

## Towards predictive models in food engineering: Parameter estimation dos and don'ts

Eva Balsa-Canto, Carlos Vilas, Ana Arias-Méndez, Míriam R. García, Antonio A. Alonso

*Bioprocess Engineering Group, IIM-CSIC, Vigo, Spain ([ebalsa@iim.csic.es](mailto:ebalsa@iim.csic.es))*

### ABSTRACT

Rigorous, physics based, modeling is at the core of computer aided food process engineering. Models often require the values of some, typically unknown, parameters (thermo-physical properties, kinetic constants, etc). Therefore, parameter estimation from experimental data is critical to achieve desired model predictive properties. Unfortunately, it must be admitted that often experiment design and modeling are fully separated tasks: experiments are not designed for the purpose of modeling and models are usually derived without paying especial attention to available experimental data or experimentation capabilities. When, at some point, the parameter estimation problem is put on the table, modelers use available experimental data to "manually" tune the unknown parameters. This results in inaccurate parameter estimates, usually experiment dependent, with the implications this has in model validation.

This work takes a new look into the parameter estimation problem in food process modeling. First the common pitfalls in parameter estimation are described. Second we present the theoretical background and the numerical techniques to define a parameter estimation protocol to iteratively improve model predictive capabilities. This protocol includes: reduced order modeling, structural and practical identifiability analyses, data fitting with global optimization methods and optimal experimental design.

And, to finish, we illustrate the performance of the proposed protocol with an example related to the thermal processing of packaged foods. The model was experimentally validated in the IIM-CSIC pilot plant.

*Keywords: Model calibration; Optimal experimental design; Parameter uncertainty; Food processes; Predictive models*

### INTRODUCTION

The food industry have increased efforts to adopt cost effective technologies that offer better quality and safe products motivated by consumer demands. The use of model based food engineering can contribute to the convergence of consumers and industry requirements and needs. In this respect the role of rigorous models in the context of food process simulation and design has been largely discussed (Bruin & Jongen, 2003).

Model validation is, however, still a challenge (Datta, 2008; Perrot et al., 2011; Trystram, 2012) mainly because of two reasons: (i) most of models available in the literature describing food processes consist of sets of partial and ordinary differential equations (PDEs and ODEs) that depend on unknown transport and kinetic properties. (ii) In general, experiments are not designed for the purpose of modelling. Often the available experimental data are used to "manually" tune the unknown parameters without the possibility of performing a quantitative analysis of the uncertainty associated with those estimates. This results in inaccurate parameter estimates, usually experiment dependent, with the implications this may have in the predictive capabilities of the model.

Alternatively, optimization based approaches have been suggested for a more systematic data based parameter estimation. The idea is to compute the parameter values that minimize the distance between the experimental data and the corresponding model predictions. Unfortunately, this can be a difficult task in the case of nonlinear PDE models due to the computational cost and the usual presence of multiple suboptimal solutions (multimodality). Reduced order modelling techniques, such as the proper orthogonal decomposition approach (Sirovich, 1987), may be used to reduce the computational burden. Multimodality may be surmounted by the use of global optimization methods (Rodriguez-Fernandez et al., 2007). In addition, a careful model based experimental design may contribute to reduce the experimental burden for the purpose of parameter estimation (Nahor et al., 2001; Van Derlinden et al., 2013).

Despite previous efforts, there are two critical questions in parametric identification that are often disregarded. The first is related to the possibility of giving unique values to the parameters, i.e. is the model structurally identifiable?. The second has to do with the uncertainty associated to the parameter estimates

and how this propagates to model predictions. Answering the first question is critical for the success of parameter estimation, while answering the second is critical for model validation. The first question is related to the mathematical structure of the model, while the second is related to the quality and quantity of experimental data.

Lack of structural identifiability, multimodality of the optimization problem or lack of appropriate experimental data for parameter estimation are usual pitfalls in food process modelling. In this scenario, it is necessary to organize available tools and methods in a protocol and to standardise data acquisition to estimate unknown transport and kinetic properties of different food materials, packages and processes by means of data fitting. This protocol should lead us to a successful model validation with minimum experimental effort.

The following protocol is presented in this work: (i) Model formulation and reduced order modelling; (ii) structural identifiability analysis; (iii) parameter estimation using global optimization methods; (iv) practical identifiability analysis to evaluate the uncertainty of the parameter estimates; (v) optimal experimental design to compute those experiments that maximize data quality for the purpose of parameter estimation; (vi) model validation. The performance of the protocol is illustrated and experimentally validated using a model that describes the thermal sterilization of solid packaged foods.

## MATERIALS & METHODS

### Common pitfalls in parameter estimation

**Lack of structural identifiability:** in some cases different parameter values lead to exactly the same fit. This may become critical if the range of possible solutions for the parameters is high or even infinite. Most of food processing models are distributed, dynamic and non-linear, thus being impossible to find an analytical solution for the observed quantities. This makes undoubtedly difficult to anticipate whether a unique value for the unknown parameters can be found attending exclusively to the mathematical structure of the model. However it is important to anticipate such a problem in order to improve model validation.

**Multimodality:** it is common practice to use standard least squares optimization methods to solve the parameter estimation problem. However, in the presence of suboptimal solutions, this type of methods does not guarantee convergence to the best solution (global optimum). Sub-optimal solutions lead to bad fits to the data. In this scenario, it will be rather difficult to distinguish between a problem in the model and suboptimal parameter estimation.

**Good fit but unsuccessful validation:** often, experiments are not designed for the purpose of modelling. Typically, factorial experimental designs are used to test the system responses to different levels of the factors influencing the process (e.g., temperature levels, processing times, etc.). Although quite intuitive and simple to implement, this technique presents a number of limitations. Specially when using the data to estimate parameters for dynamic models. It is necessary to design experiments for the purpose of modelling.

### Model identification protocol

The above mentioned pitfalls can be alleviated using appropriate tools. The identification protocol proposed in this work is suited to diagnose for possible sources of difficulties in parameter identification and to iteratively improve model predictive capabilities with the minimum experimental cost. The main steps of the proposed identification procedure are depicted in Figure 1.

At the core of the identification protocol we find the mathematical models constituted by sets of ordinary differential equations (ODEs). Spatially distributed processes, described by sets of partial differential equations (PDEs), can be numerically approximated by sets of ODEs using classical techniques like the *finite element method* (Reddy, 1993). However, such techniques are usually computationally demanding, particularly for 2D and 3D domains and for tasks that need to run the model hundreds or thousands of times, such as identifiability analysis; model calibration; optimal experimental design or process optimization. This work proposes the use of reduced order models obtained by means of the *proper orthogonal decomposition* (POD) approach (Sirovich, 1987) as an alternative to classical methods. This approach is based on capturing the most representative (slow) system dynamics while neglecting the fastest dynamics. The interested reader can find a deep insight, including a practical guidance on the implementation of the POD, in Garcia et al. (2007) and references therein.

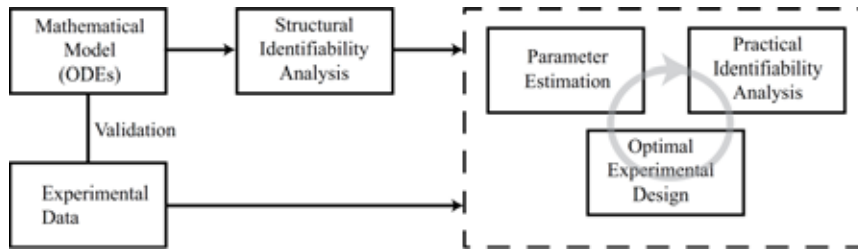


Figure 1. General scheme of the proposed identification protocol

The first step of the protocol is therefore to obtain such a reduced order model (ROM). Once the ROM is available the *structural identifiability analysis* aims to analyse whether the unknown parameters can be uniquely computed. Structural identifiability is independent of the parameter values as well as experimental data (perfect data are assumed) and it is related to model structure and type of control profiles. When the model is structurally unidentifiable, model reformulation is suggested since this problem cannot be solved by including more experimental data. Details on model structural identifiability analysis, and the methods for testing it, can be found in Chapman et al. (2003); Walter & Pronzato (1997). Once the model is confirmed to be structurally identifiable we proceed with the next step of the protocol, i.e. *parameter estimation* (Walter & Pronzato, 1997). The aim of this step is to compute the unknown parameter values that minimize the distance from experimental data to model predictions. A global optimizer eSS is suggested to facilitate convergence to the global optimum (Egea et al., 2009). In addition to the value of the parameters, it is important to compute its uncertainty since large parameter uncertainties may lead to unreliable model predictions. *Practical identifiability analysis* is then used to evaluate such uncertainty, in terms of confidence intervals, given a set of experimental data. A robust Monte Carlo based sampling method is suggested (Joshi et al., 2006). Finally, *optimal experimental design* (OED), aims at increasing parameter confidence by searching for the experimental schemes (control profiles, initial conditions, sampling times, etc) that provide the most informative data possible. The problem is formulated as a dynamic optimization problem and solved using the control vector parameterization approach together with a global optimizer (Balsa-Canto et al, 2008; Garcia et al., 2008). The optimally designed schemes are implemented in the process, the resulting data is employed to perform new parameter estimation, practical identifiability analysis and, if necessary, new OED. This procedure is repeated until parameter uncertainty is small enough to result into reliable model predictions and validation is satisfactory. All the steps of the identification protocol, except structural identifiability analysis, are implemented in the AMIGO toolbox (Balsa-Canto & Banga, 2011).

## RESULTS & DISCUSSION

In this section, the application of the proposed protocol will be illustrated using the steam sterilization of a packaged solid foodstuff as a case study. The food container is made of glass with a metallic cover. Under several assumptions (Alonso et al., 1997), the dynamics of the product temperature  $T(r,z,t)$  is described by the 2D Fourier's heat equation:

$$\rho c_p \frac{\partial T}{\partial t} = \kappa \left( \frac{1}{r} \frac{\partial}{\partial r} + \frac{\partial^2}{\partial r^2} + \frac{\partial^2}{\partial z^2} \right) T \quad (1)$$

where  $\kappa$ ,  $\rho$  and  $c_p$  are, respectively, the solid product thermal conductivity, density and specific heat.  $r$  and  $z$  stand for the spatial coordinates. Heat flux boundary conditions are considered:

$$\kappa \mathbf{n} \cdot \nabla T_g = h_g (T_r - T_g) \quad (2)$$

$$\kappa \mathbf{n} \cdot \nabla T_m = h_m (T_r - T_m) \quad (3)$$

with  $T_g$  and  $T_m$  being, respectively, the temperature in the glass and metal boundaries.  $h_g$  and  $h_m$  are the heat transfer coefficients.  $\nabla$  is the gradient operator whereas  $T_r$  is the steam retort temperature which can be measured and controlled. Finally,  $\mathbf{n}$  is a unit vector pointing outward the boundary.

The solution of the model with the FEM requires solving 462 ODEs what makes it unsuitable for tasks such as identifiability analysis, parameter estimation or OED. The POD technique was used to obtain a ROM. In the POD, the temperature field is approximated as a series of time dependent coefficients ( $m_i(t)$ ) and spatial dependent functions or basis functions ( $\varphi_i(r,z)$ ), i.e.

$$T(t, r, z) = \sum_{i=1}^n m_i(t) j_i(r, z) \quad (4)$$

The basis functions are computed *a priori* using *in silico* experimental data obtained via simulation of the FEM model. Time dependent coefficients are computed by projecting system (1)-(3) over the basis functions (Sirovich, 1987; García et al., 2007). The ROM reads as follows:

$$r c_p \frac{dm}{dt} = (kA_k + h_g A_g + h_m A_m) m + (h_g B_g + h_m B_m) T_r \quad (5)$$

where  $m$  is the vector of time coefficients.  $A_k$ ,  $A_g$ ,  $A_m$ ,  $B_g$  and  $B_m$  are given matrices. In this work,  $n = 13$  in equation (4) has been employed. This leads to solving 13 ODEs. The simulation results using the POD have been compared with the FEM under different operation conditions and using different values for the parameters. Less than 1% of difference between both techniques has been found. In the sequel, model (5) will be used to test the protocol for parameters  $\kappa$ ,  $\rho$ ,  $c_p$ ,  $h_g$  and  $h_m$ .

### Structural identifiability analysis

Inspection of model (5) tells us that: (i) parameters  $\rho$ ,  $c_p$  cannot be computed independently. Note that the same value  $\rho = \rho c_p$  can be obtained for different combinations of  $\rho$  and  $c_p$ ; (ii) quantity  $\rho = \rho c_p$  cannot be computed independently from the other model parameters. These issues are not related to the quantity or quality of the experimental measurements but only depend on the model structure. Therefore, the model must be reformulated. To that purpose, let us divide equation (5) by  $\rho c_p$ :

$$\frac{dm}{dt} = (\rho_k A_k + \rho_h A_g + \rho_m A_m) m + (\rho_g B_g + \rho_m B_m) T_r \quad (6)$$

In order to test the structural identifiability of model (6) with respect to parameters  $\rho_k$ ,  $\rho_g$ ,  $\rho_m$ , the GenSSI toolbox (Chis et al., 2011) was used. This toolbox is based on the generating series approach. The conclusion of the analysis was that using two thermocouples, i.e. measuring the temperature at two different spatial points, parameters are structurally globally identifiable. Two thermocouples were used in subsequent steps.

### Parameter estimation and identifiability analysis

The set of data required for parameter estimation was obtained from three different experiments, with constant retort temperature, carried out in a pilot plant. Experimental error was estimated using three replicates. The AMIGO toolbox was used to perform the parameter estimation, using eSS and to estimate the parameter confidence intervals via the Monte Carlo sampling method. It was observed that the objective function did not change significantly for values of  $\rho_m$  larger than  $1.5 \times 10^{-3} \text{ m s}^{-1}$ . Such a large value is equivalent to using a Dirichlet boundary condition in equation (3), i.e. the value of the temperature in the metallic cover is  $T_r$ . The results for the other two parameters are  $\rho_k = (4.74 \pm 0.29) \times 10^{-7} \text{ m}^2 \text{ s}^{-1}$  and  $\rho_g = (1.41 \pm 0.25) \times 10^{-5} \text{ m s}^{-1}$ . Model was compared against a new set of data. To improve predictive capabilities and to reduce uncertainty we proceeded with optimal experimental design.

### Model based optimal experimental design

The objective is to compute the experimental scheme (initial conditions, control profile, experiment duration, etc) that maximizes quality of the data. Zeroth-order and first order polynomials were considered to approximate the control variable. Constraints on such control were also considered to avoid package damage as a consequence of steep pressure drops. The best results were obtained using linear polynomials. The solution of the OED problem is represented in Figure 2(a). At the beginning of the experiment, the control temperature increases very rapidly (at the maximum rate allowed of  $1 \text{ }^\circ\text{C per min}$ ) and from minute 6 changes in the temperature are smoother. The optimal experiment duration was 30 min.

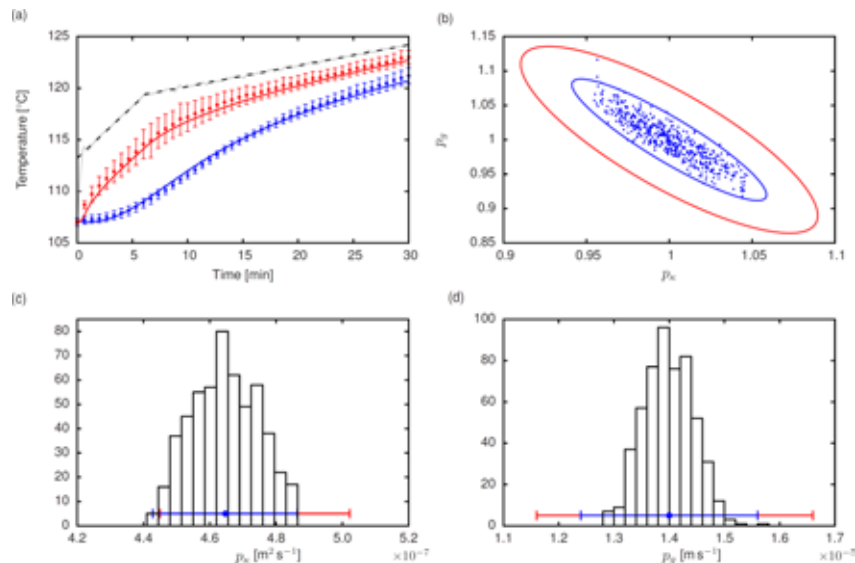


Figure 2. (a) OED profile result (dashed black line) and its implementation in the pilot plant (continuous grey line). Red and blue colours correspond to the locations close to the product top and the centre of the product. Marks are thermocouple measurements and lines are model predictions. (b) Cloud of points from the Monte Carlo sampling method and its associated hyper ellipsoid. (c) and (d) distributions of solutions for parameters  $p_k$  and  $p_g$ . The mean value of the replicates is represented by a circle whereas continuous horizontal black lines define the confidence region.

### New parameter estimation and identifiability analysis

The optimally designed experiment was implemented in the pilot plant and the recorded data were added to the previous experimental data to perform a new parameter estimation and identifiability analysis. Figure 2(a) shows how now the model (lines) is able to reproduce the pilot plant behaviour (marks). Again, a Monte Carlo sampling method was used to robustly compute the parameter confidence intervals associated to the parameters resulting in  $p_k = (4.64 \pm 0.22) \times 10^{-7} \text{ m}^2 \text{ s}^{-1}$  and  $p_g = (1.40 \pm 0.16) \times 10^{-5} \text{ m s}^{-1}$ . The cloud of points obtained from this procedure is used to compute the confidence hyper ellipsoid (see Figure 2(b)), since their volume and eccentricity are related, respectively, to identifiability and correlation issues. Figure 2(c) and (d) show the distributions of solutions obtained from the Monte Carlo method. Red lines in Figure 2(b)-(d) represent the confidence regions without the data from the OED. The new confidence intervals (continuous horizontal lines) are smaller (around a 25% for  $p_k$  and a 35% for  $p_g$ ).

### Model validation

In order to validate the model, a new experiment carried out under different conditions. The experimental data and model predictions are represented in Figure 3. As shown in the figure, the model is able to reproduce the plant behaviour, particularly for temperatures larger than 100 °C. It should be stressed that some of the assumptions made to derive the model equations are not fulfilled below 100 °C.

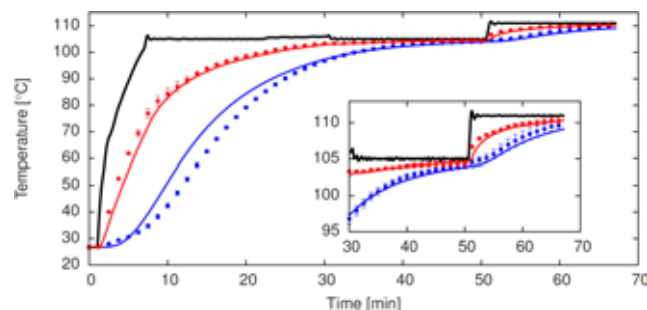


Figure 3. Model validation results. Red and blue colours refer, respectively, to a location close to the product top and the centre of the product. Marks and continuous lines correspond, respectively, to experimental data and model simulation. Black line is the retort temperature.

## CONCLUSION

This work presents a data-based protocol for model identification in food engineering. The protocol is intended to diagnose and to surmount the most common difficulties found during model parametric identification of food process models.

The protocol was experimentally validated by means of an example related to the thermal sterilization of packaged foods. Experiments were performed in our pilot plant at the IIM-CSIC. The protocol started with the definition of a reduced order model which is as accurate as a finite element based description and ended with a validation experiment. Results reveal that the protocol leads to an iterative improvement of model predictive capabilities while minimizing uncertainty.

The application of the protocol is supported by software tools available from the authors. Methods and tools are general in the sense that can be applied to any other food process model.

## REFERENCES

- Chapman M.J., Godfrey K.R., Chappell M.J. & Evans N.D. 2003. Structural identifiability for a class of non-linear compartmental systems using linear/non-linear splitting and symbolic computation. *Mathematical Biosciences*, 183, 215–215.
- Alonso A.A., Banga, J.R. & Pérez-Martín R.I. 1997. A complete dynamic model for the thermal processing of bioproducts in batch units and its application to controller design. *Chemical Engineering Science*, 52, 1307–1322.
- Balsa-Canto, E., A.A. Alonso and J.R. Banga. 2008. Computing optimal dynamic experiments for model calibration in predictive microbiology. *Journal of Food Process Engineering*, 31:186-206.
- Balsa-Canto E. & Banga J.R. 2011. AMIGO, a toolbox for advanced model identification in systems biology using global optimization. *Bioinformatics*, 27, 2311–2313
- Bruin S. & Jongen T.R.G. 2003. Food process engineering: The last 25 years and challenges ahead. *Comprehensive Reviews in Food Science and Food Safety*, 2, 42–81.
- Chis O., Banga J.R. & Balsa-Canto E. 2011. GenSSI: a software toolbox for structural identifiability analysis of biological models. *Bioinformatics*, 27, 2610–2611.
- Egea, J. A., Vazquez, E., Banga, J. R., and Marti, R. 2009. Improved scatter search for the global optimization of computationally expensive dynamic models. *Journal of Global Optimization*, 43(2-3):175-190.
- García M.R. 2008. Identification and Real Time Optimisation in the Food Processing and Biotechnology Industries. PhD thesis, University of Vigo, Spain.
- García M.R., Vilas C., Banga J.R. & Alonso, A.A. 2007. Optimal field reconstruction of distributed process systems from partial measurements. *Industrial & Engineering Chemistry Research*, 46, 530–539.
- Joshi M., Seidel-Morgenstern A. & Kremling A. 2006. Exploiting the bootstrap method for quantifying parameter confidence intervals in dynamical systems. *Metabolic Engineering*, 8, 447–455.
- Nahor H.B., Scheerlinck N., Verniest R., De Baerdemaeker J. & Nicolai B.M. 2001. Optimal experimental design for the parameter estimation of conduction heated foods. *Journal of Food Engineering*, 48, 109–119.
- Perrot N., Trelea I.C., Baudrit C., Trystram G. & Bourguin P. 2011. Modelling and analysis of complex food systems: State of the art and new trends. *Trends in Food Science and Technology*, 22, 304–314.
- Reddy J.N. 1993. *An Introduction to the Finite Element Method*. McGraw-Hill, 2nd edition.
- Rodríguez-Fernández M., Balsa-Canto E., Egea J.A. & Banga, J.R. 2007. Identifiability and robust parameter estimation in food process modeling: Application to a drying model. *Journal of Food Engineering*, 83, 374–383.
- Sirovich L. 1987. Turbulence and the dynamics of coherent structures. Part I: Coherent structures. *Quarterly of Applied Mathematics*, 45, 561–571.
- Trystram G. 2012. Modelling of food and food processes. *Journal of Food Engineering*, 110, 269–277.
- Van Derlinden E., Mertens L. & Van Impe, J. 2013. The impact of experiment design on the parameter estimation of cardinal parameter models in predictive microbiology. *Food Control*, 29, 300–308.
- Vanrolleghem P.A. & Dochain D. 1998. Bioprocess model identification. In: van Impe J.F.M., Vanrolleghem P.A., & Iserentant D.M. (Eds.). *Advanced instrumentation, data interpretation, and control of biotechnological process*. Inc. Kluwer Academic Publishers.
- Walter E. & Pronzato L. 1997. *Identification of parametric models from experimental data*. Springer.