# Edinburgh Research Explorer

# Prediction of the performance of pre-packed purification columns through machine learning

# Prediction of the performance of pre-packed purification columns through machine learning

**SCHOLARONE™**
Manuscripts

1    **Prediction of the performance of pre-packed purification columns through**

2    **machine learning**

3

4    Qihao Jiang[1], Sohan Seth[2], Theresa Scharl[3,4], Tim Schroeder[5], Alois Jungbauer[3,6],

5    Simone Dimartino[1],*

6

7    [1]Institute of Bioengineering, The School of Engineering, The University of
8    Edinburgh, Edinburgh, EH9 3DW, UK
9

10    [2]The School of Informatics, The University of Edinburgh, Edinburgh, EH9 3DW, UK
11

12    [3]Austrian Centre of Industrial Biotechnology, Vienna, Austria
13

14    [4]Institute of Statistics, University of Natural Resources and Life Sciences Vienna,
15    Vienna, Austria
16

17    [5]Atoll, Weingarten, Germany
18

19    [6]Department of Biotechnology, University of Natural Resources and Life Sciences,
20    Vienna, Austria
21

22    *Corresponding author: Dr Simone Dimartino, Faraday Building, Colin McLaurin
23    Road, The King's Buildings, Edinburgh EH9 3DW, UK, Phone: +44 131 6507305,
24    Email: simone.dimartino@ed.ac.uk
25

26

27    **Keywords:**

28    Asymmetry, Plate Height, Machine learning, Pre packed columns, Porous media

29

30

31

32

33

34

35 **Abstract**

36 Pre-packed columns have been increasingly used in process development and

37 biomanufacturing thanks to their ease of use and consistency. Traditionally, packing

38 quality is predicted through rate models, which require extensive calibration efforts

39 through independent experiments to determine relevant mass transfer and kinetic rate

40 constants. Here we propose machine learning as a complementary predictive tool for

41 column performance. A machine learning algorithm, extreme gradient boosting, was

42 applied to a large data set of packing quality (plate height and asymmetry) for pre-packed

43 columns as a function of quantitative parameters (column length, column diameter,

44 particle size) and qualitative attributes (backbone and functional mode). The machine

45 learning model offered excellent predictive capabilities for the plate height and the

46 asymmetry (90% and 93%, respectively), with packing quality strongly influenced

47 by backbone (~70% relative importance) and functional mode (~15% relative

48 importance), well above all other quantitative column parameters. The results highlight

49 the ability of machine learning to provide reliable predictions of column performance

50 from simple, generic parameters, including strategic qualitative parameters such

51 as backbone and functionality, usually excluded from quantitative considerations. Our

52 results will guide further efforts in column optimization, e.g. by focusing on

53 improvements of backbone and functional mode to obtain optimised packings.

54

55

56

57

58

59

60

61

62

63

64

## 1. Introduction

Pre-packed chromatography columns are widely employed in process development and biomanufacturing. Their biggest advantage is to take away the burden of costly and time consuming packing procedures and associated validation protocols, ultimately ensuring a consistent product [1–4] . The production of pre-packed columns should be simple, cost-effective, and robust over the long term (decades) to ensure consistent quality of columns.

The performance of pre-packed columns is assured by the manufacturer before the sale, with packing quality measured in terms of the height equivalent to a theoretical plate (HETP) and asymmetry. Both parameters are calculated from the response of the column following a pulse injection of a non-binding tracer, i.e. residence time distribution (RTD) experiments. The HETP corresponds to the column length over the number of theoretical plates (N), with efficient columns characterized by relatively large N and small HETP values. According to the general rate model, the RTD response of a "well-packed" column is a symmetrical Gaussian peak. To better assess packing quality, RTD experiments are usually run under conditions for which hydrodynamic dispersion is the dominant contribution to mass transfer (negligible intraparticle mass transfer, no adsorption). Under these conditions (reduced velocity of about 1 to 10), the HETP The minimum HETP value theoretically depends only on the properties of the tracer, the velocity of the mobile phase, and the size of the chromatographic particles [5]. However, the general rate model is unable to capture how the HETP is influenced by key factors of practical relevance such as column size (column diameter and length) or ease of packing across different chromatographic resins [6]. For example, Scharl et al. [7] qualitatively discussed the importance of material backbone on packing quality of a range of pre-packed columns. Deviations from symmetrical peaks are often observed in practice, with peak fronting or tailing associated with a number of non-idealities such as wall effects, inhomogeneous packing, inhomogeneous distribution of the solute over the bed

3

95    at the column inlet/distributor and/or at the outlet/collector, and dispersion in

96    the extra column volumes [8–12]. Such deviations are measured through the

97    asymmetry, an empirical parameter used to quantify the degree of peak skewness

98    and employed to assess packing quality in tandem to the HETP [13].

99

100   Mathematical models to predict column performance and chromatographic

101   processes, including the general rate model, are generally based on first principles.

102   In particular, they include details of mass transfer phenomena and binding

103   kinetics to describe peak profiles and breakthrough curves [14,15]. While the

104   predictive power of these models is often excellent, they require extensive

105   calibration efforts through independent experiments, e.g. to determine key model

106   parameters such as mass transfer and kinetic coefficients [16,17]. Flow non-

107   idealities such as wall effects and distribution/collection of the fluid at the column

108   inlet/outlet also require independent experiments for them to be accounted for in

109   the models. These additional experiments are specific to the chromatographic

110   system (external column volumes) and column (diameter, length) employed,

111   therefore cannot be extrapolated to different systems or different columns. Finally,

112   such models based on first principles do not take into account qualitative variables

113   such as resin backbone and functional chemistry by design.

114

115   Machine learning (ML) could represent an alternative modelling approach to

116   analyze and predict column performance. The main advantage of ML is the ability

117   to extract information from large data sets using no or only minimum assumptions,

118   eventually determining generalizable predictive patterns between multiple inputs

119   (including quantitative, qualitative and categorical parameters), and the output

120   variables [18,19]. A number of algorithms, e.g. support vector machine, decision

121   tree, gradient boosting, and deep neural networks have been developed over the

122   years, and have proved their ability in dealing with complex data problems in a

123   practical manner [20,21]. ML has been applied to chromatography systems, with

124   many successful applications e.g. in peak observation [22–24], retention

4

125    modelling [25–28], process optimization[29–31], and real-time process

126    monitoring [32,33]. The main challenge associated to the application of ML is the

127    availability of very large experimental data sets for the ML algorithm to draw

128    meaningful correlations.

129

130    In this work, we consider a large data set of around 25,000 quality assurance

131    experiments of pre-packed columns manufactured and tested under standardized

132    conditions for a period of over 10 years [7]. We first examine the time series of the

133    data set using correlation and autocorrelation analysis to ensure the data are self-

134    consistent and time-independent. We then employ ML methods to find a

135    correlation between column performance (measured in terms of HETP and

136    asymmetry) and qualitative as well as qualitative column variables, namely resin

137    backbone, functionalization chemistry, column size (length and diameter) and

138    particle size. The results are finally commented in relation to the main key

139    variables affecting column performance.

140

141    **2. Materials and Methods**

142    **2.1 Experimental Data Set**

143    The data set employed in this work is a subset of that previously employed by

144    Scharl et al consisting of 24,951 quality control runs of pre-packed small-scale

145    columns over a period of about 10 years [7]. The data contain relevant column

146    parameters (i.e. column length and diameter, particle size, backbone material,

147    functional mode, and date of testing) together with reduced HETP ($h$) and

148    asymmetry ($A_s$). Column diameter and length ranged between 5 – 11.3 mm and

149    10 – 100 mm, respectively, while particle diameter varied between 15 to 400 $\mu$m.

150    2232 experimental runs (approximately 10%) were removed from the original

151    data set as they lacked one or more column parameter inputs, reducing the data

152    set to a total of 22,359 tests. Columns with same attributes were manufactured

153    and tested more than once over the ten year monitored, with some popular types

154    examined hundreds of times (e.g. see table 1 in SI). All experiments having same

155    set of input features were treated as a single entry, with $h$ and $A_S$ averaged over

156    the available runs for that column type. This step was necessary to prevent data

157    leakage in the ML model, i.e. the use of same column type in both the training and

158    testing data sets (see 2.3), as well as to prevent overfitting of the most popular

159    column types over the ones infrequently produced. The standard error for $h$ and

160    $A_S$ was always lower than 10%, indicating that the average $h$ and $A_S$ are

161    representative output indicators of column performance for any given column

162    type. After the averaging process, the data set contained a total of 546 independent

163    runs.

164

165    All columns used to generate the data set were packed by slurry packing under

166    vibration following a standardized procedure developed by the packing company

167    (Atoll, now Repligen). The packing quality of the columns was evaluated using a

168    standardised experimental set up and experimental protocol as reported in Scharl

169    et al [7]. Briefly, the response of the column following an acetone or sodium nitrate

170    injection was measured, and the resulting chromatographic peak analysed to

171    extract $h$ and $A_S$. This simple experiment allowed to isolate the contribution to

172    band broadening associated with hydrodynamic dispersion (which in turn

173    depends on packing quality and extra column dispersion) as the tracers employed

174    are both non-retained (i.e. zero retention factor), with practically same diffusion

175    coefficients ($1.2 \times 10^{-5}$ and $1.3 \times 10^{-5}$ cm$^2$/s for acetone [34] and sodium nitrate

176    [35], respectively), and tested under reduced velocities comprised between 1 and

177    20 for which the minimum HETP is obtained [14].

178

179    **2.2 Extreme Gradient Boosting**

180    Extreme gradient boosting (XGBoost) is a scalable ML system for tree boosting

181    [36]. XGBoost is a decision-tree-based ensemble learning method [37] that

182    provides a systematic solution to a given problem by combining the predictive

183    power of several different or same ML algorithms. The algorithm used in XGBoost

184    is the Classification and Regression Tree (CART ) [38] which employs a binary tree

6

185  that can be constantly segmented by data features, thus enabling dynamic growth

186  of the tree. The characteristics of the input data will eventually fall into the leaf

187  nodes n the tree, where each leaf node corresponds to a specific score, and the sum

188  of the scores in all the leaf nodes computes the final prediction value of a certain

189  feature, e.g., $h$ or $A_s$.

190

191  Essential details of the mathematical formulation of the XGBoost model are

192  presented in the following, with additional details in the SI. For a given data set

193  with $n$ examples and $m$ features $\mathcal{D} = \{(x_i, y_i)\} | (|\mathcal{D}| = n, x_i \in \mathbb{R}^m, y_i \in \mathbb{R})$, the tree

194  ensemble model uses $K$ additive functions to predict the output.

195

196
$$\hat{y}_i = \sum_{k=1}^{K} f_k(x_i), f_k \in \mathcal{F} \#(1)$$

197

198  where $\mathcal{F} = \{f(x) = w_{q(x)}\}(q:\mathbb{R}^m \to T, w \in \mathbb{R}^T)$ is the space of the regression trees.

199  The $q$ represents the structure of each tree that maps an example to the

200  corresponding leaf index. $T$ is the number of leaves in the tree. Each $f_k$

201  corresponds to an independent tree structure $q$ and tree weights $w$. More

202  mathematical details can be found from the original XGBoost paper [36].

203

204  The regularized objective function defined for XGBoost, $\mathcal{L}$, can be written as:

205

206
$$\mathcal{L} = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) \#(2)$$

207
$$\Omega(f) = \gamma \cdot T + \frac{1}{2}\lambda \|w\|^2 \#(3)$$

208

209  Here $l$ is a differentiable convex loss function that measures the difference

210  between the prediction $\hat{y}_i$ and the target $y_i$. The second $\Omega$ term prevents

211  unnecessary large trees by penalizing the complexity of the model, in turn

7

212   avoiding overfitting. The additional regularization term $\frac{1}{2}\lambda\|w\|^2$ helps smooth the

213   final learnt weights. The shrinkage parameter $\gamma$ is an additional design to prevent

214   over-fitting. The $\gamma$ is utilized to multiply the score of each leaf node by a reduction

215   weight during the iteration, which ensures that the influence of each tree is not

216   too large, leaving more space for the tress generated later to optimize.

217

218   XGBoost is also used to determine the relative importance of the input features.

219   The definition of the relative importance is followed by the study of H. Friedman

220   [39]. For a tree model whose number of terminal nodes is $J$, the relative

221   importance of a given input feature, $I$, is calculated by the sum of the

222   corresponding empirical improvements, $i^2$, with $t$ referring to a non-terminal

223   node and $v_t$ acting as splitting variable for that node. The $i^2$ term is determined

224   from the two sub-region $R_l$ and $R_r$, where $\bar{y}_l$ and $\bar{y}_r$ are the response means,

225   respectively, and $w_l$ and $w_r$ are the corresponding sums of the weights. In Python,

226   the contribution of each input features can be automatically transferred into the

227   percentage version.

228

229
$$I_j^2(T) = \sum_{t=1}^{J-1} i_t^2(v_t = j) \#(4)$$

230
$$i^2(R_l, R_r) = \frac{w_l w_r}{w_l + w_r}(\bar{y}_l - \bar{y}_r)^2 \#(5)$$

231

232   **2.3 Data Pre-Processing and Model Implementation**

233   Functional modes and backbone are two categorical features which cannot be

234   operated by many ML algorithms directly. One-hot encoding was applied to

235   transfer them into numerical values [40], with each feature normalized between

236   0 and 1. All other numerical parameters were also normalized between 0 and 1

237   before input into the ML model as most ML algorithms perform better or converge

238   faster with features on relatively similar scale [41].

239

240 An XGBoost regression model was created in Python 3.6 combining i)

241 GridSearchCV (ten-folds) to select and determine the model's hyper-parameters

242 (e.g. learning rate, maximum tree depth, and minimum child weight) [42] and ii)

243 XGBRegressor as the main package to process our data set [43]. The whole data

244 set was then separated randomly into a training set (66.7%) and testing set

245 (33.3%), with the training set utilized for training the ML model and the testing

246 set used for inspecting the final model accuracy. Mean absolute error (MAE) [44]

247 was used as the evaluation metric during model training. The final prediction

248 precision of the model is reported by the mean absolute percentage error (MAPE)

249 between the prediction results and the testing data set. The overall model

250 prediction capability remained the same when changing initial seeding to

251 randomly generate different training and testing data sets.

252

253 **3. Results and Discussion**

254 The main goal of this study was the identification of a general relationship

255 between column parameters (column length, column diameter, particle diameter,

256 functional mode, backbone material) and chromatographic performance (reduced

257 HETP, $h$ and peak asymmetry, $A_s$) using ML algorithms as an alternative to classical

258 rate models for chromatography. Classical rate models are derived from first

259 principles and thus tend to be the preferred choice when it comes to the modelling

260 of chromatographic separations. However, some of the parameters entering rate

261 models often are either determined through empirical expressions (e.g. the

262 Wilson Geankopolis correlation for the estimation of the mass transfer coefficient

263 [45]) or simply adjusted to best-fit experimental results (e.g. diffusion or

264 dispersion coefficients [46]).

265

266 The introduction of a certain degree of empiricism in physical models is necessary

267 to capture important elements of the model hard to describe in mathematical

268 terms. For example, the three dimensional configuration of chromatographic beds

9

269  deviates from the theoretical close random packing limit [47], with the resulting

270  bed arrangement strongly influenced by attributes linked to the material and

271  column properties (e.g. Young modulus, friction factor, wall roughness) as well as

272  the packing procedure itself [47,48] . For example, Knox demonstrated that

273  hydrodynamic dispersion in columns packed with smooth non-porous glass beads

274  is smaller than those measured in columns packed with porous glass [49]. Knox

275  explained this result in terms of bed homogeneity, and speculated that smooth

276  glass particles are able to form relatively regular packings, while porous glass

277  particles are affected by greater interparticle friction forces, in turn resulting in

278  particle bridging and the formation of pockets where local mixing occurs. These

279  insights were demonstrated experimentally by Patel and coauthors [50], who

280  confirmed that the A term in the van Deemter equation is primarily associated

281  with radial heterogeneities in the bed. On the opposite front, Malkin et al. showed

282  that submicrometer silica particles tend to pack close to the limit of a face centered

283  cubic arrangement [51], resulting in reduced plate heights below 1. Khirevich et

284  al. also reported that the local microscopic disorder in packings was highly

285  correlated with eddy dispersion, directly affecting column performance [52].

286  Along the same line, Gritti et al [53] reported the outstanding performance of

287  columns packed with core-shell particles, partly attributing these results to the

288  propensity that these particles have to create homogeneous beds. More recent

289  studies on 3D printed ordered beds further confirm the advantages of perfectly

290  ordered packing, with simulated reduced plate heights below 0.1 for specific

291  arrangements (e.g. octahedral particles in simple cubic configuration) of non-

292  porous stationary phases under non-retained conditions [54].

293

294  The concept of "goodness of packing" as proposed by Knox is strongly correlated

295  to the A term of the van Deemter equation [55], with lower A values associated to

296  lower reduced plate heights and hence higher chromatographic efficiency.

297  According to the general rate model for chromatography, the A term can be

298  expressed as [16]:

10

299

$$A = 2\chi d_p \#(6)$$

301

302    or in dimensionless terms:

303

$$a = 2\chi \#(7)$$

305

306    where $d_p$ is the average particle diameter and $\chi$ is the dispersivity of the stationary

307    phase. The dispersivity is a characteristic determined by the hydrodynamics in the

308    column, in turn defined by type of particles and their packing. For a given column,

309    the dispersivity can be determined through estimation of the plate height under

310    conditions suppressing both axial diffusion (i.e. large velocity, negligible B term)

311    and mass transfer and kinetic resistances (i.e injection of a small, fast diffusing

312    non-adsorbing tracer, negligible C-term) for which the van Deemter equation

313    reduces to:

314

$$h = a = 2\chi \#(8)$$

316

317    While this equation represents a relatively rapid method to assess the

318    hydrodynamic properties of a given column, lack of correlations for the estimation

319    of the dispersivity coefficient represents a limitation to predict band broadening

320    due to axial dispersion. In particular, there exist no quantitative method to assess

321    how the dispersivity depends on different column properties such as:

322    -    backbone material and functional mode, closely related to the propensity

323        of the particles to generate regular packing;

324    -    column and particle diameters, i.e. the column to particle ratio, in turn

325        determining the importance of non-homogeneities close to the column wall

326        with respect to the rest of the column volume;

327    -    column length and column diameter, which are associated to both bed

328        compressibility [56], as well as defining the relative influence of extra-

11

329    column dispersion effects, e.g. due to non-uniformity of the velocity profile

330    resulting from non-idealities in the extra-column volumes.

331

332    Fronting or tailing deviations from the ideal symmetrical peak are often observed

333    in chromatographic practice, negatively impacting the separation performance.

334    Such deviation is often quantified through the asymmetry factor, $A_s$, defined as the

335    ratio between the width of the tailing end and of the peak front at 10% peak height

336    [57,58]. Large asymmetry factors are associated with the heterogeneity of the

337    column packing [59,60], making $A_s$ another excellent descriptor for "goodness of

338    packing". However, search for a quantitative relationship between asymmetry and

339    column parameters has been elusive so far. In this context, ML is an excellent tool

340    to extract poorly understood links between variables such as the column input

341    parameters and the asymmetry factor.

342

343    The data set of pre-packed column performance offers an opportunity to

344    quantitatively analyse the dependence of the dispersivity on a range of qualitative

345    and quantitative column attributes. The two performance parameters, $h$ and $A_s$,

346    are measured from the experimental response of an injection of a small non-

347    retained tracer (acetone or sodium nitrate). Same experimental and data analysis

348    methods were used to generate the entire data set [7]. Only resins intended to

349    separate proteins or other larger biomolecules were tested, ensuring much larger

350    pores than that of the tracers. Such conditions ensure only the hydrodynamic

351    dispersion is captured in the experiments and that the Van Deemter equation can

352    be simplified into Eq. 8.

353

354    In short, we propose here to employ ML as a powerful alternative to traditional

355    chromatographic models to investigate a correlation between the different

356    column input parameters and the output performance parameters. ML is

357    especially valuable in this context given the complexity of the problem described

358    and the qualitative nature of some of the relevant variables such as column

12

359  backbone and functional mode. ML is also able to suggest the relative importance

360  of the different inputs with respect to the outputs, thus helping the identification

361  of the key descriptors for the performance parameters.

362

363  **3.1 Time Series of Reduced Plate Height and Asymmetry**

364  Column performance can change over time due to variations in the manufacturing

365  line, e.g. improvement in the packing procedures, change of suppliers of raw

366  materials, and ageing of the production line. Scharl et al. qualitatively observed

367  that the plate height of the prepacked columns tested was stable over ten years

368  [7]. However, any interdependence between $h$ and $A_s$ with time needs to be either

369  identified or excluded in quantitative terms to avoid any input bias to the ML

370  model. In other words, it is first necessary to determine if time represents an input

371  variable to the ML model, as well as if sampling and testing of the columns changed

372  significantly over time. Autocorrelation and partial autocorrelation analysis was

373  employed onto the data set to address these two aims, respectively. In particular,

374  the autocorrelation function (acf) aims to detect cross-similarities of a signal with

375  itself at a different time (time lag) [61]. In this context, acf helps detect changes in

376  the manufacturing line and in the quality assurance protocols employed over time.

377  The partial autocorrelation function (pacf) instead aims to identify the possibility

378  of confounding variables which are correlated to both variables [62]. In this

379  instance, pacf aims to identify a correlation between time and performance

380  parameters, in turn suggesting if a specific pattern of column types was

381  manufactured over time. Additional details on acf and pacf are also provided in

382  the SI.

383

384  The $h$ and $A_s$ time series were first resampled by averaging the data set in day

385  intervals, irrespective of the other column parameters. Other than reducing noise,

386  resampling is customary when autocorrelation analysis is executed over large

387  time periods [58,61].

388

13

389    Figure 1 shows the time series of the two performance parameters, $h$ and $A_s$. Over

390    the 10 year time considered, the $h$ values varied between about 7.8 and 2.2, with

391    an average of around 4.5. Variability reduced significantly from 2011 onwards,

392    with a slight decrease of plate height in 2012–2013. The asymmetry ranged

393    between about 2 and 0.8, with average of 1.1. Similar to plate height, the scatter in

394    the asymmetry over the first five years is larger than after 2011. According to

395    Scharl et al. [7], industrial quality assurance tests require a column to have $h$

396    comprised should be smaller than 5 in industry, while the acceptable range for $A_s$

397    is between 0.8 and 1.6 [7]. The observed variability is a natural consequence of

398    industrial manufacturing, yet the columns produced were within specifications in

399    terms of both $h$ and $A_s$.

400

401    Figure 2 shows the results from the autocorrelation and partial autocorrelation

402    analysis on $h$ and $A_s$ using lag time of days up to one year. Other lag times were

403    also examined (i.e. weekly, monthly as well as over 2, 3 months) with no significant

404    difference. For both $h$ and $A_s$, almost all of the acf and pacf coefficients lie within

405    the 95% confidence interval. The low acf demonstrates that the dataset does not

406    have a specific pattern with time, quantitatively confirming that the

407    manufacturing line was stable over the ten year period here investigated [63]. In

408    addition, low pacf rules out the existence of confounding variables such as certain

409    patterns in terms of column sampling and testing over time. In other words, pacf

410    analysis confirms that column manufacture was unbiased, excluding the

411    possibility that a certain column type (e.g. having specific size and packed with a

412    specific particle) was manufactured predominantly over other columns over time.

413    Overall, acf and pacf demonstrate that all performance tests were time

414    independent, making the data set solely dependent on the five input parameters

415    of particle size, column diameter, column length, column backbone, and functional

416    mode.

417

418    **3.2 Influence of column parameters on packing quality**

14

419    XGBoost was utilized to assess the influence of the column parameters (i.e. the

420    inputs to ML algorithm: particle size, column length, column diameter, functional

421    mode, and resin backbone) on packing quality (i.e. ML outputs of $h$ and $A_s$). Other

422    ML algorithms such as artificial neural networks and decision-tree were also

423    employed in a preliminary model assessment (refer to SI for additional

424    information on ML models). XGBoost consistently provided the highest predictive

425    precision, mainly due to its regularization and shrinkage terms (Eqs. 2 and 3)

426    being capable of curbing over-fitting, the main cause of poor prediction.

427

428    Figure 3 summarizes the results obtained with the XGBoost model to predict the

429    experimental data. In particular, Figure 3a and 3b compare the predicted $h$ and $A_s$,

430    respectively, against the observed data of the testing data set. The predictions are

431    in good agreement with the experimental results, where the mean absolute

432    percentage error (MAPE) of predicted results to the observed values are 10% for

433    $h$ and 7% for $A_s$, with a few outliers in the 40% ∼ 50% range. These acceptable

434    errors confirm that the XGBoost model can be applied to this problem with good

435    prediction accuracy. Figure 3c and 3d report the contribution importance, $J$ (Eqs.

436    4 and 5), of the various input parameters to predict the model outputs.

437    Interestingly, column backbone resulted as the most important descriptor of

438    packing quality, accounting for 68.4% and 77.0% for the prediction of $h$ and $A_s$,

439    respectively. Functional mode was the second most significant descriptor for the

440    estimation of packing quality, accounting for about 15% contribution importance,

441    followed in various order by the other parameters (particle size, column diameter

442    and column length). Violin plots were employed to further analyze the correlation

443    between input features and column performance (Figure 4). A violin plot is an

444    extension of a box and whisker plot, clearly recognizable inside the "violins",

445    decorated with a curve whose width is related to the probability density.

446

447    **3.2.1 Resin backbone**

448  Resin backbone was the most influential parameter for the prediction of packing

449  quality. The material making up the resin backbone can be either inorganic ,

450  synthetic polymer, or natural polymer. The nature of the material employed

451  determines a number of properties such as surface roughness of the particles [64],

452  particle size distribution (linked to the manufacturing method)[65], occurrence of

453  microstructural defects, and other mechanical properties such as Young modulus

454  and density [63,64]. All these factors impact column packing, either directly or

455  indirectly, in turn influencing the homogeneity of the resulting chromatographic

456  bed, i.e. packing quality. Johnson et al. examined a range of resin materials

457  (agarose, cellulose, ceramic) through X-ray computed tomography (CT) and

458  focused ion beam (FIB) [66]. They highlighted clear variations in the chemical,

459  physical and mechanical properties of the different materials. Our analysis with

460  the XGBoost model also confirms that resin characteristics strongly influence

461  chromatographic performance.

462

463  Figure 4a and 4b present violin plots of $h$ and $A_s$, respectively, over the eight

464  different backbones tested. It is possible to observe that certain backbones have

465  worse performance than others as measured by both of the two packing quality

466  parameters $h$ and $A_s$. For example, polystyrene-divinilbenzene (PS/DVB),

467  inorganic support (IS) and dextran (DEX) have data widely distributed, with

468  average $h$ above 5 and average $A_s$ above 1.2. On the other hand, agarose, cellulose

469  and PVE hydrophilic (PVE) demonstrated consistent results (little data scatter)

470  with average $h$ and $A_s$ well below the arbitrary thresholds of 5 and 1.2, respectively.

471  This analysis clearly demonstrates the importance of backbone selection, e.g.

472  during process or method development.

473

474  It is worth noting that inorganic support (IS) was relatively popular in the first

475  three years of our data set, while polyvinyl-ether hydrophilic (PVE) matrices were

476  little used at first, becoming more mainstream after 2011. This change in

477  backbone population over time can partly explain the slight decrease of the

16

478    absolute value of $h$, as well as the reduced scatter of $h$ and $A_s$ observed from 2011

479    onwards (Figure 1).

480

481    **3.2.2 Functional mode**

482    Functional mode was the second most important parameter to predict packing

483    quality. Figure 4c and 4d show the relation between $h$ and $A_s$ over the different

484    functional modes. The influence of the functionalization chemistry on column

485    packing is less intuitive than for chromatographic backbone. Stickel and

486    Fotopoulos [67] reported the difference of the pressure-flow profiles between

487    sepharose and phenyl sepharose, which was associated to the differing

488    hydrophobic and electrostatic character of the resin beads. Electrostatic and

489    hydrophobic interactions might promote local or temporary bonding of two or

490    more particles into clusters, decreasing the degrees of freedom of the slurry, and

491    thus influencing column packing [68]. Also, functionalization procedures can

492    change the mechanical and surface properties of the beads, e.g. as a consequence

493    of the different solvents, chemicals and temperatures employed for ligand

494    immobilization. This in turn influences the packing process [69], ultimately

495    determining packing quality.

496

497    The possibility of a correlation between column functionality and backbone was

498    tested both qualitatively (mosaic plot in Figure 5) and statistically by employing

499    the chi-squared test. The size of the mosaic tiles in Figure 5 is proportional to the

500    number of chromatographic columns in the data set having a certain combination

501    of backbone and functional mode. Some of the tiles are predominant over the

502    others, e.g. agarose and methacrylate based materials are employed across affinity,

503    ion exchange and hydrophobic interaction chromatography (AF, AIEC, CIEC, HIC,

504    IMAC, MMC in Figure 5). Such columns are indeed ubiquitous in downstream

505    processing of biopharmaceuticals. Other backbones find use in specific application

506    domains, e.g. dextran is predominantly employed for SEC, and HCIC is purely

507    carried out with cellulosic adsorbents. In addition, a number of combinations of

17

508  functional mode and backbone are not represented in the data set, indicating some

509  resin materials do not find use for certain chromatographic modalities. A chi-

510  squared test of independence with 63 degrees of freedom, i.e. (8 backbones – 1) x

511  (10 functional modes – 1), and with a sample size of 546 tests indeed showed a

512  significant relationship between the two input variables, $\chi^2(63, N = 546)$

513  $= 693, p < 0.01$. While a correlation between resin material and

514  functionalisation is apparent, its influence in the ML model was eliminated by

515  averaging all experimental results measured under the same input conditions (see

516  section 2.1), especially important step to prevent same samples being present in

517  both the training and testing set thus overestimating the accuracy.

518

### 3.2.3 Column length

520  The influence of column length on $h$ is presented in Figure 4e. It is possible to

521  observe that the median for $h$, as well as its propensity to data scatter and

522  relatively large values ($h$ above 10) increase with column length. This observation

523  can be explained by a combination of packing consolidation and wall effects. The

524  former is relevant during column manufacture, i.e. when compression forces

525  transfer through the packing via inter-particle friction as well as friction between

526  particles and the column wall [70]. The uneven stress distribution created

527  between particles in the bulk and at the periphery of the column negatively affect

528  bed consolidation and packing homogeneity. The presence of the wall constrains

529  the resin particles to pack in configurations with higher local porosity in the

530  immediate vicinity of the column wall. The columns investigated in this work were

531  small scale purification columns (column volume about 1 and 10 ml) with

532  relatively large particle diameters (15 to 400 $\mu$m) and small column diameters (5

533  – 11.3 mm). The resulting column diameter to particle diameter ratio was in

534  general around 80, down to 20 for some columns. In this context, Maier et al. [71]

535  reported wall effect on axial dispersion can be observed even for columns with

536  column dimeter to particle diameter ratio greater than 100. Reising et al. [72] and

537  Fabrice Gritti [73] studied the dependence of fluid velocity with radial position,

18

538   and concluded that the velocity close to the column wall can be up to 2.2 times the

539   bulk velocity, significantly contributing to band broadening and early

540   breakthrough. Flow non-idealities arising from both uneven packing difficulties

541   and wall effects scale with column length, with packing quality and column

542   performance inversely related to it.

543

544   The contribution of column length on $A_s$ is reported in Figure 4f. No significant

545   difference can be observed across the data, other than a minor decrease in the

546   median asymmetry with column length. Asymmetry is heavily determined by

547   extra column band broadening, i.e. related to all flow non-idealities present in the

548   extra column volumes such as tubing, fitting, column distributor and collectors,

549   pumps, valves etc. This effect becomes more prominent for smaller columns, as

550   described by Kaltenbrunner et al. [74] who reported extra column volumes

551   accounting for more than 90% band broadening in small columns.

552

553   **3.2.4 Column diameter**

554   According to ML results the contribution of column diameter to the prediction of

555   $h$ is 5.2%, while it is only 0.9% for $A_s$ (Figure 3c), and no clear relationship can be

556   observed between column diameter and the two performance output parameters

557   (Figure 4g and 4h). All the three column diameters considered in this work fall in

558   the same order of magnitude (5, 8 and 11.3 mm), thus hiding any potential

559   correlation between column diameter and packing quality. Schweiger et al [3]

560   analyzed the band broadening arising from the extra-column and in-column

561   contributions of pre-packed columns with different column diameters, and

562   concluded that an increase in column diameter can lead to an increase in peak

563   width as caused by flow non-idealities in the flow distributor and collector.

564   Experimental data for wider columns is required to identify and eventually

565   quantify any possible relationship between column diameter and column

566   performance.

567

19

568 **3.2.5 Particle diameter**

569 The correlation between particle diameter and $h$ is reported in Figure 4i.

570 Accordingly to the reduced form of the van Deemter equation (Eq 8), the

571 magnitude of $h$ is not dependent on particle diameter. ML results indicate that the

572 importance contribution of particle diameter to $h$ is 10.7% (Figure 3c). In Figure

573 4i the median $h$ slightly drops with particle size, possibly resulting from packing

574 difficulties with smaller particles, as also reported by Scharl et al. [7]. No trend

575 between $A_s$ and particle size could be observed (Figure 4j).

576

577 **4. Conclusions**

578 Traditional statistical analysis (e.g. autocorrelation analysis, chi square analysis)

579 and machine learning were applied to a large data set (546 different combinations

580 of column features) of packing quality (reduced plate height, $h$, and asymmetry,

581 $A_s$) for pre-packed columns manufactured with different column sizes (column

582 length and column diameter) and packed with different resins (backbone,

583 functional mode, and particle diameter) over a ten year period.

584

585 Autocorrelation and partial autocorrelation provided a quantitative framework to

586 analyze column quality over time. The results indicate that packing quality was

587 indeed not correlated with time, indicating that column manufacture, sampling

588 and testing was consistent over the ten year period.

589

590 The XGBoost represented an excellent ML model to predict column performance,

591 with mean absolute percentage error (MAPE) of 10% and 7% on $h$ and $A_s$,

592 respectively. According to the ML tool employed, column backbone contributed

593 the most to its predictive capability. In other words, the resin material employed

594 had the most significant impact on column performance. A trend between column

595 length and performance was also observed, with $h$ raising slightly as the length

596 increased, consistent with a larger contribution to band broadening due to wall

597 effects and axial dispersion.

20

Overall, this work demonstrates the capability of ML to evaluate and predict column performance solely from the knowledge of some basic column characteristics (column length and diameter, particle size, backbone material, functional mode). These results could be employed to extrapolate the expected performance characteristics on new and existing columns types, help set QA protocols for new and existing manufacturing lines for pre-packed chromatography columns, or as a reference benchmark for columns packed traditionally in lab settings, especially for hard to pack columns such as PS-DVB and inorganic supports. The results presented here can guide further efforts in column optimization, e.g. informing potential inefficiencies in the packing process, and suggesting improvements of backbone and functional modes to obtain easy to pack resins prone to form ordered packing arrangements with high chromatographic performance.

More in general, ML provides a quantitative tool to describe complex problems with multiple input features, including categorical features such as resin backbone and functional mode. ML methods can also be employed in other chromatographic areas, e.g. for generating accurate retention models, resolving complex chromatography peaks and for searching column structures with improved performance.

**References**

[1]     Brenac Brochier, V., Chabre, H., Lautrette, A., Ravault, V., Couret, M. N., Didierlaurent, A., Moingeon, P., High throughput screening of mixed-mode sorbents and optimisation using pre-packed lab-scale columns for the purification of the recombinant allergen rBet v 1a. *J. Chromatogr. B Anal. Technol. Biomed. Life Sci.* 2009, 877, 2420–2427.

21

626 [2] Shukla, A. A., Gottschalk, U., Single-use disposable technologies for

627 biopharmaceutical manufacturing. *Trends Biotechnol.* 2013, 31, 147–154.

628 [3] Schweiger, S., Jungbauer, A., Scalability of pre-packed preparative

629 chromatography columns with different diameters and lengths taking into

630 account extra column effects. *J. Chromatogr. A* 2018, DOI:

631 10.1016/j.chroma.2018.01.022.

632 [4] Schweiger, S., Hinterberger, S., Jungbauer, A., Column-to-column packing

633 variation of disposable pre-packed columns for protein chromatography. *J.*

634 *Chromatogr. A* 2017, 1527, 70–79.

635 [5] Glueckauf, B. Y. E., Theory of chromatography. 1954, 34–44.

636 [6] Kaltenbrunner, O., Watler, P., Yamamoto, S., Column Qualification in

637 Process Ion-Exchange Chromatography. Elsevier Masson SAS 2000.

638 [7] Scharl, T., Jungreuthmayer, C., Dürauer, A., Schweiger, S., Schröder, T.,

639 Jungbauer, A., Trend analysis of performance parameters of pre-packed

640 columns for protein chromatography over a time span of ten years. *J.*

641 *Chromatogr. A* 2016, 1465, 63–70.

642 [8] Fornstedt, T., Zhong, G., Guiochon, G., Peak tailing and mass transfer

643 kinetics in linear chromatography. 1996, 741.

644 [9] Fornstedt, T., Zhong, G., Guiochon, G., Peak tailing and slow mass transfer

645 kinetics in nonlinear chromatography. *J. Chromatogr. A* 1996, 742, 55–68.

646 [10] Wakamatsu, A., Morimoto, K., Shimizu, M., Kudoh, S., A severe peak tailing

647 of phosphate compounds caused by interaction with stainless steel used

648        for liquid chromatography and electrospray mass spectrometry. *J. Sep. Sci.*

649        2005, 28, 1823–1830.

650    [11]  Kirkland, J. J., Yau, W. W., Stoklosa, H. J., Dilks, C. H., Sampling and extra-

651        column effects in high-performance liquid chromatography; influence of

652        peak skew on plate count calculations. *J. Chromatogr. Sci.* 1977, DOI:

653        10.1093/chromsci/15.8.303.

654    [12]  Chapel, S., Heinisch, S., Strategies to circumvent the solvent strength

655        mismatch problem in online comprehensive two-dimensional liquid

656        chromatography. *J. Sep. Sci.* 2022, 45, 7–26.

657    [13]  Mitchell, N. S., Hagel, L., Fernandez, E. J., In situ analysis of protein

658        chromotography and column efficiency using magnetic resonance imaging.

659        *J. Chromatogr. A* 1997, DOI: 10.1016/S0021-9673(97)00457-3.

660    [14]  Carta, G., Jungbauer, A., Protein Chromatography. 2017.

661    [15]  Guiochon, G., Felinger, A., Shirazi, D. G., Fundamentals of Preparative and

662        Nonlinear Chromatography. Elsevier 2006.

663    [16]  Dimartino, S., Boi, C., Sarti, G. C., A validated model for the simulation of

664        protein purification through affinity membrane chromatography. *J.*

665        *Chromatogr. A* 2011, 1218, 1677–1690.

666    [17]  Dimartino, S., Boi, C., Sarti, G. C., Scale-up of affinity membrane modules:

667        Comparison between lumped and physical models. *J. Mol. Recognit.* 2015,

668        28, 180–190.

669    [18]  Neal, R. M., Pattern Recognition and Machine Learning. *Technometrics*

23

670    2007, DOI: 10.1198/tech.2007.s518.

671    [19]   O'Rourke, J., Toussaint, G. T., Handbook of Discrete and Computational

672           Geometry, Third Edition. 2017.

673    [20]   Rasouli, K., Hsieh, W. W., Cannon, A. J., Daily streamflow forecasting by

674           machine learning methods with weather and climate inputs. *J. Hydrol.*

675           2012, 414–415, 284–293.

676    [21]   Xu, X., Zhang, Y., Zou, L., Wang, M., Li, A., A gene signature for breast cancer

677           prognosis using support vector machine. *2012 5th Int. Conf. Biomed. Eng.*

678           *Informatics, BMEI 2012* 2012, 928–931.

679    [22]   Bos, T. S., Knol, W. C., Molenaar, S. R. A., Niezen, L. E., Schoenmakers, P. J.,

680           Somsen, G. W., Pirok, B. W. J., Recent applications of chemometrics in one-

681           and two-dimensional chromatography. *J. Sep. Sci.* 2020, 43, 1678–1727.

682    [23]   Risum, A. B., Bro, R., Using deep learning to evaluate peaks in

683           chromatographic data. *Talanta* 2019, 204, 255–260.

684    [24]   Kantz, E. D., Tiwari, S., Watrous, J. D., Cheng, S., Jain, M., Deep neural

685           networks for classification of LC-MS spectral peaks. *Anal. Chem.* 2019, 91,

686           12407–12413.

687    [25]   Marengo, E., Gianotti, V., Angioi, S., Gennaro, M. C., Optimization by

688           experimental design and artificial neural networks of the ion-interaction

689           reversed-phase liquid chromatographic separation of twenty cosmetic

690           preservatives. 2004, 1029, 57–65.

691    [26]   Hervás, C., Martínez, A. C., Silva, M., Serrano, J. M., Improving the

692        quantification of highly overlapping chromatographic peaks by using

693        product unit neural networks modeled by an evolutionary algorithm. *J.*

694        *Chem. Inf. Model.* 2005, 45, 894–903.

695    [27]   Vasiljević, T., Onjia, A., Čokeša, D., Laušević, M., Optimization of artificial

696        neural network for retention modeling in high-performance liquid

697        chromatography. *Talanta* 2004, 64, 785–790.

698    [28]   Kensert, A., Collaerts, G., Efthymiadis, K., Desmet, G., Cabooter, D., Deep Q-

699        learning for the selection of optimal isocratic scouting runs in liquid

700        chromatography. *J. Chromatogr. A* 2021, 1638, 461900.

701    [29]   Ben Hameda, A., Elosta, S., Havel, J., Journal of Chromatography A. 2005.

702    [30]   Wang, G., Briskot, T., Hahn, T., Baumann, P., Hubbuch, J., Estimation of

703        adsorption isotherm and mass transfer parameters in protein

704        chromatography using artificial neural networks. *J. Chromatogr. A* 2017,

705        1487, 211–217.

706    [31]   Narayanan, H., Seidler, T., Luna, M. F., Sokolov, M. Morbidelli, M., Butté, A.,

707        Hybrid Models for the simulation and prediction of chromatographic

708        processes for protein capture. *J. Chromatogr. A* 2021, 1650, 462248.

709    [32]   Narayanan, H., Luna, M. F., Stosch, M., Cruz Bournazou, M. N., Polotti, G.,

710        Morbidelli, M., Butté, A., Sokolov, M., Bioprocessing in the Digital Age: The

711        Role of Process Models. *Biotechnol. J.* 2020, 15, 1900172.

712    [33]   Narayanan, H., Behle, L., Luna, M. F., Sokolov, M., Guillén-Gosálbez, G.,

713        Morbidelli, M., Butté, A., Hybrid-EKF: Hybrid model coupled with extended

714       Kalman filter for real-time monitoring and control of mammalian cell

715       culture. *Biotechnol. Bioeng.* 2020, 117, 2703–2714.

716   [34]   Cussler, E. L., Cussler, E. L., Diffusion: Mass Transfer in Fluid Systems.

717       Cambridge university press 2009.

718   [35]   Yeh, H. S., Wills, G. B., Diffusion coefficient of sodium nitrate in aqueous

719       solution at 25.deg. as a function of concentration from 0.1 to 1.0M. *J. Chem.*

720       *Eng. Data* 1970, 15, 187–189.

721   [36]   Chen, T., Guestrin, C., XGBoost: A scalable tree boosting system. *Proc. ACM*

722       *SIGKDD Int. Conf. Knowl. Discov. Data Min.* 2016, 13-17-Augu, 785–794.

723   [37]   Freund, Yoav,  and L. M., The alternating decision tree learning algorithm.

724       *Int. Conf. Mach. Learn.* 1999, DOI: 10.1093/jxb/ern164.

725   [38]   Lewis, R. J., Ph, D., Street, W. C., An Introduction to Classification and

726       Regression Tree ( CART ) Analysis. *2000 Annu. Meet. Soc. Acad. Emerg. Med.*

727       2000, DOI: 10.1.1.95.4103.

728   [39]   Friedman, J. H., Greedy Function Approximation: A Gradient Boosting

729       Machine. *Ann. Stat.* 2001, 29, 1189–1232.

730   [40]   Knapp, S. K., Accelerate FPGA macros with one-hot approach. *Electron. Des.*

731       1990.

732   [41]   Singh, D., Singh, B., Investigating the impact of data normalization on

733       classification performance. *Appl. Soft Comput. J.* 2019, 105524.

734   [42]   Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,

735       Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos,

736      A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, É., Scikit-learn:

737      Machine learning in Python. *J. Mach. Learn. Res.* 2011.

738  [43]  Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., Chen, K.,

739      Mitchell, R., Cano, I., Zhou, T., Mu, L., Xie, J., Lin, M., Geng, Y., Li, Y., Package

740      "Xgboost." 2019.

741  [44]  Chai, T., Draxler, R. R., Root mean square error (RMSE) or mean absolute

742      error (MAE)? -Arguments against avoiding RMSE in the literature. *Geosci.*

743      *Model Dev.* 2014, 7, 1247–1250.

744  [45]  Wilson, E. J., Geankoplis, C. J., Liquid mass transfer at very low reynolds

745      numbers in packed beds. *Ind. Eng. Chem. Fundam.* 1966, 5, 9–14.

746  [46]  Sarwar, M. S., Simon, U., Dimartino, S., Experimental investigation and

747      mass transfer modelling of 3D printed monolithic cation exchangers. *J.*

748      *Chromatogr. A* 2021, 1646, 462125.

749  [47]  Guiochon, G., Gritti, F., Shell particles, trials, tribulations and triumphs. *J.*

750      *Chromatogr. A* 2011, 1218, 1915–1938.

751  [48]  Gritti, F., Leonardis, I., Abia, J., Guiochon, G., Physical properties and

752      structure of fine core–shell particles used as packing materials for

753      chromatography. *J. Chromatogr. A* 2010, 1217, 3819–3843.

754  [49]  Knox, J. H., Band dispersion in chromatography - A universal expression

755      for the contribution from the mobile zone. *J. Chromatogr. A* 2002, 960, 7–

756      18.

757  [50]  Patel, K. D., Jerkovich, A. D., Link, J. C., Jorgenson, J. W., In-depth

758   characterization of slurry packed capillary columns with 1.0-µm

759   nonporous particles using reversed-phase isocratic ultrahigh-pressure

760   liquid chromatography. *Anal. Chem.* 2004, 76, 5777–5786.

761   [51]   Malkin, D. S., Wei, B., Fogiel, A. J., Staats, S. L., Wirth, M. J., Submicrometer

762   Plate Heights for Capillaries Packed with Silica Colloidal Crystals. *Anal.*

763   *Chem.* 2010, 82, 2175–2177.

764   [52]   Khirevich, S., Daneyko, A., Höltzel, A., Seidel-Morgenstern, A., Tallarek, U.,

765   Statistical analysis of packed beds, the origin of short-range disorder, and

766   its impact on eddy dispersion. *J. Chromatogr. A* 2010, 1217, 4713–4722.

767   [53]   Gritti, F., Leonardis, I., Shock, D., Stevenson, P., Shalliker, A., Guiochon, G.,

768   Performance of columns packed with the new shell particles, Kinetex-C18.

769   *J. Chromatogr. A* 2010, DOI: 10.1016/j.chroma.2009.12.079.

770   [54]   Dolamore, F., Dimartino, S., Fee, C. J., Numerical Elucidation of Flow and

771   Dispersion in Ordered Packed Beds: Nonspherical Polygons and the Effect

772   of Particle Overlap on Chromatographic Performance. *Anal. Chem.* 2019,

773   91, 15009–15016.

774   [55]   Knox, J. H., Band dispersion in chromatography - A new view of A-term

775   dispersion. *J. Chromatogr. A* 1999, 831, 3–15.

776   [56]   Lan, T., Gerontas, S., Smith, G. R., Langdon, J., Ward, J. M., Titchener-Hooker,

777   N. J., Investigating the use of column inserts to achieve better

778   chromatographic bed support. *Biotechnol. Prog.* 2012, 28, 1285–1291.

779   [57]   Jaulmes, A., Ignatiadis, I., Cardot, P., Vidal-Madjar, C., Characterization of

780     peak asymmetry with overloaded capillary columns. *J. Chromatogr. A*

781     1987, 395, 291–306.

782  [58]  Pápai, Z., Pap, T. L., Analysis of peak asymmetry in chromatography. *J.*

783     *Chromatogr. A* 2002, 953, 31–38.

784  [59]  Miyabe, K., Guiochon, G., Estimation of the column radial heterogeneity

785     from an analysis of the characteristics of tailing peaks in linear

786     chromatography. *J. Chromatogr. A* 1999, 830, 29–39.

787  [60]  Miyabe, K., Guiochon, G., Peak tailing and column radial heterogeneity in

788     linear chromatography. *J. Chromatogr. A* 1999, 830, 263–274.

789  [61]  Madsen, H., Time Series Analysis. 2007.

790  [62]  Palma, W., Long-Memory Time Series. John Wiley & Sons, Inc., Hoboken,

791     NJ, USA 2007.

792  [63]  Mills, T. C., ARMA Models for Stationary Time Series. *Appl. Time Ser. Anal.*

793     2019, 31–56.

794  [64]  Bendada, K., Hamdi, B., Boudriche, L., Balard, H., Calvet, R., Surface

795     characterization of reservoir rocks by inverse gas chromatography: Effect

796     of a surfactant. *Colloids Surfaces A Physicochem. Eng. Asp.* 2016, 504, 75–

797     85.

798  [65]  Bacskay, I., Sepsey, A., Felinger, A., Determination of the pore size

799     distribution of high-performance liquid chromatography stationary phases

800     via inverse size exclusion chromatography. *J. Chromatogr. A* 2014, 1339,

801     110–117.

29

802    [66]    Johnson, T. F., Bailey, J. J., Iacoviello, F., Welsh, J. H., Levison, P. R., Shearing,

803            P. R., Bracewell, D. G., Three dimensional characterisation of

804            chromatography bead internal structure using X-ray computed

805            tomography and focused ion beam microscopy. *J. Chromatogr. A* 2018,

806            1566, 79–88.

807    [67]    Stickel, J. J., Fotopoulos, A., Pressure-flow relationships for packed beds of

808            compressible chromatography media at laboratory and production scale.

809            *Biotechnol. Prog.* 2001, 17, 744–751.

810    [68]    Kawachi, Y., Ikegami, T., Takubo, H., Ikegami, Y., Miyamoto, M., Tanaka, N.,

811            Chromatographic characterization of hydrophilic interaction liquid

812            chromatography stationary phases: Hydrophilicity, charge effects,

813            structural selectivity, and separation efficiency. *J. Chromatogr. A* 2011,

814            1218, 5903–5919.

815    [69]    McCue, J. T., Cecchini, D., Hawkins, K., Dolinski, E., Use of an alternative

816            scale-down approach to predict and extend hydroxyapatite column

817            lifetimes. *J. Chromatogr. A* 2007, DOI: 10.1016/j.chroma.2007.07.053.

818    [70]    Dorn, M., Eschbach, F., Hekmat, D., Weuster-Botz, D., Influence of different

819            packing methods on the hydrodynamic stability of chromatography

820            columns. *J. Chromatogr. A* 2017, 1516, 89–101.

821    [71]    Maier, R. S., Kroll, D. M., Davis, H. T., Diameter-dependent dispersion in

822            packed cylinders. *AIChE J.* 2007, 53, 527–530.

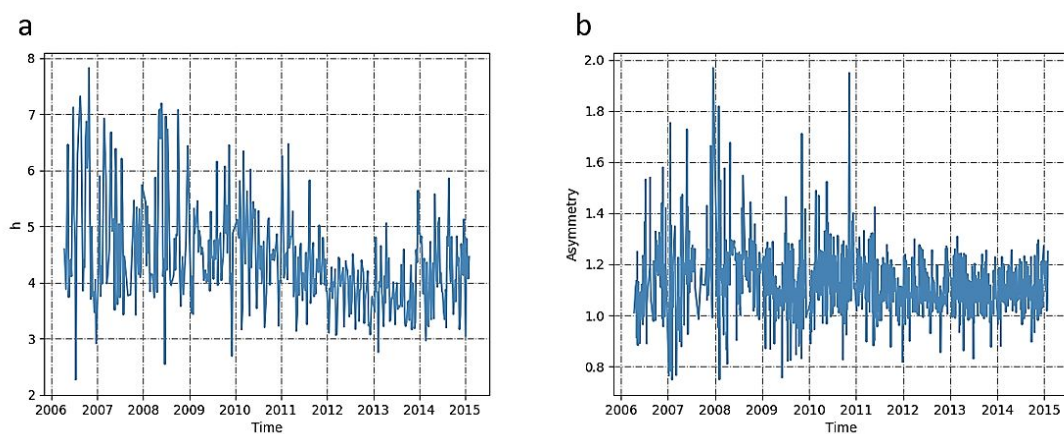823    [72]    Reising, A. E., Schlabach, S., Baranau, V., Stoeckel, D., Tallarek, U., Analysis

824   of packing microstructure and wall effects in a narrow-bore ultrahigh

825   pressure liquid chromatography column using focused ion-beam scanning

826   electron microscopy. *J. Chromatogr. A* 2017, 1513, 172–182.

827  [73] Gritti, F., On the relationship between radial structure heterogeneities and

828   efficiency of chromatographic columns. *J. Chromatogr. A* 2018, 1533, 112–

829   126.

830  [74] Kaltenbrunner, O., Jungbauer, A., Yamamoto, S., Prediction of the

831   preparative chromatography performance with a very small column. *J.*

832   *Chromatogr. A* 1997, 760, 41–53.

833

834     **Figure Captions:**

835

836     Figure 1 Time series of a) reduced plate height, $h$, and b) asymmetry, $A_s$, for pre-

837     packed purification columns manufactured over the 10 year period monitored.

838



839
840

841    Figure 2 Autocorrelation (acf) and partial autocorrelation (pacf) analysis of
842    reduced plate height, $h$, and asymmetry, $A_s$. a) Acf of $h$; b) pacf of $h$; c) acf of $A_s$; d)
843    pacf of $A_s$. The blue shaded areas correspond to 95% confidence interval.
844



845
846

847     Figure 3. XGBoost prediction results for a) $h$ and b) $A_s$ over the testing data set.

848     Variable importance contributions of c) $h$ and d) $A_s$ are reported. The importance

849     is calculated based on the improvement of the performance measured by each

850     attribute split points, weighted by the number of the observations the node is

851     responsible for. The importance contributions, named by Gain in XGBoost (refer

852     to Eq 4&5), were transferred into percentage.

853



854
855

856   Figure 4. Violin plots of $h$ and $A_s$ against input parameters (backbone, functional

857   mode, column length, column diameter, particle diameter). a) $h$ vs backbone

858   (PS/DVB: polystyrene divynilbenzene; IS: Inorganic support; MET: Methacrylate;

859   AGR: Agarose; POL: Polymer grafted; DEX: Dextran; CEL: Cellulose; PVE:

860   polyvinyl-ether hydrophilic). b) $A_s$ vs backbone. c)  $h$ vs functional mode (CIEC:

861   cation exchange chromatograph; AF: affinity chromatography; HA: hydroxyl-

862   apatite chromatography; AIEC: anion exchange chromatography; HIC:

863   hydrophobic interaction chromatography; SEC: size-exclusion  chromatography;

864   IMAC: immobilized metal affinity chromatography; MMC: mixed-mode

865   chromatography; FA: fluorophore adsorption chromatography; HCIC:

866   hydrophobic charge induction chromatography). d) $A_s$ vs functional mode. e) $h$ vs

867   column length. f) $A_s$ vs column length. g) $h$ vs column diameter. h) $A_s$ vs column

868   diameter. i) $h$ vs particle diameter. j) $A_s$ vs particle diameter.

869

35

872    Figure 5. Mosaic plot of the combinations of functional mode and backbone

873    material tested. The size of the tiles represents the relative frequency of each

874    combination. PS/DVB: polystyrene divynilbenzene; IS: Inorganic support; MET:

875    Methacrylate; AGR: Agarose; POL: Polymer grafted; DEX: Dextran; CEL: Cellulose;

876    PVE: polyvinyl-ether hydrophilic; CIEC: cation exchange chromatograph; AF:

877    affinity chromatography; HA: hydroxyl-apatite chromatography; AIEC: anion

878    exchange chromatography; HIC: hydrophobic interaction chromatography; SEC:

879    size-exclusion       chromatography;   IMAC:   immobilized    metal   affinity

880    chromatography;   MMC:   mixed-mode   chromatography;   FA:   fluorophore

881    adsorption    chromatography;   HCIC:   hydrophobic   charge   induction
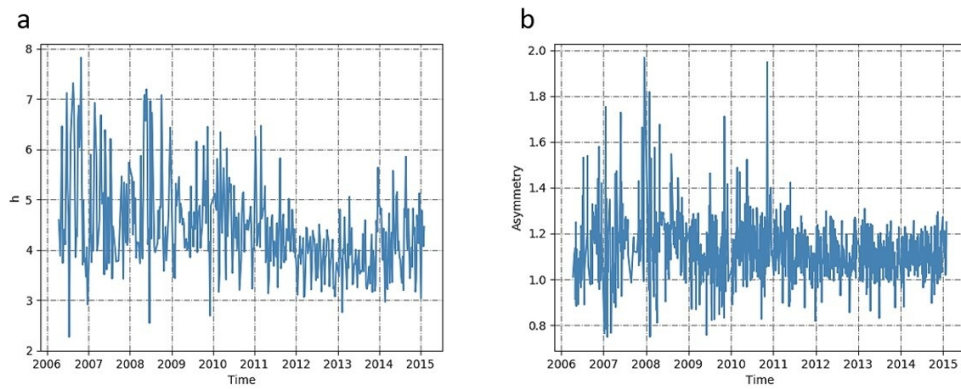
882    chromatography.

883



884

885

Figure 1 Time series of a) reduced plate height, h, and b) asymmetry, A_s.
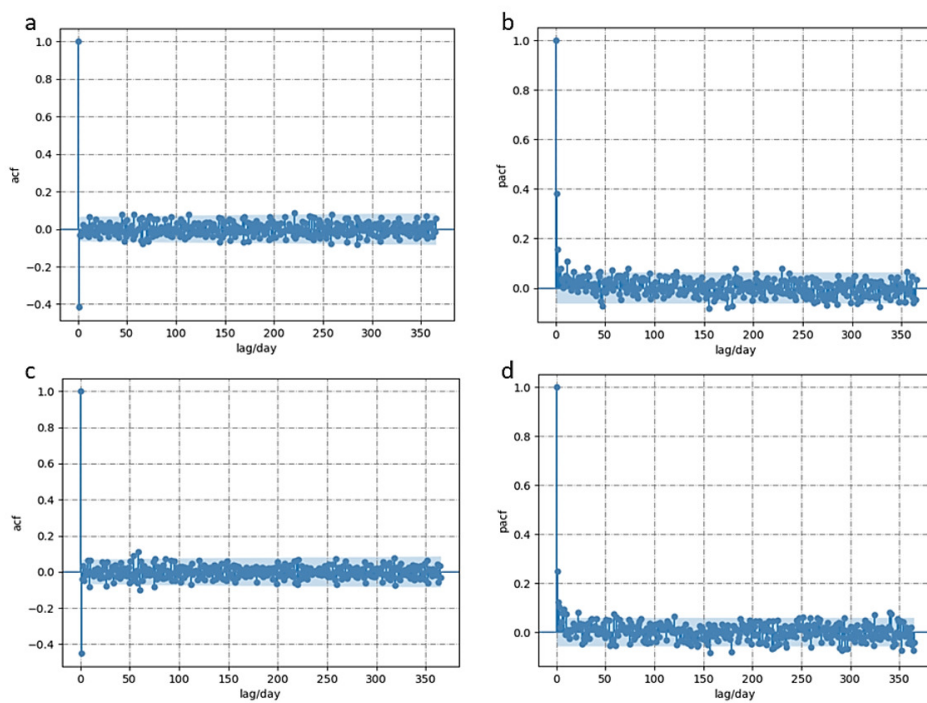
99x40mm (300 x 300 DPI)
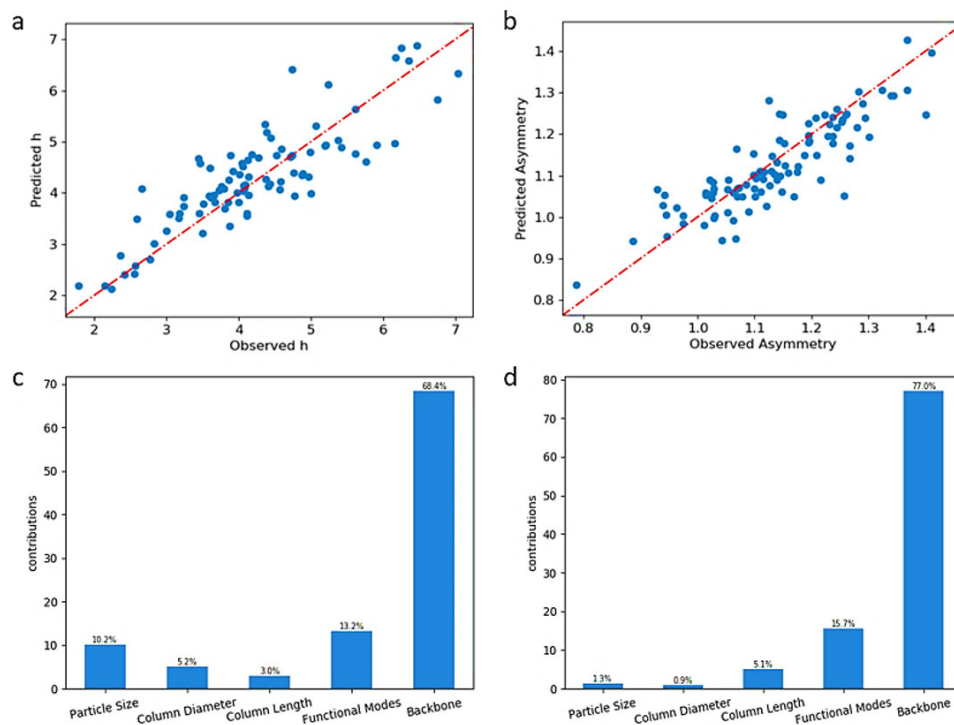
Figure 2 Autocorrelation (acf) and partial autocorrelation (pacf) analysis of reduced plate height, h, and asymmetry, A_s. a) Acf of h; b) pacf of h; c) acf of A_s; d) pacf of A_s. The blue shaded areas correspond to 95% confidence interval.

101x74mm (300 x 300 DPI)

Figure 2 Autocorrelation (acf) and partial autocorrelation (pacf) analysis of reduced plate height, h, and asymmetry, A_s. a) Acf of h; b) pacf of h; c) acf of A_s; d) pacf of A_s. The blue shaded areas correspond to 95% confidence interval.
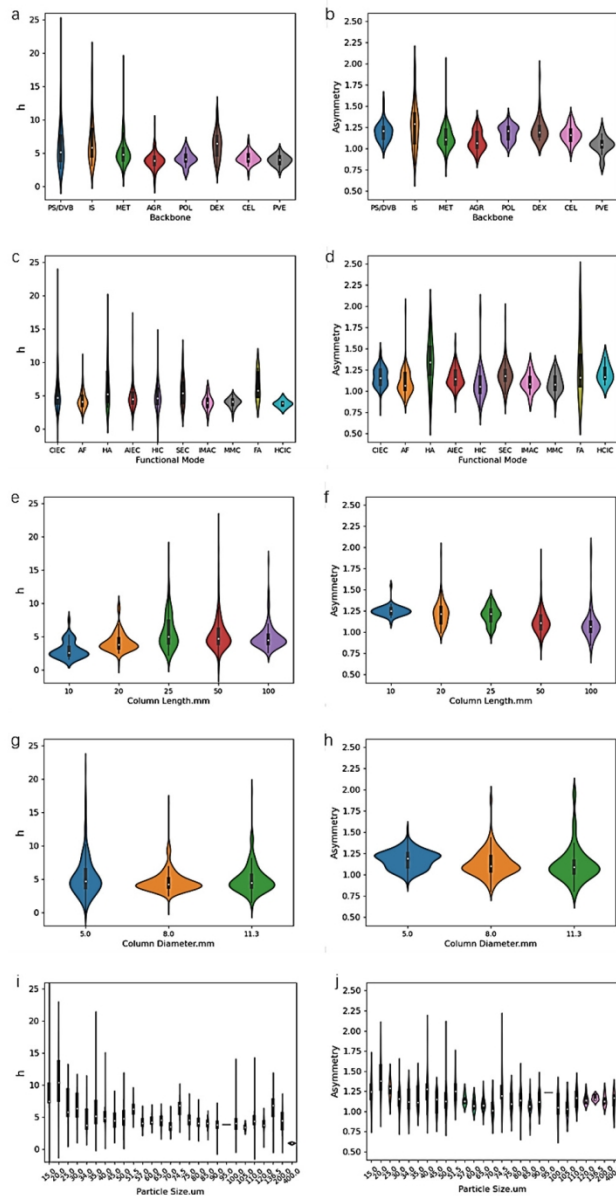
57x42mm (500 x 500 DPI)

Figure 4. Violin plots of h and As against input parameters (backbone, functional mode, column length, column diameter, particle diameter). a) h vs backbone (PS/DVB: polystyrene divynilbenzene; IS: Inorganic support; MET: Methacrylate; AGR: Agarose; POL: Polymer grafted; DEX: Dextran; CEL: Cellulose; PVE: polyvinyl-ether hydrophilic). b) A_s vs backbone. c) h vs functional mode (CIEC: cation exchange chromatograph; AF: affinity chromatography; HA: hydroxyl-apatite chromatography; AIEC: anion exchange chromatography; HIC: hydrophobic interaction chromatography; SEC: size-exclusion chromatography; IMAC: immobilized metal affinity chromatography; MMC: mixed-mode chromatography; FA: fluorophore adsorption chromatography; HCIC: hydrophobic charge induction chromatography). d) A_s vs functional mode. e) h vs column length. f) A_s vs column length. g) h vs column diameter. h) A_s vs column diameter. i) h vs particle diameter. j) A_s vs particle diameter.
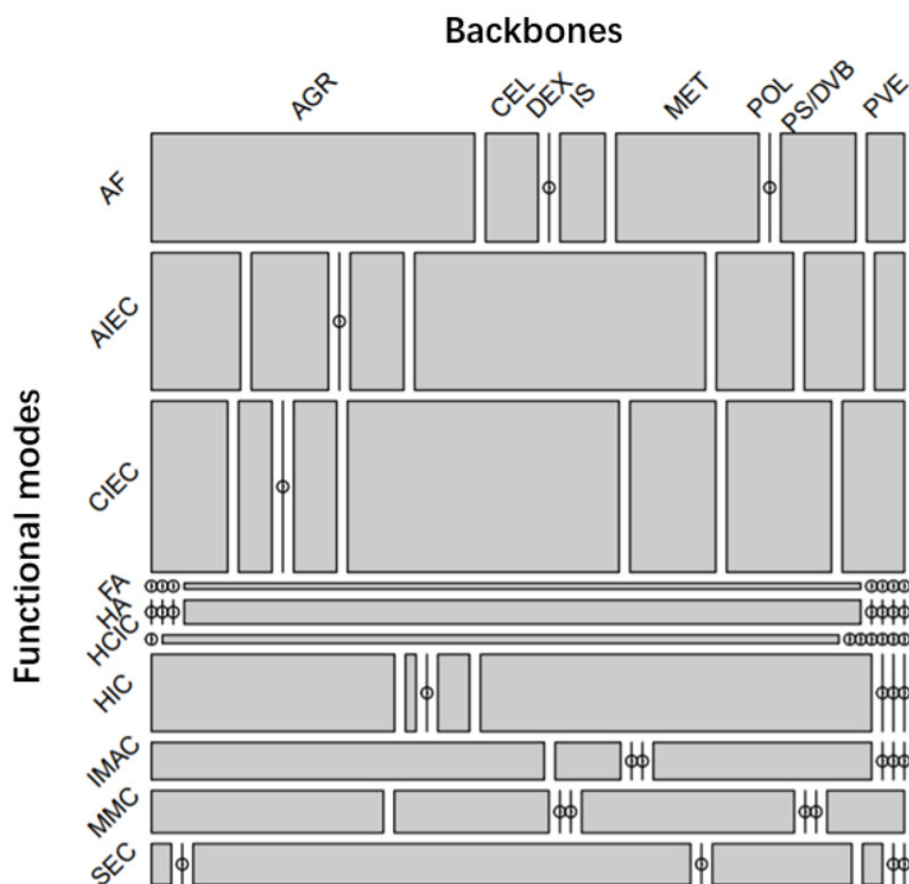
26x48mm (1000 x 1000 DPI)

Figure 5. Mosaic plot of the combinations of functional mode and backbone material tested. The size of the tiles represents the relative frequency of each combination. PS/DVB: polystyrene divynilbenzene; IS: Inorganic support; MET: Methacrylate; AGR: Agarose; POL: Polymer grafted; DEX: Dextran; CEL: Cellulose; PVE: polyvinyl-ether hydrophilic; CIEC: cation exchange chromatograph; AF: affinity chromatography; HA: hydroxyl-apatite chromatography; AIEC: anion exchange chromatography; HIC: hydrophobic interaction chromatography; SEC: size-exclusion chromatography; IMAC: immobilized metal affinity chromatography; MMC: mixed-mode chromatography; FA: fluorophore adsorption chromatography; HCIC: hydrophobic charge induction chromatography.

21x20mm (1000 x 1000 DPI)