

# Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes

Luca Mencarelli, Alexandre Pagot, Pascal Duchêne

## ► To cite this version:

Luca Mencarelli, Alexandre Pagot, Pascal Duchêne. Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes. Computers & Chemical Engineering, 2020, 135, pp.106772. 10.1016/j.compchemeng.2020.106772. hal-02553492

# HAL Id: hal-02553492 https://ifp.hal.science/hal-02553492

Submitted on 24 Apr 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes

Luca Mencarelli\*, Alexandre Pagot, Pascal Duchêne

Keywords: Surrogate modeling Principal component analysis Subset selection Catalytic reforming Light naphtha isomerization

## ABSTRACT

In this paper, we first briefly survey the main surrogate model building approaches discussed in the literature considering also design of experiments strategies and dimensionality reduction procedures: we mainly focus on sub-set approaches and sampling strategies for constrained regression problems. We delineate a systematic methodology for surrogate modelling in presence of model constraints, such as non-negativity of the model responses. The main contribution of this paper is twofold: from one side we extend the principal component analysis framework to the case of constrained regression problem, from the other we propose a novel methodology which integrates the subset selection and the previous principal component regression procedure. Finally, we apply the two novel algorithms to two fundamental chemical processes in petroleum refinery, namely catalytic reforming and light naphtha isomerization. The numerical results show the comparisons between the two algorithms in terms of computational and accuracy trade-offs.

© 2020 Elsevier Ltd. All rights reserved.

43

44

45

46

47

## 1. Introduction

The aim of chemical process synthesis engineering consists in 2 modeling, designing and optimizing complex chemical processes. 3 The possible high computational cost of the process estimation and 4 optimization can be circumvented by means of surrogate models 5 (or meta-models) that represent a systematic approximation of the 6 mathematical relationships between the degrees of freedom (in-7 put data) and the variables of interest (output data). A systematic 8 9 methodology to identify dependent and independent variables of 10 a given chemical process unit is given by Henao and Maravelias (2010, 2011). Instead of obtaining the output data via experimen-11 tal measurements, numerical simulators are often available, but in 12 several cases obtaining output values from a given input configu-13 ration is rather time consuming. 14

15 Surrogate models can be useful either if the first-principle model is too complex or time consuming to optimize or if the first-16 principle model does not exist at all. In the first case it is possible 17 18 to collect initial simulated data by choosing a sampling strategy and sampling more points later; in the latter case instead the data 19 are obtained by physical experiments. Moreover, in both cases the 20 data could be affected by noise or characterized by incomplete in-21 22 formation.

https://doi.org/10.1016/j.compchemeng.2020.106772 0098-1354/© 2020 Elsevier Ltd. All rights reserved. The objective of the present paper consists in defining a 23 methodology to derive a surrogate model starting from noiseless 24 simulated data in presence of model constraints: in our case, in 25 fact, a numerical simulator is available for the chemical processes 26 we consider, and data can be obtained easily and rather quickly 27 from the simulations. 28

From one side, the surrogate model should be sufficiently com-29 plex to catch the relationships of the given process, hence it should 30 be characterized by high accuracy, from the other one, instead, it 31 should be sufficiently simple to speed up the computational times, 32 so that its complexity is low. High accuracy and low complexity 33 are obviously conflicting targets and determine a trade-off between 34 the quality of the approximation and the quantity of computational 35 effort. Generally, in fact, surrogate models are preferred to other 36 approaches, such as rigorous models or simplified physical approx-37 imation models, when the computational time is expected to be a 38 crucial aspect Psaltis et al. (2016): rigorous or physical approxima-39 tion models could be rather time consuming since they have been 40 usually obtained by discretizing complex dynamic equations, such 41 as systems of partial differential equations. 42

In our study we introduce a novel methodology to find sufficiently accurate surrogate models and to simultaneously perform dimensionality reduction with regards to the number of model parameters. Hence, procedures for reducing the total number of experiments are out of scope of the present paper.

Many different approaches have been discussed in the literature in order to build an effective surrogate model from a set of 49

<sup>\*</sup> Corresponding author.

E-mail addresses: luca.mencarelli@ifpen.fr (L. Mencarelli), alexandre.pagot@ ifpen.fr (A. Pagot), pascal.duchene@ifpen.fr (P. Duchêne).

input and output couples of experimental or simulated data (see 50 51 for instance the surveys Queipo et al., 2005; Forrester and Keane, 2009; Vu et al., 2017; Bhosekar and Ierapetritou, 2018; McBride 52 and Sundmacher, 2019), but all the methodologies usually struc-53 ture the surrogate modeling into roughly four main steps: (i) de-54 55 sign of experiments (DOE), (ii) numerical simulations or experimental measurements, (iii) surrogate model selection and identi-56 fication, and (iv) model testing. At step (i) the design space is con-57 veniently sampled in order to define a set of input data configura-58 59 tions. At step (ii) several experiments are performed or a simulator 60 is used to obtain the output data corresponding to the input configurations. Therefore, at step (iii) a specific surrogate model is se-61 lected and trained with respect to the so-called training set, which 62 is composed of input and output data couples, i.e., the parameters 63 64 of the surrogate model are estimated. Finally, at step (iv) the per-65 formances of the model are analyzed with respect to the so-called test set. If the performances of the surrogate model in terms of 66 complexity and accuracy are not satisfactory, then the procedure 67 restarts from step (i). 68

In our applications, we adopt a one-shot approach, i.e., we consider all the sampling points at once, so that our approach is composed of three main steps: (i) we sample the design space in order to obtain a training set of input/output values which opportunely cover the entire design space, (ii) we build the surrogate model by using all the sampling data, and (iii) we test the performance of our model on the test set.

Moreover, surrogate models are also used in a posterior opti-76 77 mization step to retrieve the optimal operating conditions for the 78 chemical process for instance in a superstructure framework (we 79 refer the interested reader to the survey Mencarelli et al., 2019a). In this approach a superstructure is defined by the set of all the 80 81 possible alternative structures of a given chemical process: in sur-82 rogate driven superstructure approach the model units (reactors, distillation columns, or even entire sub-processes) are replaced by 83 84 their surrogate models (Henao and Maravelias, 2010; 2011) and the 85 superstructure is described by a general disjunctive problem (GDP) 86 or a mixed integer (non)linear problem (MI(N)LP), which is then solved to determine the optimal alternative structure. 87

88 As aforementioned, we assume in our applications that a pro-89 cess simulator is available, but too time consuming to be used directly (Mencarelli et al., 2019b). Moreover, we would like to de-90 velop an appropriate surrogate model in order to exploit its an-91 alytic expression and derivatives during the posterior optimiza-92 93 tion phase, instead of directly applying a derivative-free approach. We have already shown that mesh-adaptive algorithms, such as 94 95 NOMAD Audet et al. and Le Digabel (2011), or EGO-based methods (Jones et al., 1998) perform quite poorly in our applications 96 97 (Mencarelli et al., 2019b).

Several recent papers deal with superstructure optimization by 98 99 considering artificial neural networks (ANNs) as surrogate models (see, e.g., Altissimi et al., 1998; Nascimiento et al., 2000; Fahmi and 100 Cremaschi, 2012): in particular, Fahmi and Cremaschi (2012) pro-101 posed a superstructure optimization methodology, by combining 102 GDP and ANN, in which each process unit is replaced by an ANN 103 104 which is trained by simulated data and embedded in a GDP formulation. 105

106 Generally, several outputs are considered at the same time so 107 that we aim to build a surrogate model for each output. How-108 ever, possible constraints should be enforced: these constraints can regard a single surrogate model independently from the other 109 110 (in this case we have intra-model constraints) or contemporaneously a set of models (in this latter case we have inter-model 111 constraints). Typical examples of intra-model constraints consist in 112 non-negativity of the model responses (if the output represents 113 physical measurements, or molar or mass fractions of several com-114 pounds): in this case the non-negativity constraint is referred to 115

only one single model and each non-negativity constraint for a 116 model is independent from the others. If the outputs represent 117 compounds fractions, we should impose an equality inter-model 118 constraints, enforcing the sum of the fractions is equal to 1: in this 119 latter case the inter-model constraints depend contemporaneously 120 from a set of models responses (for instance, in the previous ex-121 ample, the sum of the responses of all models should be 1, so that 122 all the models responses appear in the corresponding inter-model 123 constraint). 124

The presence of output constraints in derivative-free context 125 has been investigated in recent papers. For instance, Conn and 126 Le Digabel (2013) illustrate a hybrid methodology which com-127 bines quadratic models as surrogate models and mesh-adaptive 128 direct search for constrained black-box problems, i.e., optimiza-129 tion problems in which the analytic expression of both the ob-130 jective and the constraints is not available or it is too com-131 plex to evaluate. Boukoulava et al. (2017) and Boukoulava and 132 Floudas (2017) introduce a data-driven methodology that com-133 bines surrogate modeling approach and deterministic global op-134 timization algorithm, by extending the parallel AlgoRithms for 135 Global Optimization of coNstrAined grey-box compUTational prob-136 lems (p-ARGONAUT) algorithm (Beykal et al., 2018b). Moreover, 137 p-ARGONAUT has been extended to multi-objective optimization 138 problems in Beykal et al. (2018a). For a review in constrained 139 derivative-free optimization we refer the interested reader to the 140 survey (Boukoulava et al., 2016). 141

Therefore, the contribution of this paper is threefold: (i) we 142 briefly summarize the main (recent) approaches to surrogate mod-143 eling together with sampling strategies and dimensionality reduc-144 tion techniques focusing on constrained regression problems, (ii) 145 we extend a well-known dimensionality reduction technique such 146 as principal component analysis (PCA) to the case of regression 147 problem with constraints for the response, and finally (iii) we pro-148 pose a novel methodology by combining the previous version of 149 PCA with subset selection (SS) in order to further reduce the di-150 mensionality of the surrogate model, i.e., the number of parame-151 ters of the model. In fact, since the modelling step is usually fol-152 lowed by an optimization phase where the first-principle models 153 are substituted by the surrogate model and the resulting optimiza-154 tion problem is solved, considering a model with a low number of 155 parameters is crucial to solve the final optimization problem in a 156 reasonable amount of time. Obviously, the dimensionality reduc-157 tion procedure should maintain an acceptable quality of the surro-158 gate model in terms of accuracy in data representation. 159

The rest of the paper is organized as follows. In Section 2 we 160 briefly overview the sampling methods proposed in the litera-161 ture. Section 3 is devoted to discuss the surrogate building pro-162 cedures. In Sections 4.1 and 4.2 we deal with two kinds of possi-163 ble additional constraints for the surrogate model, namely inter-164 model constraints and intra-model equality constraints, respec-165 tively. Then, we propose a two-phase PCA method (Section 5) and 166 hybrid algorithm obtained by combining the two-phase PCA with 167 SS (Section 6). The proposed methodologies are then applied to 168 two relevant chemical processes in petroleum refinery framework, 169 namely catalytic reforming in Section 7.1 and light naphtha iso-170 merization in Section 7.2, respectively, which this report exten-171 sively discusses the numerical results of the computational exper-172 iments. Finally, conclusions and future work perspectives follow in 173 Section 8. 174

In the rest of the paper we adopt the following notation. Given 175 a positive integer scalar  $N \in \mathbb{N}$ , we indicate  $[N] := \{1, ..., N\}$ . Moreover, given a set of N vectors  $a_n \in \mathbb{R}^m$   $(n \in [N])$ ,  $a_n^m$  represents the *m*-th entry of the *n*-th vector in the set. In particular, we define a set of N data  $x_n \in \mathbb{R}^K$   $(n \in [N])$  per the *k*-th input  $(k \in [K])$  and a set  $z_n \in \mathbb{R}^M$   $(n \in [N])$  of N data per the *m*-th output  $(m \in [M])$ . 180

## 181 2. Sampling step

In this section, we review the main techniques adopted for the sampling step. In particular, we report the papers which deal with sampling strategies for constrained problems: in our applications, in fact, the design space is described by a set (linear) constraints (see Section 4).

Two main sampling approaches have been exploited in the literature: (i) geometrical designs, and (ii) statistical designs. To the best of our knowledge, this distinction is introduced for the first time in Vu et al. (2017). In geometrical designs the DOE is defined by taking into account the geometrical shape of the design space; while in statistical designs the response of the surrogate model is assimilated to a realization of a random process.

194 Among geometrical designs the most adopted ones are: (i) full factorial design (FFD) (Box et al., 2005; Forrester et al., 2008), and 195 (ii) latin hypercube design (LHD) (McKay et al., 1979). In both cases 196 197 the design space is uniformly divided into regular cells with same dimensionality: in FFD the centre and the extreme points of each 198 199 cell are selected, while in LHD we keep only a proper subset of the centres of the cells such that there is no couple of points shar-200 201 ing the same coordinate. It is worth to notice, in fact, that the FFD guarantees the design space is sampled in a uniform way; on the 202 203 contrary, LHD does not guarantee the design space is sampled uni-204 formly (see Fig. 2 in Viana (2016)). In order to avoid such a situa-205 tion, several methods have been proposed in literature by choosing LHD according to space-filling criteria (Johnson et al., 1990; Pron-206 207 zato, 2017), obtaining the so-called minimax LHD (van Dam, 2008), which minimize the covering radius, and maximin LHD (Morris 208 209 and Mitchell, 1995; van Dam et al., 2007; Joseph and Hung, 2008; Husslage et al., 2011), which maximize the minimal pairwise dis-210 tance between sampled points. Petelet et al. Petelet et al. (2010) in-211 212 troduce a sampling approach to deal with constrained LHD, i.e., LHD with inequality constraints, based on permutation technique 213 214 applied to (unconstrained) LHD.

215 Among statistical designs the most common ones are (i) D-216 optimal, which aim to find the design which maximizes the determinant of the correlation matrix of the data, and (ii) I-217 218 optimal, whose designs minimize the average prediction variance 219 (Goos et al., 2016) (for an insight about statistical designs see, e.g., Cornell (2002) and Smith (2005)). D-optimal and I-optimal designs 220 have been extended to the case of linear constrained regression, in 221 222 which (linear) additional constraints are presented, by Coetzer and 223 Haines (2017).

Recently, an adaptive method for the sampling phase is intro-224 225 duced in Garud et al. (2017b): the idea consists in initially sampling the whole region according to a given sampling strategy de-226 227 scribed above and then iteratively placing the new sampling points 228 in order to sample the original function as far as possible from the 229 already placed points and in the region where the quality of the 230 approximation is poor. Mixed adaptive sampling strategy for surrogate models represented by ANNs has been proposed by Eason and 231 232 Cremaschi (2014). Straus and Skogestad (2019) have proposed another sampling algorithm relying on a termination criterion based 233 234 on partial least squares regression (PLSR).

Sampling techniques and size choices are compared in Davis et al. (2018) with respect to different surrogate model building approaches. For an exhaustive discussion about the different strategies for sampling phase we refer the interested reader to survey (Garud et al., 2017a).

## 240 3. Surrogate model building

In this section, we survey the main techniques employed to develop surrogate models from input and output data. In statistical regression approaches, we consider a set *B* of basis functions  $f_j(x) : \mathbb{R}^K \to \mathbb{R}$   $(j \in B)$  over the input variables, whose (lin-244 ear) combinations give the responses of the surrogate model. The 245 basis functions could be polynomials, transcendental and trigono-246 metric functions, or even radial basis functions (RBFs). 247

RBFs arise from the seminal paper by Broomhead and 248 Lowe (1988) and have been originally use to smoothly interpolate multivariable functions by Hardy (1971): they are universal approximators for functions over a finite number of real variables 251 (Park and Sandberg, 1993). For a complete insight into RBF topics 252 see the excellent survey by Buhmann (2000). 253

Other approaches include, for instance, Kriging methodol-254 ogy and support vector regression (SVR). Kriging dates back to 255 the papers (Krige, 1952; Matheron, 1963) and have been ap-256 plied to the design and analysis of computational experiments 257 by Sacks et al. (1989). For a detailed analysis of the Krig-258 ing technique we refer the interested reader to the survey by 259 Kleijnen (2009). Caballero and Grossmann (2008) propose a Krig-260 ing approach for the flowsheet optimization problem. Recently, 261 Bouhlel et al. (2016) and Gaspar et al. (2017) have proposed two 262 hybrid approaches by combining Kriging techniques with PLSR 263 and with trust region method, respectively. On the contrary, SVR 264 has been introduced by Vapnik (1995) (see also Vapnik et al., 265 1997) which extends the support vector machine technique to ap-266 proximate nonlinear functions (for an introduction to SVR see the 267 tutorial (Smola and Schölkopf, 2004)). Papers (Li et al., 2009; Ivan-268 ciuc, 2007) present several chemical applications for SVR. 269

Moreover, several papers are then devoted to compare the 270 different approaches (see, e.g., Clarke et al., 2004; Amouzgar 271 and Strömberg, 2017; Jin et al., 2001; Bhosekar and Ierapetritou, 272 2018 and references therein): however, there is no definitive un-273 derstanding about the dominance relationships of one type of sur-274 rogate model with respect to the others in terms of accuracy in 275 data representation. Müller and Shoemaker (2015) systematically 276 address the influence of the surrogate model choice and the sam-277 pling method selection on the accuracy of the resulting model. 278 In our study we restrict ourselves to polynomial surrogate mod-279 els since we have already shown in Mencarelli et al. (2019b) that 280 quadratic polynomial models perform sufficiently well with respect 281 to the application we consider in the present paper; however the 282 approaches we will describe can be extended to other types of ba-283 sis functions. 284

Let  $\mathcal{A} \subseteq \mathbb{R}^{|\mathcal{B}|}$  be the set of *a priori* constraints on the regression 285 coefficients  $\beta \in \mathbb{R}^{|B|}$ , such as, e.g., non-negativity constraints. In re-286 gression we aim to minimize an objective function  $g(\beta): \mathbb{R}^{|\mathcal{B}|} \to \mathbb{R}^{|\mathcal{B}|}$ 287  $\mathbb{R}$  representing the distance between the simulated observations 288 or measurements (output data) and the model responses. If the 289 coefficients appear linearly (resp. nonlinearly) in the correspond-290 ing model, then we define the problem as linear regression (LR) 291 (resp. nonlinear regression (NLR)). In our applications, we consider 292 LR in which  $f(x; \beta) := \sum_{j} \beta_{j} f_{j}(x)$ , where  $f_{j}(x) : \mathbb{R}^{K} \to \mathbb{R}$  may be gen-293 eral functions of x. 294

Typical objective functions for regression problems are the sum 295 of the absolute distances: 296

$$\mathbf{g}(\boldsymbol{\beta}) := \sum_{n \in [N]} \left| z_n - \sum_{j \in \mathcal{B}} \beta_j f_j(\mathbf{x}_n) \right|, \qquad (LAD)$$

which gives rise to the so-called least absolute deviation (LAD) cri-297 terion; or the sum of the squares of the residuals: 298

$$g(\beta) := \sum_{n \in [N]} \left( z_n - \sum_{j \in B} \beta_j f_j(x_n) \right)^2, \qquad (ORL)$$

which defines the so-called ordinary least square regression (OLR) 299 problem. For a detailed discussion about advantages and disadvantages of the two previous choices for the objective function, see 301

Chapter 1 in Huber (1981). Regression problems with (LAD) as ob-302 303 jective function can be equivalently reformulated as constrained linear problems (LPs): it is sufficient to replace the absolute values 304 305 with new variables  $t_n$  ( $n \in [N]$ ) and consider the following inequalities for all  $n \in [N]$ : 306

$$t_n \ge z_n - \sum_{j \in \mathcal{B}} \beta_j f_j(x_n)$$
  
$$t_n \ge \sum_{i \in \mathcal{B}} \beta_j f_j(x_n) - z_n.$$
 (1)

In most of the cases, however, the previous approach results 307 into a highly dense surrogate model, i.e., with a large number of 308 309 non-zero coefficients, so that the resulting surrogate model could 310 be difficult to analyze and to optimize. In the surrogate-based superstructure approach, the curse of dimensionality is particularly 311 critical, since the surrogate models are used as blocks in a more 312 313 complex optimization framework.

314 Hence, in surrogate building process dimensionality reduction 315 is therefore a key step to induce sparsity, i.e. to reduce the to-316 tal number of non-zero coefficients: classical approaches include PCA (Cunningham, 2007; Fodor, 2002), random projections (RP) 317 (Krahmer and Ward, 2016), and subset selection (SS) (Miller, 2002). 318 319 Sparsity in the surrogate model can be induced also by adding a 320 regularization parameter to the objective function such as in LASSO approach (Hastie et al., 2015): in this paper we consider only SS 321 approaches since, in this latter case, we have a direct control on 322 the number of non-zero coefficients in the surrogate models. More 323 324 recently, Straus and Skogestad introduce two novel dimensionality 325 reduction methods based on PLSR (Straus and Skogestad, 2017a; 2017b) and self-optimizing control (Straus and Skogestad, 2018), 326 327 respectively.

Dimensionality reduction with application to water-flooding 328 329 production optimization has been recently addressed by Sorek et al. (2017) by introducing the functional control method 330 (FCM) and the interpolation control method (ICM) relying on 331 332 polynomial approximation and piecewise polynomial interpolation controls, respectively (see also Beykal et al., 2018b). 333

In particular, in SS regression only a given subset of dimension 334  $T \leq |\mathcal{B}|$  of basis functions is considered. Hence the following prob-335 336 lem is defined:

$$\begin{array}{ll}
\min_{\mathbf{y}\in\mathbb{R}^{|\mathcal{B}|},\boldsymbol{\beta}\in\mathcal{A}} & g(\boldsymbol{\beta}) \\
\text{s.t.} & \sum_{j\in\mathcal{B}} \mathbf{y}_j = T \\
\underline{\beta}\,\mathbf{y}_j \leq \boldsymbol{\beta}_j \leq \overline{\boldsymbol{\beta}}\,\mathbf{y}_j & \forall j \in B \\
\mathbf{y}_j \in \{0,1\} & \forall j \in B.
\end{array}$$
(SS)

In (SS) |B| binary variables  $y_i$  are introduced to switch on/off 337 338 the corresponding regression parameters  $\beta_i$  ( $j \in B$ ). Depending on the linearity (resp. nonlinearity) of the objective function, SS is a 339 340 MILP (resp. MINLP). Coefficients  $\beta$  and  $\overline{\beta}$  can be estimated in the preprocessing phase using any feasible value for  $\beta$  computed, e.g., 341 by means of the (N)LR approach. In practical implementation we 342 followed the procedure which is suggested in Cozad et al. (2014), 343 i.e., we sum up the absolute values of  $\beta$  found for unconstrained 344 345 regression and we set the obtained numerical value as the upper bound  $\beta$ : then for the lower bound we simply set  $\beta := -\beta$ . 346

There exists several heuristic procedures to argue an opportune 347 numerical value for T: the most adopted ones consist in forward-348 349 and backward-stepwise regressions. Forward-stepwise regression incrementally builds surrogate models by increasing T starting 350 form  $\mathcal{B} = \emptyset$  until a given information criterion, which includes the 351 complexity and the accuracy of the model, is worsen. A possible 352

information criterion is the correct Akaike criterion (AIC<sub>c</sub>):

$$AIC_{c}(T, \beta) := N \log \left( \frac{1}{N} \sum_{n \in [N]} \left( z_{n} - \sum_{j \in B} \beta_{j} f_{j}(x_{n}) \right)^{2} \right)$$
  
+ 2T +  $\frac{2T(T+1)}{N-T-1}$ , (AIC<sub>c</sub>)

which is constituted by a weighted sum of the accuracy of the 354 model, represented by the squares of the model residuals given 355 by the distance between the output data and the surrogate model 356 responses, and the relative complexity of the model, which takes 357 into account the number of basis functions and the total num-358 ber of observations. The flowchart of the SS algorithm is shown 359 in Fig. 1, where  $(\beta^*, y^*)$  is the optimal solution for problem (SS). 360 For other information criteria we refer the interested reader to 361 the paper (Wilson and Sahinidis, 2017). The backward-stepwise re-362 gression approach, on the contrary, initially considers all the basis 363 functions and progressively removes the less significant ones. 364

A comparison between different SS regression strategies is per-365 formed in Kim and Boukoulava (2019). Cozad et al. (2014) intro-366 duce a procedure, the automated learning of algebraic models for 367 optimization (ALAMO), to solve (SS) with a forward-stepwise phi-368 losophy. A comprehensive description of ALAMO with applications 369 to chemical problems is given by Wilson and Sahinidis (2017). 370 Other software packages for surrogate building process are de-371 scribed in Bhosekar and Ierapetritou (2018). 372

## 4. Additional constraints

In practical applications several additional constraints on the 374 responses of the surrogate models might be present. We divide 375 them into two classes: (i) intra-model and (ii) inter-model con-376 straints. Intra-model constraints regard the response of a single 377 surrogate model, while inter-model constraints concern the re-378 sponses of a subset of the models. The presence of inter-model 379 constraints forces the procedure to address the corresponding sub-380 set of surrogate models at the same time. 381

## 4.1. Intra-model constraints

We consider a set of M outputs so that we have one problem 383 of (SS) -type per output. In the notation the variables and the pa-384 rameters of each output model are identified by the superscript 385  $m \in [M]$ . 386

We consider non-negativity constraints and we treat them by 387 means of the approach introduced by Cozad et al. (2015). They pro-388 pose a two-phase procedure: in the first step they build the surro-389 gate model with respect to a given finite set of observations, while 390 in the second one the points corresponding to the maximum vi-391 olation with respect to the constraints for the resulting surrogate 392 configuration are found and added to the set of the first step. The 393 algorithm stops when no violated point is found. 394

In particular, in the first phase we solve a master problem with 395 positivity constraints restricted to a (finite) subset  $\mathcal{X}$  of the closed 396 set  $\mathcal{D} \subset \mathbb{R}^K$  describing the design space. Since in our case studies 397 variables x represent mass fractions (see Section 7), we consider 398 design spaces such that  $\mathcal{D} := \{x \in \mathbb{R}^K : x \in [\underline{x}, \overline{x}] \land \sum_{k \in [K]} x^k = 1\}.$ 399 400

The restricted master problem is:

a ( Om)

$$\min_{\substack{y^m \in \mathbb{R}^{|\mathcal{B}|}, \beta^m \in \mathcal{A}}} g(\beta^m)$$
s.t.
$$\sum_{j \in \mathcal{B}^m} y_j^m = T^m$$

$$\frac{\beta^m y_j^m \le \beta_j^m \le \overline{\beta}^m y_j^m \forall j \in \mathcal{B}^m}{y_j^m \in \{0, 1\} \forall j \in \mathcal{B}^m}$$

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772

J

353

373



Fig. 2. Flowchart of the two-phase algorithm for intra-model constraints

$$\sum_{j\in\mathcal{B}^m}\beta_j^mf_j(x)\geq 0 \forall x\in\mathcal{X}. \tag{$P_1^m$}$$

401 In the second phase, given is a feasible solution  $(\hat{y}^m, \hat{\beta}^m)$  of the 402 problem  $(P_1^m)$ , a non-negative scalar  $\varepsilon_f \in \mathbb{R}_+$ , representing the fea-403 sibility tolerance, and we solve the following optimization problem 404 identifying the maximum violation:

$$\begin{array}{ll} \min_{\mathbf{x}\in\mathcal{D}} & \sum_{j\in\mathcal{B}^m} \hat{\beta}_j^m f_j(\mathbf{x}) \\ \text{s.t.} & \sum_{i\in\mathcal{B}^m} \hat{\beta}_j^m f_j(\mathbf{x}) \leq -\varepsilon_f. \end{array} \tag{P_2}$$

405 A positive small value for  $\varepsilon_f$  enforces strictly positivity for the 406 violation (Cozad et al., 2015), or, in other words, we consider as feasible points the ones for which the corresponding violation is 407 strictly less than  $\varepsilon_{\rm f}$ . 408

5

The optimal solution  $x^*$  of problem (P<sub>2</sub>) is then added to the 409 set  $\mathcal{X}$  and the first phase is performed again. As suggested in 410 Cozad et al. (2015), in order to speed up the algorithm, instead 411 of the optimal solution  $x^*$ , every set of (isolated) feasible solutions 412 found for problem  $(P_2)$  by a state-of-the-art optimization solver, 413 such as BARON (Tawarmalani and Sahinidis, 2005), can be added 414 to  $\mathcal{X}$ . The procedure alternates the two phases until the problem 415 (P<sub>2</sub>) becomes infeasible. We note the previous procedure converges 416 since at each iteration the amount of the violation of the cur-417 rent solution decreases because an increasing number of feasible 418 points is taken into account in the first phase. Problems  $(P_1^m)$  can 419 be solved separately for each output. 420

We initially set  $\mathcal{X}$  equal to the set of all sampled points. As we 421 422 said before, we use all the sampling points in one shot to build the surrogate model. Moreover, we observe that the objective function 423 424 is always evaluated over the same set of initial points for which we know also the real outputs. 425

#### 4.2. Equality inter-model constraints 426

Let  $c \in \mathbb{R}^{K}$  and  $d \in \mathbb{R}^{M}$  be L given vectors of opportune dimen-427 428 sions. We consider additional equality constraints linking the responses of several surrogate models, as follows: 429

$$\sum_{k \in [K]} c_{\ell}^{k} x^{k} = \sum_{m \in [M]} d_{\ell}^{m} \sum_{j \in \mathcal{B}^{m}} \beta_{j}^{m} f_{j}(x) \quad \forall x \in \mathcal{D} \land \forall \ell \in [L].$$
(2)

In particular the previous relationship must hold for  $x_n$  ( $n \in [N]$ ). 430 The resulting problem is semi-infinite since it has an infinite num-431 ber of constraints: we have one constraint for each design configu-432 ration  $x \in \mathcal{D}$ . In our computational experiments we practically con-433 sider problems with hundreds of constraints (see Section 7). In or-434 der to solve the resulting semi-infinite problem, in this case, we 435 should consider all the M surrogate models at once by choosing 436 the objective function  $\sum_{m \in [M]} \omega^m g(\beta^m)$ , which is the weighted sum 437 438 of the objective functions of the models. In the implementation, we simply set  $\omega^m := 1$  for all  $m \in [M]$ ; however, the choice of the 439 model weights  $\omega^m$  constitutes a degree of freedom which can be 440 further explored. In our approach we use the constraints (2) to ex-441 press L variables as functions of the other ones for all the observa-442 443 tions. In this way the equality constraint is automatically satisfied by definition. 444

To be more precise, we model only K - L outputs and we de-445 rive the others from the equality constraints. We note that possible 446 intra-model constraints, such as, e.g., non-negativity constraints, 447 448 should be enforced for all the outputs in order to guarantee possible infeasible solutions are not generated. To the best of our 449 knowledge, the previous technique for equality inter-model con-450 straints is novel and can be applied to all the chemical balance 451 452 constraints regarding, for instance, mass or energy balance.

#### 5. Two-phase PCA 453

454 In order to reduce the number of basis functions in the previous procedure we propose an integration between the two-phase ap-455 456 proach with PCA regression technique. In PCA regression approach a PCA step is performed before the regression procedure. We im-457 plement a PCA step over the basis functions, by decomposing in 458 principal components the (Pearson) correlation matrix  $C \in \mathbb{R}^{|B| \times |B|}$ 459 of the basis functions evaluated over the initial input data. Then 460 we consider only the eigenvectors corresponding to the first largest 461 eigenvalues: hence, we derive new basis functions by projecting 462 463 the original functions onto the subspace generated by the eigenvectors corresponding to the selected eigenvalues. 464

465 We note that performing the PCA over the correlation matrix 466 can be seen as a standardization of the data in order to have the 467 same variation data scale, since in order to define the (Pearson) correlation coefficients we subtract the means and we divide per 468 469 the standard deviations. The (Pearson) correlation coefficients are defined as 470

$$C_{j_{1},j_{2}} = \frac{\sum_{n \in [N]} \left( f_{j_{1}}(x_{n}) - \overline{f}_{j_{1}} \right) \left( f_{j_{2}}(x_{n}) - \overline{f}_{j_{2}} \right)}{\sqrt{\sum_{n \in [N]} \left( f_{j_{1}}(x_{n}) - \overline{f}_{j_{1}} \right)^{2}} \sqrt{\sum_{n \in [N]} \left( f_{j_{2}}(x_{n}) - \overline{f}_{j_{2}} \right)^{2}}},$$
  
$$\forall (j_{1}, j_{2}) \in \mathcal{B} \times \mathcal{B},$$
(3)

471

where  $\overline{f}_{j_1} = (\sum_{n \in [N]} f_{j_1}(x_n))/N$  and  $\overline{f}_{j_2} = (\sum_{n \in [N]} f_{j_2}(x_n))/N$ . In the two-phase approach we solve the first phase consider-477 ing the new basis functions (reducing the dimensionality of the 473 corresponding surrogate building problem) obtained by projecting 474

the original basis function onto the space defined by the princi-475 pal components of the correlation matrix. In the second phase we 476 consider the original design space: the value of the function in the 477 new point is then obtained by projecting the result of the second 478 phase onto the new space. 479

The correlation matrix is decomposed as  $C = \Lambda \Sigma \Lambda^{T}$ , where  $\Lambda$ 480 is the matrix whose columns are the orthonormal eigenvectors of 481 the correlation data matrix and  $\Sigma$  is the diagonal matrix whose 482 diagonal entries are the eigenvalues of matrix C sorted in non-483 increasing order. Let  $\Lambda_{i'}$  be the matrix whose columns are the first 484  $|\mathcal{B}'|$  eigenvectors of the correlation data matrix and F(x) be the ma-485 trix whose rows are the basis function  $f_i(x)$   $(j \in B)$ , we define the 486 new projected matrix of the basis functions as  $F'(x) := \Lambda_{ij}^T F(x)$ . 487 The rows of the matrix F'(x) give the new projected basis func-488 tions  $f'_{i'}(x)$   $(j' \in \mathcal{B}')$ . The problems solved in the first phase read 489 as follows: 490

$$\min_{\substack{\beta^m \in \mathcal{A}'}} g'(\beta^m)$$
s.t.  $\sum_{j' \in \mathcal{B}'} \beta_{j'}^m f_{j'}'(x) \ge 0 \quad \forall x \in \mathcal{X},$  (P<sub>1,PCA</sub>)

where for instance we set  $g'(\beta^m) := \sum_{n \in [N]} |z_n - \sum_{j' \in B'} \beta_{j'} f'_{j'}(x_n)|$ . 491 The set A of a priori constraints over the surrogate parameters is 492 replaced by its projected version  $\mathcal{A}'$ . Problem (P<sub>1.PCA</sub>) is a (N)LP and 493 depends on the number of principal components selected in the 494 PCA step. We note that each principal component corresponds to a 495 basis function. Therefore, in the two-phase PCA approach the selec-496 tion of the basis functions is driven by the value of the eigenvalues 497 of the correlation data matrix of the basis functions calculated over 498 the input data. 499

Moreover, we observe that the PCA step is independent from 500 the output data and can be performed in a preprocessing phase if 501 different outputs should be considered at the same time. This is 502 the case for example when different process alternative configu-503 rations should be evaluated: for instance, if a given chemical pro-504 cess can be realized with a different number of reactors (for a case 505 study see Section 7.1), different surrogate models can be calculated 506 for each possible number of reactors starting from the same in-507 put data opportunely generated according to a DOE strategy (see 508 Section 21 509

510

517

The problem addressed in the second phase is instead

$$\min_{\mathbf{x}\in\mathcal{D}} \qquad \sum_{j'\in\mathcal{B}'} \hat{\beta}_{j'}^m f_{j'}'(\mathbf{x})$$
s.t. 
$$\sum_{j'\in\mathcal{B}'} \hat{\beta}_{j'}^m f_{j'}'(\mathbf{x}) \leq -\varepsilon_f. \qquad (\mathsf{P}_{2,\mathsf{PCA}})$$

In the objective function and in the inter-model constraints we 511 project the original basis function onto the new space defined by 512 the selected eigenvectors. The typology of the two problems and 513 the stopping criterion follow the same philosophy as the two-514 phase approach sketched in the Section 4.1. The flowchart of the 515 two-phase PCA algorithm is given in Fig. 2. 516

## 6. Hybrid approach

In this section we describe a hybrid approach obtained by com-518 bining SS philosophy and PCA regression. In particular, in the hy-519 brid algorithm a SS step is performed once at the beginning to de-520 termine a lower number of representative basis functions and then 521 the two-phase PCA procedure described in the previous section is 522 applied to further reduce the dimensionality of the problem. The 523 complete flowchart of the hybrid algorithm is given in Fig. 3. 524

To be more precise, we first implement the SS algorithm tak-525 ing into account only the intra-model constraints: in this way we 526 can solve a SS-type problem separately for each modelled output 527





528 by decomposing the original problem into *M* simpler independent 529 subproblems with the same structure. We note a parallel imple-530 mentation setting can be exploited in this context (in order to 531 fairly compare the different approaches we consider only pure se-532 quential implementations).

In particular, in presence of equality inter-model constraints, 533 we can avoid to solve the SS-type problem for the L outputs ob-534 tained as functions of the other K - L outputs through the equality 535 constraints. Then, we sum up the values of the binary variables 536 537 introducing to switch on/off the basis functions over the outputs 538 and the values obtained per input are multiplied with the original functions in order to weight them. The functions selected for mul-539 540 tiple outputs in the SS step will have a larger weight and will have, hence, a larger probability to be selected in the PCA step. We ob-541 542 serve that the selected basis functions are employed to model only K-L outputs, while L outputs are still obtained from the equality 543 544 inter-model constraints.

This approach combines the efficacy of the SS strategy to find 545 546 representative basis functions and the computational speed of the 547 PCA regression. The numerical results show, in fact, that a lower number of principal components should be considered to obtain 548 549 (in average) the same accuracy of the surrogate models in case a preliminary SS step is performed. Moreover, we note that in 550 551 the case of pure SS approach a MI(N)LP should be solved at 552 each iteration, while in the two-phase PCA only a (N)LP should be addressed: the (non)linearity of the problem depends on the 553 (non)linearity of the objective function  $g(\beta)$ . Furthermore, per-554 forming a separate SS step per output allows to better capture the 555 complexity of each surrogate model by dealing with different num-556 ber  $T^m$  of basis function per the *m*-th output  $(m \in [M])$ . 557

## 558 7. Computational experiments

559 In our computational setting we divided the input variables into two classes, namely process variables  $p_i$  ( $i \in [J]$ ), which represent 560 561 the control variables in the posterior optimization phase, and composition variables  $c_i$  ( $i \in [I]$ ), which coincide with the mass or molar 562 563 composition of the compounds. We chose polynomial models because we are looking for simple surrogate models, since we would 564 optimize them in a posterior phase in order to retrieve the best 565 operational conditions for the analyzed process. We consider two 566 567 types of surrogate models: (i) polynomial quadratic models and 568 (ii) polynomial cubic models. In particular, in case (i) we consider models with the composition and the process variables occurring 569 linearly and with bilinear interactions, i.e., bilinear mixed products, 570 between process variables and composition variables: such models 571 572 can be expressed in the form

$$\beta_0 + \sum_{i \in [I]} \beta_{1,i} c_i + \sum_{i \in [I]} \sum_{j \in [J]} \beta_{2,ij} c_i p_j + \sum_{i \in [I]} \beta_{3,i} c_i^2$$
(4)

on the contrary in case (ii) we have models where the composition
variables appear linearly and the process variables appear quadratically, i.e., such that can they be expressed in the form

$$\beta_{0} + \sum_{i \in [I]} \beta_{1,i} c_{i} + \sum_{i \in [I]} \sum_{j \in [J]} \beta_{2,ij} c_{i} p_{j} + \sum_{i \in [I]} \beta_{3,i} c_{i}^{2} + \sum_{i \in [I]} \sum_{j' \in [J]} \sum_{j'' \in [J]} \beta_{4,ij'j''} c_{i} p_{j'} p_{j''}.$$
(5)

The number  $n_p$  of parameters in case (i) is given by  $n_p = 1 + |I| + |I||I|$ ; in case (ii) we have instead  $n_p = 1 + |I| + |I||I| + |I||I|^2$ .

After preliminary computational tests with polynomial functions by considering the full cubic model with all the interactions, we decided to restrict ourselves to models resulting from the multiplication of the compositions variables appearing linearly and the process variables up to the quadratic terms, i.e. surrogate models (5) (considering other terms do not add much in terms of model 583 accuracy for the case studies). 584

We choose function (LAD) as objective function for the re-585 gression step (first phase). Each MILP is solved by means 586 of IBM ILOG CPLEX 12.8 IBM ILOG CPLEX with option 587 numericalemphasis and scaind activated and parallelization 588 setting enable (up to 2 threads). The MINLPs are solved through 589 BARON 18.5.8 Tawarmalani and Sahinidis (2005). All the codes 590 are implemented in GAMS 25.1.2 McCarl et al. (2017) on a Dell 591 machine with Intel(R) Xeon(R) CPU E5-1620 v3 at 3.50 GHz with 4 592 GB RAM. 593

We set a time limit of 100 CPU seconds for the first phase and 594 15 CPU seconds for the second phase. Preliminary computational 595 experiments have shown that either the second phase finds a fea-596 sible solution relatively quickly or no solution is found within a 597 larger CPU time, declaring hence the problem as infeasible. More-598 over, in order to speed up the algorithm for the MILP solved in 599 the SS step we stop at the feasible solution found at the root node 600 in the Branch-and-Bound tree (we have observed that the feasible 601 solution found at the root node is not so far from the optimum 602 solution: however, a relatively large amount of time is needed to 603 certify its optimality). 604

For the second phase, we keep the first feasible solution found 605 by BARON and we add it to the restricted master problem. We 606 set  $\varepsilon_f := 0.2$ . Generally, the choice of the numerical value for  $\varepsilon_f$ 607 depends on the order of magnitude of the involved functions: in 608 this case for the second phase we are considering as acceptable 609 points the ones for which the corresponding response is greater 610 than -0.2, which represents a reasonable value in our case stud-611 ies. The computation of the eigenvalues and of the eigenvectors 612 of the correlation data matrix is performed by means of the alge-613 braic tools available in GAMS based on the LAPACK DSYEV rou-614 tine Anderson et al. (1999). 615

## 7.1. Case study 1: Catalytic reforming

In this section we present a real-world application in petroleum 617 refinery industry, namely the catalytic reforming process. Cat-618 alytic reforming (CR) is a chemical refinery process transform-619 ing raw naphtha into high octane gasoline called reformate con-620 taining aromatic hydrocarbons and iso-alcanes (Turaga and Ra-**62**1 manathan, 2003; Gjervan et al., 2004; Lapinski et al., 2014) 622 (for an historical perspective on the studies about CR, see also 623 Rahimpour et al. (2013)). 624

The catalytic reforming process was originally introduced in 625 the 1940s by the Charles Stark Dreape laureate Vladimir Haensel<sup>1</sup>, 626 who proposed the so-called Platforming process, adopting a cata-627 lyst containing platinum. The CR unit is the most important pro-628 cess in refinery industry to produce lead-free automobile fuel and 629 hydrogen. CR unit is composed of a sequence of reactors (usually 630 from 3 to 5) characterized by operating conditions (temperature, 631 pressure, molar hydrogen-to-hydrocarbon ratio, and feed composi-632 tion) and equipped with catalyst (typically with platinum). 633

The main chemical reactions in CR are, in fact, dehydrogenation 634 and hydroisomerization of naphtenes transforming them into aro-635 matics, isomerization and dehydrocyclization of alcanes convert-636 ing them into aromatics and iso-alcanes, hydrocraking of alcanes 637 into smaller components, hydrogenolysis, and coke formation. Coke 638 formation is a relatively slow process and it represents a damag-639 ing reaction since coke reduces the performances of the catalyst. 640 Typical operating conditions are high temperature (450–500°C), 641

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772

8

<sup>&</sup>lt;sup>1</sup> see patents Alumina-platinum-halogen catalyst and preparation thereof. 1949, August 16. U.S. Patent No. 2,479,109 and Process of reforming a gasoline with an alumina-platinum-halogen catalyst. 1949, August 16. U.S. Patent No. 2,479,110.

able 1 cronyms.	
Acronyms (al	phabetical order)
AICc	correct Akaike information criterion
ALAMO	Automated learning of algebraic models for optimization
ANN	Artificial neural network
ARGONAUT	Algorithms for global optimization of constrained grey-box computational problems
ATOUT	Advanced tools for optimization and uncertainty treatment
CR	Catalytic reforming
DOE	Design of experiments
FCM	Functional control method
FFD	Fully factorial design
GDP	General disjunctive problem
ICM	Interpolation control method
iP	iso-alcanes
IS	Isomerization
IAD	Least absolute deviation
LHD	Latin hypercube design
LP	Linear problem
LR	Linear regression
MILP	Mixed integer linear problem
MINLP	Mixed integer nonlinear problem
NLR	Nonlinear regression
nP	n-alcanes
OLR	Ordinary least square regression
PCA	Principal component analysis
PLSR	Partial least square regression
RBF	Radial basis function
RON	Research octane number
RP	Random projections
SIP	Semi-infinite problem
SS	Subset selection
SVR	Support vector regression
WHSV	Weight hourly space velocity

Table 2 Input compounds of the catalytic reforming process.

Input compounds	
nP <sub>7</sub>	heptane
iP7	isoheptane
N7	naphthenes with 7 atoms of carbons
A7	toluene
nP <sub>8</sub>	octane
iPa	isoctane
N <sub>8</sub>	naphthenes with 8 atoms of carbons
Ag	ethyl-benzene + xylenes

medium level of pressure (3-35 atm) and molar hydrogen-tohydrocarbon (H<sub>2</sub>/HC) ratio between 3 and 8 Ancheyta-Juárez and Villafuerte-Macías (2001).

645 In particular, for illustration we consider the C7–C8 cut. The in-646 puts of the model are the research octane number (RON) and the 647 mass percentage of the hydrocarbon compounds occurring in the 648 CR process (see Tables 2–3). The pressure and the temperature of 649 the chemical reactions are treated as constants. Hence, we have 650 one process variable (RON) and eight composition variables (mass 651 percentages).

In this case, since the model outputs represent mass fractions, 652 653 non-negativity and summing up to 1 constraint mush be enforced. Moreover, in CR we have one equality constraint (2) for each *l*-654 655 th chemical element ( $\ell \in [L]$ ) (hydrogen and carbon in hydrocarbon 656 compounds) and *n*-th observation ( $n \in [N]$ ), expressing the equivalence between the total number  $\xi_{\ell,n}^{in}$  of moles of the  $\ell$ -th element 657 in the process inflow and the total number  $\xi_{\ell,n}^{out}$  of the moles of the 658  $\ell$ -th element in the process outflow in the *n*-th observation, i.e., 659  $\xi_{\ell,n}^{in} = \xi_{\ell,n}^{out}$  for all  $\ell \in [L]$  and  $n \in [N]$ . The number of moles for the 660 *e*-th element is given by the weighted sum of the molar percent-661 age of all the chemical compounds in the stream, whose weights 662 are the number of moles of hydrogen and carbon occurring in the 663

## Table 3

Output compounds of the catalytic reforming process.

Output compounds	
H <sub>2</sub>	hydrogen
P <sub>1</sub>	methane
P <sub>2</sub>	ethane
P <sub>3</sub>	propane
P4	butane
P <sub>5</sub>	pentane
6P <sub>6</sub>	n-hexane
5P6	simple branched alcanes with 6 atoms of carbons
4P <sub>6</sub>	double branched alcanes with 6 atoms of carbons
N <sub>6</sub>	naphthenes with 6 atoms of carbons
A <sub>6</sub>	benzene
7P <sub>7</sub>	n-heptane
6P7	simple branched alcanes with 7 atoms of carbons
5P7	double branched alcanes with 7 atoms of carbons
N7	naphthenes with 7 atoms of carbons
A7	toluene
8P8	n-octane
7P <sub>8</sub>	simple branched alcanes with 8 atoms of carbons
6P8	double branched alcanes with 8 atoms of carbons
N <sub>8</sub>	naphthenes with 8 atoms of carbons
A <sub>8</sub>	ethyl-benzene + xylenes

corresponding compound, i.e.,

$$\xi_{\ell,n}^{in} = \sum_{k \in [K]} \frac{\operatorname{mol}_{\ell,k}^{in}}{\operatorname{mass}_{k}^{in}} c_{n}^{k}, \tag{6}$$

where  $mol_{\ell,k}^{in}$  and  $mass_k^{in}$  are the number of moles of the  $\ell$ -th element in the k-th input compounds and the molar mass of the k-th input compounds, respectively, and 667

$$\xi_{\ell,n}^{out} = \sum_{m \in [M]} \frac{\operatorname{mol}_{\ell,m}^{out}}{\operatorname{mass}_{m}^{out}} z_{n}^{m},$$
(7)

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772

9

### Table 4

<b>R</b> 2	and weighted	$\mathbb{R}^2$ (w $\mathbb{R}^2$ ) for	r PCA and SS+PC	A per number (	of principal	components
<b>n</b> -		V_ I MVV_ I IO				CONTRACTOR

	PCA				SS+PCA			
	R <sup>2</sup>		wR <sup>2</sup>		R <sup>2</sup>		WR <sup>2</sup>	
No. comp.	Train	Test	Train	Test	Train	Test	Train	Test
30	0.85	0.78	0.94	0.90	0.87	0.81	0.95	0.92
31	0.85	0.78	0.94	0.90	0.88	0.87	0.96	0.95
32	0.88	0.87	0.96	0.95	0.88	0.87	0.96	0.95
33	0.88	0.87	0.96	0.95	0.88	0.87	0.96	0.95
34	0.88	0.87	0.96	0.96	0.88	0.87	0.96	0.95
35	0.88	0.87	0.96	0.95	0.89	0.87	0.96	0.96
36	0.89	0.88	0.96	0.96	0.89	0.86	0.96	0.96
37	0.89	0.88	0.96	0.96	0.88	0.88	0.96	0.96
38	0.89	0.89	0.96	0.96	0.89	0.88	0.96	0.96
39	0.89	0.89	0.96	0.96	0.90	0.89	0.97	0.95
40	0.90	0.91	0.97	0.97		-	-	-
avg	0.88	0.86	0.96	0.95	0.88	0.87	0.96	0.95

where  $mol_{\ell,m}^{(ndt)}$  and  $mass_m^{(ndt)}$  are the number of moles of the  $\ell$ -th element in the *m*-th output compounds and the molar mass of the *m*-th output compounds, respectively.

Note that satisfying the balance constraints on the number of moles implies all the mass fractions sum up to 1 and also the mass balance holds. Hence, these constraints have not been enforced in the first phase problem.

675 We have generated the plan of experiments by means of the 676 software package Design Expert 10 Design Expert In particular, we adopt the D-optimal design for the training set and the 677 I-optimal design for the test set for model (5). We chose statistical 678 designs since for this kind of DOE Design Expert lets us to consider 679 680 other additional constraints (such as that the sum of all the mass fractions for a given observation is equal to 1) when the plan of 681 682 experiments is developed: D-optimal and I-optimal design are referred to the constraint case for the cubic model (5). In particular, 683 684 we define a training set with 226 samples and a test set with 200 685 samples.

The software tool used to numerically simulate the catalytic reforming process is OSCAR 1.1, a numerical software developed at IFP Energies nouvelles.

We build individual surrogate models for each output component: as stated in the previous sections, we model only m - 2 outputs since we have 2 equations for the molar balance of carbon and hydrogen, and we obtain the remaining 2 outputs from the equality inter-model constraints.

In order to evaluate the performance of a surrogate model, in 694 695 addition to average coefficient of determination (R<sup>2</sup>) over the outputs, we consider a weighted  $R^2$  (w $R^2$ ) which consists in the aver-696 697 age of the R<sup>2</sup> for the single models weighted with the average output mass fraction computed over the observations. In Table 4 we 698 report the values of R<sup>2</sup> and wR<sup>2</sup> per number of principal compo-699 nents for the two-phase PCA approach and for the hybrid approach 700 701 (SS+PCA).

The maximum number of principal components coincides with the number of considered basis functions: in the surrogate model for the catalytic reforming we are considering only the composition variables appearing linearly and the interaction between the process variable and the compositions variables up to the quadratic term, i.e., model (5) (see Section 7). Before dimensionality reduction the cubic model (5) has 40 parameters.

For the hybrid algorithm the maximum number of the principal components is given by the number of basis functions found in the SS step. It is worth observing that in the SS step the number of basis functions is already reduced from 40 to 39. The gain is relatively small: in this case we are not using the full cubic models with all the parameters, but we select *a priori* a subset of terms for

Table 5
CPU times (in seconds) for PCA and SS+PCA per
number of principal components.

PCA	SS+PCA	Δ
457.89	543.34	1 <b>5.73%</b>
<b>501.58</b>	480.36	-4.42%
452.73	504.94	1 <b>0.34%</b>
476.43	504.58	5.58%
508.24	568.73	10.64%
497.81	566.78	1 <b>2.17%</b>
566.99	574.13	1.24%
587.06	591.95	0.83%
583.76	628.23	7.08%
603.74	625.85	3.53%
550.84	1	-
526.10	558.89	6.27%
	PCA 457.89 501.58 452.73 476.43 508.24 497.81 566.99 587.06 583.76 603.74 550.84 550.84 526.10	PCA         SS+PCA           457.89         543.34           501.58         480.36           452.73         504.94           476.43         504.58           508.24         568.73           497.81         566.78           566.99         574.13           587.06         591.95           583.76         628.23           603.74         625.85           550.84         -           526.10         558.89

the interactions between variables. Preliminary computational ex-715 periments with the full cubic model with all the possible mixed 716 bilinear products have shown the SS step allows us to significantly 717 reduce the number of parameters in the surrogate models. The to-718 tal number of parameters for a polynomial model of degree d in 719 *n* variables is  $\binom{n+d}{d}$ : a full cubic model has 220 parameters. In the 720 case of full cubic model, the SS step let us to decrease the num-721 ber of basis functions up to 18, obtaining a  $R^2$  and a  $wR^2$  indices 722 for the training set equal to 0.87 and 0.96, respectively and a  $R^2$ 723 and a  $wR^2$  indices for the test set equal to 0.89 and 0.96, respec-724 tively. In general, larger the number of basis functions the SS can 725 chose among, better the performances of the SS step in terms of 726 selection of the basis functions. 727

Table 4 shows the hybrid approach achieves the same perfor-728mances as the two-phase PCA approach with a smaller number of729principal components, both for the training and the test set. More-730over, we note the values of  $wR^2$  is always larger than the ones of731 $R^2$ , because the surrogate model in presence of equality constraints732tends to better estimate the components characterized by a higher733percentage concentration.734

From Table 5, which reports the computational times for the 735 two proposed algorithms and the percentage increase  $\Delta$  between 736 the computational time of SS+PCA and the computational time of 737 PCA. It is clear that the SS step has an important impact on the 738 CPU times: except for the cases of 31 principal components, in 739 which the CPU time of the hybrid method is smaller than the one 740 of the two-phase PCA, in all the other cases the CPU time of the 741 hybrid method is the largest one. The average increase of the CPU 742 time is 6.27% with a maximum of 15.73% for the case of 30 prin-743 cipal components. Higher the increase in the relative performance 744 of the surrogate model, higher the increase of the CPU time. There-745



Fig. 4. H<sub>2</sub>, predicted vs. actual plots. The values given by quadratic model (4) in green (QM); and the values given by the cubic model (5) in blue (CM).



Fig. 5. H<sub>2</sub>, predicted vs. actual plots, cubic model (5). The values given by the unconstrained regression in green (UR); and the values given by constrained regression in blue (CR).

fore, the results underline a trade-off between the relative reduc-tion on the number of basis functions and the relative increase ofcomputational effort.

Figs. 4a-7b show the scatter parity plots for the H<sub>2</sub> for the 749 training set (figures with caption (a)) and the test set (figures with 750 caption (b)). In the figures the x-axis reports the values of the sim-751 752 ulated outputs, while the y-axis reports the values of the response of the surrogate models. In the limit case in which the response of 753 754 the models coincides with the simulated output, the corresponding point lies on the bisector of the first quadrant. Smaller the distance 755 756 between the value and the bisector of the first quadrant, better the performance of the surrogate model. 757

758 The quality of the surrogate models strongly depends on the 759 grade of the mixed products between the process and the com-760 position variables as shown in Fig. 4a and b which compare the quadratic model, i.e., Eq. (4), and the cubic model, i.e., Eq. (5), and 761 on the fact that we are obliged to consider constrained regres-762 763 sion (see Fig. 5a and b which report the parity plots obtained by constrained regression over all the basis functions and the uncon-764 strained LAD for cubic model (5)). 765

A key driver of the performance of the surrogate models is clearly the number of principal components and hence of basis 
 Table 6

 Relative R<sup>2</sup> and wR<sup>2</sup> increases for the 30/35 and 35/40 principal components.

Index	30/35	35/40
R <sup>2</sup> train	3.41%	2.22%
R <sup>2</sup> test	10.34%	4.40%
wR <sup>2</sup> train	2.08%	1.03%
wR <sup>2</sup> mest	5.26%	2.06%

functions considered in the estimation process: Fig. 6a and b show 768 the predicted outputs for the two-phase PCA by varying the num-769 ber of principal components. From the parity plot, it is possible to 770 observe that the relative increase of estimation quality between 30 771 and 35 principal components is more accentuated than the one ob-772 tained by passing from 35 to 40 principal components (the relative 773  $R^2$  and  $wR^2$  increases for the 30/35 principal components and for 774 the 35/40 principal components are reported in Table 6). 775

Fig. 7a and b report the comparison between the two-phase 776 PCA and the hybrid approach for 35 principal components, show-777 ing the second method is slightly better than the first one in terms 778 of distance between simulated and predicted outputs. 779

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772



Fig. 6. H<sub>2</sub>, predicted vs. actual plots, cubic model (5). The values given by PCA with 30 principal components in green (PCA 30); the values given by PCA with 35 principal components in blue (PCA 35); and the values given by PCA with 40 principal components in orange (PCA 40).



Fig. 7. H<sub>2</sub>, predicted vs. actual plots, cubic model (5). The values given by the PCA with 35 principal components in green (PCA); and the values given by SS+PCA with 35 principal components in blue (SS+PCA).

780 Finally, Fig. 8a and b show the residual gap, calculated as

$$\frac{\sum_{j\in\mathcal{B}^m}\beta_j^m f_j(x) - z_n}{\sum_{j\in\mathcal{B}^m}\beta_j^m f_j(x)} \quad \forall n \in [N],$$
(8)

sorted according to the RON input value (in the figures we report only the values for H<sub>2</sub>). It is worth noting that the surrogate models reproduced sufficiently the behaviour of the simulated values with regards to the process variable, which is indeed the control variable in the operating phase of the considered process: the values of the residual gaps for H<sub>2</sub> are in absolute value less than 0.25 for the training set and 0.15 for the test set.

In conclusion, in the CR case study the performances of the hybrid method are generally slightly better than the ones of the two phase PCA.

## 791 7.2. Case study 2: Isomerization

Isomerization (IS) is a chemical process which increases RON
index of light hydrocarbon (C5–C6) by transforming n-alcanes into
branched iso-alcanes with higher octane index (Valavarasu and
Sairam, 2013; Sullivan et al., 2014). IS has been introduced in 1930s
by Vladimir Ipatieff, who proposed a new chemical process to
transform butane into isobutane, and used during the World War

II to produce high octane aviation gasoline (for an historical anal-798 ysis of IS process, see Sullivan et al. (2014)). Nowadays IS is fun-799 damental to produce high octane fuel and reduce the level of ben-800 zene, aromatics and olefins in gasoline. IS unit is usually composed 801 of a single reactors operating at relatively low temperatures (110-802 150°C) Valavarasu and Sairam (2013). In theory, in fact, hydrocar-803 bons with C6 content could be treated via the catalytic reforming, 804 but the constraint over the benzene content in the gasoline makes 805 the process infeasible. 806

Low operating temperatures are necessary in order to minimize 807 cracking of the hydrocarbons, and imply the chemical reactions are 808 relatively slow: this effect is balanced, for instance, by means of 809 highly active catalysts. 810

Table 7 reports the input and output compounds in the IS pro-811cess. We consider the inverse of the weight hourly space veloc-812ity (WHSV), temperature, pressure, and mass fractions of the in-813put compounds as input data and the mass fractions of the output814compounds as output data.815

As in the previous case study, we have one balance constraints 816 for each chemical element (carbon and hydrogen) which should 817 be satisfied by all the input configurations belonging to the design space. A latin hypercube design of experiments is gener-



Fig. 8. H2, RON vs. residual gaps, cubic model (5). The residual gap given by the PCA with 35 principal components in violet (PCA); and the residual gap given by SS+PCA with 35 principal components in green (SS+PCA).

Table	7			
Input	and	output	compounds	o
the is	omeri	ization c	process.	

Input compounds					
nP <sub>5</sub>	pentane				
iP5	isopentane				
nP <sub>6</sub>	hexane				
2iP <sub>6</sub>	2-methylpentane				
3iP6	3-methylpentane				
22iP6	2,2-dimethylhexane				
23iP6	2,3-dimethylhexane				
Output	compounds				
nP4	butane				
iP4	isobutane				
nP <sub>5</sub>	pentane				
iP5	isopentane				
nPe	hexane				
2iP5	2-methylhexane				
3iP5	3-methylhexane				
22iP <sub>6</sub>	2.2-dimethylhexane				
23iP <sub>6</sub>	2,3-dimethylhexane				

Table 8  $R^2$ , weighted  $R^2$  (w $R^2$ ), and time (in seconds) for PCA per number of principal components.

	R <sup>2</sup>		wR <sup>2</sup>		
No. comp.	Train	Test	Train	Test	Time
50	< 0	< 0	< 0	< 0	599.80
90	72.36	<b>72.3</b> 1	98.43	98.65	800.49
91	71.89	72.37	98.46	98.72	837.86
92	71.97	72.38	98.50	98.72	<b>964.9</b> 1
93	72.03	<b>72.4</b> 1	98.54	98.82	882.16
94	72.05	72.50	98.54	98.78	884.78
95	72.01	72.45	98.54	98.79	810.49
96	72.02	<b>72.5</b> 1	98.55	98.84	853.40
97	72.34	72.79	98.72	98.91	933.40
98	72.74	73.16	98.93	98.13	835.45
avg	<b>72</b> .16	72.54	98.58	98.71	866.99

Table 9

R<sup>2</sup>, weighted R<sup>2</sup> (wR<sup>2</sup>), and time (in seconds) for SS+PCA per number of principal components.

-	R <sup>2</sup>		wR <sup>2</sup>		
No. comp.	Train	Test	Train	Test	Time
47	<b>66.2</b> 1	< 0	93.44	48.03	488.01
48	65.80	< 0	93.46	50.15	487.16
49	72.05	68.38	97.91	94.05	467.78
50	72.37	73.24	98.40	99.42	283.67
avg	<b>72.2</b> 1	70.81	<b>98</b> .16	96.7 <b>4</b>	375.73

tational time for the hybrid method to reach a satisfactory grade of accuracy is approximately half of the time needed for the two-843 phase PCA time to achieve the same quality level.

Fig. 9a and b show the scatter parity plots for the 2iP<sub>6</sub> for the training set (figures with caption (a)) and the test set (figures with caption (b)). As for the CR process case study, the quality of the surrogate models clearly depends on the grade of the 848 polynomial models (see Fig. 9a and b) and on the presence of the 849 850 additional non-negativity and molar conservation constraints (see Fig. 10a and b). The previous plots refer to models (4) and (5) with 851 the *a priori* selection of the interactions (bilinear terms) between 852 process and composition variables. 853

Fig. 11a and b report the parity plot for the two-phase PCA with 854 90, 95, and 98 principal components, showing the light quality im-855 provement of the surrogate models by considering an increasing 856 number of principal components. 857

ated by means of the software package Advanced tools for opti-820 mization and uncertainty treatment (ATOUT 1.1) developed and 821 maintained by IFP Energies nouvelles. In particular, we define a 822 823 training set with 200 samples and a test set with 5000 samples. We use the simulation software developed at IFP Energies nou-824 velles in order to simulate the industrial process performances. The 825 total number of basis functions in the model before dimensionality 826 827 reduction is 98 for the cubic model (5).

The same numerical trends observed for the CR process are 828 829 valid also for the IS. Tables 8 and 9 show the  $R^2$ ,  $wR^2$  and the computational time (in seconds) for the two-phase PCA and the 830 hybrid algorithm, respectively (the average referred to the principal 831 components with positive value of R<sup>2</sup> indices). We note the two-832 833 phase PCA method is able to reduce the number of basis functions 834 by maintaining however a comparable model accuracy. The hybrid 835 method achieves the same grade of model accuracy with less basis functions: in particular, already with 49 basis function the hybrid 836 algorithm is rather capable to capture the informations of the (sim-837 838 ulated) data, while in order to achieve the same model quality the two-phase PCA method approach needs 90 basis functions. 839

840 For the computational time related to the surrogate model generation, the lower number of basis functions implies the compu-841

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772



Fig. 9. 2iP<sub>6</sub>, predicted vs. actual plots. The values given by quadratic model (4) in green (QM); and the values given by the cubic model (5) in blue (CM).



Fig. 10. 2iP<sub>6</sub>, predicted vs. actual plots, cubic model (5). The values given by the unconstrained regression in green (UR); and the values given by constrained regression in blue (CR).



Fig. 11. 2iP<sub>5</sub>, predicted vs. actual plots, cubic model (5). The values given by PCA with 90 principal components in green (PCA 90); the values given by PCA with 95 principal components in blue (PCA 95); and the values given by PCA with 98 principal components in orange (PCA 98).



Fig. 12. 21P6, predicted vs. actual plots, cubic model (5). The values given by the PCA with 98 principal components in green (PCA 98); and the values given by SS+PCA with 50 principal components in blue (SS+PCA 50).

Finally, Fig. 12a and b report the comparison between the two-858 859 phase PCA and the hybrid approach, graphically showing that the 860 latter algorithm is characterized by a model accuracy comparable to the one of the two-phase PCA method, by considering, however, 861 a lower number of basis functions. 862

In conclusion, the IS case study highlights the benefits of the 863 864 hybrid approach, which is able to obtain the same model accuracy with about half the number of basis functions of the two-phase 865 866 PCA approach.

#### 8. Conclusions 867

868 In this paper, we have designed a systematic methodology to 869 define and compute surrogate models for a given black-box pro-870 cess. In particular, we have discussed the DOE strategies and the main approaches for the identification of the surrogate model, fo-871 cusing on the SS perspective. Moreover, we illustrate how to deal 872 873 with possible intra-model and equality inter-model constraints. We 874 have introduced a new two-phase PCA procedure for constrained regression problems by combining the two-phase approach in 875 Cozad et al. (2015) with a PCA regression strategy. Therefore, we 876 have discussed a possible hybrid strategy to combine the SS ap-877 878 proach with the introduced two-phase PCA procedure.

The methods are evaluated and compared with respect to two 879 880 case studies in petroleum refinery process, namely catalytic reforming and isomerization. For both case studies we have shown 881 that the two-phase PCA method is able to reduce the number of 882 basis functions required to obtain a satisfactory model accuracy. 883 and the hybrid algorithm (the two-phase PCA preceded by a SS 884 885 step) achieves a satisfactory model quality with a lower number of basis functions than the simple two-phase PCA methodology. The 886 887 reduction in the number of basis function is significant in the isomerization case study, where the hybrid approach is able to obtain 888 889 a satisfactory model accuracy with half the number of basis functions of the two-phase PCA. 890

891 In future work we would like to extend our methodology to 892 consider also PLSR approaches. Moreover, our approaches can be 893 easily extended to consider other functions than polynomials as basis functions: hence, analyzing the performances of these sur-894 rogate models in the context of PCA and SS+PCA could be another 895 896 interesting future research axis. In our study, we have considered all the sampling points at once: consequently incrementally adding 897 the sampling points could be an interesting strategy in order to ob-898 tain an accurate surrogate model with a lower number of sampling 899

points. We have considered only noiseless data, we would like to 900 test the two-phase procedure for noisy data relative to physical ex-901 periments in a further study. 902

Then, we are also interested in embedding the surrogate model-903 ing approach into an optimization framework, where the surrogate 904 model replaces the physical model to retrieve the optimal config-905 uration and operating conditions of a given chemical process. The 906 surrogate model is sufficiently accurate for this, and as its form 907 is simple, it is fast to compute, and it allows the use of power-908 ful global optimizers, that can fully exploit its analytic expression. 909 Moreover, we are confident that the two-phase dimension reduc-910 tion approaches described in the paper could be applied to other 911 chemical processes. If needed, other basis functions than polyno-912 mial terms could be used as basis functions. 913

## Acknowledgments

We are grateful to Master student Amina Bougueroua for 915 preliminary computational experiments concerning the numerical 916 characterization of the design of experiments for the isomerization 917 process. We would thank also two anonymous referees whose sug-918 gestions significantly improve the quality of this paper. 919

## References

- Altissimi, R., Brambilla, A., Deidda, A., Semino, D., 1998. Optimal operation of a sepa-921 ration plant using artificial neural networks. Computers & Chemical Engineering 922 22 (Supplement 1), S939-S942. 923
- Amouzgar, K., Strömberg, N., 2017. Radial basis functions as surrogate models with a 924 priori bias in comparison with a posteriori bias. Structural and Multidisciplinary 925 Optimization 55 (4), 1453-1469
- Ancheyta-Juárez, J., Villafuerte-Macías, E., 2001. Experimental validation of a kinetic 927 model for naphtha reforming. Studies in Surface Science and Catalysis 133, 615-928 618 929
- Anderson, E., Bai, Z., Bischof, C., Blackford, S., Demmel, J., Dongarra, J.J., Du Cruz, J., Greenbaum, A., Hammarling, S., McKenney, A., Sorensen, D., 1999, LAPACK User's Guide. Society for Industrial and Applied Mathematics (SIAM).
- Audet, C., Le Digabel, S., Tribes, C., Rochon Montplaisir, V.,. The NOMAD project. Software available at https://www.gerad.ca/nomad/.
- Beykal, B., Boukoulava, F., Floudas, C.A., Pistikopoulos, E.N., 2018. Optimal design of 935 936 energy systems using constrained grey-box multi-objective optimization. Computers & Chemical Engineering 116, 488-502. 937 938
- Beykal, B., Boukoulava, F., Floudas, C.A., Sorek, N., Zalavadia, H., Gildin, E., 2018. Global optimization of grey-box computational systems using surrogate functions and application to highly constrained oil-field operations. Computers & Chemical Engineering 114, 99-110.
- Bhosekar, A., Ierapetritou, M., 2018, Advances in surrogate based modeling, feasi-942 bility analysis, and optimization: A review. Computers & Chemical Engineering 943 108, 250-267 944

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772

920

926

930

931

932

933

934

939

940

941

914

947

961

962

963

964

965

966

967

968

985

987

991

992

993

995

997

1001

- Bouhlel, M.A., Bartoli, N., Otsmane, A., Morlier, J., 2016. Improving kriging surrogates 945 of high-dimensional design models by Partial Least Squares dimension reduc-946 tion. Structural and Multidisciplinary Optimization 53 (5), 935-952.
- Boukoulava, F., Floudas, C.A., 2017. ARGONAUT: AlgoRithms for Global Optimization 948 949 of coNstrAined grey-box compUTational problems. Optimization Letters 11 (5), 950 895-913.
- 951 Boukoulava, F., Hasan, M.M.F., Floudas, C.A., 2017. Global optimization of general 952 constrained grey-box models: new method and its application to constrained PDEs for pressure swing adsorption. Journal of Global Optimization 67 (1-2), 953 954
- 955 Boukoulava, F., Misener, R., Floudas, C.A., 2016. Global optimization advances in 956 mixed-integer nonlinear programming, MINLP, and constrained derivative-free 957 optimization, CDFO, European Journal of Operational Research 255 (3), 701-727.
- 958 Box, G.E.P., Hunter, J.S., Hunter, W.G., 2005. Statistics for experimenters: Design, innovation, and discovery. John Wiley & Sons, Hoboken, NJ. 959
- Broomhead, D.S., Lowe, D., 1988. Radial basis functions, multi-variable functional 960 interpolation and adaptive networks. 4148. Royal Signal & Radar Establishment. Buhmann, M.D., 2000. Radial basis functions. Acta Numerica 9, 1-38.
  - Caballero, J.A., Grossmann, I.E., 2008. An algorithm for the use of surrogate models in modular flowsheet optimization. AIChE Journal 54 (10), 2633-2650.
  - Clarke, S.M., Griebsch, J.H., Simpson, T.W., 2004. Analysis of support vector regression for approximation of complex engineering analyses. Journal of Mechanical Design 127 (6), 1077-1087.
- Coetzer, R., Haines, L.M., 2017. The construction of D- and I-optimal designs for mix-969 ture experiments with linear constraints on the components. Chemometrics and 970 Intelligent Laboratory Systems 171, 112-124.
- Conn. A.R., Le Digabel, S., 2013, Use of quadratic models with mesh-adaptive di-971 rect search for constrained black box optimization. Optimization Methods and 972 973 Software 28 (1), 139-158.
- Cornell, J.A., 2002. Experiments with mixtures: Designs, models, and the analysis of 974 975 mixture data. John Wiley & Sons, New York.
- 976 Cozad, A., Sahinidis, N.V., Miller, D.C., 2014. Learning surrogate models for simulation-based optimization. AIChE Journal 60 (6), 2211-2227, 977
- Cozad, A., Sahinidis, N.V., Miller, D.C., 2015. A combined first-principles and data-978 979 driven approach to model building. Computers & Chemical Engineering 73, 116-980 127.
- 981 Cunningham, P., 2007. Dimension reduction. UCD-CSI-2007-7. University College Dublin. 982
- Davis, S.E., Cremaschi, S., Eden, M.R., 2018. Efficient surrogate model development: 983 984 Impact of sample size and underlying model dimensions. Computer Aided Chemical Engineering 44, 979-984.
- van Dam, E.R., 2008. Two-dimensional minimax Latin hypercube designs. Discrete 986 Applied Mathematics 158 (18), 3483-3493.
- 988 van Dam, E.R., Husslage, B.G.M., den Hertog, D., Melissen, H., 2007. Maximin Latin 989 hypercube design in two dimensions. Operations Research 55 (1), 158–169. 990 Design Expert., https://www.statease.com/software.html.
  - Eason, J., Cremaschi, S., 2014. Adaptive sequential sampling for surrogate model generation with artificial neural networks. Computers & Chemical Engineering 68 (4), 220-232,
- 994 Fahmi, I., Cremaschi, S., 2012. Process synthesis of biodiesel production plant using artificial neural networks as the surrogate models. Computers & Chemical 996 Engineering 46, 105–123.
- Fodor, I.K., 2002. A survey of dimension reduction techniques. Technical Report. 998 Center for Applied Scientific Computing, Lawrence Livermore National Labora-999 tory.
- 1000 Forrester, A.I.J., Keane, A.J., 2009. Recent advances in surrogate-based optimization.
- Progress in Aerospace Sciences 45 (1-3), 50-79. Forrester, A.I.J., Sobester, A., Keane, A.J., 2008. Engineering design via surrogate mod-elling: A practical guide. John Wiley & Sons, Hoboken, NJ. 1002 1003
- 1004 Garud, S.S., Karimi, I.A., Kraft, M., 2017. Design of computer experiments: A review. Computers & Chemical Engineering 106, 71-95. 1005
- 1006 Garud, S.S., Karimi, I.A., Kraft, M., 2017. Smart sampling algorithm for surrogate 1007 model development. Computers & Chemical Engineering 96, 103-114.
- 1008 Gaspar, B., Teixeira, A.P., Guedes Soares, C., 2017. Adaptive surrogate model with 1009 active refinement combining Kriging and a trust region method. Reliability Engineering & System Safety 165, 277-291. 1010
- Gjervan, T., Prestvik, R., Holmen, A., 2004. Catalytic reforming. In: Baerns, M. (Ed.), 1011 1012 Basic Principles in Applied Catalysis. Springer, Berlin, Heidelberg, pp. 125-158. 1013 Goos, P., Jones, B., Syafitri, U., 2016. I-optimal design of mixture experiments. Journal
- 1014 of the American Statistical Association 111 (514), 899–911. Hardy, R.L., 1971. Multiquadratic equations of topography and other irregular sur-1015 1016 faces. Journal of Geophysical Research 76 (8), 1905-1915.
- 1017 Hastie, T., Tibshirani, R., Wainwright, M., 2015. Statistical learning with sparsity: The 1018 Lasso and generalizations. CRC Press.
- 1019 Henao, C.A., Maravelias, C.T., 2010. Surrogate-based process synthesis. Computer 1020 Aided Chemical Engineering 28, 1129–1134.
- Henao, C.A., Maravelias, C.T., 2011. Surrogate-based superstructure optimization framework. AIChE Journal 57 (5), 1216–1232. 1021 1022
- Huber, P.J., 1981. Robust statistics. John Wiley & Sons, Hoboken, NJ. 1023
- 1024 Husslage, B.G.M., Rennen, G., van Dam, E.R., den Hertog, D., 2011. Space-filling Latin 1025 hypercube designs for computer experiments. Optimization and Engineering 12 (4), 611-630 1026
- IBM ILOG CPLEX, https://www.ibm.com/analytics/cplex-optimizer. 1027
- 1028 Ivanciuc, O., 2007. Application of support vector machine in chemistry. In: Lipkowitz, K., Cundari, T., Boyd, D. (Eds.), Reviews in Computational Chemistry, 23. 1029 1030 John Wiley & Sons, Hoboken, NJ.

- lin, R., Chen, W., Simpson, T.W., 2001, Comparative studies of metamodelling tech-1031 niques under multiple modelling criteria. Structural and Multidisciplinary Opti-1032 mization 23 (1), 1-13. 1033
- Johnson, M.E., Moore, L.M., Ylvisaker, D., 1990. Minimax and maximin distance de-1034 signs. Journal of Statitical Planning and Inference 26 (2), 131-148. 1035 Jones, D.R., Schonlau, M., Welch, W.J., 1998. Efficient global optimization of expen-1036

1037

1038

1071

1**072** 

1**073** 

1074

1075

1076

1085

1100

1104

1106

1107

- sive black-box functions. Journal of Global Optimization 13 (4), 455-492. Joseph, V.R., Hung, Y., 2008. Orthogonal-maximin latin hypercube designs. Statistica Sinica 18 (1), 171–186.
- 1039 Kim, S.H., Boukoulava, F., 2019. Machine learning-based surrogate modeling for 1040
- data-driven optimization: A comparison of subset selection for regression tech-1**04**1 nique. Optimization Letters 1-22 1042 1043
- Kleijnen, J.P.C., 2009. Kriging metamodeling in simulation: A review. European Journal of Operational Research 192 (3), 707-716. 1044
- Krahmer, F., Ward, R., 2016. A unified framework for linear dimensionality reduction 1045 in 11. Results in Mathematics 70 (1-2), 209-231. 1046
- Krige, D.G., 1952. A statistical approach to some basic mine valuation problems on 1047 the Witwatersrand. Journal of the Southern African Institute of Mining and Met-1048 allurgy 52 (9), 201-203. Lapinski, M.P., Metro, S., Pujadó, P.R., Moser, M., 2014. Catalytic reforming in 1049
- 1050 petroleum processing. In: Treese, S., Jones, D., Pujadó, P. (Eds.), Handbook of 1051 Petroleum Processing. Springer, Cham, pp. 1-25. 1052
- algorithm. ACM Transactions on Mathematical Software 37 (4), 44:1-44:15.
- istry. Chemometrics and Intelligent Laboratory Systems 95 (2), 188-198.
- Matheron, G., 1963. Principles of Geostatistics. Economic Geology 58 (8), 1246–1266. McBride, K., Sundmacher, K., 2019. Overview of surrogate modeling in chemical pro-
- 1059 McCarl, B. A., Meeraus, A., van der Eijk, P., Bussieck, M., Dirkse, S., Nelissen, F., 1060 2017. McCarl Expanded GAMS User Guide, GAMS Release 24.6.GAMS Develop-1061 1062
- 1063 selecting values of input variables in the analysis of output from a computer 1064 code. Technometrics 21 (2), 239-245. 1065
- Mencarelli, L., Chen, Q., Pagot, A., Grossmann, I.E., 2019. A review on superstructure 1066 optimization approaches in process system engineering. Technical Report. IFP 1067 Energies nouvelles, Solaize, France, and Department of Chemical Engineering, 1068 Carnegie Mellon University, Pittsburgh. 1069 1070

Mencarelli, L., Duchêne, P., Pagot, A., 2019. Optimization approaches to the integrated system of catalytic reforming and isomerization processes in petroleum refinery. Technical Report. IFP Energies nouvelles, Solaize, Prance.

- Miller, A., 2002. Subset selection in regression. Chapman & Hall/CRC, Boca Raton, Florida.
- Morris, M.D., Mitchell, T.J., 1995. Exploratory designs for computational experiments. Journal of Statitical Planning and Inference 43 (3), 381-402. Müller, J., Shoemaker, C.A., 2015. Influence of ensemble surrogate models and sam-
- 1077 pling strategy on the solution quality of algorithms for computationally expen-1078 sive black-box global optimization problems. Journal of Global Optimization 60 1079 (2), 123-144. 1080
- Nascimiento, C.A.O., Giudici, R., Guardani, R., 2000, Neural network based approach 1081 for optimization of industrial chemical processes. Computers & Chemical Engi-1082 neering 24 (9-10), 2303-2314. 1083 1084
- Park, J., Sandberg, I.W., 1993. Approximation and radial-basis-function networks. Neural Computation 5 (2), 305-316.

Petelet, M., Iooss, B., Asserin, O., Loredo, A., 2010. Latin hypercube sampling with 1086 inequality constraints. AStA Advances in Statistical Analysis 9 (4), 325-339. 1087

Pronzato, L. 2017. Minimax and maximin space-filling designs: Some properties and 1088 methods for construction. Journal de la Société Française de Statistique 158 (1), 1089 1090 7-36.

- Psaltis, A., Sinoquet, D., Pagot, A., 2016. Systematic optimization methodology for 1091 heat exchanger network and simultaneous process design. Computers & Chem-1092 ical Engineering 95, 146-160. 1093
- Queipo, N.V., Hafka, R.T., Shyy, W., Goel, T., Vaidyanathan, R., Tucker, P.K., 2005. 1094 Surrogate-based analysis and optimization. Progress in Aerospace Sciences 41 1095 1096 (1), 1-28.
- Rahimpour, M.R., Jafari, M., Iranshavi, D., 2013. Progress in cambytic naphtha reform-1097 ing process: A review. Applied Energy 109, 79-93. 1098 1099
- Sacks, J., Welch, W.J., Mitchell, T.J., Wynn, H.P., 1989. Design and analysis of computer experiments. Statistical Science 4 (4), 409–485. Smith, W.F., 2005. Experimental design for formulation. ASA-SIAM Series on Statis-
- 1101 tics and Applied Probability, 15. SIAM, Philadelphia. 1102 1103
- Smola, A.J., Schölkopf, B., 2004. A tutorial on support vector regression. Statistics and Computing 14 (3), 199-222.
- Sorek, N., Gildin, E., Boukoulava, F., Beykal, B., Floudas, C.A., 2017. Dimensionality 1105 reduction for production optimization using polynomial approximations. Computational Geosciences 21 (2), 247-266. Straus, J., Skogestad, S., 2017. Use of latent variables to reduce the dimension of
- 1108 surrogate models. Computer Aided Chemical Engineering 40, 445-450. 1109 1110
- Straus, J., Skogestad, S., 2017. Variable reduction for surrogate modelling. In: Pro seeding of Foundations of Computer-Aided Process Operations, Tucson, AZ, USA, 1111 8-12 January 2017. Straus, J., Skogestad, S., 2018. Surrogate model generation using self-optimizing vari-1112
- 1113 ables. Computers & Chemical Engineering 119, 143-151. 1114

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772

Le Digabel, S., 2011. Algorithm 909: NOMAD: Nonlinear optimization with the MADS 1053 1054 Li, H., Lianh, Y., Xu, Q., 2009. Support vector machines and in applications in chem-1055 1056 1057 1058

cess engineering. Chemie Ingenieur Technik 91 (3), 228-239.

ment Corporation, Washington, DC, USA. McKay, M.D., Beckman, R.J., Conover, W.J., 1979. A comparison of three methods for

1139

- Straus, J., Skogestad, S., 2019. A new termination criterion for sampling for surrogate model generation using partial least squares regression. Computers & Chemical 1115 1116 1117 Engineering 121, 75-85.
- 1118 Sullivan, D., Metro, S., Pujadó, P.R., 2014. Isomerization in Petroleum Processing. 1119 In: Treese, S., Jones, D., Pujadó, P. (Eds.), Handbook of Petroleum Processing. 1120 Springer, Charn, pp. 1-15.
- Tawarmalani, M., Sahinidis, N.V., 2005. A polyhedral branch-and-cut approach to 1121 1122
- global optimization. Mathematical Programming 103 (2), 225-249. Turaga, U.T., Ramanathan, R., 2003. Catalytic naphtha reforming: Revisiting its im-1123 1124 portance in the modern refinery. Journal of Scientific and Industrial Research 62 1125 (10), 963–978.
- 1**126** Valavarasu, G., Sairam, B., 2013. Light naphtha isomerization process: A review. 1127 Petroleum Science and Technology 31 (6), 580-595.
- Vapnik, V., 1995. The nature of statistical learning theory. Springer-Verlag, NY. Vapnik, V., Golowich, S., Smola, A.J., 1997. Support vector method for function ap-
- 1129 proximation, regression estimation, and signal processing. In: Mozer, M., Jor-1130 dan, M., Petsche, T. (Eds.), Advances in Neural Information Processing Systems 1131 9. MIT Press, Cambridge, MA, pp. 281-287. 1132
- Viana, F.A.C., 2016. A tutorial on latin hypercube design of experiments. Quality and Reliability Engineering International 32 (5), 1975–1985. Vu, K.K., D'Ambrosio, C., Hamadi, Y., Liberti, L., 2017. Surrogate-based methods for
- 1135 black-box optimization. International Transactions in Operational Research 24 1136 (3), 393-424. 1137 1138
- Wilson, Z.T., Sahinidis, N.V., 2017. The ALAMO approach to machine learning. Computers & Chemical Engineering 106, 785-795.

Please cite this article as: L. Mencarelli, A. Pagot and P. Duchêne, Surrogate-based modeling techniques with application to catalytic reforming and isomerization processes, Computers and Chemical Engineering, https://doi.org/10.1016/j.compchemeng.2020.106772