature measurements of dataset D and the relevant model predictions is shown in Fig. 7.12.

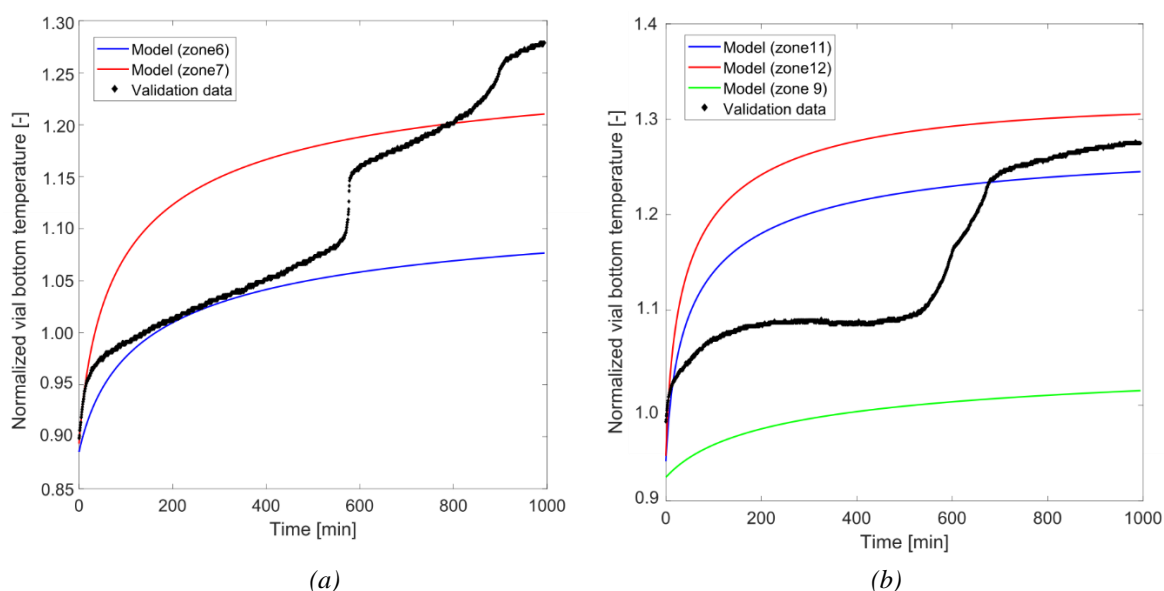


Figure 7.12. Model validation for (a) thermocouple #5 and (b) thermocouple #6 of the validation dataset D.

The first important remark to be made is on the experimental profile of the vial bottom temperature. This profile (black dots) is characterized by an abrupt inflection at around half of the total duration of primary drying, due to the transition of the frozen layer through the bottom of the vial. In principle, if the thermocouple touches the bottom of the vial, the inflection should correspond to the end of the ice sublimation process. Starting from this instant, the heating flow from the shelf is used to significantly increase the product temperature, but it is still not large enough to promote a significant water desorption (Pikal, 2000). From a practical point of view, it is difficult to guarantee a complete adherence between the thermocouple and the vial bottom, and the very presence of the thermocouple affects the vial behavior during ice sublimation. For

these reasons, the inflection on the product temperature profile can only be considered as a qualitative indication that the sublimation process is close to the end.

The mono-dimensional model (7.1)-(7.10) only accounts for the ice sublimation process, and is therefore not suitable to describe the vial behavior when ice sublimation ends. Therefore, the behavior of the product temperature after the inflection cannot be captured by the available model. This behavior can be described by introducing a desorption model in the original formulation (Fissore *et al.*, 2012).

Due to the different geometry configurations between the calibration and validation datasets, it is not possible to describe the experimental profile with the model calibrated for a specific zone. However, Fig 7.12 shows that the model predictions of zones adjacent to the one where a temperature measurement is available are able to bracket the relevant behavior of the system during ice sublimation²¹. A better representation of the system could be obtained by developing an integrated multi-vial model accounting for the different thermal behaviors between the vials. Such a model is still an open research area that will be investigated in the future.

Summarizing, the following conclusions to this preliminary analysis can be drawn:

1. the model (7.1)-(7.10) is structurally consistent with the experimental observations as ice sublimation progresses. However, the different behaviors of the vials placed at different locations on the shelf cannot be captured by this model. The model can only be used to capture the behavior of each zone (assuming that all the vials of the same zone behave in the same way), and therefore its applicability to different configurations is limited.
2. In the experimental conditions used for this study, the number of parameters that need to be identified by MBD_{oE} is only two: one involving the heat transfer coefficient (C_1), the other involving the mass transfer resistance (R_1). These parameters can be identified from historical data with high precision (high t-values), but their accuracy is not guaranteed. However, if a different freeze-dryer geometry configuration, different vial, different formulation or different loading system (e.g., in frames or not) is used, optimally designed experiments are still needed to maximize their information in order to assist the parameter estimation activity.

The latter issue is thoroughly discussed in the next section.

7.5 Model-based design of experiments for primary drying

The preliminary analysis described above showed that the two most influential model parameters C_1 and R_1 can be estimated from historical data with good precision. However, their values are affected by factors (such as system geometry and type of formulation) that are not

²¹ In absolute values, the prediction of the vial bottom temperature is bracketed within a $\pm 2^\circ\text{C}$ accuracy by the model predictions.

taken into account explicitly in the model structure. Therefore, a new estimation of these parameters is required whenever the system configuration changes or a different formulate is used. In order to maximize the information that can be obtained from a new identification experiment, MBDoe techniques can be used to design optimal experiments based on the mono-dimensional model (7.1)-(7.10).

The objective of this section is to design a new experiment in order to maximize the precision that can be obtained for the estimation of C_1 and R_1 . The same system configuration as the previous preliminary analysis will be kept, and a comparison between the precision on the estimated values of C_1 and R_1 that can be obtained with or without the new designed experiment will be discussed. Note that, in practical applications where both the system configuration and the type of formulation change, a reasonable estimate of the other model parameters (C_2, C_3, R_0, R_2) is required. These estimates can be obtained with the same considerations described in the previous section, with the only difference that the value of R_0 may not be available in the literature. In this scenario, R_0 should be taken into account in the MBDoe activity as an additional parameter that must be estimated.

7.5.1 Problem statement

The design of a new optimal experiment, based on the model (7.1)-(7.10), requires the following information:

1. initial conditions of the system, namely initial thickness of the frozen layer;
2. assignment of the control variables whose profile need to be optimized, namely total chamber pressure and heating fluid temperature;
3. time profile of the control (i.e., manipulated) inputs, which can be assumed to be time-invariant, piecewise constant or piecewise linear. If a time-varying profile is chosen, the number of intervals and the upper and lower bounds of the control variables for each interval must be specified. For the system considered, a piecewise linear behavior for T_{fluid} and a piecewise constant behavior for p_c have been chosen;
4. assignment of the measurements that have to be collected and sampling frequency. For the system under investigation, the only variable that can be measured is the vial bottom temperature T_B at specific locations within the freeze dryer, and a sampling interval of 30 s was set;
5. constraints on the state variables that need to be satisfied over the entire experiment or at the end of the experiment. As discussed in section 7.1.2, the two constraints that need to be satisfied during primary drying are (i) product temperature below the glass transition/eutectic temperature of the product, and (ii) sublimation flux lower than the maximum allowable flux to avoid choked flow. Both constraints were assumed to be time-invariant interior-point constraints to be satisfied over the entire drying operation;

6. model parameters that need to be estimated from the designed experiment, namely C_1 and A .

The results of the MBDoe activity are:

1. total duration of the experiment. Note that this variable can also be set to a fixed value depending on the type of experiment desired;
2. length of the switching intervals for the manipulated inputs;
3. profile (for T_{fluid}) or value (for p_c) of the control inputs within each switching interval.

The MBDoe activity was performed based on the following assumptions and operational constraints for the experimental equipment:

- a. the model originally calibrated for zone #6 (central vials) was used as the reference model for the MBDoe. This is because central vials represent the majority of the vials within the freeze dryer, and therefore the experiment is optimized based on the behavior of the majority of the vials within the system. Note that this may cause the edge vials to not fulfill the constraint on the maximum allowable product temperature;
- b. a maximum normalized product temperature $\bar{T}_B^{max} = 1.32$ and a maximum normalized sublimation flux $\bar{J}_w^{max} = 1.26$ were set based on the type of formulation and machine considered.
- c. the maximum allowable number of intervals was set equal to 13 due to limitations in the recipe implementation on the machine.

The results of the MBDoe activity are reported in the next section.

7.5.2 Designed experiment: “in silico” results

The experiment was designed using the D-optimal criterion (section 7.3) and the total duration of the experiment was set to 990 min (16.5 h). The optimal profiles of the control inputs are shown graphically in Fig. 7.13 and numerically in Table 7.14.

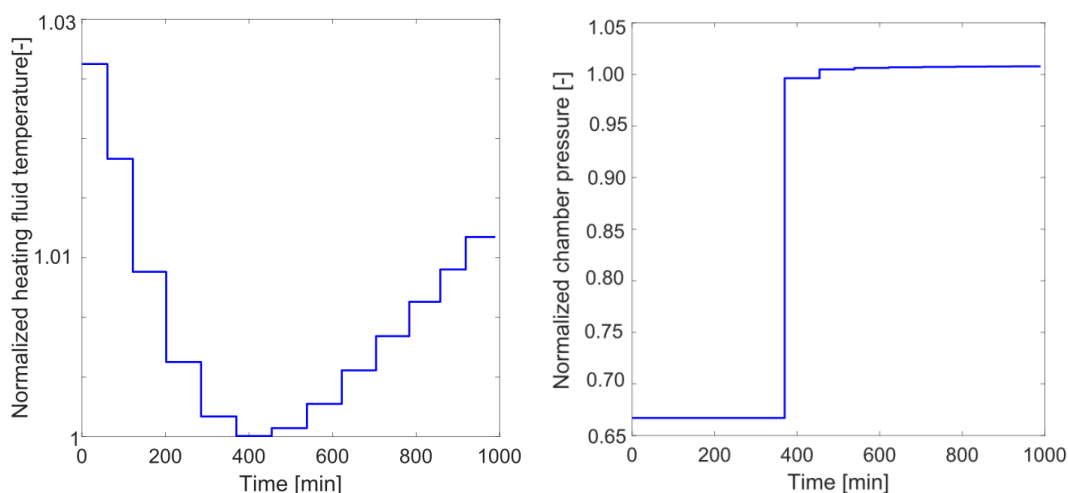


Figure 7.13. Control input profiles for the optimal experiment designed using MBDoe.

Notice that, although a piecewise linear profile was looked for, the heating fluid temperature is kept constant in each interval, and is characterized by an initial decrease followed by a slight increase of temperature towards the end of the experiment. With respect to the chamber pressure, only one meaningful switch is required from ~ 50 mT to 75 mT after ~ 380 min.

The designed experiment can be simulated *in silico* in order to study how the precision in the estimation of parameters C_1 and A improves with respect to the historical (non-designed) experiments.

The results of the parameter estimation activity using the *in-silico* designed experiment are summarized in Table 7.15. For the designed experiment, two scenarios are analyzed, namely one where both the designed experiment and the historical datasets are used to estimate the parameters, and one where *only* the designed experiment is used to this purpose. It is worth noticing that whereas the designed experiment collects measurements of vial bottom temperature only, the historical data collect measurements of vial bottom temperature *and* cumulative amount of water sublimated. As can be seen, in both designed scenarios the precision of the estimation increases (much larger t-values are obtained) than when only historical experiments are used.

Table 7.14. Numerical values of the control inputs for each interval of the optimally designed experiment of Figure 7.13.

Interval #	Interval duration [min]	Normalized heating fluid temperature [-]	Normalized chamber pressure [-]
1	61.0	1.025	0.667
2	61.0	1.019	0.667
3	79.8	1.011	0.667
4	83.5	1.005	0.667
5	84.5	1.001	0.667
6	84.6	1	0.667
7	84.3	1.001	0.667
8	83.5	1.002	1
9	82.0	1.004	1
10	79.5	1.007	1
11	74.3	1.009	1
12	60.9	1.011	1
13	70.7	1.014	1

Table 7.15. Final values and estimate precision for parameters C_1 and A using historical and optimal (*in silico*) experiments

Parameter	95% t-value (only historical data)	95% t-value (historical data + designed experiment)	95% t-value (only designed experiment)	Reference t-value
C_1	73.62	111.7	81.12	1.645
A	59.6	72.6	63.21	1.645

It is worth noticing, however, that the two parameters are strongly correlated between, as confirmed by the following correlation matrices (index 1 = C_1 , index 2 = R_1):

$$\text{Only historical data:} \quad \begin{pmatrix} 1 & -0.932 \\ -0.932 & 1 \end{pmatrix}, \quad (7.21)$$

$$\text{Historical data+ designed exp.:} \quad \begin{pmatrix} 1 & -0.993 \\ -0.993 & 1 \end{pmatrix}, \quad (7.22)$$

$$\text{Only designed exp.:} \quad \begin{pmatrix} 1 & -0.974 \\ -0.974 & 1 \end{pmatrix}. \quad (7.23)$$

This may have a potential effect on the *accuracy* of the estimates of the two parameters, since their values are jointly dependent on each other. This issue is still under investigation and is not covered in this Dissertation

7.6 Final remarks and future work

In this Chapter, a methodology to design a new experiment for the primary drying stage of a pharmaceutical freeze dryer has been proposed.

First, a step-by-step approach has been developed in order to assess the structural consistency and the most influential parameters of a mono-dimensional model available in the literature. The methodology allowed to reduce the number of parameters to be estimated to only two parameters, related to heat and mass transfer. Moreover, a systematic approach to identify the nominal values of the other parameters has been proposed.

Secondly, a new optimal experiment has been designed using MBDoE according to a D-optimal criterion. The designed experiment has been simulated *in silico* and promising results have been obtained for the parameter estimation activity, with a significant increase of the precision of the parameter estimates over the one that could be achieved using historical (non-designed) experiments. A significant correlation between the two parameters has also been identified at this stage.

Future work involves the implementation of the designed experiment on the real equipment. Moreover, the issue related to the high correlation of the two parameters will be tackled by considering re-parameterization techniques and different parameter estimation activities. Finally, further improvements to the available model will be implemented in order to capture the behavior of the partial pressure of water in the chamber in a more detailed fashion.

Conclusions and future perspectives

The description of the design space (DS) of a new pharmaceutical product is one of the most important activities that has been introduced within the *Quality by Design* (QbD) framework. An approved DS can represent a competitive advantage for the applicant company, since the process can be run anywhere within the DS without requiring any additional regulatory approvals. This translates into an enhanced process flexibility, which in turn translates into an increase of profitability.

Different activities related to DS description are involved during the lifecycle of a pharmaceutical product. In pharmaceutical development, design space determination and uncertainty quantification (i.e., quantification of “assurance” of quality for the final product) represent two fundamental steps that companies should implement in order to prepare a submission for a new drug approval to the relevant regulatory bodies (e.g. Food and Drug Administration, FDA). During the technology transfer phase, design space *scale-up* and design space *scale-out* allows transferring the DS obtained at the laboratory scale to the commercial plant scale, or across different manufacturing sites. Finally, during commercial manufacturing, design space *maintenance* allows obtaining a continuous update of the DS as more information about the manufacturing process is gained through plant operation.

Model-based approaches play a key role to tackle all the aforementioned activities related to DS description. Models can be used to link the raw material properties and critical process parameters of the manufacturing process to the critical quality attributes of the final product, thus allowing determination of the DS of the product. Moreover, models can be used to assist technology transfer activities related to DS description, as well as to assist DS maintenance during plant operation.

In this Dissertation, advanced modeling strategies for the description and maintenance of the design space of new (or existing) pharmaceutical products have been developed. The proposed approaches aim at utterly fulfilling the definitions and requirements defined by the regulatory documents, in terms of both *demonstration* of the DS (i.e., rigorous scientific explanation of the methodology adopted for DS determination) and of quantification of *assurance* of quality (i.e., probability that the product will meet its quality specifications). Specifically, the following four different research areas were investigated:

1. development of rigorous methodologies and metrics for the quantification of the assurance of quality for the final product;

2. development of computationally feasible and user-friendly (to be intended as ease of interpreting the results) methodologies for design space description;
3. development of a methodology to adapt and maintain the model-based representation of a DS as process operation progresses;
4. exploitation of model-based design of experiments (MBDoE) techniques for the identification of first-principles models to be used for DS description.

Table C.1 summarizes the main achievements of this Dissertation, with indication of example applications, type of data used, and references where the results have been discussed.

In Chapter 1, a thorough review of the methodologies that have been proposed in the literature to assist the description of a design space during product/process development, technology transfer and commercial manufacturing has been given.

With reference to the issue of **the quantification of “assurance” of quality**, in Chapter 3 and Chapter 4 two methodologies have been proposed. The methodology presented in Chapter 3 is derived with concepts of frequentist statistics and suited for situations where a projection onto latent structures (PLS) model is to be used for DS description. The methodology presented in Chapter 4 is built upon concepts of Bayesian statistics and jointly exploits PLS modeling and Bayesian multivariate linear regression model to obtain a probabilistic representation of the design space.

In more detail, in Chapter 3 two frequentist models were developed to back-propagate the uncertainty from the space of product quality (*output space* from a modeling perspective) to the space of raw material properties and process parameters (*input space*) of a pharmaceutical process. The backpropagation models were obtained for situations where a PLS model (built on historical manufacturing data of products “similar” to the one under development) is used to relate the raw material properties and process parameters to the product quality characteristics. It was shown how PLS model inversion, together with uncertainty back-propagation, allows identifying a restricted portion of the original input domain (called *experiment space*) where the design space of the new product is expected to lie with a given degree of confidence. The experiment space can then be used to tailor an experimental campaign within this restricted set of input combinations, with advantages in terms of reduction of the duration and cost of the experimental campaign. The effectiveness of the methodology was tested on five case studies involving typical unit operations of the pharmaceutical industry, such as dry/wet granulation and roll compaction, involving both simulated and real experimental data. For all the case studies considered, the predicted experiment space effectively bracketed the actual design space of the pharmaceutical product under investigation. In Chapter 4, a methodology was proposed to exploit Bayesian multivariate regression together with PLS modeling in order to quantify the assurance of quality for the final product. The metric used to this purpose is the *probability* (to be intended according to its Bayesian interpretation) that the product will meet its quality specification, given the historical data available. The

proposed methodology combines: (i) a PLS model to reduce the dimensionality of the original input space for an easier discretization and interpretation; (ii) Markov-chain Monte-Carlo techniques to reconstruct the posterior probability distribution (PPD) of the product quality attributes for each input combination considered. The PPD obtained for all the combinations of the input domain allows identifying a probabilistic (or Bayesian) design space, defined as the set of input combinations for which the probability that the product will meet its specifications is greater than an assigned threshold. The DS obtained according to this approach is fully consistent with the regulatory requirements, since it contains a *quantitative* metric of the risk that the product may not reach its quality targets. The final outcome is a step-by-step methodology that companies may decide to use in design space submission to the regulatory agencies. The effectiveness of the proposed approach was tested on three case studies involving real experimental data (taken from the literature) of pharmaceutical unit operations.

With respect to the issue of **design space determination**, a methodology was proposed in Chapter 5 to obtain a computationally feasible and user-friendly representation of the DS. The methodology was obtained by combining two different modeling strategies: (i) a PLS model to reduce the dimensionality of the input space; (ii) surrogate-based feasibility analysis to identify the DS on the latent representation of the input space obtained with the PLS model. It was shown how a cubic radial-basis function surrogate can be used to obtain a computationally cheap approximation of the original process model, and how an adaptive sampling strategy can be used to determine the DS boundary on the latent space identified by the PLS model. The methodology effectiveness was tested on a continuous line for pharmaceutical tablet manufacturing, and the computational benefits (~80% reduction of the overall computational time) were presented together with the improvements in terms of ease of DS representation (a single two-dimensional latent plot can be used to represent a six-dimensional design space). Moreover, the robustness of the methodology was tested with other strongly nonlinear mathematical examples.

The problem of **design space maintenance** during plant operation was discussed in Chapter 6, where a methodology was proposed to obtain a continuous adaptive refinement of the representation of the DS that can be obtained using a first-principles model. Given the available measurements coming from plant sensors, the methodology jointly exploits (i) a dynamic state estimator and (ii) surrogate-based feasibility analysis to perform continuous adaption of the DS as process operation progresses. It was shown how the state estimator can be used to adjust in real time a small subset of the model parameters in order to compensate for process/model mismatch, and obtain an up-to-date representation of the state of the system. Feasibility analysis were then used to identify the DS boundary based on the up-to-date model returned by the state estimator

Table C.1. Summary of the main achievements of this Dissertation, with indication of their relevant applications, the original of the data used and the related references.

Chapter	Main achievement	Application	Data origin	Reference
Chapter 1	Review on methodologies for design space description	(-)	(-)	Bano, G., Facco, P., Bezzo, F., Barolo M. (2018) Design space description in pharmaceutical development and manufacturing: a review. <i>In preparation.</i>
Chapter 3	Development of two uncertainty back-propagation models for PLS model inversion in pharmaceutical product development	(-) <ul style="list-style-type: none"> • Roll compaction • High-shear wet granulation • Dry granulation 	Simulated and laboratory	<p>Bano, G., Meneghetti, N., Facco, P., Bezzo, F., Barolo M. (2017) Uncertainty back-propagation in PLS model inversion for design space determination in pharmaceutical product development. <i>Comput. Chem. Eng.</i> 101, 110-124.</p> <p>Facco, P., Bano, G., Bezzo, F., Barolo M. (2017) Multivariate statistical approaches to aid pharmaceutical processes: product development and process monitoring in a Quality-by-design perspective. EuroPACT 2017. May 10-12, Potsdam (Germany).</p> <p>Bano, G., Meneghetti, N., Facco, P., Bezzo, F., Barolo, M. (2017) Determinazione dello spazio di progetto di un nuovo prodotto farmaceutico: caratterizzazione dell'incertezza. <i>Convegno GRICU 2017: gli orizzonti dell'ingegneria chimica</i>. September 24-29, Anacapri (NA; Italy).</p>
Chapter 4	Joint Bayesian/latent variable approach for the quantification of the “assurance” of quality of a new pharmaceutical product	<ul style="list-style-type: none"> • Roll compaction • High-shear wet granulation • Dry granulation 	Simulated and laboratory	<p>Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018) Probabilistic design space determination in pharmaceutical development: a Bayesian/latent variable approach. <i>AIChE J.</i> 64, 2438-2449.</p> <p>Bano, G., Facco, P., Bezzo, F., Barolo, M (2017) Handling parametric and measurement uncertainty for design space determination in pharmaceutical product development: a Bayesian approach. Presented</p>

				at: <i>World Congress of Chemical Engineering WCCE10</i> . October 1-5, Barcelona (Spain).
<i>Chapter 5</i>	General methodology to describe the design space of a pharmaceutical product through a joint PLS/feasibility analysis approach	<ul style="list-style-type: none"> • Continuous direct compaction manufacturing line 	Simulated	<p>Bano, G., Wang, Z., Facco, P., Bezzo, F., Barolo, M., Ierapetritou, M. (2018) A novel and systematic approach to identify the design space of pharmaceutical processes. <i>Comput. Chem. Eng.</i> 115, 309-322.</p> <p>Bano, G., Z. Wang, P. Facco, F. Bezzo, M. Barolo, M. Ierapetritou (2018). Dimensionality reduction in feasibility analysis by latent variable modeling. In: <i>Computer-Aided Chemical Engineering 44, Proc. of the 13th International Symposium on Process Systems Engineering – PSE 2018, July 1-5 2018</i> (M.R. Eden, M.G. Ierapetritou, G.P. Towler, Eds.), Elsevier, Amsterdam (The Netherlands), 1477-1482.</p>
<i>Chapter 6</i>	Systematic methodology for online maintenance of the design space	<ul style="list-style-type: none"> • Dry granulation • Penicillin fermentation 	Simulated	Bano, G., Facco, P., Ierapetritou, M., Bezzo, F., Barolo, M. (2018) Design space maintenance by online model adaptation in pharmaceutical manufacturing. Submitted to: <i>Comput. Chem. Eng.</i>
<i>Chapter 7</i>	Design of optimal experiments for pharmaceutical freeze drying	<ul style="list-style-type: none"> • Pharmaceutical freeze drying 	Industrial	(-)

The methodology was tested with two case studies, the former involving the granulation of a pharmaceutical formulate, the latter involving the fermentation of penicillin in pilot scale bioreactor. In both situations, the ability of the proposed approach to timely track the DS was proven. The methodology was tested on two case studies, the former involving the granulation of a pharmaceutical formulate, the latter involving the fermentation of penicillin in pilot scale bioreactor. In both situations, the ability of the proposed approach to continuously track the DS was proven.

With respect to the issue of the **design of optimal experiments** for the identification of first-principles models to be used for DS description, in Chapter 7 it was shown how model-based design of experiments (MBD_{oE}) techniques can be exploited to reach this target. MBD_{oE} was used to design a new experiment in order to extract the maximum amount of information for the identification of the most influential parameters of a mechanistic model of an experimental pharmaceutical freeze dryer. A preliminary analysis was first performed based on historical data in order to assess the structural consistency of the model and to determine its most influential parameters (with respect to the key performance indicators of the system). A new experiment was then designed using MBD_{oE} in order to extract the maximum amount of information for the identification of the two most influential model parameters. Significant improvements were obtained in terms of precision in the parameter estimates, both *in silico* and in the experiment performed on the real equipment.

Some **areas of further investigation** can be discussed at this point.

- Validation of the methodologies presented in Chapter 3 through Chapter 6 was carried out using laboratory or simulated data. It would be very useful if industrial data could be used instead. Chapter 7 has started considering the use of industrial datasets, but more work needs to be done in this direction.
- The methodology proposed in Chapter 3 could be extended to situations where multivariate quality specifications are imposed on the product quality. The current methodology can only handle univariate quality specifications, due to a lack of proper prediction uncertainty models for PLS. The identification of such models is still an open research area that could be investigated and integrated with the uncertainty back-propagation approach presented in Chapter 3.
- The methodology presented in Chapter 4 could be extended to nonlinear models (data-driven or mechanistic) in order to handle strong process nonlinearities, where the proposed linear approach is deemed to fail in most practical situations. The application of Bayesian methodologies on complex nonlinear models poses some (still unresolved) issues in terms of both computational demand and choice of the prior distributions of

the parameters. The solution to these problems could be beneficial in terms of probabilistic DS description for complex process models.

- The methodology presented in Chapter 5 could be extended to surrogates other than the cubic radial basis function one, and a comparison between the performance of the different surrogates could be performed. Moreover, there is still a need to develop new metrics to link the projection error of the PLS model to the accuracy metrics of the surrogate-based feasibility analysis step.
- Areas of future research for the methodology presented in Chapter 6 could be the extension to other state estimation algorithms (e.g., particle filters), as well as the development of more sophisticated methodologies to identify the parameters to be adjusted in real time during plant operation. Moreover, a validation with real plant data could be beneficial to test the robustness of the methodology.
- The work presented in Chapter 7 is still ongoing and different areas of investigation are currently being considered. First, the possibility of performing some improvements (from a structural point of view) to the mechanistic model are being studied. These improvements include a better representation of the dynamics of water partial pressure within the drying chamber, as well as a better representation of the drying heterogeneity within the vials. With respect to the MBD_{oE} activity, other experiments considering different types of operational constraints and different experimental settings are currently being designed.

Appendix A

Penicillin fermentation model

This Appendix collects a detailed description of the penicillin fermentation model proposed by Birol *et al.* (2002) that has been used in Chapter 6 to simulate the pilot plant behavior. The relevant model equations, as well as all the relevant model parameters, are reported.

A.1 Model equations

The model equations are reported below. These equations were implemented in MATLAB™ 2017b.

Biomass growth:

$$\frac{dB}{dt} = \mu B - \left(\frac{B}{S_F V}\right) U \quad (\text{A.1})$$

with:

$$\mu = \left[\frac{\mu_x}{1 + \frac{K_1}{[H^+]} + \frac{K_2}{[H^+]}} \right] \frac{S}{K_b B + S} \frac{C_L}{K_{Ox} B + C_L} \left\{ \left[k_g \exp\left(-\frac{E_g}{RT}\right) \right] - \left[k_d \exp\left(-\frac{E_d}{RT}\right) \right] \right\}. \quad (\text{A.2})$$

Substrate utilization:

$$\frac{dS}{dt} = -\mu \left(\frac{B}{Y_{B/S}}\right) - \mu_{pp} \left(\frac{B}{Y_{P/S}}\right) - m_B B + \frac{US_F}{V} - \frac{S}{V} \frac{dV}{dt} \quad (\text{A.3})$$

Effect of pH:

$$\frac{d[H^+]}{dt} = \gamma \left(\mu B - \frac{UB}{V}\right) + \frac{\left[\frac{-\delta + \sqrt{(\delta^2 + 4 \cdot 10^{-14})}}{2} - [H^+] \right]}{\Delta t} \quad (\text{A.4})$$

with:

$$\delta = \left[\frac{10^{-14}}{[H^+]} - [H^+] \right] V - \frac{C_{a/b}(F_a + F_b)\Delta t}{V + (F_a + F_b)\Delta t}. \quad (\text{A.5})$$

Dissolved oxygen material balance:

$$\frac{dC_L}{dt} = \frac{\mu}{Y_{B/o}} B - \frac{\mu_{pp}}{Y_{P/o}} B - m_o B + K_{la}(C_L^* - C_L) - \frac{C_L}{V} \frac{dV}{dt} \quad (\text{A.6})$$

with

$$K_{la} = \alpha \sqrt{f_g} \left(\frac{P_w}{V} \right)^\beta. \quad (\text{A.7})$$

Penicillin production:

$$\frac{dP}{dt} = \mu_{pp} B - KP - \frac{P}{V} \frac{dV}{dt} \quad (\text{A.8})$$

with

$$\mu_{pp} = \mu_{pp}^{max} \frac{S}{K_p + S + \frac{S^2}{K_1}} \frac{C_L}{K_{op} B + C_L}. \quad (\text{A.9})$$

Volume change:

$$\frac{dV}{dt} = \frac{U}{s_F} + F_{a/b} - F_{loss} \quad (\text{A.10})$$

with

$$F_{loss} = V \lambda \left(e^{\frac{5(T-T_0)}{T_v-T_0}} - 1 \right). \quad (\text{A.11})$$

Heat generation:

$$\frac{dQ_r}{dt} = r_{q1} \frac{dB}{dt} V + r_{q2} B V. \quad (\text{A.12})$$

Energy balance:

$$\frac{dT}{dt} = \frac{U}{s_F} (T_f - T) + \frac{1}{V \rho c_p} \left[Q_r - \frac{\alpha F_C^{b+1}}{F_C + \left(\frac{\alpha F_C^b}{2 \rho c_p} \right)} \right]. \quad (\text{A.13})$$

CO₂ evolution:

$$\frac{d[\text{CO}_2]}{dt} = \alpha_1 \frac{dB}{dt} + \alpha_2 B + \alpha_3. \quad (\text{A.14})$$

The state variables involved in the model (A.1)-(A.14) are collected in Table A.1. A full list of the initial conditions that have been used for the simulations presented in Chapter 6, as well as a full list of all the model parameters with their respective values are reported in Table A.2 and Table A.3 respectively.

Table A.1. List of state variables involved in the process of Birol

Symbol	Definition
B	Biomass concentration
S	Substrate concentration
P	Penicillin concentration
C_L	Dissolved oxygen concentration
$[H^+]$	Concentration of hydrogen ions
V	Volume
Q_r	Heat generation power
T	Temperature
$[CO_2]$	Concentration of CO_2 produced

Table A.2. List of initial conditions for the state variables

State variable	Initial value	Units
Biomass concentration	0.1	$[g_B/L]$
Substrate concentration	15	$[g_S/L]$
Penicillin concentration	0	$[g_P/L]$
Dissolved oxygen concentration	1.16	$[g_O/L]$
Volume	100	$[L]$
CO_2 concentration	0.5	$[mmol/L]$
Hydrogen ion concentration	$10^{-5.1}$	$[mol/L]$
Temperature	297	$[K]$
Heat generation	0	$[cal]$

Table A.3. List of parameters for the penicillin fermentation model of Birol et al. (2002).

Parameter	Definition	Units	Value
μ_x	Maximum specific biomass growth rate	[h ⁻¹]	0.092
K_x	Contois saturation constant	[g/L]	0.15
$Y_{B/S}$	Yield constant	[g _B /g _S]	0.45
$Y_{B/o}$	Yield constant	[g _B /g _o]	0.04
$Y_{P/S}$	Yield constant	[g _P /g _S]	0.90
$Y_{P/o}$	Yield constant	[g _P /g _o]	0.20
K_1	Constant for μ	[mol/L]	10 ⁻¹⁰
K_2	Constant for μ #2	[mol/L]	7 × 10 ⁻⁵
m_x	Maintenance coefficient on substrate	[h ⁻¹]	0.014
m_o	Maintenance coefficient on oxygen	[h ⁻¹]	0.467
α_1	Constant relating CO ₂ to growth	[mmol CO ₂ /g _B]	0.143
α_2	Constant relating CO ₂ to maintenance energy	[mmol CO ₂ /(g _B h)]	4 × 10 ⁻⁷
α_3	Constant relating CO ₂ to penicillin production	[mmol CO ₂ /(L h)]	10 ⁻⁴
K_{ox}, K_{op}	Oxygen limitations constants (no limitation)	[-]	0
K_{ox}, K_{op}	Oxygen limitations constants (with limitation)	[-]	2 × 10 ⁻² , 5 × 10 ⁻⁴
μ_{pp}^{max}	Specific rate of penicillin production	[h ⁻¹]	0.005
K_p	Inhibition constant	[g/L]	0.0002
K_l	Inhibition constant for product formation	[g/L]	0.10
p	Constant	[-]	3
K	Penicillin hydrolysis rate constant	[h ⁻¹]	0.04
k_g	Arrhenius constant for growth	[-]	7 × 10 ⁻³
E_g	Activation constant for growth	[cal/mol]	5100
k_d	Arrhenius constant for cell death	[-]	10 ³³
E_d	Activation energy for cell death	[cal/mol]	50000
ρc_p	Density × heat capacity of the medium	[cal/(°C L)]	1/1500
ρc_{pc}	Density × heat capacity of the cooling liquid	[cal/(°C L)]	1/2000
r_{q1}	Yield of heat generation	[cal/g _B]	60
r_{q2}	Constant in heat generation	[cal/g _B h]	1.6783 × 10 ⁻⁴
a	Heat transfer coefficient	[cal/h °C]	1000
b	Constant	[-]	0.60
α	Constant for K_{la}	[-]	70
β	Constant for K_{la}	[-]	0.4
λ	Constant in F_{loss}	[h ⁻¹]	2.5 × 10 ⁻⁴
γ	Proportionality constant	[mol [H ⁺]/g _B]	10 ⁻⁵
T_f	Feed temperature of substrate	[K]	298

Appendix B

Additional numerical results for the freeze-drying process

In this Appendix, the results that have been obtained for dataset C (siliconized vials) of the process described in Chapter 7 are presented. The results are shown following the same structure that has been used for dataset A (non-siliconized vials) in Chapter 7.

B.1 Gravimetric estimation of K_v for dataset C

Table B.1 – B-3 collect the estimates of the heat transfer coefficient obtained with the gravimetric method. Similar results to the ones for dataset A have been obtained, thus suggesting that the heat transfer coefficient is not sensibly affected by the presence of the silicone layer on the internal surface of the vials.

Table B.1. Mean values of the heat transfer coefficient and standard deviations for dataset C at $p_c/\bar{p}_c = 0.667$.

Zone #	Normalized mean value of K_v [-]	Standard deviation	Vial type
1	1.413	0.129	Edge
2	1.081	0.082	Edge
3	1.225	0.121	Edge
4	1.041	0.133	Central
5	1.187	0.145	Edge
6	0.972	0.122	Central
7	1.019	0.154	Central
8	1.346	0.124	Edge
9	0.868	0.122	Central
10	1.611	0.296	Edge
11	1.299	0.201	Edge
12	2.212	0.287	Edge (window)

Table B.2. Mean values of the heat transfer coefficient and standard deviations for dataset A at $p_c/\bar{p}_c = 1$.

Zone #	Normalized mean value of K_v [-]	Standard deviation	Vial type
1	1.715	0.094	Edge
2	1.213	0.086	Edge
3	1.455	0.135	Edge
4	1.124	0.095	Central
5	1.341	0.083	Edge
6	1.119	0.087	Central
7	1.201	0.085	Central
8	1.513	0.112	Edge
9	1.050	0.122	Central
10	1.746	0.264	Edge
11	1.422	0.101	Edge
12	2.974	0.196	Edge (window)

Table B.3. Mean values of the heat transfer coefficient and standard deviations for dataset A at $p_c/\bar{p}_c = 1.333$.

Zone #0	Normalized mean value of K_v [-]	Standard deviation	Vial type
1	1.642	0.113	Edge
2	1.310	0.125	Edge
3	1.542	0.145	Edge
4	1.055	0.055	Central
5	1.412	0.142	Edge
6	1.200	0.133	Central
7	1.322	0.194	Central
8	1.848	0.368	Edge
9	1.136	0.341	Central
10	1.697	0.145	Edge
11	1.655	0.193	Edge
12	2.412	0.166	Edge (window)

The “global” heat transfer coefficient that would be obtained by averaging all the values within the drying chamber is reported in Table B.4, together with its standard deviation.

Table B.4. Mean values and standard deviations of the “global” (i.e. computed for all the 12 zones) heat transfer coefficient at the three different values of chamber pressure considered.

Normalized chamber pressure [-]	Normalized mean value of K_v [-]	Standard deviation
0.667	1.341	0.793
1	1.552	0.892
1.333	1.601	0.889

B.2 Maximum likelihood parameter estimation for dataset C

Table B.5 collects the results of the MLE activity for dataset C. The results suggest a strong effect of the silicone layer on the mass transfer resistance, while only a slight effect on the heat transfer coefficient (as already suggested by the gravimetric estimates).

Table B.5. Maximum likelihood parameter estimation for the 12 zones considered. Dataset C (siliconized vials).

Parameter	% nominal value	95% t-value	Reference t-value	Parameter	% nominal value	95% t-value	Reference t-value
<i>Zone 1</i>				<i>Zone 7</i>			
C_1	103.44	74.41	1.65	C_1	103.10	68.11	1.65
R_1	112.11	69.18	1.65	R_1	121.79	46.22	1.65
<i>Zone 2</i>				<i>Zone 8</i>			
C_1	109.00	72.11	1.65	C_1	108.40	65.41	1.65
R_1	109.12	57.14	1.65	R_1	107.92	50.12	1.65
<i>Zone 3</i>				<i>Zone 9</i>			
C_1	101.16	75.82	1.65	C_1	111.61	71.46	1.65
R_1	118.32	62.12	1.65	R_1	101.12	48.57	1.65
<i>Zone 4</i>				<i>Zone 10</i>			
C_1	106.21	64.59	1.65	C_1	103.12	68.92	1.65
R_1	107.56	52.13	1.65	R_1	107.11	65.16	1.65
<i>Zone 5</i>				<i>Zone 11</i>			
C_1	112.14	68.26	1.65	C_1	112.30	71.25	1.65
R_1	113.12	60.12	1.65	R_1	108.19	56.25	1.65
<i>Zone 6</i>				<i>Zone 12</i>			
C_1	121.41	68.59	1.65	C_1	109.10	61.12	1.65
R_1	119.71	61.92	1.65	R_1	119.51	44.32	1.65

References

- Aamir, E., Z.K. Nagy and C.D. Rielly (2010). Optimal seed recipe design for crystal size distribution control for batch cooling crystallization processes. *Chem. Eng. Sci.*, **65**, 3602-3614.
- Adi, V. S. K., and Chang, C. T. (2013). A mathematical programming formulation for temporal flexibility analysis. *Comput. Chem. Eng.*, **57**:151-158.
- Adi, V. S. K., Laxmidewi, R., & Chang, C. T. (2016) An effective computation strategy for assessing operational flexibility of high-dimensional systems with complicated feasible regions *Chem. Eng. Sci.* **147**: 137-149.
- Agatonovic-Kustrin, S., & Beresford, R. (2000). Basic concepts of artificial neural network (ANN) modeling and its application in pharmaceutical research. *Journal of pharmaceutical and biomedical analysis*, **22**, 717-727.
- Akkisetty, P.K.; Lee, U.; Reklaitis, G.V.; Venkatasubramanian, V. (2010) Population balance model-based hybrid neural network for a pharmaceutical milling process. *J. Pharm. Innov.* **5**, 161–168.
- am Ende, D., Bronk, K. S., Mustakis, J., O'Connor, G., Santa Maria, C. L., Nosal, R., & Watson, T. J. (2007). API quality by design example from the Torcetrapib manufacturing process. *J. Pharm. Innov.*, **2**, 71-86.
- am Ende, D. J. (2010). Chemical engineering in the pharmaceutical industry: an introduction. *Chemical Engineering in the Pharmaceutical Industry: R&D to Manufacturing*, 1-20.
- Anand, A., Curtis, J. S., Wassgren, C. R., Hancock, B. C., & Ketterhagen, W. R. (2008). Predicting discharge dynamics from a rectangular hopper using the discrete element method (DEM). *Chem. Eng. Sci.*, **63**, 5821-5830.
- Anderson M.J. and Whitcomb. P.J. Find the Most Favorable Formulations (1998). *Chem. Eng. Progr.*, **94**(4), 63–67.
- Andersson, M., Josefson, M., Langkilde, F. W. and Wahlund, K. G. (1999). Monitoring of a film coating process for tablets using near infrared reflectance spectrometry. *J. Pharm. Biomed. Anal.*, **20**, 27-37.
- Asprey, S. P., Macchietto, S. (2000). Statistical tools for optimal dynamic model building. *Comput. Chem. Eng.*, **24**, 1261-1267.
- Bajpai, R. K., and Reuss, M. (1980) A mechanistic model for penicillin production. *J. Chem. Technol. Biotechnol.* **30**: 332-344.
- Banerjee, I., & Ierapetritou, M. G. (2002). Design optimization under parameter uncertainty for general black-box models. *Ind. Eng. Chem. Res.*, **41**, 6687-6697.
- Banerjee, I., & Ierapetritou, M. G. (2004). Model independent parametric decision making. *Annals of Operations Research*, **132**, 135-155.

- Banerjee, I., Pal, S., & Maiti, S. (2010). Computationally efficient black-box modeling for feasibility analysis. *Comput. Chem. Eng.*, **34**, 1515-1521.
- Bano G., Facco P., Meneghetti N., Bezzo F., Barolo M. (2017). Uncertainty back-propagation in PLS model inversion for design space determination in pharmaceutical product development. *Comput. Chem. Eng.* **101**: 110-124.
- Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018a) Probabilistic design space determination in pharmaceutical product development: a Bayesian/latent variable approach. *AIChE J.* **64**: 2438:4779.
- Bano, G., Facco, P., Bezzo, F., Barolo, M. (2018b) A novel and systematic approach to identify the design space of pharmaceutical processes. *Comput. Chem. Eng.*, **115**, 309-322.
- Barrasso, D., & Ramachandran, R. (2015). Multi-scale modeling of granulation processes: bi-directional coupling of PBM with DEM via collision frequencies. *Chemical Engineering Research and Design*, **93**, 304-317.
- Barrasso, D., El Hagrasy, A., Litster, J. D., & Ramachandran, R. (2015). Multi-dimensional population balance model development and validation for a twin screw granulation process. *Powder Technology*, **270**, 612-621.
- Bayes, T.R. (1763) An essay towards solving a problem in the doctrine of chances *Phil. Trans. Roy. Soc. London*, 53:370–418, 1763.
- Bhosekar, A., & Ierapetritou, M. (2017). Advances in surrogate based modeling, feasibility analysis and optimization: A review. *Comput. Chem. Eng.* **108**, 250-267.
- Bilgili, E. and B. Scarlett (2005). Population balance modeling of non-linear effects in milling processes. *Powder Tech.*, **153**, 59-71.
- Biol, G., Ündey, C., Cinar, A. (2002) A modular simulation package for fed-batch fermentation: penicillin production. *Comput. Chem. Eng.*, **26**: 1553-1565.
- Björkman, M., & Holmström, K. (2000). Global optimization of costly nonconvex functions using radial basis functions. *Optim. Eng.*, **1**: 373-397.
- Blanco, M., M. Alcalá, J.M. Gonzales and E. Torras (2006). Near infrared spectroscopy in the study of polymorphic transformation. *Anal. Chim. Acta*, **567**, 262-268.
- Bosca, S., Barresi, A. A., & Fissore, D. (2013). Use of a soft sensor for the fast estimation of dried cake resistance during a freeze-drying cycle. *Int. J. Pharm.*, **451**, 23-33.
- Bosca, S., Fissore, D., & Demichela, M. (2015). Risk-based design of a freeze-drying cycle for pharmaceuticals. *Ind. Eng. Chem. Res.*, **54**, 12928-12936.
- Boukouvala, F., Muzzio, F. J., & Ierapetritou, M. G. (2010). Design space of pharmaceutical processes using data-driven-based methods. *J. Pharm. Inn.*, **5**: 119-137.
- Boukouvala, F., A. Dubey, A. Vanarase, R. Ramachandran, F.J. Muzzio and M. Ierapetritou (2011a). Computational approaches for studying the granular dynamics of continuous blending processes, 2 – Population balance and data-based method. *Macromol. Mat. Eng.*, **297**, 9-19.

- Boukouvala, F., Muzzio, F. J., & Ierapetritou, M. G. (2011b). Dynamic data-driven modeling of pharmaceutical processes. *Industrial & Engineering Chemistry Research*, *50*(11), 6743-6754.
- Boukouvala, F., & Ierapetritou, M. G. (2012). Feasibility analysis of black-box processes using an adaptive sampling Kriging-based method. *Comput. Chem. Eng.*, **36**, 358-368.
- Boukouvala, F., & Ierapetritou, M. G. (2013). Surrogate-based optimization of expensive flowsheet modeling for continuous pharmaceutical manufacturing. *J. Pharm. Innov.*, **8**, 131-145.
- Boukouvala, F., & Ierapetritou, M. G. (2014). Derivative-free optimization for expensive constrained problems using a novel expected improvement objective function. *AIChE J.*, **60**: 2462-2474.
- Box, G. E. P. and Wilson, K.B. (1951) On the Experimental Attainment of Optimum Conditions (with discussion). *J. R. Stat. Soc. Series B* **13**, 1-45.
- Box, G. E., & Wilson, K. B. (1992). On the experimental attainment of optimum conditions. In *Breakthroughs in statistics* (pp. 270-310). Springer, New York, NY.
- Brueggemeier, S. B., Reiff, E. A., Lyngberg, O. K., Hobson, L. A., & Tabora, J. E. (2012). Modeling-based approach towards quality by design for the ibipinabant API step. *Org. Proc. Res. Dev.*, **16**, 567-576.
- Brülls, M., & Rasmuson, A. (2002). Heat transfer in vial lyophilization. *Int. J. Pharm.*, **246**, 1-16.
- Bu, Dongsheng, Boyong Wan, and Gary McGeorge (2013). A Discussion on the Use of Prediction Uncertainty Estimation of NIR Data in Partial Least Squares for Quantitative Pharmaceutical Tablet Assay Methods. *Chemom. Intell. Lab. Syst.*, **120**, 84–91.
- Bück, A., Peglow, M., Naumann, M., & Tsotsas, E. (2012). Population balance model for drying of droplets containing aggregating nanoparticles. *AIChE J.*, **58**, 3318-3328.
- Burt, J.L., A.D. Braem, A. Ramirez, B. Mudryk, L. Rossano, S. Tummala, (2011). Model-guided design space development for a drug substance manufacturing process. *J. Pharm. Innov.* **6**, 181–192.
- Cahyadi, C., Karande, A. D., Chan, L. W., and Heng, P. W. S. (2010). Comparative study of non-destructive methods to quantify thickness of tablet coatings. *Int. J. Pharm.*, **398**, 39-49.
- Cassidy, D. T., & Reid, J. (1982). Atmospheric pressure monitoring of trace gases using tunable diode lasers. *Applied Optics*, **21**, 1185-1190.
- Castilho, L. R., & Anspach, F. B. (2003). CFD-aided design of a dynamic filter for mammalian cell separation. *Biotechnology and bioengineering*, **83**, 514-524.
- Chalus, P., Roggo, Y., Walter, S., & Ulmschneider, M. (2005). Near-infrared determination of active substance content in intact low-dosage tablets. *Talanta*, **66**, 1294-1302.

- Charoo, N. A., Shamsheer, A. A., Zidan, A. S., & Rahman, Z. (2012). Quality by design approach for formulation development: a case study of dispersible tablets. *Int. J. Pharm.*, **423**, 167-178.
- Chatfield, M. A. K. B. M., Karmarkar, E. H. P. J. S., Andy, A. M. A. M. P., Trone, R. D. S. M. D., & Zhao, Q. W. Z. W. Y. (2017). Evaluating Progress in Analytical Quality by Design. Available at: <http://www.pharmtech.com/evaluating-progress-analytical-quality-design> (last accessed: 23/09/2018)
- Chatzizacharia, K. A., and Hatzivramidis, D. T. (2014) Design Space Approach for Pharmaceutical Tablet Development. *Ind. Eng. Chem. Res.*, **53**, 12003.
- Chaudhury, A., Barrasso, D., Pandey, P., Wu, H., & Ramachandran, R. (2014). Population balance model development, validation, and prediction of CQAs of a high-shear wet granulation process: towards QbD in drug product pharmaceutical manufacturing. *J. Pharm. Innov.*, **9**, 53-64.
- Chen, W., Wang, J., Desail S., Chang, I., Kiang, S., Lyngberg, O. (2017) A strategy for tablet active film coating formulation development using a content uniformity model and Quality by Design principles. In *Comprehensive Quality by design for pharmaceutical product development and manufacturing* (editors Reklatis *et al.*).
- Chen, H., Rustagi, S., Diep, E., Langrish, T. A., & Glasser, B. J. (2018). Scale-up of fluidized bed drying: Impact of process and design parameters. *Powder Technology*, **339**, 8-16.
- Cheng, Y. S., López-Isunza, F., Mongkhonsi, T., & Kershenbaum, L. (1993). Estimation of catalyst activity profiles in fixed-bed reactors with decaying catalysts. *Appl. Catal., A*, **106**, 193-199.
- Cheng, Y.S., Mongkhonsi, T., Kershenbaum, L.S (1997) Sequential estimation for nonlinear differential and algebraic systems- theoretical development and application. *Comput. Chem. Eng.* **21**:1051-1067.
- Chong, I.G., and C.H. Jun (2005). Performance of some variable selection methods when multicollinearity is present. *Chemom. Intell. Lab. Syst.*, **78**, 103-112.
- Christodoulou, C., Sorensen, E., García-Muñoz, S., & Mazzei, L. (2018). Mathematical modelling of water absorption and evaporation in a pharmaceutical tablet during film coating. *Chem. Eng. Sci.*, **175**, 40-55.
- Cioppa, T. M., & Lucas, T. W. (2007). Efficient nearly orthogonal and space-filling Latin hypercubes. *Technometrics*, **49**, 45-55.
- Close, E. J., Jeffrey R. S., Bracewell D.G., & Sorensen E. (2014). A Model Based Approach for Identifying Robust Operating Conditions for Industrial Chromatography with Process Variability. *Chem. Eng. Sci.*, **116**: 284–95.
- Coleman, M. C., & Block, D. E. (2006). Bayesian parameter estimation with informative priors for nonlinear systems. *AIChE J.*, **52**, 651-667.

- Collins, P. C. (2018). Chemical engineering and the culmination of quality by design in pharmaceuticals. *AIChE J*, **64**, 1502-1510.
- Cook, J., Cruaños, M. T., Gupta, M., Riley, S., & Crison, J. (2014). Quality-by-design: are we there yet? *AAPS PharmSciTech*, **15**, 140-148.
- Cruz, J., & Perkins, W. (1964). A new approach to the sensitivity problem in multivariable feedback system design. *IEEE Transactions on Automatic Control*, **9**, 216-223.
- Daraoui, N., Dufour, P., Hammouri, H., & Hottot, A. (2010). Model predictive control during the primary drying stage of lyophilisation. *Control Eng. Pract.*, **18**, 483-494.
- de Jong, S. (1993). SIMPLS: an alternative approach to partial least squares regression. *Chemom. Intell. Lab. Syst.*, **18**, 251-263.
- De Vries, S, and Ter Braak C.J.F. (1995). Prediction Error in Partial Least Squares Regression: A Critique on the Deviation Used in the Unscrambler. *Chemom. Intell. Lab. Syst.*, **30**, 239-45.
- Debrus, B., Lebrun, P., Kindenge, J. M., Lecomte, F., Ceccato, A., Caliaro, G., & Hubert, P. (2011). Innovative high-performance liquid chromatography method development for the screening of 19 antimalarial drugs based on a generic approach, using design of experiments, independent component analysis and design space. *J. Chrom.*, **1218**, 5205-5215.
- Deloitte (2017). *A new future for R&D? Measuring the return from pharmaceutical innovation 2017*. Available at: <https://www2.deloitte.com/content/dam/Deloitte/uk/Documents/life-sciences-health-care/deloitte-uk-measuring-roi-pharma.pdf> (last accessed: 23/09/2018)
- Deloitte (2018). *2018 Global Life Sciences Outlook*. Available at: <https://www2.deloitte.com/global/en/pages/life-sciences-and-healthcare/articles/global-life-sciences-sector-outlook.html> (last accessed: 23/09/2018)
- Denham, M. C. (1997). Prediction intervals in partial least squares. *J. Chemom.*, **11**, 39-52.
- Dienemann, E., & Osifchin, R. (2000). The role of chemical engineering in process development and optimization. *Current opinion in drug discovery & development*, **3**, 690-698.
- Dimitriadis, V.D., Pistikopoulos E.N. (1995) Flexibility analysis of dynamic systems. *Ind. Eng. Chem. Res.* **34**: 4451-4462.
- Ding, B. (2018). Pharma Industry 4.0: Literature review and research opportunities in sustainable pharmaceutical supply chains. *Process Saf. Environ.*, **119**, 115-130.
- Duralliu, A., Matejtschuk, P., & Williams, D. R. (2018). Humidity induced collapse in freeze dried cakes: A direct visualization study using DVS. *Eur. J. Pharm. Biopharm.*, **127**, 29-36.
- EFPIA (2017) *The pharmaceutical industries in figures. Key data 2017*. Annual report from the European Federation of Pharmaceutical Industries and Associations. Available at: https://www.efpia.eu/media/219735/efpia-pharmafigures2017_statisticbroch_v04-final.pdf (last accessed: 23/09/2018)

- Eitzlmayr, A., & Khinast, J. (2015a). Co-rotating twin-screw extruders: detailed analysis of conveying elements based on smoothed particle hydrodynamics. Part 1: hydrodynamics. *Chem. Eng. Sci.*, **134**, 861-879.
- Eitzlmayr, A., & Khinast, J. (2015b). Co-rotating twin-screw extruders: Detailed analysis of conveying elements based on smoothed particle hydrodynamics. Part 2: Mixing. *Chem. Eng. Sci.*, **134**, 880-886.
- Elkeshen, S. A., Badawi, S. S., & Badawi, A. A. (1996). Optimization of a reconstitutable suspension of rifampicin using 24 factorial design. *Drug development and industrial pharmacy*, **22**, 623-630.
- Eriksson, L., E. Johansson, N. Kettaneh-Wold, J. Trygg, C. Wikström and S. Wold (2006). *Multi- and megavariate data analysis. Part I. Basic principles and applications*. Umetrics AB, Umeå (Sweden).
- Escotet-Espinoza, M. S., Foster, C. J., & Ierapetritou, M. (2018). Discrete Element Modeling (DEM) for mixing of cohesive solids in rotating cylinders. *Powder Technology*, **335**, 124-136.
- EvaluatePharma (2018). *World preview 2018. Outlook to 2024*. Available at: <http://info.evaluategroup.com/rs/607-YGS-364/images/WP2018.pdf> (last accessed: 23/09/2018)
- Faber, K., and Kowalski B.R. (1996) Prediction error in least squares regression: Further critique on the deviation used in The Unscrambler. *Chemom. Intell. Lab. Syst.*, **34**, 283-292.
- Faber, K., and Kowalski B.R. (1997). Propagation of Measurement Errors for the Validation of Predictions Obtained by Principal Component Regression and Partial Least Squares. *J. Chemom.*, **11**, 181-238.
- Faber, K. (2002). Uncertainty Estimation for Multivariate Regression Coefficients. *Chemom. Intell. Lab. Syst.* **64**, 169-79.
- Facco, P., Dal Pasto F., Meneghetti N., Bezzo F., and Barolo M. (2015). Bracketing the Design Space within the Knowledge Space in Pharmaceutical Product Development. *Ind. Eng. Chem. Res.*, **54**, 5128-38.
- Facco, P., Santomaso, A. C., & Barolo, M. (2017). Artificial vision system for particle size characterization from bulk materials. *Chem. Eng. Sc.*, **164**: 246-257.
- FDA (2004a). Innovation or stagnation. Challenge and opportunity on the critical path to new medical products. *U.S. Department of Health and Human Services. U.S. Food and Drug Administration*.
- FDA (2004b). Pharmaceutical CGMPs for the 21st century – A risk based approach. Final report. *U.S. Department of Health and Human Services. U.S. Food and Drug Administration*.
- FDA (2004c). Guidance for industry. PAT - A framework for innovative pharmaceutical development, manufacturing and quality assurance. *Center for Drug Evaluation and Research, U.S. Food and Drug Administration, Rockville (MD), USA*.

- FDA (2006). Guidance for Industry - Quality Systems Approach to Pharmaceutical CGMP Regulations. *U.S. Department of Health and Human Services. U.S. Food and Drug Administration.*
- Fisher, A. C., Lee, S. L., Harris, D. P., Buhse, L., Kozlowski, S., Yu, L., ... & Woodcock, J. (2016). Advancing pharmaceutical quality: an overview of science and research in the US FDA's Office of Pharmaceutical Quality. *Int. J. Pharm.*, **515**, 390-402.
- Fissore, D., Pisano, R., & Barresi, A. A. (2010). On the methods based on the Pressure Rise Test for monitoring a freeze-drying process. *Drying Technol.*, **29**, 73-90
- Fissore, D., & Barresi, A. A. (2011b). Scale-up and process transfer of freeze-drying recipes. *Drying Technol.* **29**, 1673-1684.
- Fissore, D., Pisano, R., & Barresi, A. A. (2011a). Advanced approach to build the design space for the primary drying of a pharmaceutical freeze-drying process. *J. Pharm. Sci.*, **100**, 4922-4933.
- Fissore, D., Pisano, R., & Barresi, A. A. (2012). A model-based framework to optimize pharmaceuticals freeze drying. *Drying Technol.*, **30**, 946-958.
- Flaten G.R. (2018) Model maintenance. In: *Multivariate analysis in the pharmaceutical industry* (Ferreira, A. P., Menezes, J. C., & Toba, M. (Eds.)). Academic Press.
- Forrester, A. I., & Keane, A. J. (2009). Recent advances in surrogate-based optimization. *Prog. Aerosp. Sci.*, **45**, 50-79.
- Franceschini, G., & Macchietto, S. (2008). Model-based design of experiments for parameter precision: State of the art. *Chem. Eng. Sci.*, **63**, 4846-4872.
- Franks, F. (1998). Freeze-drying of bioproducts: putting principles into practice. *Europ. J. Pharm. Biopharm.*, **45**, 221-229.
- Gabrielsson, J., Lindberg, N. O., & Lundstedt, T. (2002). Multivariate methods in pharmaceutical applications. *Journal of Chemometrics: A Journal of the Chemometrics Society*, **16**, 141-160.
- Galvanin, F. (2010) Optimal model-based design of experiments in dynamic systems: novel techniques and unconventional applications. *PhD Thesis*. University of Padova.
- Ganguly, A., Nail, S. L., & Alexeenko, A. (2013). Experimental determination of the key heat transfer mechanisms in pharmaceutical freeze-drying. *J. Pharm. Sci.*, **102**, 1610-1625.
- García-Muñoz, S., T. Kourti and J.F. MacGregor (2004). Model Predictive Monitoring for Batch Processes. *Ind. Eng. Chem. Res.*, **43**, 5929-5941.
- García-Muñoz, S., J.F. MacGregor and T. Kourti (2005). Product transfer between sites using Joint-Y PLS. *Chemom. Intell. Lab. Syst.*, **79**, 101-114.
- García-Muñoz, S., Kourti T., MacGregor J.F., Apruzzese F., and Champagne M. (2006). Optimization of Batch Operating Policies. Part I. Handling Multiple Solutions. *Ind. Eng. Chem. Res.*, **45**, 7856-66.

- García-Muñoz, S. (2009). Establishing multivariate specifications for incoming materials using data from multiple scales. *Chemom. Intell. Lab. Syst.*, **98** 51-57.
- Garcia-Munoz, S., Luciani, C. V., Vaidyaraman, S., & Seibert, K. D. (2015). Definition of design spaces using mechanistic models and geometric projections of probability maps. *Org. Proc. Res. Dev.*, **19**, 1012-1023.
- García-Muñoz, S. and C.A. Oksanen (2010). Process modeling and control in drug development and manufacturing. *Comput. Chem. Eng.*, **34**, 1007-1008.
- García-Muñoz, S., Butterbaugh, A., Leavesley, I., Manley, L. F., Slade, D., & Bermingham, S. (2018). A flowsheet model for the development of a continuous process for pharmaceutical tablets: An industrial perspective. *AIChE J.*, **64**, 511-525.
- Garud S.S., Karimi I.A., Kraft M. (2017) Design of computer experiments: A review. *Comput. Chem. Eng.*, **106**: 71-95.
- Geladi, P., and Kowalski B.R. (1986). Partial Least-Squares Regression: A Tutorial. *Anal. Chim. Acta*, **185**: 1-17.
- Geman, S., & Geman, D. (1993). Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *J. Appl. Stat.*, **20**, 25-62.
- Georgakis, C. (2013). Design of dynamic experiments: A data-driven methodology for the optimization of time-varying processes. *Ind. Eng. Chem. Res.* **52**:12369-12382.
- Gernaey, K. V., Cervera-Padrell, A. E., & Woodley, J. M. (2012). A perspective on PSE in pharmaceutical process development and innovation. *Comput. Chem. Eng.*, **42**, 15-29.
- Gieseler, H., Kramer, T., & Pikal, M. J. (2007). Use of manometric temperature measurement (MTM) and SMART™ freeze dryer technology for development of an optimized freeze-drying cycle. *J. Pharm. Sci.*, **96**, 3402-3418.
- Global Regulatory Authority Websites. Available at : <https://www.pda.org/scientific-and-regulatory-affairs/regulatory-resources/global-regulatory-authority-websites> (last accessed: 23/09/2018)
- Goel, T., Haftka, R. T., Shyy, W., & Queipo, N. V. (2007). Ensemble of surrogates. *Struct. Multidiscip. Optim.*, **33**: 199-216.
- Goff, J. A. (1946). Low-pressure properties of water-from 160 to 212° F. *Trans. Am. Heat. Vent. Eng.*, **52**, 95-121.
- Goyal, V., & Ierapetritou, M. G. (2002). Determination of operability limits using simplicial approximation. *AIChE J.*, **48**: 2902-2909.
- Grossmann, I. E., Calfa, B. A., & Garcia-Herreros, P. (2014). Evolution of concepts and models for quantifying resiliency and flexibility of chemical processes. *Comp. Chem. Eng.*, **70**: 22-34.
- Gujral, B., Stanfield, F., & Rufino, D. (2007). Monte Carlo simulations for risk analysis in pharmaceutical product design. In *Crystal Ball User Conference*.

- Guo, Y., Kafui, K. D., Wu, C. Y., Thornton, C., & Seville, J. P. K. (2009). A coupled DEM/CFD analysis of the effect of air on powder flow during die filling. *AIChE J.*, **55**, 49-62.
- Guo, Y., Wu, C. Y., Kafui, K. D., & Thornton, C. (2011). 3D DEM/CFD analysis of size-induced segregation during die filling. *Powder Technology*, **206**, 177-188.
- Gutmann, H. M. (2001). A radial basis function method for global optimization. *J. Glob. Optim.*, **19**, 201-227.
- Halemane, K. P., & Grossmann, I. E. (1983). Optimal process design under uncertainty. *AIChE J.*, **29**, 425-433.
- Harinath, E., Foguth, L. C., and Braatz, R. D. (2015). Robust optimal control for the maximization of design space. In *2015 American Control Conference (ACC)*, 3886-3891.
- Haseltine, E. L., & Rawlings, J. B. (2005). Critical evaluation of extended Kalman filtering and moving-horizon estimation. *Ind. Eng. Chem. Res.*, **44**, 2451-2460.
- Hastings, W. K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, **57**, 97-109.
- Helton, J. C., & Davis, F. J. (2003). Latin hypercube sampling and the propagation of uncertainty in analyses of complex systems. *Reliability Engineering & System Safety*, **81**, 23-69.
- Herwig, C., Wölbeling, C., Zimmer, T. (2017) A holistic approach to production control from industry 4.0 to pharma 4.0. *Pharmaceutical engineering*. **37**, 44-49.
- Horibe, M., Sonoda, R., & Watano, S. (2018). Scale-Up of Lubricant Mixing Process by Using V-Type Blender Based on Discrete Element Method. *Chemical and Pharmaceutical Bulletin*, **66**, 548-553.
- Höskuldsson, A. (1988). PLS regression methods. *J. Chemom.*, **2**, 211-228.
- Hotelling, H. (1931). The generalization of Student's ratio. *Ann. Math. Stat.*, **2**, 360-378.
- Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *J. Educ. Psychol.*, **24**, 417-441.
- Hottot, A., Peczkalski, R., Vessot, S., & Andrieu, J. (2006). Freeze-drying of pharmaceutical proteins in vials: Modeling of freezing and sublimation steps. *Drying Technol.*, **24**, 561-570.
- Høy, M., Steen K., and Martens H. (1998). Review of Partial Least Squares Regression Prediction Error in Unscrambler. *Chemom. Intell. Lab. Syst.*, **44**, 123-33.
- Hsu, S.H., Reklaitis, G.V., Venkatasubramanian, V. (2010a) Modeling and control of roller compaction for pharmaceutical manufacturing. Part I: process dynamics and control framework, *J. Pharm. Innov.*, **5**: 14-23.
- Hsu, S.H., Reklaitis, G.V., Venkatasubramanian, V. (2010b) Modeling and control of roller compaction for pharmaceutical manufacturing. Part II: control system design. *J. Pharm. Innov.*, **5**: 24-36.
- Huang, D., Allen, T. T., Notz, W. I., & Zeng, N. (2006). Global optimization of stochastic black-box systems via sequential Kriging meta-models. *J. Glob. Optim.*, **34**: 441-466.

- Huang, J., Kaul, G., Cai, C., Chatlapalli, R., Hernandez-Abad, P., Ghosh, K., & Nagi, A. (2009). Quality by design case study: an integrated multivariate approach to drug product and process development. *Int. J. Pharm.*, **382**, 23-32.
- IBM Business Consulting Services (2005). Transforming industrialization. A new paradigm for pharmaceutical development. Available at: <http://www-935.ibm.com/services/uk/igs/pdf/ge510-3997-transforming-industrialization.pdf> (last accessed: 23/09/2018).
- ICH (1999). ICH harmonised tripartite guide. Specifications: test procedures and acceptance criteria for new drug substances and new drug products: chemical substances. Q6A. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2005). ICH Harmonized Tripartite Guideline, Quality Risk Management Q9.
- ICH (2006). ICH harmonised tripartite guide. Quality risk management Q9. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2009). ICH Harmonised Tripartite Guideline, Guidance for Industry, Pharmaceutical Development Q8(R2). International conference on harmonisation of technical requirements for registration of pharmaceuticals for human use. Silver Spring, MD: ICH, 2009.
- ICH (2009a). ICH harmonised tripartite guide. Pharmaceutical quality system Q10. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2009b). ICH harmonised tripartite guide. Pharmaceutical development Q8(R2). Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2010). Quality implementation working group on Q8, Q9 and Q10. Questions & Answers (R4). Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2012a). ICH quality implementation working group. Points to consider (R2). ICH-endorsed guide for ICH Q8/Q9/Q10 implementation. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2012b). Q11- Development and manufacture of drug substances. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2016). Q7- Good manufacturing practice guidance for active pharmaceutical ingredients. Guidance for industry. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)

- ICH (2018a). Q11- Development of manufacture of drug substances – Questions and answers. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2018b). Q7- Good manufacturing practice guidance for active pharmaceutical industry. Guidance for industry- Questions and answers. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: 23/09/2018)
- ICH (2018c) Q12- Technical and regulatory considerations for pharmaceutical product lifecycle management. Core guideline for industry. Available at: <http://www.ich.org/products/guidelines/quality/article/quality-guidelines.html> (last accessed: xxx)
- IFPMA (2017) *The pharmaceutical industry and global health. Facts and figures*. Annual report from the International Federation of Pharmaceutical Manufacturers & Associations. Available at: <https://www.ifpma.org/resource-centre/the-pharmaceutical-industry-and-global-health/> (last accessed: 23/09/2018)
- Immel, B. K. (2001). A brief history of the GMPs for pharmaceuticals. *Pharm. Technol.*, **25**, 44-53.
- Iooss, B., & Lemaître, P. (2015). A review on global sensitivity analysis methods. In *Uncertainty management in simulation-optimization of complex systems* (pp. 101-122). Springer, Boston, MA.
- J. Pharm. Innov.*, **2**, 71-86.
- Jackson, J. E (1991). *A user's guide to principal components*. John Wiley & Sons, Inc., New York (U.S.A.).
- Jaekle, C.M., and MacGregor J.F. (1998). Product Design Through Multivariate Statistical Analysis of Process Data. *AIChE J.*, **44**, 1105–18.
- Jaekle, C. M. and MacGregor, J. F. (2000). Industrial Applications of Product Design Through the Inversion of Latent Variable Models. *Chemom. Intell. Lab. Syst.*, **50**, 199–210.
- Jaekle, C.M. and J.F. MacGregor (2000). Product transfer between plants using historical process data. *AIChE J.*, **46**, 1989-1997.
- Jang, J., & Arastoopour, H. (2014). CFD simulation of a pharmaceutical bubbling bed drying process at three different scales. *Powder Technology*, **263**, 14-25.
- Jerier, J. F., Hathong, B., Richefeu, V., Chareyre, B., Imbault, D., Donze, F. V., & Doremus, P. (2011). Study of cold powder compaction by using the discrete element method. *Powder Technology*, **208**, 537-541.
- Jia, Z., Davis, E., Muzzio, F. J., & Ierapetritou, M. G. (2009). Predictive modeling for pharmaceutical processes using kriging and response surface. *J. Pharm. Innov.*, **4**, 174-186.
- Johanson, J.R. (1965) A rolling theory for granular solids. *J. Appl. Mech. B.*, **32**: 842-848.

- Jones, D. R., Schonlau, M., & Welch, W. J. (1998). Efficient global optimization of expensive black-box functions. *J. Glob. Opt.*, **13**: 455-492.
- Junod, S.W. (2004) “God, motherhood and the flag”: implementing the first pharmaceutical current good manufacturing practices (cGMP’s) regulations. Available at: <https://www.fda.gov/downloads/AboutFDA/History/ProductRegulation/UCM593504.pdf> (last accessed: 23/09/2018)
- Kager, J., Herwig, C., and Stelzer, I. V. (2018). State estimation for a penicillin fed-batch process combining particle filtering methods with online and time delayed offline measurements. *Chem. Eng. Sci.*, **177**: 234-244.
- Kaiser, Henry F. (1960). The Application of Electronic Computers to Factor Analysis. *Educ. Psychol. Meas.*
- Kagermann, H., Lukas W-D., Wahlster, W. (2011) *Industrie 4.0: Mit dem Internet der Dinge auf dem Weg zur 4. industriellen Revolution*. Available at: <https://www.vdi-nachrichten.com/Technik-Gesellschaft/Industrie-40-Mit-Internet-Dinge-Weg-4-industriellen-Revolution> (last accessed: 23/09/2018)
- Kaiser, Henry F. (1960). The Application of Electronic Computers to Factor Analysis. *Educ. Psychol. Meas.*, **20**, 141-151.
- Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, **82**: 35-45
- Kauffman, J. F., & Geoffroy, J. M. (2008). Propagation of uncertainty in process model predictions by Monte Carlo simulation. *Amer Pharm Rev*, 75-79.
- Kennard R.W, Stone L.A. (1969). Computer Aided Design of Experiments, *Technometrics*, **11**, 137-148.
- Ketterhagen, W. R., am Ende, M. T., & Hancock, B. C. (2009). Process modeling in the pharmaceutical industry using the discrete element method. *J. Phar. Sci.*, **98**, 442-470.
- Kim, M., H. Chung, Y. Woo and M.S. Kempes (2007). A new non-invasive, quantitative Raman technique for the determination of an active ingredient in pharmaceutical liquids by direct measurement through a plastic bottle. *Anal. Chim. Acta*, **587**, 200-207.
- Kishida, M., and Braatz R.D. (2012) A model-based approach for the construction of design spaces in quality-by-design. In *2012 American Control Conference (ACC)*.
- Kishida, M., and Braatz R.D. (2014) Skewed structured singular value-based approach for the construction of design spaces: theory and applications. *IET Control Theory A.*, **8**, 1321-1327.
- Kleijnen, J. P. (2009). Kriging metamodeling in simulation: A review. *Eur. J. Oper. Res.*, **192**, 707-716.
- Kopcha, M. (2017) *Continuous manufacturing- common guiding principles can help ensure progress*. Available at: <https://blogs.fda.gov/fdavoce/index.php/2017/09/continuous-manufacturing-common-guiding-principles-can-help-ensure-progress/> (last accessed: 23/09/2018).

- Kourti, T., & Davis, B. (2012). The business benefits of quality by design (QbD). *Pharm Eng*, **32**, 1-10.
- Koziel, S., & Michalewicz, Z. (1999). Evolutionary algorithms, homomorphous mappings, and constrained parameter optimization. *Evol. Comp.*, **7**: 19-44.
- KPMG (2018). *Pharma outlook 2030: from evolution to revolution*. Available at: <https://assets.kpmg.com/content/dam/kpmg/xx/pdf/2017/02/pharma-outlook-2030-from-evolution-to-revolution.pdf> (last accessed: 23/09/2018)
- Kremer, D. M., & Hancock, B. C. (2006). Process simulation in the pharmaceutical industry: a review of some basic physical models. *J. Pharm. Sci.*, **95**, 517-529.
- Krige, D. G. (1951). *A statistical approach to some mine valuation and allied problems on the Witwatersrand: By DG Krige* (Doctoral dissertation, University of the Witwatersrand).
- Kumar, S., Gokhale, R., & Burgess, D. J. (2014). Quality by design approach to spray drying processing of crystalline nanosuspensions. *Int. J. Pharm.*, **464**, 234-242.
- Lainiotis, D. (1971). Optimal adaptive estimation: Structure and parameter adaption. *IEEE Trans. Automat. Contr.*, **16**, 160-170.
- Lebrun, P., Boulanger, B., Debrus, B., Lambert, P., and Hubert, P. (2013). A Bayesian design space for analytical methods based on multivariate models and predictions. *J. Biopharm. Stat.*, **23**, 1330-1351.
- Lebrun, P., Krier, F., Mantanus, J., Grohganz, H., Yang, M., Rozet, E., & Hubert, P. (2012). Design space approach in the optimization of the spray-drying process. *European Journal of Pharmaceutics and Biopharmaceutics*, **80**, 226-234.
- Lee, S. L., O'Connor, T. F., Yang, X., Cruz, C. N., Chatterjee, S., Madurawe, R. D., ... & Woodcock, J. (2015). Modernizing pharmaceutical manufacturing: from batch to continuous production. *J. Pharm. Innov.*, **10**, 191-199.
- Lenhoff, A. M., & Morari, M. (1982). Design of resilient processing plants—I Process design under consideration of dynamic aspects. *Chem. Eng. Sci.*, **37**, 245-258.
- Lenk, P. J., & DeSarbo, W. S. (2000). Bayesian inference for finite mixtures of generalized linear models with random effects. *Psychometrika*, **65**, 93-119.
- Lepore, J., and Spavins J. (2008) PQLI design space. *J. Pharm. Innov*, **3**, 79-87.
- Li, G., Rosenthal, C., & Rabitz, H. (2001). High dimensional model representations. *J. Phys. Chem. A*, **105**, 7765-7777.
- Liapis, A. I., & Bruttini, R. (1995). Freeze-drying of pharmaceutical crystalline and amorphous solutes in vials: Dynamic multi-dimensional models of the primary and secondary drying stages and qualitative features of the moving interface. *Drying Technol.*, **13**, 43-72.
- Liapis, A. I., Pim, M. L., & Bruttini, R. (1996). Research and development needs and opportunities in freeze drying. *Drying Technol.*, **14**, 1265-1300.

- Lima, F. V., Jia, Z., Ierapetritou, M., & Georgakis, C. (2010). Similarities and differences between the concepts of operability and flexibility: The steady-state case. *AIChE J.*, **56**, 702-716.
- Lindenberg, C., & Mazzotti, M. (2009). Experimental characterization and multi-scale modeling of mixing in static mixers. Part 2. Effect of viscosity and scale-up. *Chem. Eng. Sci.*, **64**, 4286-4294.
- Liu, Z., M.-J. Bruwer, J.F. MacGregor, S.S.S. Rathore, D.E. Reed and M.J. Champagne (2011a). Modeling and optimization of a tablet manufacturing line. *J. Pharm. Innov.*, **6**, 170-180.
- Liu, Z., M.-J. Bruwer, J.F. MacGregor, S.S.S. Rathore, D.E. Reed and M.J. Champagne (2011b). Scale-Up of a pharmaceutical roller compaction process using a Joint-Y partial least squares model. *Ind. Eng. Chem. Res.*, **50**, 10696-10706.
- Lundsberg-Nielsen, L., & Bruce, D. (2017) Does Quality by Design have the “X-factor”? Available at: https://www.nsf.org/newsroom_pdf/pb_does_quality_design_have_x_factor_wp.pdf (last accessed: 23/09/2018)
- MacGregor, J. F, and Bruwer MJ. (2008). A Framework for the Development of Design and Control Spaces. *J. Pharm. Innov* , **3**, 15–22.
- Mantanus, J., E. Ziémons, P. Lebrun, E. Rozet, R. Klinkenberg, B. Streeel, B. Evrard and P. Hubert (2010). Active content determination of non-coated pharmaceutical pellets by near-infrared spectroscopy: method development, validation and reliability evaluation. *Talanta*, **80**, 1750-1757.
- Mardia, K.V., J.T. Kent and J.M. Bibby (1979). *Multivariate analysis*. Academic Press Limited, London (U.K.).
- Marigo, M., Davies, M., Leadbeater, T., Cairns, D. L., Ingram, A., & Stitt, E. H. (2013). Application of Positron Emission Particle Tracking (PEPT) to validate a Discrete Element Method (DEM) model of granular flow and mixing in the Turbula mixer. *Int. J. Pharm.*, **446**, 46-58
- Martino, L., Read, J., & Luengo, D. (2015). Independent doubly adaptive rejection Metropolis sampling within Gibbs sampling. *IEEE Trans. Signal Process.*, **63**: 3123-3138.
- Mascarenhas, W. J., Akay, H. U., & Pikal, M. J. (1997). A computational model for finite element analysis of the freeze-drying process. *Comput. Methods Appl. Mech. Eng.*, **148**, 105-124.
- Matero, S., S. Poutiainen, J. Leskinen, K. Järvinen, J. Ketolainen, A. Poso and S.P. Reinikainen (2010). Estimation of granule size distribution for batch fluidized bed granulation process using acoustic emission and N-way PLS. *J. Chemom.*, **24**, 464-471.
- Mattsson, S. E., and Söderlind, G. (1993). Index reduction in differential-algebraic equations using dummy derivatives. *SIAM J. Sci. Comput.* **14**: 677-692.

- McCarthy, J. J., Jasti, V., Marinack, M., & Higgs, C. F. (2010). Quantitative validation of the discrete element method using an annular shear cell. *Powder Technology*, **203**, 70-77.
- McKinsey & Company. (2009). Challenges to Quality by Design. Available at: http://www.ipqpubs.com/wp-content/uploads/2012/01/McKinsey_report_QbD.pdf (last accessed: 23/09/2018)
- McKinsey&Company (2018). *Outlook on pharma operations*. Available at: https://www.mckinsey.com/~media/mckinsey/dotcom/client_service/operations/pdfs/outlook_on_pharma_operations.ashx (last accessed: 23/09/2018)
- Metta, N., Ierapetritou, M., & Ramachandran, R. (2018a). A multiscale DEM-PBM approach for a continuous comilling process using a mechanistically developed breakage kernel. *Chem. Eng. Sci.*, **178**, 211-221.
- Metta, N., Verstraeten, M., Ghijs, M., Kumar, A., Schafer, E., Singh, R., & Ierapetritou, M. (2018b). Model development and prediction of particle size distribution, density and friability of a comilling operation in a continuous pharmaceutical manufacturing process. *Int. J. Pharm.*, **549**, 271-282.
- Mevik, B. H., Segtnan, V. H., & Næs, T. (2004). Ensemble methods and partial least squares regression. *J. Chemom.*, **18**: 498-507.
- Millman, M. J., Liapis, A. I., & Marchello, J. M. (1985). An analysis of the lyophilization process using a sorption-sublimation model and various operational policies. *AIChE J.*, **31**, 1594-1604.
- Morris, M. D. (1991). Factorial sampling plans for preliminary computational experiments. *Technometrics*, **33**, 161-174.
- Mortier, S. T. F., Gernaey, K. V., De Beer, T., & Nopens, I. (2013). Development of a population balance model of a pharmaceutical drying process and testing of solution methods. *Comput. Chem. Eng.*, **50**, 39-53.
- Mortier, S. T. F., Van Bockstal, P. J., Corver, J., Nopens, I., Gernaey, K. V., & De Beer, T. (2016). Uncertainty analysis as essential step in the establishment of the dynamic Design Space of primary drying during freeze-drying. *Eur. J. Pharm. Biopharm.*, **103**, 71-83.
- Mortier, S.T.F.C., T. De Beer, K.V. Gernaey, J.P. Remon, C. Vervaet and I. Nopens (2011). Mechanistic modelling of fluidized bed drying processes of wet porous granules: a review. *Europ. J. Pharm. Biopharm.*, **79**, 205-225.
- Mukherjee, R., Mao, C., Chattoraj, S., & Chaudhuri, B. (2018). DEM based computational model to predict moisture induced cohesion in pharmaceutical powders. *Int. J. Pharm.*, **536**, 301-309.
- Naelapää, K., Allesø, M., Kristensen, H. G., Bro, R., Rantanen, J., & Bertelsen, P. (2009). Increasing process understanding by analyzing complex interactions in experimental data. *J. Pharm. Sci.*, **98**, 1852-1861.

- Nagy, Z.K., M. Fujiwara and R.D. Braatz (2008). Modeling and control of combined cooling and antisolvent crystallization processes. *J. Process Control*, **18**, 856-864.
- Nasr, M. M., Krumme, M., Matsuda, Y., Trout, B. L., Badman, C., Mascia, S., & Konstantinov, K. (2017). Regulatory Perspectives on Continuous Pharmaceutical Manufacturing: Moving From Theory to Practice. *J. Pharm. Sci.*, **106**, 3199-3206.
- Ngo, T. H., Vertommen, J., & Kinget, R. (1997). Formulation of artemisinin tablets. *Int. J. Pharm.*, **146**, 271-274.
- Ogunnaike, B. A. (1994). On-line modelling and predictive control of an industrial terpolymerization reactor. *Int. J. Control*, **59**, 711-729.
- Oka, S., Kašpar O., Tokárová V., Sowrirajan K., Wu H., Khan M., Muzzio F., Štěpánek F., and Ramachandran R. (2015). A Quantitative Study of the Effect of Process Parameters on Key Granule Characteristics in a High Shear Wet Granulation Process Involving a Two Component Pharmaceutical Blend. *Adv. Powder Technol.*, **26**, 315–22.
- Ottoboni, S., Steven, C., Meehan, E., Barton, A., Firth, P., Mitchell, A., & Tahir, F. (2018). Development of a novel continuous filtration unit for pharmaceutical process development and manufacturing. *J. Pharm. Sci. In press*. DOI: <https://doi.org/10.1016/j.xphs.2018.07.005>
- Pandey, P., Song, Y., Kayihan, F., & Turton, R. (2006). Simulation of particle movement in a pan coating device using discrete element modeling and its comparison with video-imaging experiments. *Powder Technology*, **161**, 79-88.
- Pantelides C.C., Shah N., Adjiman C.S. Design Space, Models, and Model Uncertainty. *Comprehensive Quality by Design in Pharmaceutical Development and Manufacture*. AIChE 2012 Annual Meeting: Nashville, TN; paper 417f
- Pantelides, C. C. (1988). The consistent initialization of differential-algebraic systems. *SIAM Journal on Scientific and Statistical Computing*, **9**: 213-231.
- Pantelides, C. C., Renfro, J. G. (2013) The online use of first-principles models in process operations: Review, current status and future needs. *Comput. Chem. Eng.* **51**: 136-148.
- Pantelides, C.C., Bano, G., Cheng Y.S., Spatenka S., Matzopoulos M., *Model-based real time monitoring of ethylene cracking furnaces*, 2016 AIChE Spring Meeting & 12th Global congress on process safety: Houston, TX; paper 444934
- Pantelides, C.C., Shah, N., Adjiman, C. S. (2009) Design Space, Models, and Model Uncertainty. *Comprehensive Quality by Design in Pharmaceutical Development and Manufacture*; AIChE Annual Meeting: Nashville, TN; paper 417f
- Parker, J., LaMarche, K., Chen, W., Williams, K., Stamato, H., & Thibault, S. (2013). CFD simulations for prediction of scaling effects in pharmaceutical fluidized bed processors at three scales. *Powder technology*, **235**, 115-120.
- Patwardhan, S. C., Narasimhan, S., Jagadeesan, P., Gopaluni, B., Shah, S. L. (2012) Nonlinear Bayesian state estimation: A review of recent developments. *Control Eng. Pract.* **20**: 933-953.

- Paul, E. L., & Rosas, C. B. (1990). Challenges for chemical engineers in the pharmaceutical industry. *Chem. Eng. Progr.* **86**, 17-25.
- Peinado, A., J. Hammond and A. Scott (2011). Development, validation and transfer of a near infrared method to determine in-line the end point of a fluidised drying process for commercial production batches of an approved oral solid dose pharmaceutical product. *J. Pharm. Biomed. Anal.*, **54**, 13-20.
- Peña, R., Burcham, C. L., Jarmer, D. J., Ramkrishna, D., & Nagy, Z. K. (2017). Modeling and optimization of spherical agglomeration in suspension through a coupled population balance model. *Chem. Eng. Sci.*, **167**, 66-77.
- Persson, A. S., & Frenning, G. (2012). An experimental evaluation of the accuracy to simulate granule bed compression using the discrete element method. *Powder technology*, **219**, 249-256.
- Peterson J.J. (2004). A Posterior Predictive Approach to Multiple Response Surface Optimization.
- Peterson JJ (2008). A Bayesian approach to the ICH Q8 definition of design space. *J. Biopharm. Stat.*, **18**: 959-75
- Peterson JJ, Yahyah M (2009). A Bayesian design space approach to robustness and system suitability for pharmaceutical assays and other processes. *Stat. Biopharm. Res.* **1**: 441-449.
- Peterson, J. J., and Lief, K. (2010). The ICH Q8 definition of design space: A comparison of the overlapping means and the Bayesian predictive approaches. *Stat. Biopharm. Res.*, **2**, 249-259..
- Phatak, A., Reilly, P. M., & Penlidis, A. (1993). An approach to interval estimation in partial least squares regression. *Anal. Chim. Act.*, **277**, 495-501.
- Picheny, V., Ginsbourger, D., Richet, Y., & Caplin, G. (2013). Quantile-based optimization of noisy computer experiments with tunable precision. *Technometrics*, **55**: 2-13
- Pierna, JAF., Jin L., Wahl F., Faber NM., and Massart DL. (2003). Estimation of Partial Least Squares Regression Prediction Uncertainty When the Reference Values Carry a Sizeable Measurement Error. *Chemom. Intell. Lab. Syst.*, **65**, 281-91.
- Pikal, M. J. (1985). Use of laboratory data in freeze drying process design: heat and mass transfer coefficients and the computer simulation of freeze drying. *PDA J. Pharm. Sci. Technol.*, **39**, 115-139.
- Pikal, M. J. (2000). Heat and mass transfer in low pressure gases: applications to freeze drying. *Drugs and the Pharmaceutical Sciences*, **102**, 611-686.
- Pikal, M. J., Bogner, R., Mudhivarthi, V., Sharma, P., & Sane, P. (2016). Freeze-drying process development and scale-up: scale-up of edge vial versus center vial heat transfer coefficients, Kv. *J. Pharm. Sci.*, **105**, 3333-3343.
- Pisano, R., Fissore, D., & Barresi, A. A. (2011). Freeze-drying cycle optimization using model predictive control techniques. *Ind. Eng. Chem. Res.*, **50**, 7363-7379.

- Pisano, R., Fissore, D., Barresi, A. A., & Rastelli, M. (2013). Quality by design: scale-up of freeze-drying cycles in pharmaceutical industry. *AAPS PharmSciTech*, **14**, 1137-1149.
- Pistikopoulos, E. N., & Mazzuchi, T. A. (1990). A novel flexibility analysis approach for processes with stochastic parameters. *Comput. Chem. Eng.*, **14**, 991-1000.
- Pistikopoulos, E. N., & Ierapetritou, M. G. (1995). Novel approach for optimal process design under uncertainty. *Comput. Chem. Eng.*, **19**, 1089-1110.
- Polizzi, M. A., & García-Muñoz, S. (2011). A framework for in-silico formulation design using multivariate latent variable regression methods. *Int. J. Pharm.*, **418**, 235-242.
- Pöllänen, K., A. Hakkinen, S.P. Reinikainen, J. Rantanen, M. Karjalainen, M. Louhi-Kultanen and L. Nystrom (2005). IR spectroscopy together with multivariate data analysis as a process analytical tool for in-line monitoring of crystallization process and solid-state analysis of crystalline product. *J. Pharm. Biomed. Anal.*, **38**, 275-284.
- Poon, J.M., R. Ramachandran, C.F.W. Sanders, T. Glaser, C.D. Immanuel and F.J. Doyle III (2009). Experimental validation studies on a multi-dimensional and multi-scale population balance model of batch granulation. *Chem. Eng. Sci.*, **64**, 775-786.
- Poozesh, S., Lu, K., & Marsac, P. J. (2018). On the particle formation in spray drying process for bio-pharmaceutical applications: Interrogating a new model via computational fluid dynamics. *International Journal of Heat and Mass Transfer*, **122**, 863-876.
- Pordal, H. S., Matice, C. J., & Fry, T. J. (2002). Computational fluid dynamics in the pharmaceutical industry. *Pharm. Technol*, **26**, 72-79.
- Process Systems Enterprise Ltd. (2014) gSOLIDS[®], VERSION 4.1.0; Process Systems Enterprise Ltd: London, UK
- Process Systems Enterprise Ltd. (2014) gPROMS Model Builder; Process Systems Enterprise Ltd: London, UK
- Prpich, A., am Ende, M. T., Katschner, T., Lubczyk, V., Weyhers, H., & Bernhard, G. (2010). Drug product modeling predictions for scale-up of tablet film coating—a quality by design approach. *Comput. Chem. Eng.*, **34**, 1092-1097.
- PwC (2018) Pharma 2020: from vision to decision. Available at: <https://www.pwc.com/gx/en/pharma-life-sciences/pharma2020/assets/pwc-pharma-success-strategies.pdf> (last accessed: 23/09/2018)
- QuantileIMS Institute (2016). *Outlook for Global Medicines Through 2021: Balancing Cost and Value*. Available at: <https://www.iqvia.com/newsroom/2016/quintilesims-institute-forecast> (last accessed: 23/09/2018)
- Ramachandran, R., C.D. Immanuel, F. Stepanek, J.D. Litster and F.J. Doyle III (2009). A mechanistic model for breakeage in population balances of granulation: theoretical kernel development and experimental validation. *Chem. Eng. Res. Des.*, **87**, 598-614.

- Rambhatla, S., Ramot, R., Bhugra, C., & Pikal, M. J. (2004). Heat and mass transfer scale-up issues during freeze drying: II. Control and characterization of the degree of supercooling. *Aaps Pharmscitech*, **5**, 54-62.
- Ramšak, M., Ravnik, J., Zadavec, M., Hriberšek, M., & Iljaž, J. (2017). Freeze-drying modeling of vial using BEM. *Engineering Analysis with Boundary Elements*, **77**, 145-156.
- Rao, C. V., Rawlings, J. B., & Lee, J. H. (2001). Constrained linear state estimation—a moving horizon approach. *Automatica*, **37**, 1619-1628.
- Rauschnabel, J. (2018). *2018 Pharma Industry Outlook*. Available at: https://www.contractpharma.com/issues/2018-01-01/view_features/pharma-industry-outlook (last accessed: 23/09/2018)
- Raw, A. S., Lionberger, R., & Lawrence, X. Y. (2011). Pharmaceutical equivalence by design for generic drugs: modified-release products. *Pharmaceutical research*, **28**, 1445-1453.
- Ray, W. (1981). *Advanced Process Control*. McGraw-Hill- Chemical Engineering Series.
- Rehrl, J., Karttunen, A. P., Nicolai, N., Hörmann, T., Horn, M., Korhonen, O., ... & Khinast, J. G. (2018). Control of three different continuous pharmaceutical manufacturing processes: Use of soft sensors. *Int. J. Pharm.* **543**, 60-72.
- Reis, M., and Saraiva PM. (2005). Integration of Data Uncertainty in Linear Regression and Process Optimization. *AIChE J.*, **51**, 3007–19.
- Reklaitis, G. V., Garcia-Munoz, S., & Seymour, C. (Eds.). (2017). *Comprehensive Quality by Design for Pharmaceutical Product Development and Manufacture*. John Wiley & Sons.
- Reuters (2018). *New drug approvals hit 21-year high in 2017*. January 2, 2018. Available at: <https://uk.reuters.com/article/us-pharmaceuticals-approvals/new-drug-approvals-hit-21-year-high-in-2017-idUKKBN1ER0P7> (last accessed: 23/09/2018)
- Riascos, C. A., Pinto, J. M. (2004) Optimal control of bioreactors: a simultaneous approach for complex systems. *Chem. Eng. J.* **99**: 23-34.
- Riccati, J. (1724). Animadversiones in aequationes differentiales secundi gradus. *Actorum Eruditorum Supplementa*, **8**, 66-73.
- Ricker, N. L., and Lee, J. H. (1995). Nonlinear modeling and state estimation for the Tennessee Eastman challenge process. *Comput. Chem. Eng.*, **19**: 983-1005.
- Rogers, A., Ierapetritou, M. (2015), Feasibility analysis of black-box processes. Part 1: surrogate-based feasibility analysis. *Chem. Eng. Sci.* **137**: 986-1004.
- Rogers, A. and M. Ierapetritou (2015). Challenges and opportunities in modeling pharmaceutical manufacturing processes, *Comput. Chem. Eng.*, **81**, 32-39.
- Rogers, A. J., Hashemi, A., & Ierapetritou, M. G. (2013). Modeling of particulate processes for the continuous manufacture of solid-based pharmaceutical dosage forms. *Processes*, **1**, 67-127.
- Rogers, A., and Ierapetritou, M. (2015). Feasibility and flexibility analysis of black-box processes Part 1: Surrogate-based feasibility analysis. *Chem. Eng. Sc.*, **137**: 986-1004.

- Sadikoglu, H., & Liapis, A. I. (1997). Mathematical modelling of the primary and secondary drying stages of bulk solution freeze-drying in trays: Parameter estimation and model discrimination by comparison of theoretical results with experimental data. *Drying Technol.*, **15**, 791-810.
- Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Gatelli, D., & Tarantola, S. (2008). *Global sensitivity analysis: the primer*. John Wiley & Sons.
- Saltelli, A., Annoni, P., Azzini, I., Campolongo, F., Ratto, M., & Tarantola, S. (2010). Variance based sensitivity analysis of model output. Design and estimator for the total sensitivity index. *Comp. Phys. Comm.*, **181**: 259-270.
- Sampat, C., Bettencourt, F., Baranwal, Y., Paraskevakos, I., Chaturbedi, A., Karkala, S., ... & Ierapetritou, M. (2018). A parallel unidirectional coupled DEM-PBM model for the efficient simulation of computationally intensive particulate process systems. *Comput, Chem. Eng. In press*. DOI: <https://doi.org/10.1016/j.compchemeng.2018.08.006>
- Schneider, R., Georgakis, C. (2012) How to NOT make the Extended Kalman Filter fail. *Ind. Eng. Chem. Res.* **52**, 3354-3362.
- Schölkopf, B., Smola, A., & Müller, K. R. (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput.*, **10**: 1299-1319.
- Scutellà, B., Plana-Fattori, A., Passot, S., Bourlès, E., Fonseca, F., Flick, D., & Trelea, I. C. (2017a). 3D mathematical modelling to understand atypical heat transfer observed in vial freeze-drying. *Appl. Therm. Eng.*, **126**, 226-236.
- Scutellà, B., Passot, S., Bourlès, E., Fonseca, F., & Trélea, I. C. (2017b). How vial geometry variability influences heat transfer and product temperature during freeze-drying. *J. Pharm. Sci.*, **106**, 770-778.
- Searles, J. (2004). Observation and implications of sonic water vapor flow during freeze-drying. *Am. Pharmaceut. Rev.*, **7**, 58-69.
- Sen, M., Dubey, A., Singh, R., & Ramachandran, R. (2013). Mathematical development and comparison of a hybrid PBM-DEM description of a continuous powder mixing process. *J. Powder Technology*.
- Sen, M., Barrasso, D., Singh, R., & Ramachandran, R. (2014). A multi-scale hybrid CFD-DEM-PBM description of a fluid-bed granulation process. *Processes*, **2**, 89-111.
- Serneels, S., Lemberge, P., & Van Espen, P. J. (2004). Calculation of PLS prediction intervals using efficient recursive relations for the Jacobian matrix. *J. Chemom.*, **18**, 76-80.
- Shah, R.B., M.A. Tawakkul and M.A. Khan (2007). Process analytical technology: chemometric analysis of Raman and near infra-red spectroscopic data for predicting physical properties of extended release matrix tablets. *J. Pharm. Sci.*, **96**, 1356-1365.
- Shan, S., & Wang, G. G. (2010). Survey of modeling and optimization strategies to solve high-dimensional design problems with computationally-expensive black-box functions. *Struct. Multidiscipl. Optim.*, **41**: 219-241.

- Sheehan, P., & Liapis, A. I. (1998). Modeling of the primary and secondary drying stages of the freeze drying of pharmaceutical products in vials: Numerical results obtained from the solution of a dynamic and spatially multi-dimensional lyophilization model for different operational policies. *Biotechnology and Bioengineering*, **60**, 712-728.
- Shirazian, S., Darwish, S., Kuhs, M., Croker, D. M., & Walker, G. M. (2018). Regime-separated approach for population balance modelling of continuous wet granulation of pharmaceutical formulations. *Powder Technology*, **325**, 420-428.
- Siiriä, S. M., Antikainen, O., Heinämäki, J., & Yliruusi, J. (2011). 3D simulation of internal tablet strength during tableting. *Aaps Pharmscitech*, **12**, 593-603.
- Sin, G., Ödman, P., Petersen, N., Lantz, A. E., & Gernaey, K. V. (2008). Matrix notation for efficient development of first-principles models within PAT applications: Integrated modeling of antibiotic production with *Streptomyces coelicolor*. *Biotech. Bioeng.*, **101**, 153-171.
- Smith, P. J., Thornhill, G. D., Dance, S. L., Lawless, A. S., Mason, D. C., and Nichols, N. K. (2013). Data assimilation for state and parameter estimation: application to morphodynamic modelling. *Q. J. Royal Meteorol. Soc.* **139**: 314-327.
- Sobol, I. M. (1993). Sensitivity estimates for nonlinear mathematical models. *Mathematical modelling and computational experiments*, **1**, 407-414.
- Sobol, I. M. (2001). Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. *Mathematics and computers in simulation*, *55*(1-3), 271-280.
- Souih, N., Reynolds G., Tajarobi P., Wikström H., Haeffler G., Josefson M., and Trygg J. (2015). Roll Compaction Process Modeling: Transfer Between Equipment and Impact of Process Parameters. *Int. J. Pharm.*, **484**, 192–206.
- Stockdale, G. W., and Cheng, A. (2009). Finding design space and a reliable operating region using a multivariate Bayesian approach with experimental design. *Qual. Technol. Quant. Manag.*, **6**, 391-408.
- Straub, D. A., & Grossmann, I. E. (1993). Design optimization of stochastic flexibility. *Comput. Chem. Eng.*, **17**, 339-354.
- Swann, J. P. (1999). The 1941 sulfathiazole disaster and the birth of good manufacturing practices. *Pharmacy in history*, **41**, 16-25.
- Thirunahari, S., Chow, P. S., & Tan, R. B. (2011). Quality by design (QbD)-based crystallization process development for the polymorphic drug tolbutamide. *Crystal Growth & Design*, **11**, 3027-3038.
- Tomba, E., Barolo, M., & García-Muñoz, S. (2012). General framework for latent variable model inversion for the design and manufacturing of new products. *Ind. Eng. Chem. Res.*, **51**, 12886-12900.

- Tomba, E., Facco, P., Bezzo, F., and Barolo, M. (2013). Latent variable modeling to assist the implementation of Quality-by-Design paradigms in pharmaceutical development and manufacturing: a review. *Int. J. Pharm.*, **457**, 283-297.
- Trelea, I. C., Passot, S., Fonseca, F., & Marin, M. (2007). An interactive tool for the optimization of freeze-drying cycles based on quality criteria. *Drying Technol.*, **25**, 741-751.
- Troup, G. M., & Georgakis, C. (2013). Process systems engineering tools in the pharmaceutical industry. *Comput. Chem. Eng.* **51**, 157-171.
- Tsinontides, S. C., Rajniak, P., Pham, D., Hunke, W. A., Placek, J., & Reynolds, S. D. (2004). Freeze drying—principles and practice for successful scale-up to manufacturing. *Int. J. Pharm.*, **280**, 1-16.
- Uztürk, D., & Georgakis, C. (2002). Inherent dynamic operability of processes: General definitions and analysis of SISO cases. *Ind. Eng. Chem. Res.*, **41**, 421-432.
- van der Voet, H. (1999). Pseudo-Degrees of Freedom for Complex Predictive Models: The Example of Partial Least Squares. *J. Chemom.*, **13**, 195–208.
- Vanarase, A.U., M. Alcala, J.I.J. Rozo, F.J. Muzzio and R.J. Romanach (2010). Real-time monitoring of drug concentration in a continuous powder mixing process using NIR spectroscopy. *Chem. Eng. Sci.*, **65**, 5728-5733.
- Vanlaer, J., Gins, G., Van Impe, J. F. M. (2013) Quality Assessment of a Variance Estimator for Partial Least Squares Prediction of Batch-End Quality. *Comput. Chem. Eng.*, **52**, 230-239.
- Velardi, S. A., & Barresi, A. A. (2008a). Development of simplified models for the freeze-drying process and investigation of the optimal operating conditions. *Chem. Eng. Res. Des.*, **86**, 9-22.
- Velardi, S. A., Rasetto, V., & Barresi, A. A. (2008b). Dynamic parameters estimation method: advanced manometric temperature measurement approach for freeze-drying monitoring of pharmaceutical solutions. *Ind. Eng. Chem. Res.*, **47**, 8445-8457.
- Vemavarapu, C., Surapaneni M., Hussain M., and Badawy S.. (2009). Role of Drug Substance Material Properties in the Processibility and Performance of a Wet Granulated Product. *Int. J. Pharm.*, **374**, 96–105.
- Verma, S., Lan, Y., Gokhale, R., & Burgess, D. J. (2009). Quality by design approach to understand the process of nanosuspension preparation. *Int. J. Pharm.*, **377**, 185-198.
- Vetter, T., Iggland, M., Ochsenbein, D. R., Hänseler, F. S., & Mazzotti, M. (2013). Modeling nucleation, growth, and Ostwald ripening in crystallization processes: a comparison between population balance and kinetic rate equation. *Crystal Growth & Design*, **13**, 4890-4905.
- Vinson, D. R., & Georgakis, C. (2000). A new measure of process output controllability. *J. Process Control*, **10**, 185-194.
- Wan, E. A., & Van Der Merwe, R. (2000). The unscented Kalman filter for nonlinear estimation. In *Adaptive Systems for Signal Processing, Communications, and Control Symposium 2000. AS-SPCC. The IEEE 2000* (pp. 153-158).

- Wang, D. Q., Hey, J. M., & Nail, S. L. (2004). Effect of collapse on the stability of freeze-dried recombinant factor VIII and α -amylase. *J. Pharm. Sci.*, **93**, 1253-1263.
- Wang, Z., and Ierapetritou, M. (2017a). A novel feasibility analysis method for black-box processes using a radial basis function adaptive sampling approach. *AIChE J.* **63**:532-550.
- Wang, Z., and Ierapetritou, M. (2017b). A Novel Surrogate-Based Optimization Method for Black-Box Simulation with Heteroscedastic Noise. *Ind. Eng. Chem. Res.*, **56**:10720-10732.
- Wang, Z., Escotet-Espinoza, M. S., & Ierapetritou, M. (2017). Process analysis and optimization of continuous pharmaceutical manufacturing using flowsheet models. *Comp. Chem. Eng.* **107**: 77-91.
- Wassgren, C. and J.S. Curtis (2006). The application of computational modeling to pharmaceutical materials science. *MRS Bulletin*, **31**, 900-904.
- Watson, T. J., Nosal, R., Lepore, J., & Montgomery, F. (2018). Misunderstanding Design Space: a Robust Drug Product Control Strategy Is the Key to Quality Assurance. *J. Pharm. Innov.*, 1-3.
- WifOR (2016). *The economic footprint of selected pharmaceutical companies in Europe*. Available at: https://www.wifor.com/tl_files/wifor/PDF_Publikationen/161219_Efpia_EF_report_WifOR_updated.pdf (last accessed : 23/09/2018)
- Winkle, H.N. (2007). Implementing quality by design. *Proc. of PDA/FDA joint regulatory conference*, September 24-28 2007, Washington D.C., U.S.A.
- Wise, B. M., & Roginski, R. T. (2015). A calibration model maintenance roadmap. *IFAC-PapersOnLine*, **48**(8), 260-265.
- Wishart, J. (1928). The generalised product moment distribution in samples from a normal multivariate population. *Biometrika*, 32-52.
- Wold, S. (1978). Cross-validatory estimation of the number of components in factor and principal components models. *Technometrics*, **20**(4), 397-405.
- Wold, S., H. Martens and H. Wold (1983). The multivariate calibration problem in chemistry solved by the PLS method. *Lecture Notes in Math.*, **973**, 286-293.
- Wold, S., Sjöström, M., and Eriksson, L. (2001). PLS-Regression: A Basic Tool of Chemometrics. *Chem. Intell. Lab. Syst.*, **58**, 109-30.
- Woo, X. Y., Tan, R. B., & Braatz, R. D. (2008). Modeling and computational fluid dynamics–population balance equation–micromixing simulation of impinging jet crystallizers. *Crystal Growth and Design*, **9**, 156-164.
- Woodcock J., (2013). FDA Check Up: Drug Development and Manufacturing Challenges. Available at: <http://www.fda.gov/NewsEvents/Testimony/ucm378343.html> (last accessed: 26/09/2018)
- Wu, C. Y., & Cocks, A. C. (2006). Numerical and experimental investigations of the flow of powder into a confined space. *Mechanics of Materials*, **38**, 304-324.

- Wu, C. Y. (2008). DEM simulations of die filling during pharmaceutical tableting. *Particuology*, **6**, 412-418.
- Wu, C. Y., & Guo, Y. (2012). Numerical modelling of suction filling using DEM/CFD. *Chem. Eng. Sci.*, **73**, 231-238.
- Xie, L. I. N., Wu, H., Shen, M., Augsburger, L. L., Lyon, R. C., Khan, M. A., & Hoag, S. W. (2008). Quality-by-design (QbD): effects of testing parameters and formulation variables on the segregation tendency of pharmaceutical powder measured by the ASTM D 6940-04 segregation tester. *J. Pharm. Sci.*, **97**, 4485-4497.
- Yacoub, F. and J.F. MacGregor (2011). Robust processes through latent variable modeling and optimization. *AIChE J.*, **57**, 1278-1287.
- Yi, L. Y., Dong, K. J., Zou, R. P., & Yu, A. B. (2011). Coordination number of the packing of ternary mixtures of spheres: DEM simulations versus measurements. *Ind. Eng. Chem. Res.*, **50**, 8773-8785.
- Yoshino, H., Hara, Y., Dohi, M., Yamashita, K., Hakomori, T., Kimura, S. I., & Itai, S. (2018). A Scale-up Approach for Film Coating Process Based on Surface Roughness as the Critical Quality Attribute. *AAPS PharmSciTech*, **19**, 1243-1253.
- Yu L.X., Amidon, G., Khan, M. A., Hoag, S. W., Polli, J., Raju, G. K., & Woodcock, J. (2014). Understanding pharmaceutical quality by design. *The AAPS journal*, **16**, 771-783.
- Yu L.X., & Woodcock, J. (2015). FDA pharmaceutical quality oversight. *Int. J. Pharm.*, **491**, 2-7.
- Yu L.X., & Kopcha, M. (2017). The future of pharmaceutical quality and the path to get there. *Int. J. Pharm.*, **528**, 354-359.
- Yu, X., Hounslow, M. J., Reynolds, G. K., Rasmuson, A., Niklasson Björn, I., & Abrahamsson, P. J. (2017). A compartmental CFD-PBM model of high shear wet granulation. *AIChE J.*, **63**, 438-458.
- Zacour, B.M., J.K. Drennen III and C.A. Andreson (2012). Development of a fluid bed granulation design space using
- Zhai, S., Su, H., Taylor, R., & Slater, N. K. (2005). Pure ice sublimation within vials in a laboratory lyophiliser; comparison of theory with experiment. *Chem. Eng. Sci.*, **60**, 1167-1176.
- Zhang, L., and Garcia-Munoz, S. (2009). A Comparison of Different Methods to Estimate Prediction Uncertainty Using Partial Least Squares (PLS): A Practitioner's Perspective. *Chem. Intell. Lab. Syst.*, **97**, 152-58.
- Zhang, Y., & Fearn, T. (2015). A linearization method for partial least squares regression prediction uncertainty. *Chem. Intell. Lab. Syst.*, **140**, 133-140. am Ende, D.J., K. S. Bronk, J. Mustakis, G. O'Connor, C.L. Santa Maria, R. Nosal and T.J.N.

- Zhu, H. P., Zhou, Z. Y., Yang, R. Y., & Yu, A. B. (2008). Discrete particle simulation of particulate systems: a review of major applications and findings. *Chem. Eng. Sci.*, **63**, 5728-5770.
- Zhu, T., Moussa, E. M., Witting, M., Zhou, D., Sinha, K., Hirth, M., & Alexeenko, A. (2018). Predictive models of lyophilization process for development, scale-up/tech transfer and manufacturing. *European Journal of Pharmaceutics and Biopharmaceutics*, **128**, 363-378.
- Ziémons, E., J. Mantanus, P. Lebrun, E. Rozet, B. Evrard and P. Hubert (2010). Acetaminophen determination in low-dose pharmaceutical syrup by NIR spectroscopy. *J. Pharm. Biomed. Anal.*, **53**, 510-516.

Acknowledgments

Finally the acknowledgments!

There are many people who gave an important contribution to my journey as a PhD student that I would like to thank with deep gratitude.

First of all, I would like to thank my supervisor, prof. Massimiliano Barolo, as well as prof. Fabrizio Bezzo and dr. Pierantonio Facco, for their guidance and invaluable technical support throughout my entire career as a PhD student. Thank you for the several fruitful discussions and for the scientific and professional advice that you gave me during these three years. The most important lesson that I learned from you is that in science, as well as in life, attention to details is everything. I will treasure this teaching in my future professional career.

Thanks to prof. Marianthi Ierapetritou for accepting me as a guest during my stay at Rutgers University, as well as to all the friends and colleagues that I had the pleasure to meet there: Zilong, Lisia, Ou, Nirupa, Sebastian, Atharv, Parham. Thanks to all the friends of the AirBorne futsal team: Jose, Chris (x3), Eddie, Rosendo, Jason....we had a lot of fun!

Grazie ai miei genitori, Riccardo e Clarice, per avermi insegnato cosa significhi spirito di sacrificio e per avermi sempre aiutato e supportato nei miei studi. Siete i due esempi più importanti nella mia vita di ogni giorno. Grazie ai miei fratelli, Emanuele e Milena, per l'affetto che mi date. Grazie nonna Maria e zio Esto, che mi guardate da lassù: spero di avervi resi orgogliosi.

Grazie a tutti gli amici, vecchi e nuovi, che ho avuto l'onore di incontrare al CAPE-Lab: Pierantonio (che non è stato solo il mio mentore, ma anche un amico), Natascia, Filippo, Andrea "*Il Maestro*", Riccardo "*Number one*", Christopher "*Cooking water*", Myriam, Federico, Francesco e tutti i laureandi che ho avuto il piacere di incontrare durante i tre anni.

Last, but not least, thanks to Kimberly. Thank you for changing my life in such a short amount of time and for the love and support that you give me every day. Words are not enough to describe my gratitude towards you: thank you for making me understand that love has no distance.