

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

Analysis of the performance of a Crude-oil Desalting System based on historical data

Ehsan Ranaee¹, Hamzeh Ghorbani², Sajjad Keshavarzian¹, Pejman Ghazaeipour Abarghoei³,
Monica Riva⁴, Fabio Inzoli^{1*}, Alberto Guadagnini⁴

¹Dipartimento di Energia, Politecnico di Milano, Via Lambruschini 4, 20156 Milano, Italy

²Islamic Azad University, Ahvaz branch, Young Researchers Club, Iran

³Iran University of Science and Technology, Narmak, Tehran, Iran

⁴Dipartimento di Ingegneria Civile e Ambientale, Politecnico di Milano, Piazza L. Da Vinci 32,
20133 Milano, Italy

* Corresponding author. Tel. +39 02 2399 3883. Fax. +39 02 2399 3913

E-mail address: fabio.inzoli@polimi.it

Abstract

Our study is keyed to the development of a methodological approach to assess the workflow and performance associated with the operation of a crude-oil desalting/demulsification system. Our analysis is data-driven and relies on the combined use of (a) Global Sensitivity Analysis (GSA), (b) machine learning, and (c) rigorous model discrimination/identification criteria. We leverage on an extensive and unique data-set comprising observations collected at a daily rate across a three-year period at an industrial plant where crude oil is treated through a combination of demulsification/desalting processes. Results from GSA enable us to quantify the system variables which are most influential to the overall performance of the industrial plant. Machine learning is then applied to formulate a set of candidate models whose relative skill to represent the system behavior is quantified upon relying on model identification criteria. The integrated approach we propose can then effectively assist to (a) modern and reliable interpretation of data associated with performances of the crude oil desalting process and (b) robust evaluation of future performance scenarios, as informed by historical data.

Keywords: Crude oil plant assessment, Sensitivity Analysis, Uncertainty Quantification, Principal Component Analysis, Machine Learning.

1. Introduction

Crude oil extracted from reservoirs generally contains water with soluble salts, solid particles, and/or dissolved chemicals, including metals (e.g., vanadium, nickel) [1-4]. The presence of water reduces the total volume of the oil delivered to refineries and increases operational costs. Soluble salts can cause negative impacts on refinery facilities, in the form of, e.g., corrosion, and/or reduced effectiveness of the catalysts that are used for crude oil refining processes. Such phenomena, in turn, lead to increased maintenance/service expenses, in terms of intensive energy use, reduction of refinement system performance, and uncontrolled changes in oil properties, which can then hamper the effectiveness of the entire energy cycle.

Previous experiences (e.g., [5-6]) document that corrosive properties of water-soluble salts can cause significant damage to operating equipment (e.g., pipes, valves, pumps, tanks, and tanker vessels) as well as internal parts of refinery distillation towers. Impurities may cause fouling on interior surfaces of the equipment, thus potentially leading to clogging of the oil heater pipes and increase of temperature and pressure therein, with critical implications on energy use. Laboratory experiments performed by Hamadi et al. [7] document that impurities in crude oil lead to increased values of density and viscosity for crude oil emulsion. This, in turn, can reduce the API gravity degree of the crude oil emulsion as compared to pure oil, thus yielding an oil of reduced marketable quality.

In this setting, crude oil is typically processed in desalting/demulsification units before being conveyed to refineries. The latter may be operated on the basis of various techniques (e.g., gravity-based, thermal, chemical, electrical, or through approaches grounded on clean water injection, membrane use, ultrasonic waves, and microwave) or a combination of diverse technologies.

Demulsification is typically designed as a pre-treatment to break water/oil emulsions and remove impurities and salty water from the oil. It takes place according to two main stages, i.e., (i) emulsification and (ii) separation of salt/water from oil. Emulsification relies on mixing fresh water with a salt-water-oil emulsion. Studies have been performed to assess the behavior of this process,

59 including the evaluation of the size and distribution of the water droplets in crude oil emulsion and
60 droplet breakage and coalescence under turbulent flow conditions [8]. After mixing fresh water with
61 crude oil, salt-water should be separated from the organic phases to complete the demulsification
62 process.

63 After demulsification, oil is collected from the various desalting units and a blend is formed to
64 comply with the characteristics of the target refinery. Properties of crude oils extracted from diverse
65 sedimentary basins (with differing petrophysical properties) can exhibit marked variabilities,
66 including differences of salinity and impurities content. It is then important to verify if oil batches
67 processed in diverse desalting units (and then collected into a single refining unit) are associated with
68 similar characteristics. Laboratory tests on samples of crude oil are typically performed and the type
69 and amount of demulsifiers to be employed is suggested on the basis of the crude-oil properties [9].

70 Demulsification may also be implemented jointly with other types of desalting techniques. It is
71 then challenging to have a robust appraisal of the extent at which laboratory tests can be scaled to
72 field conditions to cope with, e.g., (i) interaction of a mixture of different types of demulsifiers, and/
73 or (ii) interaction of the demulsification with other desalting techniques (possibly) implemented in
74 the field.

75 Less et al. [10] study the performance of an electrostatic coalescer by analyzing a set-up based
76 on electrical desalting technique alone and in conjunction with the injection of demulsifying materials
77 into the crude oil. Meidanshahi et al. [11] analyze the feedback between salt concentration and water
78 droplet size when demulsification is performed through electrical desalting. These results suggest that
79 simultaneous use of electrical and chemical techniques for desalting show a markedly higher potential
80 than relying on each technique as a standalone process.

81 In this broad context, it is then emphasized that optimization of oil treatment processes through
82 a firm characterization of the system functioning stands as an important step towards enhancing
83 sustainability within oil and gas industries with the aim of maximizing general efficiency with a

84 reduction of energy consumptions. Along these lines, we consider here an industrial plant comprising
85 a set of desalting units where a combination of desalting techniques is implemented in a coherent and
86 structured way (see details in Section 2). We then rely on a modern theoretical approach to the
87 analysis of experimental data collected in this industrial plant to allow evidencing the relative
88 importance and relationship between performance of the industrial plant and monitored parameters
89 driving the system behavior.

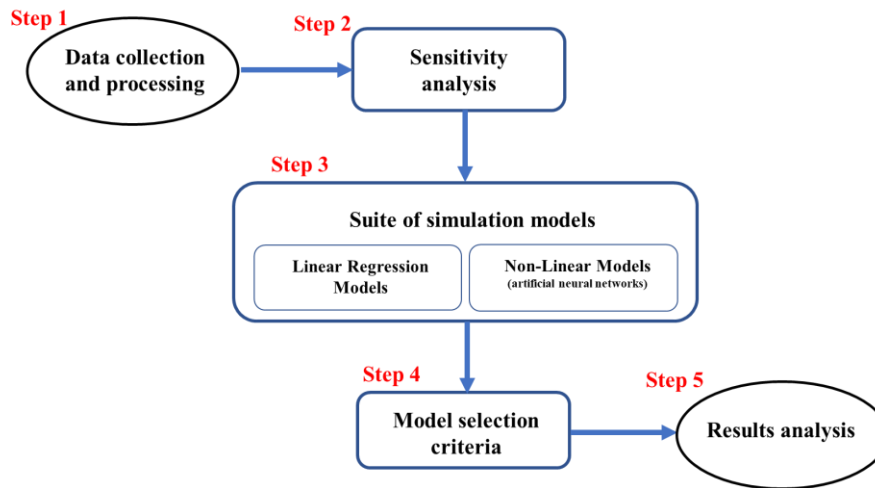
90 Dai et al. [12] recently proposed a data-driven approach for the optimization of a crude oil
91 blending system and improvement of the operating efficiency of a refinery. These authors rely on a
92 case study to illustrate the effectiveness of the proposed methodology. Multivariate probabilistic
93 analyses have been employed to diagnose and manage large datasets associated with a variety of
94 applications in the petroleum industry, and in more general energy industry sector. Qin and Chiang
95 [13] provide a comprehensive review of recent developments and applications of machine learning
96 and artificial intelligence (AI) approaches, as applied to big commercial data sets. These authors
97 remark that uncertainties and possible lack of critical measurements are major aspects to be
98 considered when applying these approaches to industrial settings. Siena et al. [14] rely on Bayesian
99 hierarchical clustering (BHC) and Principal Components Analysis (PCA) techniques to select the
100 most appropriate Enhanced Oil Recovery (EOR) technique amongst a set of available options for a
101 given reservoir. Sad et al. [15] perform a set of laboratory experiments on oil samples produced from
102 various wells in the same field to study compliance of physic-chemical properties of blends of crude
103 oil with the characteristics of a target refinery. These authors implement a hierarchical cluster analysis
104 to classify laboratory data associated with samples collected from different production wells.
105 Nonlinear methods such as artificial neural networks (ANNs) are also applied in a multivariate
106 analysis context to optimize the efficiency of a desalting unit (e.g., [16]). Mohammad et al. [17]
107 develop and apply ANN models to evaluate the performance of a liquid desiccant dehumidifier in
108 terms of the water condensation rate and dehumidifier effectiveness.

109 Some applications of machine learning in crude oil refinery planning and management can be
110 found in the literature. Mouret et al. [18] introduce a Lagrangian decomposition scheme to
111 simultaneously optimize (a) a crude-oil operations scheduling model and (b) a coarse refinery
112 planning model. More recently, Gao et al. [19] apply machine learning concepts to introduce a
113 piecewise linear modeling approach to assist refinery scheduling, with an application to the operations
114 associated with a refinery in China. Ochoa-Estopier et al. [20] and Gueddar and Dua [21] consider
115 ANNs for the optimization of heat-integrated crude oil distillation systems and refinery units,
116 respectively.

117 In this framework, Sensitivity Analysis (SA) may serve as a general framework within which
118 one can assess the way a given quantity of interest is influenced by diverse controlling (or driving)
119 elements. Al-Qahtani and Elkamel [22] consider uncertainty in process parameters to analyze
120 performance of a refinery plan as well as an integration network of petroleum refineries. These
121 authors document that the refinery model analyzed is most sensitive to variations in the pricing of
122 imported crude oil and exported final products as opposed to changes in product demand. Typical
123 applications of SA are concerned with the diagnosis of the behavior of a given model of a system
124 functioning. As such, the relationship between a target quantity (or modeling goal) and the related
125 controlling variables is filtered by the structure of the model considered. A major purpose of a SA is
126 then to ranking and screening the control variables which are most influential to the target quantities
127 according to some predefined sensitivity metrics. When one is interested in the analysis of the way
128 multiple controlling variables influence a set of target quantities that jointly contribute to define the
129 performance of the system considered, performing a SA considering each of the target variables
130 separately might be influenced by the occurrence of possible correlations among these. Lamboni et
131 al. [23] proposed to circumvent this issue by relying on a combined use of (i) a principal component
132 analysis (PCA) of the set of multivariate target variables and (ii) variance-based sensitivity indices to

133 characterize the contribution of each controlling quantity on the set of identified principal
134 components.

135 As such, a distinctive objective of our analysis is the assessment of a workflow based on a
136 suitable theoretical framework leading to the identification of the key quantities which are typically
137 monitored in a desalting plant that can facilitate the control of the optimal performance of the plant.
138 We do so by providing an approach which enables us to effectively employ a combination of SA
139 tools in the absence of a clear model structure relating controlling and target variables. We rely on a
140 unique (in terms of quality and quantity of data) dataset to evaluate the influence of system variables
141 on some desired system state quantity (which would represent a modeling goal). Thus, our analysis
142 is observation-driven and rests on a multivariate analysis framework to assess information collected
143 from field settings involving the use of different demulsifiers in combination with a set of desalting
144 techniques to discuss the way the ensuing results can assist estimation of quantities of interest (i.e.,
145 production of wastewater, processed oil and salt from the desalting unit) representing the
146 effectiveness/performance of the system. We do so by first identifying a set of control variables
147 associated with ambient/processing conditions of the joint crude oil demulsification/desalting system
148 considered. We then (i) consider jointly various SA approaches to assess the relative influence of the
149 system control variables on the performance of the demulsification/desalting system, as quantified
150 through key target quantities, (ii) formulate a suite of linear/ nonlinear models to interpret
151 observations of such target variables, and (iii) rely on formal model selection criteria to evaluate the
152 relative skills of the considered models to characterize the system behavior.



153

154 **Fig. 1** Flowchart of the analysis framework.

155 Our analyses are structured according to the flowchart depicted in Fig. 1 and comprising the
 156 following five sequential steps:

- 157 • *step 1*: acquisition of field data associated with a demulsification/desalting industrial plant;
- 158 • *step 2*: evaluation (through SA approaches) of the relative influence of the main controlling
 159 variables on the selected target quantities, with possible reduction (based on SA results) of the
 160 dimensionality of the space identified by the controlling variables;
- 161 • *step 3*: application of a collection of linear/nonlinear simulation models to interpret observations
 162 of the selected target quantities;
- 163 • *step 4*: assessment of the quality of the considered simulation models through model selection
 164 criteria;
- 165 • *step 5*: analysis of the results to summarize and highlight key relationships between controlling
 166 variables and the target quantities of interest considered to evaluate the performance of the
 167 demulsification/desalting system.

168 With reference to the simulation models, we rely on a typical linear regression modeling
 169 framework and an artificial neural networks- (ANN-) based non-linear approach (e.g., [24-26]).

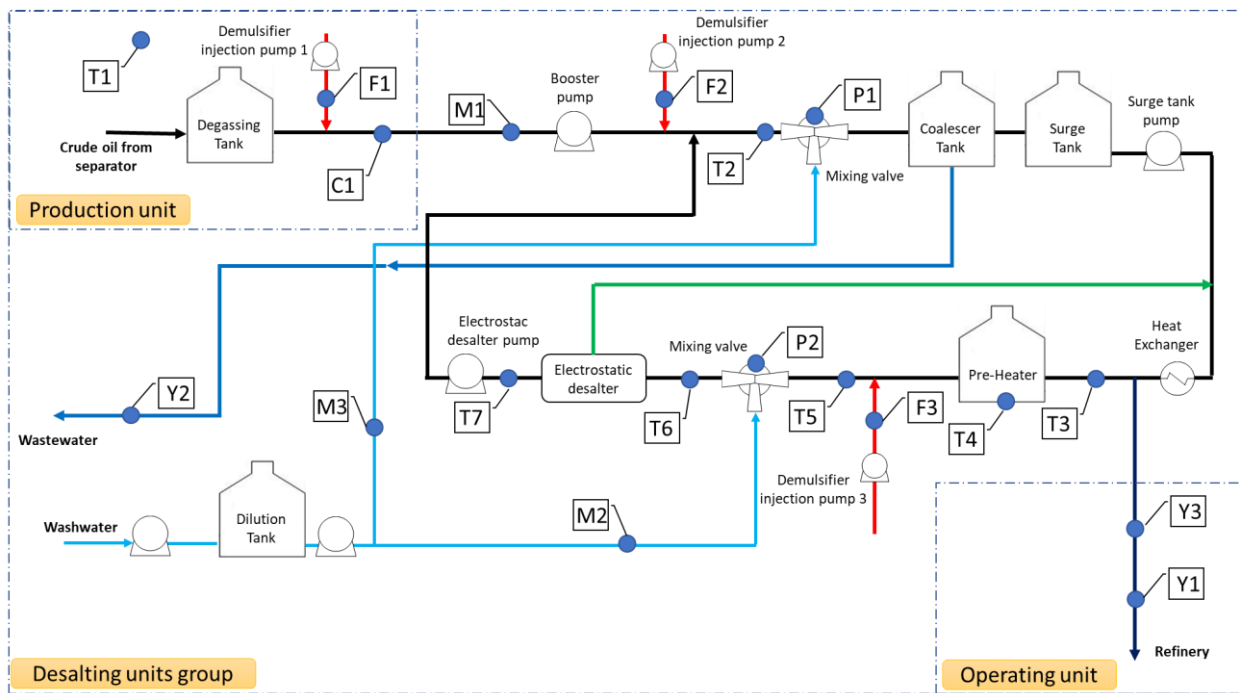
170 The study is organized as follows. We briefly describe in Section 2.1 the desalting plant
 171 analyzed and the collected data. Section 2.2 is devoted to the process of on-site data collection. We

172 then provide a synthesis of the key theoretical elements underpinning sensitivity analysis and
 173 dimensionality reduction technique (Section 2.3) as well as simulation models (Section 2.4). Model
 174 selection criteria are introduced in Section 2.5. Implementation of the analysis workflow and the
 175 ensuing results are discussed in Section 3.

176 2. Methodology and application

177 2.1. Plant analysis and data collection

178 The investigated plant comprises 3 different modules, serving the same oil field. Fig. 2 depicts
 179 a schematic representation of a desalting module, including the types of data collected. These plants
 180 are logically structured according to 3 main units, i.e., production, desalting, and operating units.



181

Collected data					
ID	Ref	description	ID	Ref	description
X ₁	M1	Oil volume flow rate @ Production unit	X ₁₁	T4	Fluid temperature @ Pre-heater bath
X ₂	M2	Fresh water volume flow rate @ Electrostatic desalter	X ₁₂	T5	Fluid temperature @ outlet Pre-heater
X ₃	M3	Fresh water volume flow rate @ Coalescer tank	X ₁₃	T6	Fluid temperature @ inlet Electrostatic desalter
X ₄	F1	Demulsifier injection @ Production Unit	X ₁₄	T7	Fluid temperature @ inlet Electrostatic desalter
X ₅	F2	Demulsifier injection @ Coalescer tank	X ₁₅	P1	Pressure drop @ mixing valve (Coalescer tank)
X ₆	F3	Demulsifier injection @ Electrostatic desalter	X ₁₆	P2	Pressure drop @ mixing valve (Electrostatic desalter)
X ₇	C1	Salt amount @ Production unit			
X ₈	T1	Ambient temperature	Y ₁	Y1	Oil volume flow rate @ Operating unit
X ₉	T2	Fluid temperature @ inlet Coalescer tank	Y ₂	Y2	Wastewater volume flow rate
X ₁₀	T3	Fluid temperature @ inlet Pre-heater	Y ₃	Y3	Salt amount @ Operating unit

182

183 **Fig. 2** Schematic drawing of the analyzed crude oil desalting module. Solid blue circles correspond

184 *to locations of probes for data collections. Bottom table lists parameter identifiers (ID) and*
185 *references to the plant schematics (Ref).*

186 At the early stage, crude oil enters the production unit and is conveyed to the degassing tank
187 for gas separation. A demulsifier is injected into the oil exiting from the degassing tank, before it
188 enters the desalting unit. The salt content of the emulsion is measured at the entrance of the desalting
189 unit. In the next stage, crude oil is pumped to the coalescer tank by means of a crude booster pump.
190 The pressure drop across a mixing valve located before the coalescer tank is recorded, fluid
191 temperature being monitored at the entrance of the coalescer tank. The emulsion remains in the
192 coalescer tank for the time needed to separate part of the wastewater from oil. The internal structure
193 of the coalescer tank enables crude oil to be subject to a rotating motion within the tank (typically for
194 24 hours, depending on the amount of crude oil) and provides sufficient time for water droplets to be
195 separated from the oil. Water droplet separation in the coalescer tank is expedited by injection of
196 fresh water from the dilution tank and by addition of a demulsifier to break the water-in-oil emulsion.
197 The resulting separated water is then conveyed to the wastewater section, the oil entering the surge
198 tank. Oil exiting from the coalescer tank typically contains around 0.5% to 2% water, depending on
199 process conditions. Oil accumulates in the surge tank and is then pumped at a uniform flow rate to
200 the heat exchanger. Afterwards, oil enters the pre-heater and, after injection of the demulsifier, enters
201 the electrostatic desalter. Fluid temperature is also recorded before, inside and after the pre-heater
202 unit. Pressure drop at a mixing valve positioned before the electrostatic desalter is recorded. A small
203 amount (about 3 - 10%) of fresh water is injected into the mixing valve upstream the electrostatic
204 desalter to decrease salt concentration in the crude oil. Relying on an electrical field and on emulsion
205 breakage increases the rate of deposition of water droplets. Accordingly, the remaining droplets in
206 the salt water are separated from the oil in the electrostatic desalting tank. Fluid temperature data are
207 also collected before and after electrostatic desalter. Finally, the output oil from the electrostatic
208 desalter enters again the heat exchanger to be cooled down, before being conveyed to the operating

209 unit. Water extracted from the crude oil in the electrostatic desalting tank is pumped (by means of an
 210 electrostatic desalter pump) to the coalescer tank to (i) increase the fluid temperature and (ii) reduce
 211 concentration of the salt at the entrance of the coalescer tank. The rate of the produced oil and the
 212 amount of separated salt is measured before sending the processed oil to the refinery. The plant
 213 performance is then assessed through the evaluation of (i) the rate of production oil (Y_1), (ii) the rate
 214 of wastewater at the plant outlet (Y_2), and (iii) the amount of separated salt detected in the oil
 215 production stream (Y_3).

216 We collect across a three-year period onsite measurements of the 19 variables in Fig. 2, from 3
 217 different desalting units (in the same oil field) where demulsification is performed through injection
 218 of 5 different demulsifiers in conjunction with (i) the gravity force desalting technique, (ii) a thermal
 219 desalting technique, and (iii) the washing water injection approach. A total number of $N = 1820$
 220 observations are available for each variable during the three-year period.

2.2. Data analysis methodology

222 Monitored values of the quantities listed in Fig. 2 are collected in matrices $\underline{\underline{\mathbf{X}}^*}$ and $\underline{\underline{\mathbf{Y}}^*}$:

$$223 \underline{\underline{\mathbf{X}}^*} = \underline{\underline{\mathbf{X}}^*} (\mathbf{X}_1^*, \mathbf{X}_2^*, \dots, \mathbf{X}_{16}^*); \quad \underline{\underline{\mathbf{Y}}^*} = \underline{\underline{\mathbf{Y}}^*} (\mathbf{Y}_1^*, \mathbf{Y}_2^*, \mathbf{Y}_3^*) \quad (1)$$

224 where the entries of each vector \mathbf{X}_i^* ($i = 1, 2, \dots, 16$) and \mathbf{Y}_k^* ($k = 1, 2, 3$) correspond to N available
 225 data.

226 The demulsification process can be described through the following functional relationship

$$227 \mathbf{Y} = f(\mathbf{X}) \quad (2)$$

228 where vectors $\mathbf{Y} = (Y_1, Y_2, Y_3)$ and $\mathbf{X} = (X_1, X_2, \dots, X_{16})$ include all (generally unknown) variables listed
 229 in Table 1 whose observed values are collected in $\underline{\underline{\mathbf{X}}^*}$ and $\underline{\underline{\mathbf{Y}}^*}$, respectively.

230 Our strategy rests on the study of the way the variability of a set of target state variables Y_k
 231 (collected in \mathbf{Y}) is influenced by the remaining system variables (also termed as controlling variables,

232 collected in \mathbf{X}). We do so through a unique application of a combination of Global Sensitivity
233 Analysis (GSA) approaches (see Section 2.3).

234 Upon relying on GSA, we assess the possibility of performing a dimensionality reduction of
235 the space identified by the system controlling variables. Accordingly, the functional relationship (2)
236 becomes

$$237 \quad f_1(\mathbf{Y}, \mathbf{X}_R) = 0 \quad (3)$$

238 where \mathbf{X}_R includes the set of R controlling variables associated with the reduced dimensionality
239 space which is eventually identified. The functional relationship (3) is then characterized through
240 linear regression and nonlinear ANN-based modeling techniques (Section 2.4). Parameters of the
241 constructed simulation models are evaluated upon relying on the available data (Section 2.1). We
242 recall that the latter are collected from 3 different desalting units (in the same field) and 5 types of
243 demulsifiers have been used across the considered 3 years temporal frame. As such, the development
244 of an interpretive model should also consider the possibility of including information about (i) the
245 specific desalting module from which each measurement is collected and (ii) the type of demulsifiers
246 used for a given record. We take this type of information into account in our regression models by
247 following a strategy proposed in prior studies (e.g., [27]) and mapping information on the location of
248 data sampling and the type of employed demulsifiers into a set dummy (categorical) variables
249 collected, respectively, in vectors \mathbf{U} and \mathbf{Q} . We set $U_1 = 1$ (while $U_2 = U_3 = 0$) when a measurement
250 record is collected in the first desalting units, corresponding notations being employed to identify
251 data sampled at the second or third desalting units. The same approach is implemented to define \mathbf{Q} ,
252 whose entries are dummy variables Q_l ($l = 1, 2, \dots, 5$) quantifying information on the type of
253 demulsifier used. Variables $U_1 - U_3$ and $Q_1 - Q_5$ can then be included in \mathbf{X}_R to characterize the
254 ensuing simulation model in a such way that:

$$255 \quad \mathbf{X}_R = \mathbf{X}_R(\mathbf{X}, \mathbf{Q}, \mathbf{U}) \quad (4)$$

2.3. Sensitivity Analysis

256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280

We rely on a Sensitivity Analysis (SA) approach to identify the system controlling variables which (amongst those listed in Fig. 2) are most influential to the effectiveness of the demulsification process. As stated in Section 1, our SA is completely data-driven and relies on the available dataset described in Sections 2.1 and 2.2. As such, our results enable us to evaluate the influence of system variables on some desired system state quantity without the filter of an interpretive model. Outcomes of the SA are then employed to drive the construction of a collection of candidate models (see Section 2.4) whose relative interpretive skills are evaluated according to the approach illustrated in Section 2.5.

Our SA starts by considering in (3) the performance of the system as rendered through the following 3 target quantities, which are measured as outputs from the Desalting Unit and form the entries of \mathbf{Y} , i.e., (i) Y_1 , oil volume; (ii) Y_2 , wastewater volume, and (iii) Y_3 , salt mass.

Various approaches to SA are proposed in the literature to diagnose the behavior of a given model, where a target modeling goal depends on a number of uncertain model parameters. These are essentially framed within the context of (i) Local Sensitivity Analysis (LSA) or (ii) Global Sensitivity Analysis (GSA). LSA approaches are typically based on the evaluation of local derivatives (also termed elementary effects or local sensitivities) of modeling goals with respect to a control variable at a given location in Γ . Such an approach allows quantifying the sensitivity of a given (dependent) target variable Y_k to perturbations of X_i and information about the importance of X_i is limited to the location where a given elementary effect is assessed [28]. These results can (in principle) be extended to the entire space where model control variables are defined when, e.g., (i) the target system state variables (Y_k) display a linear behavior with respect to the control variables (X_i), and (ii) an adequate number of trajectories is sampled within Γ for the evaluation of the sensitivities. While the elementary effects are evaluated at a set of points in the parameter space, our study rests on a GSA framework.

281 Each state variables Y_k (with $k = 1, \dots, 3$) is considered to depend on quantities X_i (with $i = 1,$
 282 \dots, I). We assess the strength of the dependence between Y_k and X_i through SA techniques which are
 283 implemented upon relying on the parameter space $\Gamma = \Gamma_{X_1} \times \dots \times \Gamma_{X_I}$, $\Gamma_{X_i} = [X_{i,\min}, X_{i,\max}]$
 284 being the support (range of variability) of X_i . This approach is based on the interpretation of each
 285 quantity X_i as a random quantity whose variability influences the output Y_k .

286 Note that we perform our GSA before model formulation, ensuing calibration and performance
 287 evaluation. As such, our GSA addresses the possibility of (i) improving our understanding of the
 288 relevance of each controlling variable on the target quantities, and (ii) identifying variables which
 289 might be unimportant in the context of a subsequent model formulation and calibration (e.g., [29-31]
 290 and references therein). We rest on the joint use of multiple global sensitivity indices, each providing
 291 a particular insight to a given aspect of sensitivity. This strategy yields a comprehensive depiction of
 292 the way model controlling parameters influence target state variables and minimizes the risk of
 293 classifying as unimportant some variables which might have a non-negligible impact on selected
 294 features of the output of interest. We illustrate some details of the sensitivity approaches and metrics
 295 we employ in Section 2.3.1, 2.3.2, and 2.3.3.

296 *2.3.1. Global Sensitivity Analysis based on Morris Indices*

297 We rely here on the evaluation of the Morris indices [32], which enable us to extend the concept
 298 of derivative-based local sensitivity to the entire space Γ . These indices are evaluated as

$$299 \mu_{X_i}^{Y_k} = E(|d_{X_i}^{Y_k}|); \quad \text{with} \quad d_{X_i}^{Y_k} = \frac{Y_k(X_1, \dots, X_i + \delta_i, \dots, X_I) - Y_k(X_1, \dots, X_i, \dots, X_I)}{\delta_i} \quad (4)$$

300 Here, $\mu_{X_i}^{Y_k}$ is the Morris Index for Y_k associated with variable X_i ; $E(\bullet)$ denotes expected value; $d_{X_i}^{Y_k}$
 301 is an approximation of the partial derivative of Y_k with respect to X_i and evaluated at a given point in
 302 Γ ; and δ_i is a perturbation to X_i .

2.3.2. Moment-based Sensitivity Indices

303
 304 Dell’Oca et al. [33] propose estimating global sensitivity metrics based on the (statistical)
 305 moments of the probability density function (*pdf*) of target variables. This technique provides richer
 306 information than the commonly used variance-based GSA [34]. Here, we consider GSA as rendered
 307 through the assessment of the *AMA* metrics (termed after the Authors’ initials [33])

$$308 \quad AMAE_{X_i}^{Y_k} = \begin{cases} \frac{1}{|Y_k^0|} E\left[|Y_k^0 - E[Y_k | X_i]|\right] & \text{if } Y_k^0 \neq 0 \\ E\left[|E[Y_k | X_i]|\right] & \text{if } Y_k^0 = 0 \end{cases} \quad (6)$$

$$309 \quad AMAV_{X_i}^{Y_k} = \frac{E\left[|V[Y_k] - V[Y_k | X_i]|\right]}{V[Y_k]} \quad (7)$$

310 where $AMA E_{X_i}^{Y_k}$ and $AMA V_{X_i}^{Y_k}$ represent the sensitivity indices associated with the first and the
 311 second moment of Y_k , respectively, as linked to variability of X_i ; the quantity Y_k^0 in (6) indicates the
 312 (unconditional) mean of Y_k , as computed over Γ ; $V[\bullet]$ denotes variance; the symbol $Y_k | X_i$ indicates
 313 conditioning of Y_k to a known value of X_i within Γ ; $E(\bullet)$ denotes expected value. Essentially,
 314 indices $AMA E_{X_i}^{Y_k}$ and $AMA V_{X_i}^{Y_k}$ quantify the expected change of the mean and variance of Y_k due to
 315 our knowledge of X_i (see [33], [35]).

2.3.3. PCA-based Sensitivity Indices

317 As stated in Section 1, Lamboni et al. [23] perform a GSA of a set of (multivariate) model
 318 outputs by combining a principal component analysis (PCA) of these with the analysis of variance
 319 (ANOVA). As such, these authors introduce a set of sensitivity indices quantifying the influence of
 320 selected quantities controlling the model behavior on the principal components associated with the
 321 collection of (multivariate) model outputs of interest.

322 We start by performing PCA through the eigenvalue decomposition of the covariance matrix
 323 of the considered set of system outputs. This enables us to transform the original outputs \mathbf{Y} into a set

324 of new orthogonal variables, which are sorted according to their contribution to the overall system
 325 variance (see [34] for additional details).

326 We consider here the first order sensitivity indices

$$327 \quad SI_{i,k} = \frac{V_{i,k}}{V_k} \quad (8)$$

328 representing the influence of X_i on the k -th principal component (P_k) of the multivariate vector \mathbf{Y} .

329 Here, $V_{i,k}$ is the variance of P_k associated with variability of X_i and V_k is the total variance of P_k .

330 The sensitivity indices (8) focus on the k -th principal component of model outputs rather than
 331 on the set of multivariate outputs as a whole. To quantify the contribution of each control variable X_i
 332 to the total variance of \mathbf{Y} , Lamboni et al. [23] introduce a generalized first order sensitivity index
 333 expressed as

$$334 \quad GSI_i = \frac{\sum_{k=1}^K \lambda_k SI_{i,k}}{\sum_{k=1}^K \lambda_k} \quad (9)$$

335 where λ_k is the k -th (with $k = \{1, \dots, K\}$) eigenvalue of the covariance matrix of \mathbf{Y} . One can see that
 336 the generalized sensitivity index (9) corresponds to the weighted average of the sensitivity indices on
 337 the principal components, the weight of each term in (9) being proportional to the eigenvalue
 338 associated with the corresponding principal component. As such, this quantity expresses the unique
 339 contribution of control variable X_i on the total variance of the multivariate outputs encapsulated in \mathbf{Y} .

340 **2.4. Candidate Simulation Models**

341 The application of the GSA techniques illustrated in Section 2.3 can yield a dimensionally
 342 reduced space formed by R control variables. We then model the dependence of the target variables
 343 \mathbf{Y} on the reduced set of controlling variables (\mathbf{X}_R) through implementation of a suite of
 344 linear/nonlinear models via (3).

345 We follow Takane and Hunter [36-37] to characterize (4) through a multivariate linear
346 regression technique (e.g., [38-39]). We also test the (possibly) nonlinear nature of relationship (3)
347 by constructing an Artificial Neural Network-based simulation model. We do so upon considering a
348 Multi-Layer Perceptron (MLP) Artificial Neural Network (ANN) and training it to reproduce the
349 observations included in $\underline{\underline{\mathbf{Y}}}$ through the unsupervised learning rule of Oja [40, 41]. Hidden layer(s)
350 contain nonlinear mapping functions. Training of the ANN is performed through serving sample
351 batches (taken from the measurements) to the model and optimizing the weights, to minimize mean
352 squared error (*MSE*) between the network responses and related measurement records of target
353 variables.

354 The use of formal model selection criteria to discriminate amongst the set of candidate models
355 which are formulated to interpret the system behavior is illustrated in Section 2.5.

356 2.5. Model Selection Criteria

357 The quality of model calibration (either linear regression or nonlinear ANNs in our case) is
358 typically appraised through the evaluation of the mean squared error, *MSE*, as

$$359 \quad MSE = \frac{J}{O \times N} \quad \text{with} \quad J = \sum_{k=1}^O \sum_{i=1}^N (\hat{Y}_{ki} - Y_{ki}^*)^2 \quad (10)$$

360 where \hat{Y}_{ki} is a model-based output corresponding to observation Y_{ki}^* (i.e., i -th sample record of the
361 k -th target quantity, O being the total number of target quantities) in matrix $\underline{\underline{\mathbf{Y}}}$.

362 We rely here on formal model selection criteria to evaluate (in a relative sense) the skill of each
363 of the models we construct to interpret the available information. Among the various model selection
364 criteria proposed in the literature to discriminate amongst models [42], we rest here on the commonly
365 used information criteria *AIC_c* ([43]) and *BIC* ([44]).

366 When multiple models are considered to interpret a physical scenario of interest, one may
367 calibrate each of them by minimizing the negative log likelihood criterion, *NLL* [45]

368
$$NLL = \frac{J}{\sigma_{\mathbf{Y}}^2} + (O \times N) \ln(2\pi\sigma_{\mathbf{Y}}^2) \quad (11)$$

369 where $\sigma_{\mathbf{Y}}^2$ corresponds to measurement error variance, its Maximum Likelihood (*ML*) estimate
 370 being given by $\hat{\sigma}_{\mathbf{Y}}^2 = \frac{J_{\min}}{(O \times N)}$ (see [46] and references therein for more details). Once the
 371 parameters associated with each model are estimated, the alternative model formulations can be
 372 ranked by evaluating the following model selection (or discrimination) criteria

373
$$AIC_c = NLL + 2m + \frac{2m(m+1)}{(O \times N) - m - 1} \quad (12)$$

374
$$BIC = NLL + m \ln(O \times N) \quad (13)$$

375 where m represents the number of model parameters. Model discrimination criteria can also be
 376 employed to determine posterior model weight (for AIC_c) or posterior model probability (for BIC),

377 $p(M_\alpha | \underline{\mathbf{Y}}^*)$, as ([47, 48] and references therein)

378
$$p(M_\alpha | \underline{\mathbf{Y}}^*) = \frac{\exp(-\frac{1}{2} \Delta IC_\alpha) p(M_\alpha)}{\sum_{l=1}^L \exp(-\frac{1}{2} \Delta IC_l) p(M_l)} \quad (14)$$

379 Here, $\Delta IC_\alpha = IC_\alpha - IC_{\min}$, IC_α being a given model selection criterion (12)-(13); IC_{\min} is the
 380 minimum value of IC_α across the L candidate models; and $p(M_\alpha)$ is the prior probability of model
 381 M_α .

382 **3. Results**

383 We start our analyses by providing basic statistics of the observations acquired at the
 384 investigated plant. These are listed in Table 1 in terms of range of variability (as expressed through

385 the lower and upper values detected for each quantity), mean, standard deviation (*std*), and the
 386 coefficient of variation (*CV*).

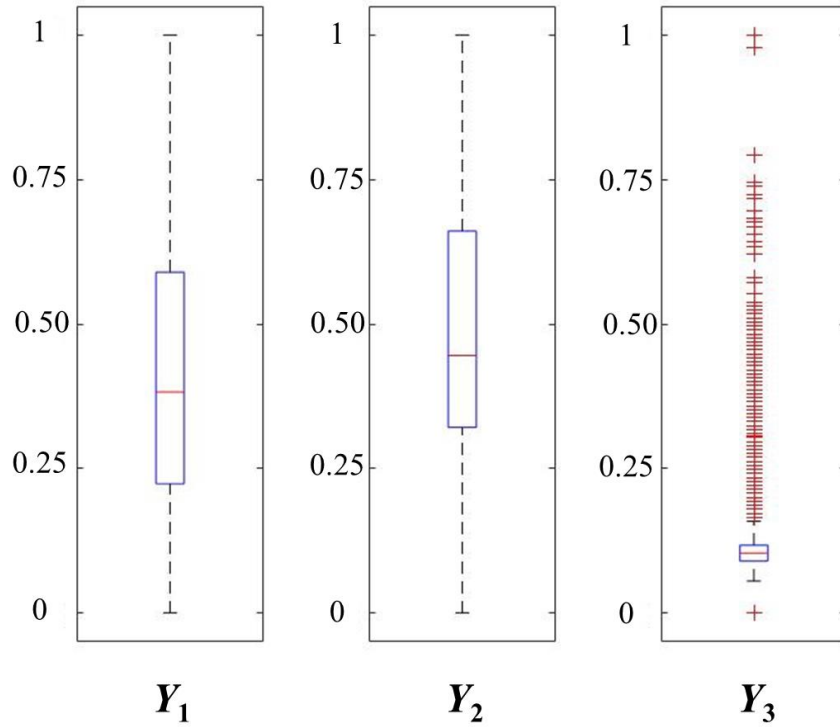
387 **Table 1** Range of variability (identified in terms of the lower, *Min*, and upper, *Max*, values detected),
 388 mean, standard deviation (*std*), and coefficient of variation (*CV*) of each sampled variable.

Data acquisition							
ID	description	units	Min	Max	Mean	Std	CV
X ₁	Oil volume flow rate @ Production unit	m ³ /day	3703	19070	10611	3505	33%
X ₂	Fresh water volume flow rate @ Electrostatic desalter	m ³ /day	0	556	249	154	62%
X ₃	Fresh water volume flow rate @ Coalescer tank	m ³ /day	0	1479	376	330	88%
X ₄	Demulsifier injection @ Production Unit	ppm	0	69	25	11	44%
X ₅	Demulsifier injection @ Coalescer tank	ppm	0	40	6	10	167%
X ₆	Demulsifier injection @ Electrostatic desalter	ppm	0	40	14	7	50%
X ₇	Salt amount @ Production unit	gr/m ³	450	21200	6759	4786	71%
X ₈	Ambient temperature	°C	0	44	28	10	36%
X ₉	Fluid temperature @ inlet Coalescer tank	°C	5	65	34	15	44%
X ₁₀	Fluid temperature @ inlet Pre-heater	°C	16	56	42	6	14%
X ₁₁	Fluid temperature @ Pre-heater bath	°C	60	793	431	230	53%
X ₁₂	Fluid temperature @ outlet Pre-heater	°C	23	65	59	4	7%
X ₁₃	Fluid temperature @ inlet Electrostatic desalter	°C	14	64	60	1.4	2%
X ₁₄	Fluid temperature @ inlet Electrostatic desalter	°C	14	53	42	1.3	3%
X ₁₅	Pressure drop @ mixing valve (Coalescer tank)	bar	0.5	12	11	0.26	2%
X ₁₆	Pressure drop @ mixing valve (Electrostatic desalter)	bar	1.2	4.5	1.37	0.22	16%
Y ₁	Oil volume flow rate @ Operating unit	m ³ /day	3231	17640	9662	3405	35%
Y ₂	Wastewater volume flow rate	m ³ /day	1.59	2432	1146	469	41%
Y ₃	Salt amount @ Operating unit	gr/m ³	4	149	24	15	63%

389 All data are then rescaled and normalized to the unit interval [0, 1], for ease of comparison.

390 Fig. 3 depicts box-plot representations of the variability of target state variables, i.e., Y₁, Y₂, and Y₃.

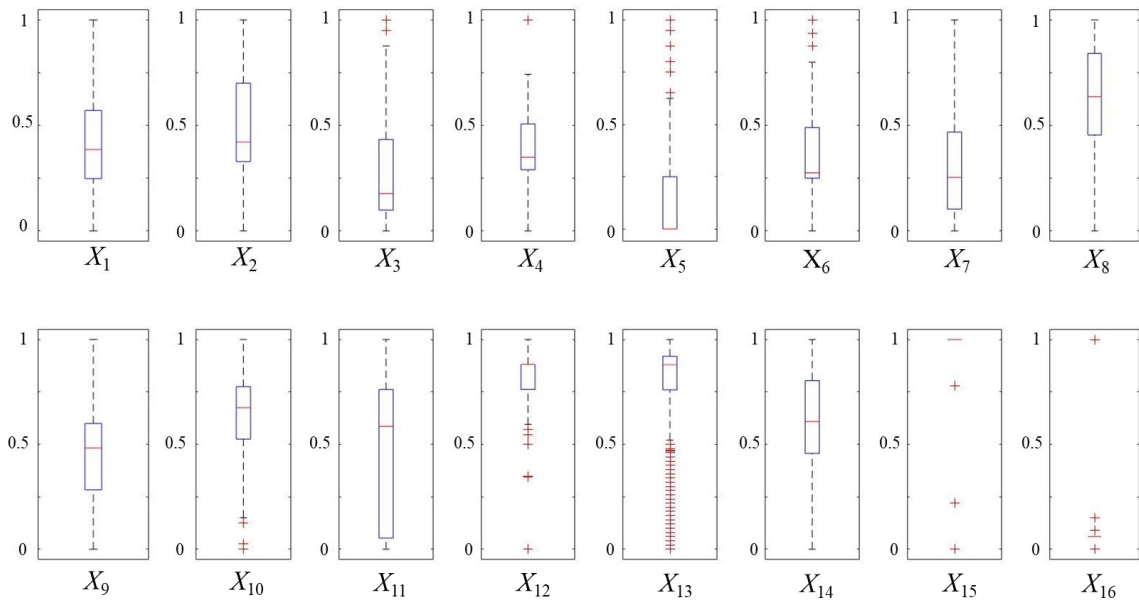
391 Box-plot representations of system controlling variables X_i are depicted in Fig. 4.



393

394 **Fig. 3** Box-plot representation of the variability of target variables (Y_1 , Y_2 , Y_3). Each box-plot
 395 identifies 25th, 50th (median), and 75th percentile scores. Whiskers correspond to the lowest
 396 observation comprised within 1.5 the interquartile range (IQR) of the lower quartile and to the
 397 highest observation still within 1.5 IQR of the upper quartile; crosses represent outliers.

398 Most of the observations of the system performance variables Y_1 and Y_2 are roughly
 399 concentrated between the 25th and the 75th percentile of the corresponding distributions. Observed
 400 values of the produced amount of salt (Y_3) exhibits a notable number of (high valued) outliers, the
 401 highest (normalized) value of the observation still within 1.5 the interquartile range (IQR) of the
 402 upper quartile being 0.16.



403

404 **Fig. 4** Box-plot representation of the variability of controlling variables ($X_1 - X_{16}$). Each box-plot
 405 identifies 25th, 50th (median), and 75th percentile scores. Whiskers correspond to the lowest
 406 observation comprised within 1.5 the interquartile range (IQR) of the lower quartile and to the
 407 highest observation still within 1.5 IQR of the upper quartile; crosses represent outliers.

408 Most of the controlling variables are concentrated between the 25th and the 75th percentile of
 409 the corresponding distributions, with few notable exceptions: (i) X_{15} and X_{16} (i.e., pressure drop at the
 410 mixing valves of the Coalescer Tank and Electrostatic Desalter) are typically set to some predefined
 411 values (by the operator onsite) and are usually clustered around these across the entire temporal
 412 observation window, with only few (high valued) outliers; (ii) the temperature monitored after the
 413 pre-heater (X_{12}) and before the electrostatic desalter (X_{13}) exhibits several (low valued) outliers; (iii)
 414 X_5 (i.e., the demulsifier injected at the coalescer tank) shows some high valued outliers; (iv) X_6 (i.e.,
 415 demulsifier injected at electrostatic desalter) shows only a few high valued outliers.

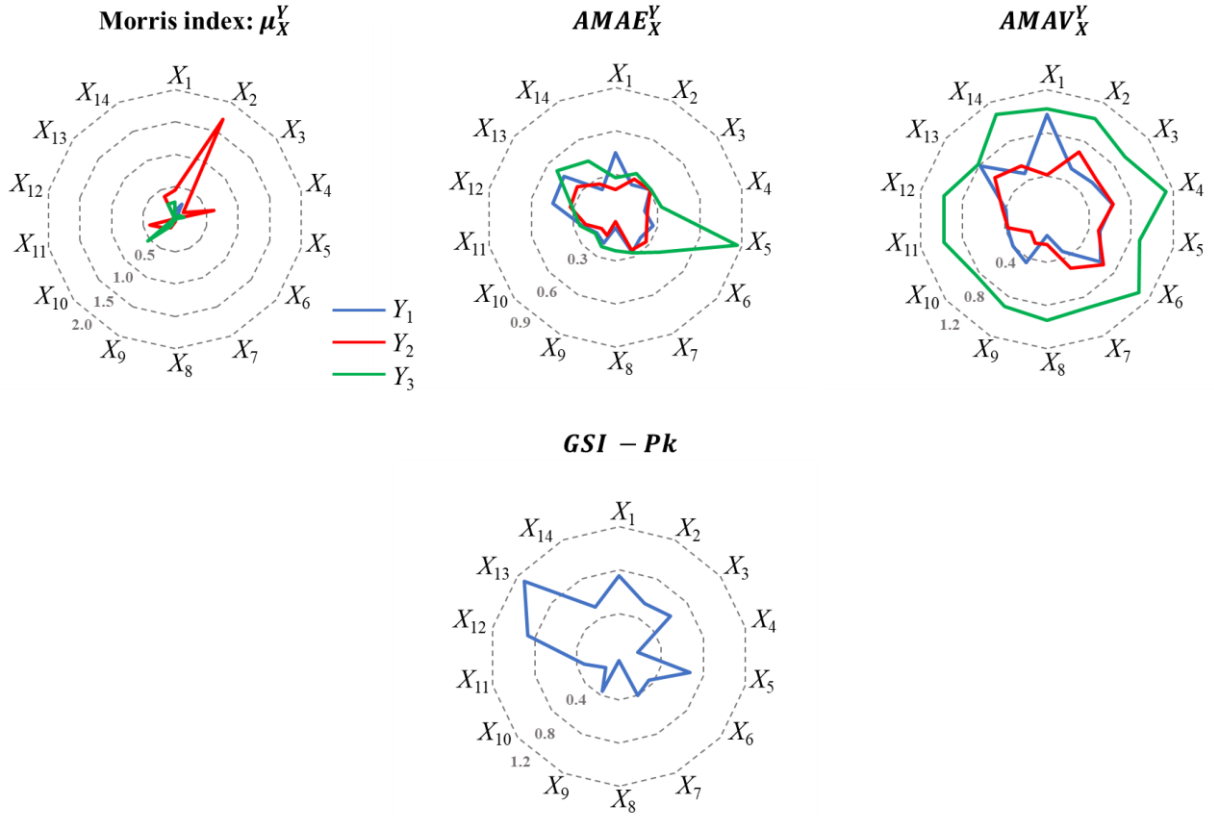
416 The analysis of the bivariate correlation among the system performance variables suggests
 417 that (i) some degree of correlation can be noted between Y_2 (i.e., waste water volume) and both Y_1
 418 (i.e., oil volume) and Y_3 (i.e., amount of salt extracted from the crude oil), as denoted by values of the

419 correlation coefficient equal to 0.30 and 0.36, respectively; and (ii) Y_1 and Y_3 appear to lack
420 correlation (with a low value of the correlation coefficient, equal to -0.06). These documented
421 correlations suggest the opportunity of considering also a PCA-based sensitivity analysis of the kind
422 illustrated in Section 2.3.3 to analyze the sensitivity of the set of multivariate target variables Y_1 , Y_2 ,
423 and Y_3 to the state variables $X_1 - X_{16}$.

424 **3.1. Global Sensitivity Analysis**

425 Following the approach outlined in Section 2.3, we quantify sensitivity of the target state
426 variables Y_k to each input variable X_i through evaluation of Morris indices $\mu_{X_i}^{Y_k}$ (5), the moment-based
427 $AMAE_{X_i}^{Y_k}$ and $AMAV_{X_i}^{Y_k}$ indices (6)-(7), and the PCA-based sensitivity index (9). A depiction of the
428 ensuing results is offered in Fig. 5.

429 We evaluate $\mu_{X_i}^{Y_k}$ upon relying on the available dataset and identifying trajectories in the
430 parameter space Γ along which all values of the controlling variables are (approximately) constant
431 with the exception of X_i . Fig. 5a illustrates $\mu_{X_i}^{Y_k}$ evaluated for the available measurement records. We
432 note that in some cases there is only a limited number of trajectories that can be extracted from the
433 measurement records to satisfy the conditions required for the evaluation of the local derivatives at
434 given observation points in Γ . As such, Morris indices cannot unambiguously evaluate sensitivity of
435 Y_k to all X_i variables and we complement this set of results through the $AMAE_{X_i}^{Y_k}$ (Fig. 5b) and
436 $AMAV_{X_i}^{Y_k}$ (Fig 5c) indices.



437

438 **Fig. 5** Radar diagram representations of: (a) Morris ($\mu_{X_i}^{Y_k}$), (b) $AMAE_{X_i}^{Y_k}$, and (c) $AMAV_{X_i}^{Y_k}$ indices
 439 of Y_k with respect to X_i ; (d) PCA-based sensitivity indices, GSI_i , quantifying the effects of X_i on the
 440 target multivariate set of quantities.

441 When assessed jointly on the basis of these three indices, none of the target state variables Y_k
 442 appears to be influenced by X_{15} and X_{16} . Thus, we omit hereafter any additional comments on possible
 443 relationships between Y_k and these controlling variables.

444 Comparing Fig. 5b - 5c, it is clear that all variables Y_k show non negligible sensitivities (in
 445 terms of $AMAE_{X_i}^{Y_k}$ and $AMAV_{X_i}^{Y_k}$) to $X_1 - X_{14}$. In other words, the mean and variance of each Y_k are in
 446 general influenced, albeit with differing extents, by all state variables $X_1 - X_{14}$. We also note that (i)
 447 Y_1 and Y_2 do not show significant sensitivity to X_8 (i.e., ambient temperature) which is however
 448 influential to Y_3 ; (ii) the variability of Y_1 is mainly affected by X_1 (i.e., volume of crude oil served to
 449 the system) and X_{13} (i.e., temperature of the mixture fluid monitored before the Electrostatic Desalter),

450 (ii) the average change of the first two statistical moments of Y_2 is mainly affected by knowledge on
451 X_2 - X_6 (i.e., injected volume of fresh water and demulsifier amount) and X_{13} , while (iii) the mean and
452 variance of Y_3 is almost uniformly influenced by all controlling variables.

453 We complete our sensitivity analysis by exploring the results stemming from the PCA-based
454 GSA described in Section 2.2.3. While a general criterion to establish if a given set of principal
455 components provides a satisfactory representation of the full variability of a dataset is still lacking, a
456 percentage larger than 75% of the total variance described by (a given number of) principal
457 components is typically deemed sufficient to grasp the main features encapsulated in the original
458 data. Results of our PCA (details not shown) suggest that almost 60.5% of the variance of $\underline{\mathbf{Y}}^*$ is
459 explained by considering the first principal component, P_1 , the second (P_2) and third (P_3) principal
460 components explaining almost 31.5% and 8.25% of the total variance of $\underline{\mathbf{Y}}^*$, respectively. These
461 results indicate that relying on the first two principal components (P_1 and P_2) enables us to retain a
462 significant amount of information about the set of target variables. Considering Fig. 5d, variable X_8
463 is characterized by negligible value of the associated PCA-based sensitivity index, suggesting that
464 the variability of this quantity is not significantly influential to the target multivariate set of state
465 variables.

466 Juxtaposing all of these results suggests that we can reduce the dimensionality of the space of
467 the controlling variables, \mathbf{X} . Accordingly, we consider model development and testing in Section 3.2
468 by relying on the reduced parameter space formed by variables $X_1 - X_{14}$ and include these in vector
469 \mathbf{X}_R .

470 **3.2. Candidate Interpretive Models and Model Selection Criteria**

471 We construct a set of 12 candidate models (listed in Table 2) to be potentially employed for (a)
472 the interpretation of the data associated with the state variables considered and representing
473 effectiveness of the crude oil desalting/demulsification process and (b) assessment of future

474 performance scenarios of the plant. Four of these models ($M_1 - M_4$) are constructed through linear
475 regression, the formulation of the remaining eight nonlinear models ($M_5 - M_{12}$) being grounded on
476 training of ANNs (see Section 2.3). Models (i) M_1 , M_5 , and M_9 are constructed upon relying on the
477 reduced set of variables included in \mathbf{X}_R , following the sensitivity analysis outlined in Section 3.1
478 (i.e., for a total of 14 controlling variables); (ii) M_2 , M_6 , and M_{10} are based on \mathbf{X}_R and the 3 dummy
479 variables $U_1 - U_3$ described in Section 2.2 (resulting in a total of 17 controlling variables) with the
480 aim of including information about the desalting unit from which observation records are collected
481 from; (iii) M_3 , M_7 , and M_{11} are based on \mathbf{X}_R and the 5 dummy variables $Q_1 - Q_5$ described in Section
482 2.2 (resulting in a total of 19 controlling variables), with the specific aim of including information
483 about the kind of demulsifier injected to the system while being agnostic about the type of desalting
484 unit associated with collected information records; and finally (iv) M_4 , M_8 , and M_{12} are grounded on
485 \mathbf{X}_R and the 8 dummy variables introduced in Section 2.2 (i.e., $U_1 - U_3$ and $Q_1 - Q_5$) (for a total of 22
486 controlling variables). We structure the ANNs either according to a single (i.e., models M_5 , M_6 , M_7 ,
487 and M_8) or double (i.e., models M_9 , M_{10} , M_{11} , and M_{12}) hidden layer(s). Quantities H_1 and H_2 (in Table
488 2 indicate the number of neurons in the first and second hidden layer of the ANNs, respectively. The
489 number (m) of parameters to be estimated for each candidate model is also listed in Table 2. Model
490 assessment is performed upon relying on normalized (to the unit interval $[0, 1]$) values of the data.

491 Characterization of the linear models is straightforward and relies on estimating model
492 parameters through minimization of (10) (or, equivalently, (11)). The training process of the ANNs
493 includes feeding sample batches (taken randomly from the available measurement records) as input
494 vectors to the network, evaluating prediction errors (in terms of mean square error, MSE (10), or,
495 equivalently, NLL (11)) with respect to measured values of the desired model output, and then tuning
496 the weights of the network to minimize such errors [49].

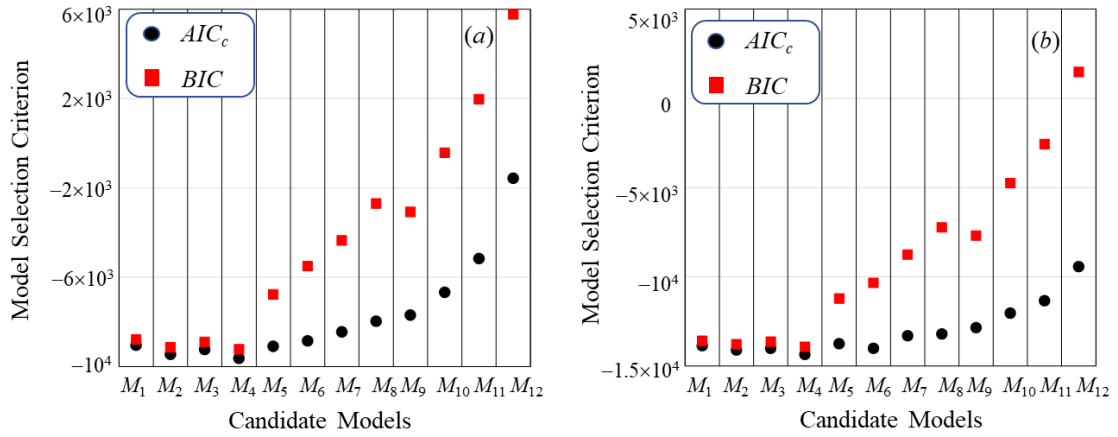
497 **Table. 2** Characteristics of the linear regression models ($M_1 - M_4$) and ANN-based nonlinear ($M_5 -$
498 M_{12}) candidate models to represent desalting process. Quantities H_1 , and H_2 indicate the number of
499 neurons in the first and second hidden layer of the ANNs, respectively; R and m denote the number
500 of controlling variables included in a given candidate model and the associated number of
501 parameters to be estimated, respectively.

Candidate Model	Linear regressions				Nonlinear ANNs							
	M_1	M_2	M_3	M_4	M_5	M_6	M_7	M_8	M_9	M_{10}	M_{11}	M_{12}
R	14	17	19	22	14	17	19	22	14	17	19	22
H_1	-	-	-	-	28	34	38	44	28	34	38	44
H_2	-	-	-	-	0	0	0	0	14	17	19	22
M	42	51	57	66	392	578	722	968	826	1207	1501	2002
$MSE (N=1200)$	0.0046	0.0041	0.0043	0.0039	0.0037	0.0034	0.0034	0.0030	0.0038	0.0033	0.0033	0.0030
$MSE (N=1820)$	0.0046	0.0043	0.0044	0.0041	0.0041	0.0036	0.0038	0.0034	0.0039	0.0037	0.0034	0.0033

502
503 Table 2 lists MSE values (10) associated with the set of 12 candidate models employed to
504 represent the desalting process and subject to calibration on the available N data. In order to assess
505 the robustness of the models and of the calibration procedure, all models have been calibrated by (i)
506 considering all N = 1820 available data and (ii) considering only a subset of N = 1200 data, randomly
507 selected within the collection of 1820 samples.

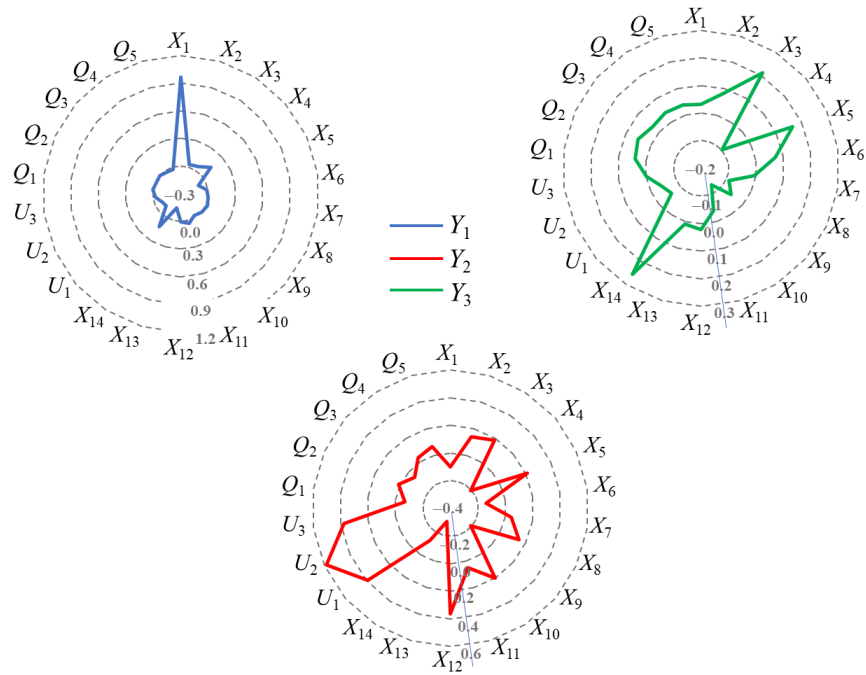
508 As expected, the value of MSE is generally lower when information about the type of desalting
509 unit or demulsifier is included in the model or when a non-linear model structure is considered. Our
510 results reveal that models M_8 and M_{12} are characterized by an identical value of $MSE(N = 1200)$ and
511 very close values of $MSE(N = 1820)$. This suggests a weak improvement of the results (in terms of
512 MSE) by inclusion of a second hidden layer in the construction of the ANN models. In general,
513 knowledge about the desalting unit from which observations are collected is more relevant (for
514 reducing MSE) than information about the injected demulsifier.

515 The ability of the calibrated models to interpret the process considered is assessed upon relying
516 on the model selection criteria (12) - (13). Fig. 6 illustrates AIC_c , and BIC evaluated for all calibrated
517 candidate models ($M_1 - M_{12}$) and upon relying on the selected subset or on the complete set of
518 measurement records. Results in Fig. 6 clearly show that the performance of the model (in terms of
519 both BIC and AIC_c) decreases when model complexity is increased, i.e., transitioning from a linear to
520 a non-linear model. The value of AIC_c is less dependent on the degree of model complexity (as
521 evaluated in terms of linear vs. non-linear structure and/or number of model parameters) than BIC .
522 Such a discrepancy between BIC and AIC_c is linked to the different weighting of the number m of
523 model parameters and of the number of data ($O \times N$) appearing in formulations (12) and (13). Our
524 results show that model M_4 (linear model considering reduced input parameter space while including
525 whole set of dummy variables \mathbf{U} and \mathbf{Q}) provides the best performance in terms of BIC and AIC_c ,
526 with a value of posterior model probability ($p(M_4 | \underline{\mathbf{Y}}^*)$, evaluated through (14)) larger than 98%.
527 On the basis of these results, we select M_4 for the additional analyses described in the following, as
528 it provides (i) adequate capability of reproducing measurement records, and (ii) simplicity of model
529 structure.



530
531 **Fig. 6** Values of model selection criteria (AIC_c and BIC) evaluated for the candidate models ($M_1 -$
532 M_{12} introduced in Table 3) calibrated on $N = (a) 1200$ and (b) 1820 measurement records.

533 The estimated values of the (linear regression) coefficients associated with the terms
534 appearing in model M_4 when using the subset of 1200 data are depicted in Fig. 7. Similar results have
535 been obtained when model construction relies on the complete data set. Our results show that variable
536 X_1 is associated with the highest coefficient value (i.e., = 0.96), this result being consistent with the
537 observation that the amount of processed oil (Y_1) should closely correspond to the volume of crude
538 oil served to the desalting unit. Evaluation of the amount of produced wastewater (Y_2) is mainly driven
539 by the weights evaluated for X_{12} and X_{13} (i.e., fluid temperature at the Pre-heater and Electrostatic
540 Desalter), as well as information on the particular desalting unit employed in the plant (as rendered
541 through U_1 , U_2 , and U_3). The negative sign of the coefficient associated with X_1 indicates that an
542 increase of the amount of crude oil served to the system correspond to a decrease (for fixed values of
543 $X_2 - X_{14}$) of extracted wastewater (Y_2) from the crude oil. The separated amount of salt (Y_3) at the
544 output of the unit is rendered by this model via a high weight of (a) the volume of the fresh water
545 (X_3) injected at Coalescer tank to the system, (b) the fluid temperature at the Electrostatic Desalter
546 (X_{14}), and (c) the amount of demulsifier (X_5) injected to the system at the Coalescer Tank. Considering
547 that the model is assessed as a best estimate of the system behavior, this result is consistent with the
548 sensitivity displayed by these parameters to the mean of Y_3 , as expressed through $AMAE_{X_i}^{Y_3}$ (see Fig.
549 5b). The positive sign and the value of the model coefficients associated with X_3 , X_5 and X_{14} indicate
550 a beneficial influence of combining (i) injection of fresh water, (ii) injecting demulsifier and (iii)
551 desalting at the Electrostatic Desalter for extraction of salt from crude oil.



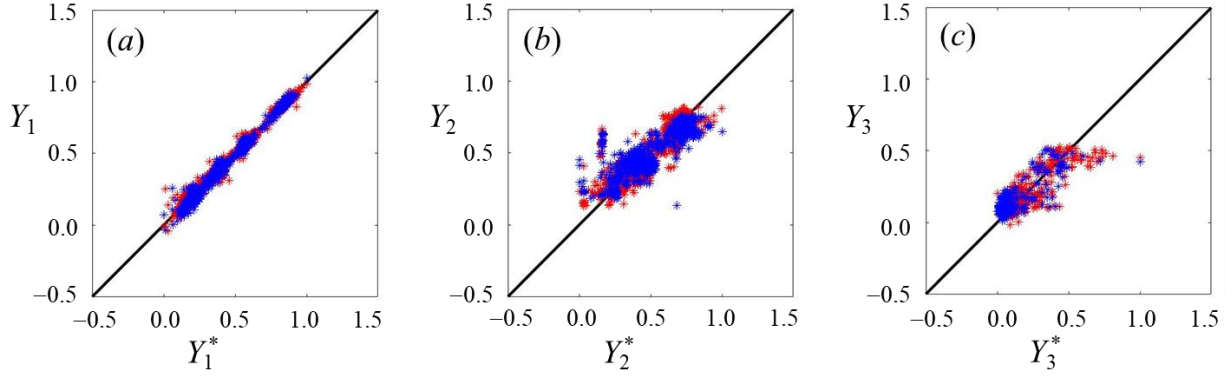
552

553 **Fig. 7** Estimated values of the linear regression coefficients associated with model M_4 and controlling
 554 variables (i.e., $X_1 - X_{14}$, $U_1 - U_3$, and $Q_1 - Q_5$) after calibration and referring to target variables: Y_1
 555 (produced oil), Y_2 (produced wastewater), and Y_3 (amount of extracted salt from the analyzed system).

556 The inclusion of available information on the desalting units from which measurements are
 557 collected (as embedded in $U_1 - U_3$) and the type of demulsifier injected to the system (i.e., $Q_1 - Q_5$) is
 558 not effective to represent the behavior of Y_1 and Y_3 . Otherwise, the impact of information content
 559 associated with $U_1 - U_3$ is notable to model Y_2 . One can also note that (i) the amount of injected fresh
 560 water (X_2) and demulsifier (X_6) at the Electrostatic Desalter; and (ii) ambient temperature (X_8), fluid
 561 temperature before the Coalescer Tank (X_9), and fluid temperature before the Pre-heater bath (X_{10})
 562 exert some degree of secondary influence to estimate Y_1 and Y_3 . Information on salt content of the
 563 fluid at the entrance of desalting unit (X_7) is only weakly influential to Y_1 , Y_2 , and Y_3 . In this context,
 564 all of the results embedded in Fig. 7 are consistent with the GSA findings included in Fig. 5.

565 Fig. 8 depicts scatterplots of observed and modeled values of the target variables, as obtained
 566 with M_4 after parameter estimation relying on the subset of 1200 observations (red symbols). The

567 remaining 620 available observations (not employed during the calibration) are also depicted in the
 568 scatterplot (blue symbols) to provide an appraisal of the robustness of the model. These results show
 569 the remarkable capability of the model to reproduce observed values of Y_1 . Otherwise, model
 570 calibration results show acceptable correspondence with data when considering Y_2 and Y_3 .



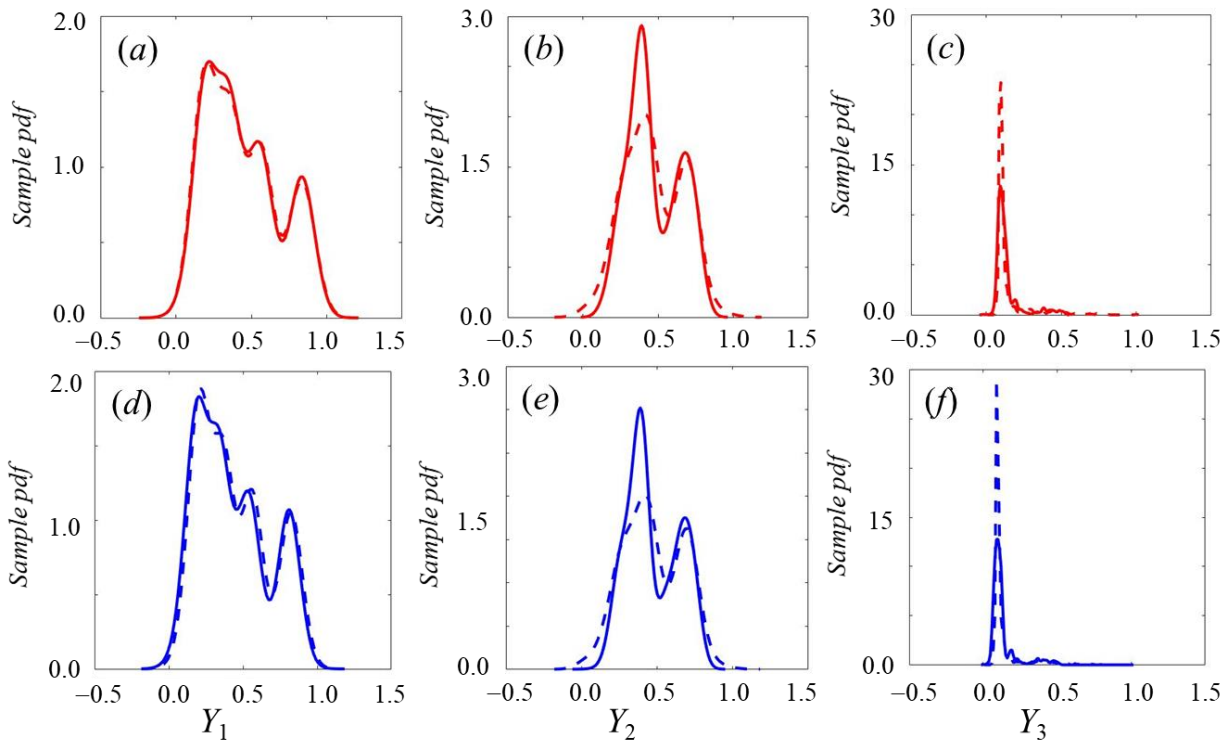
571

572 **Fig. 8** Scatterplot of observed and predicted values of Y_k obtained through the linear regression
 573 model M_4 . Results are depicted for model calibration based on the subset comprising 1200
 574 measurements (red symbols); results associated with the remaining 620 available observations (not
 575 employed during calibration) are also reported (blue symbols).

576 Finally, Fig. 9 juxtaposes empirical distributions (i.e., *sample pdf*) of the observed values of
 577 Y_k and those rendered by the corresponding values evaluated through model M_4 . The Kullback-
 578 Leibler divergence (*KLD*) [50] is used here as a metric to quantify the difference between the two
 579 distributions. Results obtained by considering model calibration against the subset of 1200
 580 observation records (blu lines) and the complete set are depicted (red lines). Results for Y_1 can be
 581 considered of good quality ($KLD = 0.02$ and 0.06 for 1200 or 1820 measurement records,
 582 respectively). Otherwise, moderate values of discrepancy can be observed for Y_2 ($KLD = 0.11$ and
 583 0.35 , for 1200 or 1820 measurement records, respectively), Y_3 being characterized by the largest
 584 differences between the observed and model-based distributions ($KLD = 5.1$ and 1.8 , for 1200 or
 585 1820 measurement records, respectively). Note that the high values of *KLD* associated with Y_3 are

586 linked to the higher peak displayed by the *sample pdf* of Y_3^* , as compared to the corresponding *pdf* of
 587 \hat{Y}_3 (see Fig. 9c and 9f). The small values of *KLD* obtained for Y_1 are consistent with the results
 588 depicted in Fig. 8.

589 Additionally, we note that multiple peaks of the *sample pdf* for Y_1 and Y_2 in Fig. 9 are related
 590 to the fact that data are collected from differing desalting units. Otherwise, measurement records of
 591 the extracted salt (Y_3) collected from different desalting units are associated with a unique pronounced
 592 peak, a very small secondary peak being barely visible (Fig 9c and 9f).



593

594 **Fig. 9** Empirical distributions (i.e., *sample pdf*) of the observed (dashed curves) and modeled
 595 (continuous curves) values of Y_k . Results are obtained by considering model calibration against the
 596 subset of 1200 (a - c) observations and the complete set of 1820 (d - f) records.

4. Conclusions

598 We analyze the performance of a system of desalting/demulsification units to process crude oil
599 before conveying it to a refinery upon relying on an original combination of (a) various Global
600 Sensitivity Analysis (GSA) approaches, (b) machine learning techniques, and (c) model
601 discrimination criteria. Our analysis is completely data-driven, as it rests on a unique data base
602 collected across a three-year period through daily measurements from three desalting units (in the
603 same oil field) where demulsification is performed through injection of 5 different demulsifiers in
604 conjunction with three desalting techniques. The overall performance of the system is evaluated in
605 terms of three main quantities, expressing (i) the output volume of the oil from Desalting Unit (Y_1),
606 (ii) the volume of the wastewater extracted from Desalting Unit (Y_2), and (iii) the separated amount
607 of salt at the output of the unit (Y_3). These are seen to depend on a total of 16 controlling variables
608 (see Table 1), representing ambient and operational conditions in the plant. Our study leads to the
609 following major conclusions.

- 610 1. Joint application of GSA techniques based on the Morris [32], Moment-based [33], as well as
611 PCA-based [40] sensitivity indices enables us to clearly identify a reduced set of 14 variables
612 that exert the highest influence on the overall performance of the industrial plant. All of these
613 indices are evaluated numerically, upon relying on the available extensive database. We find
614 that the variability of Y_1 (i.e., output volume of the oil from Desalting Unit) is mainly affected
615 by the volume of crude oil served to the system, while Y_2 (i.e., volume of the wastewater
616 extracted from Desalting Unit) is mainly influenced by the volume of injected water to the
617 system and the amount of injected demulsifier. These quantities do not show significant
618 sensitivity to the variation of the ambient temperature, which is in turn influential to the
619 extracted amount of salt from crude oil at the desalting unit (Y_3).

- 620 2. A suite of 12 linear and nonlinear regression models are derived and considered as candidates
621 for (a) the interpretation of the data associated with the performance of the crude oil desalting
622 process and (b) the evaluation of future performance scenarios of the plant. The nonlinear
623 models are grounded on training of Artificial Neural Networks (ANNs) on the available
624 observations of the reduced set of controlling variables identified through the GSA. The
625 relative skill of each of these models is evaluated via formal model identification criteria. We
626 find that a simple linear regression model shows remarkable capabilities (with posterior model
627 probability higher than 98%) to evaluate the performance of the crude oil demulsification
628 process. Results from the combination of modeling techniques employed suggest that: (a) the
629 volume of processed oil is mainly driven by the volume of crude oil served to the desalting
630 unit; (b) the volume of extracted wastewater is mainly controlled by fluid temperature at Pre-
631 eater and Electrostatic Desalter; (c) the amount of extracted salt is mainly controlled by the
632 volume of the fresh water injected to the system, amount of demulsifier injected to the system
633 at the Coalescer Tank, and the fluid temperature at the Electrostatic Desalter; and (d) the
634 benefit of including in the modeling effort available information of the desalting units from
635 which measurements are collected is notable to model the amount of extracted wastewater
636 while being less marked to represent the volume of processed oil and extracted salt.
- 637 3. Our results confirm the effectiveness of a combination of demulsification and desalting
638 techniques for the extraction of salt from crude oil. The methodological and operational
639 approach we present is general and ready to be engineered and transferred to industrial plants
640 where identification of the key variables governing the quality of the system behavior enables
641 optimal design and control through an appropriate monitoring campaign.

642

References

- 643 [1] Hu G, Li J, Zeng G. Recent development in the treatment of oily sludge from petroleum
644 industry: a review. J Hazard Mater 2013;261:470–90. doi:10.1016/J.JHAZMAT.2013.07.069.

- 645 [2] Olajire AA. The petroleum industry and environmental challenges. *J Pet Env Biotechnol*
646 2014;5:2157–7463. doi:10.4172/2157-7463.1000186.
- 647 [3] Adeniyi OD, Adediran AA, Adeniyi MI, Yahya MD, Afolabi EA, Adebayo IT, et al.
648 Evaluation of the impact of Kaduna refinery effluent on river Romi. *Niger J Technol Res*
649 2017;12(1):17–21. doi:10.4314/njtr.v12i1.4.
- 650 [4] JRC. Best available techniques (BAT) reference document for the refining of mineral oil and
651 gas industrial emissions : industrial emissions directive 2010/75/EU (integrated pollution
652 prevention and control). Publications Office of European Union; 2015.
- 653 [5] Manning FS, Thompson RE. Oilfield processing of petroleum: crude oil. Pennwell books;
654 1995.
- 655 [6] Farrokhi F, Jafari Nasr MR, Rahimpour MR, Arjmand M, Vaziri SA. An investigation on
656 simultaneous effects of several parameters on the demulsification efficiency of various crude
657 oils. *Asia-Pacific J Chem Eng* 2017;12(6):1012-1022. doi:10.1002/apj.2142.
- 658 [7] Hamadi AS, Mahmood LH. Demulsifiers for Simulated Basrah Crude Oil. *Eng Technol J*
659 2009;28(1):54–64.
- 660 [8] Becher P. *Emulsions: Theory and Practice*. Oxford University Press, New York, NY (2001).
- 661 [9] Fortuny M, Oliveira CBZ, Melo RLFV, Nele M, Coutinho RCC, Santos AF. Effect of Salinity,
662 Temperature, Water Content, and pH on the Microwave Demulsification of Crude Oil
663 Emulsions. *Energy and Fuels* 2007;21(3) 1358–1364. doi:10.1021/EF0603885.
- 664 [10] Less S, Hannisdal A, Bjørklund E, Sjöblom J. Electrostatic destabilization of water-in-crude
665 oil emulsions: Application to a real case and evaluation of the Aibel VIEC technology. *Fuel*
666 2008;87:2572–81. doi:10.1016/j.fuel.2008.03.004.
- 667 [11] Meidanshahi V, Jahanmiri A, Rahimpour MR. Modeling and optimization of two stage AC
668 electrostatic desalter. *Sep Sci Technol* 2012;47:30–42. doi:10.1080/01496395.2011.614316.
- 669 [12] Dai X, Wang X, He R, Du W, Zhong W, Zhao L, Qian F. Data-driven robust optimization for

- 670 crude oil blending under uncertainty. *Comp. and Chem. Eng.* 2019.
671 doi:10.1016/j.compchemeng.2019.106595
- 672 [13] Qin SJ, Chiang LH. Advances and opportunities in machine learning for process data analytics.
673 *Comp. and Chem. Eng.* 2019;126:465-473. doi:10.1016/j.compchemeng.2019.04.003
- 674 [14] Siena M, Guadagnini A, Della Rossa E, Lamberti A, Masserano F, Rotondi M. A novel
675 enhanced-oil-recovery screening approach based on Bayesian clustering and Principal-
676 Component Analysis. *SPE Reserv Eval Eng* 2016;19:382–90. doi:10.2118/174315-PA.
- 677 [15] Sad CMS, Da Silva M, Dos Santos FD, Pereira LB, Corona RRB, Silva SRC, et al. Multivariate
678 data analysis applied in the evaluation of crude oil blends. *Fuel* 2019;239:421–428.
679 doi:10.1016/j.fuel.2018.11.045.
- 680 [16] Gueddar T, and Dua V. Novel model reduction techniques for refinery-wide energy
681 optimisation. *Applied Energy* 2012;89(1) 117-126. doi:10.1016/j.apenergy.2011.05.056
- 682 [17] Mohammad AT, Mat SB, Sulaiman MY, Sopian K, Al-abidi AA. Implementation and validation
683 of an artificial neural network for predicting the performance of a liquid desiccant dehumidifier
684 *Energy Convers. Manage.*, 2013; 67: 240-250. doi/10.1016/j.enconman.2012.12.005
- 685 [18] Mouret S, Grossmann IE, Pestiaux P. A new Lagrangian decomposition approach applied to the
686 integration of refinery planning and crude-oil scheduling. *Comp. and Chem. Eng.*
687 2011;35:2750-2766. doi:10.1016/j.compchemeng.2011.03.026
- 688 [19] Gao X, Jiang Y, Chen T, Huang D. Optimizing scheduling of refinery operations based on
689 piecewise linear models. *Comp. and Chem. Eng.* 2015;75:105-119.
690 doi:10.1016/j.compchemeng.2015.01.022
- 691 [20] Ochoa-Estopier LM, Jobson M, Smith R. Operational optimization of crude oil distillation
692 systems using artificial neural networks. *Comp. and Chem. Eng.* 2013;59:178-185.
693 doi:10.1016/j.compchemeng.2013.05.030
- 694 [21] Gueddar T, Dua V. Disaggregation-aggregation based model reduction for refinery-wide

- 695 optimization. *Comp. and Chem. Eng.* 2011;35(9):1838-1856.
696 doi:10.1016/j.compchemeng.2011.04.016
- 697 [22] Al-Qahtani K, Elkamel A. Robust planning of multisite refinery networks: Optimization under
698 uncertainty. *Comp. and Chem. Eng.* 2010;34(6):985-995.
699 doi:10.1016/j.compchemeng.2010.02.032
- 700 [23] Lamboni, M. Monod, H. and, Makowski, D. [2011] Multivariate sensitivity analysis to measure
701 global contribution of input factors in dynamic models. *Reliability Engineering and System
702 Safety* 2011;96(4):450–459. doi:10.1016/j.res.2010.12.002
- 703 [24] Vapnik VN. *The nature of statistical learning theory*. Springer press, New York, 2000.
704 doi:10.1007/978-1-4757-3264-1.
- 705 [25] Choubineh A, Ghorbani H, Wood DA, Moosavi SR, Khalafi E, Sadatshojaei E, Improved
706 predictions of wellhead choke liquid critical-flow rates: Modelling based on hybrid neural
707 network training learning based optimization. *Fuel*, 2017;207:547-560.
708 doi:10.1016/j.fuel.2017.06.131.
- 709 [26] Ghorbani H, Wood DA, Choubineh A, Tatar A, Ghazaeipour Abarghoyi P, Madani M,
710 Mohamadian N. Prediction of oil flow rate through an orifice flow meter: Artificial intelligence
711 alternatives compared. *Petroleum* 2018: In press. doi:10.1016/j.petlm.2018.09.003.
- 712 [27] Bingöl D, Xiyili H, Elevli S, Kılıç E, Çetintaş S. Comparison of multiple regression analysis
713 using dummy variables and a NARX network model: an example of a heavy metal adsorption
714 process. *Water and Environmental Journal* 2018;32(2): 186/196 doi: 10.1111/wej.12314
- 715 [28] Saltelli A, Ratto M, Andres T, Campolongo F, Cariboni J, Gatelli D, Salsana M, Tarantola S.
716 *Global sensitivity analysis - The primer*. Wiley, 2008.
- 717 [29] Liu H, Chen W, Sudjianto A. Relative entropy based method for probabilistic sensitivity analysis
718 in engineering design. *ASME Journal of Mechanical Design* 2006; 128(2):326–336.
719 doi:10.1115/1.2159025.

- 720 [30] Hucheson RS, McAdams DA. A hybrid sensitivity analysis for use in early design. *Journal of*
721 *Mechanical Design* 2010;132(11). doi:10.1115/DETC2009-87120.
- 722 [31] Dell’Oca A, Riva M, Guadagnini A. Global sensitivity analysis for multiple interpretive models
723 with uncertain parameters. *Water Resour Res* 2020;55(2).
724 <https://doi.org/10.1029/2019WR025754>
- 725 [32] Morris MD. Factorial sampling plans for preliminary computational experiments.
726 *technometrics* 1991;33(2): 161–174
- 727 [33] Dell’Oca A, Riva M, Guadagnini A. Moment-based metrics for global sensitivity analysis of
728 hydrological systems. *Hydrol Earth Syst Sci* 2017;21:6219–34. doi:10.5194/hess-21-6219-
729 2017.
- 730 [34] Sobol’ IM. Sensitivity estimates for nonlinear mathematical models. *MMCE* 1993; 1(4):407–
731 414.
- 732 [35] Bianchi Janetti E, Guadagnini L, Riva M, Guadagnini A. Global sensitivity analyses of multiple
733 conceptual models with uncertain parameters driving groundwater flow in a regional-scale
734 sedimentary aquifer. *J. Hydrol.* 2019;574:544–556. doi:10.1016/j.jhydrol.2019.04.035
- 735 [36] Takane Y, Hunter MA. A new family of constrained principal component analysis (CPCA).
736 *Linear Algebra Appl* 2011;434:2539–2555. doi:10.1016/J.LAA.2011.01.002.
- 737 [37] Reinsel GC, Velu RP. *Multivariate reduced-rank regression*. 136. Springer New York; 1998.
738 doi:10.1007/978-1-4757-2853-8.
- 739 [38] Scholz M, Vigário R. Nonlinear PCA: a new hierarchical approach. Presented at 10th Euroean
740 Symposium on Artificial Neural Networks, Bruges, Belgium, April 24-26, 2002.
- 741 [39] Ranaee E, Porta G, Riva M, Guadagnini A. Investigation of saturation dependency of oil
742 relative permeability during WAG process through linear and non-linear PCA. Presented at
743 ECMOR XIV–14th European Conference on the Mathematics of Oil Recovery, Catania, Italy,
744 September 2014. doi:10.3997/2214-4609.20141800.

- 745 [40] Oja E. A simplified neuron model as a principal component analyzer. Springer-Verlag 1982.
- 746 [41] Neurosolution version 5.07 (2007) Neuro dimension. Gainesville Inc., Florida, (USA). 2007.
- 747 [42] Ranaee E, Riva M, Porta GM, Guadagnini A. Comparative assessment of three-phase oil
748 relative permeability models. *Water Resour Res* 2016;52:5341–56.
749 doi:10.1002/2016WR018872.
- 750 [43] Hurvich CM, Tsai CL. Regression and time series model selection in small samples.
751 *Biometrika* 1989;76(2):297–307. doi:10.1093/biomet/76.2.297
- 752 [44] Schwarz G. Estimating the dimension of a model. *Ann Stat* 1978;6:461–464.
753 doi:10.1214/aos/1176344136.
- 754 [45] Carrera J, Neuman SP. Estimation of aquifer parameters under transient and steady state
755 conditions: 1. Maximum likelihood method incorporating prior information. *Water Resour.*
756 *Res.* 1986;22(2):199-210. doi:10.1029/WR022i002p00199.
- 757 [46] Ranaee E, Porta GM, Riva M, Blunt MJ, Guadagnini A. Three-phase oil relative permeability
758 prediction through a sigmoid-based model, *Journal of Petroleum Science and Engineering*
759 2015;26:190–200. doi:10.1016/j.petrol.2014.11.034.
- 760 [47] Ranaee E, Moghadasi L, Inzoli F, Riva M, Guadagnini A. Identifiability of parameters of three-
761 phase oil relative permeability models under simultaneous water and gas (SWAG) injection.
762 *Petroleum Science and Engineering* 2017;159:942–951. doi: 10.1016/j.petrol.2017.09.062.
- 763 [48] Ye M, Meyer PD, Neuman SP. On model selection criteria in multimodel analysis. *Water*
764 *Resour. Res.* 2008;44(3):W03428, 10.1029/2008WR006803
- 765 [49] Moghadasi L, Ranaee E, Inzoli F, Guadagnini A. Petrophysical well Log analysis through
766 intelligent methods. Presented at SPE Bergen One Day Seminar, Bergen, Norway, April 2017.
767 doi:10.2118/185922-MS.
- 768 [50] Kullback S, Leibler RA. On information and sufficiency. *Ann. Math. Statistics* 1951;22(1): 79-
769 86. doi:10.1214/aoms/1177729694.