1   **Prediction of outlet dissolved oxygen in micro-irrigation sand media**

2   **filters using a Gaussian process regression**

3   Paulino J. García–Nieto[a,*], Esperanza García–Gonzalo[a], Jaume Puig–Bargués[b], Miquel

4   Duran–Ros[b], Francisco Ramírez de Cartagena[b], Gerard Arbat[b]

5   [a]Department of Mathematics, Faculty of Sciences, University of Oviedo, 33007 Oviedo, Spain

6   [b]Department of Chemical and Agricultural Engineering and Technology, University of Girona, 17003

7   Girona, Catalonia, Spain

8

9   **Abstract**

10   Sand media filters are a key component of micro-irrigation systems since they help

11   preventing emitter clogging, which greatly affects the system performance. Dissolved

12   oxygen is an irrigation water quality parameter related to organic matter loading. Low

13   values of dissolved oxygen can cause crop root hypoxia and, therefore, agronomic

14   problems. Thus, an accurate prediction of dissolved oxygen values could be of great

15   interest, especially if effluents are used in micro-irrigation systems. The aim of this

16   study was to obtain a predictive model able to forecast the dissolved oxygen values at

17   the outlets of sand media filters. In this study, a Gaussian process regression (GPR)

18   model was used for predicting the output dissolved oxygen ($DO_o$) from data

19   corresponding to 547 filtration cycles of different sand filters using reclaimed effluent.

20   This optimisation technique involves kernel parameter setting in the GPR training

21   procedure, which significantly influences the regression accuracy. To this end, the

22   height of the filter bed, filtration velocity and filter inlet values of the electrical

23   conductivity, dissolved oxygen, pH, turbidity and water temperature were monitored

---

*Corresponding author. Tel.: +34-985103417; fax: +34-985103354.
*E-mail address*: lato@orion.ciencias.uniovi.es (P.J. García–Nieto).

24 and analysed. The significance of each variable on filtration performance is presented

25 and a model for forecasting the outlet dissolved oxygen obtained. Regression with

26 optimal hyperparameters was performed and a coefficient of determination of 0.90 for

27 $DO_o$ was obtained when this new predictive GPR–based model was applied to the

28 experimental dataset. Agreement between experimental data and the model confirmed

29 the good performance of the latter.

30

31 *Keywords:* Gaussian process regression; Bayesian statistics; Machine learning

32 techniques; Drip irrigation; Clogging; Effluents

33

34 **Nomenclature**

35 Abbreviations

| ANN | Artificial neural network |
|---|---|
| DE | Differential evolution |
| DO | Dissolved oxygen |
| GPR | Gaussian process regression |
| GEP | Gene expression programming |
| $R^2$ | Coefficient of determination |
| RBF | Radial basis function |
| SCADA | Supervisory control and data acquisition |
| SE | Squared-exponential |
| SVM | Support vector machine |
| $v$ | Filtration velocity, m h$^{-1}$ |

Symbols

| $DO_i$ | Dissolved oxygen at filter inlet, mg l$^{-1}$ |
|---|---|

| $DO_o$ | Dissolved oxygen at filter outlet, mg l$^{-1}$ |
| --- | --- |
| $\delta_{ij}$ | Kronecker delta function |
| $\varepsilon$ | Additive white noise |
| $\ell$ | Length-scale for the RBF kernel |
| $\sigma_f^2$ | Variance for the RBF kernel |
| $\sigma_n^2$ | Gaussian noise variance |

36

## 1. Introduction

38 The substitution of conventional irrigation water by reclaimed effluents in areas of low

39 water availability is a common management strategy despite of its potential pollution

40 and health hazards (Ait-Mouheb et al., 2018). Among the different irrigation techniques

41 used, micro-irrigation shows several environmental and health advantages related

42 mainly to the reduced effluent exposure to humans and plants. However, one of the

43 most important disadvantages of applying effluents with micro-irrigation is emitter

44 clogging which can cause irrigation nonuniformity and system failure (Trooien & Hills,

45 2007). In order to avoid emitter clogging, micro-irrigation systems require effective

46 filtration (Nakayama, Boman, & Pitts, 2007) and sand media filters are the standard for

47 the protection of micro-irrigation systems using effluents (Trooien & Hills, 2017).

48

49 The level of dissolved oxygen (DO) decreases with the increased organic matter,

50 commonly present in wastewaters. So, DO, which can be determined easier and quicker

51 using sensors, is an indicator of irrigation water quality. Low DO values in the irrigation

52 water cause root oxygen deficiency, leading to low yields (Bhattarai, Midmore, &

53 Pendergast, 2008) and low quality (Zhou, Zhou, Xu, Muhammad, & Li, 2019). Usually,

DO increases through micro-irrigation systems, especially when water is released by the emitters (Maestre–Valero & Martínez-Álvarez, 2010). The DO increase is slight in sand media filters but it is considerably affected by the filter performance (Elbana, Ramírez de Cartagena, & Puig-Bargués, 2012; Solé–Torres, Puig–Bargués, Duran–Ros, Arbat, Pujol, & Ramírez de Cartagena, 2019b). Thus, the development of accurate models for forecasting DO at filter outlets can be very useful for the appropriate management of both sand filter performance and irrigation water quality. Optimal efficiency of drip irrigation systems is required for implementing smart irrigation techniques which aim to provide optimum use of the water resources (Canales-Ide, Zubelzu & Rodríguez-Sinobas, 2019).

In this regard, advanced techniques such as artificial neural networks (ANN) (Puig–Bargués, Duran-Ros, Arbat, Barragán, & Ramírez de Cartagena, 2012), gene expression programming (GEP) (Martí et al., 2013) and support vector machines (SVM) (García–Nieto, García–Gonzalo, Arbat, Duran–Ros, Ramírez de Cartagena, & Puig–Bargués, 2016) have been used for predicting the filtered volume and the value of dissolved oxygen at sand media filter outlets. Recently, other machine learning techniques such as gradient boosted regression have been applied to different aspects of the filter operation (García–Nieto et al. 2017, 2018).

Thus, the application of an innovative methodology that combines a Gaussian process regression (GPR) approach (Rasmussen, 2003; Kuhn & Johnson, 2018; Ebden, 2015) with a metaheuristic optimisation algorithm Differential Evolution (DE) (Storn & Price, 1997; Price, Storn, & Lampinen, 2005; Feoktistov, 2006; Chakraborty, 2008; Simon,

4

2013) to foretell the outlet dissolved oxygen in sand media filters used in microirrigation systems could be an interesting approach since this issue has not yet been yet addressed in previous investigations. GPR is a machine learning method developed on the basis of statistical and Bayesian theory. As a nonparametric regression method it can be considered a complex model with capability to model nonlinearities and variable interactions (Rasmussen, 2003; Ebden, 2015). When GPR is compared with other machine learning techniques, it has several advantages (Rasmussen & Williams, 2006): (1) it has an important generalisation capacity; (2) the hyperparameters in GPR can be self-adaptively calculated; and (3) the GPR outputs have clear probabilistic meaning. In this study, the DE method is applied to optimise the GPR hyperparameters. Previous researches show that GPR is an effective tool in many fields, such as irrigation mapping (Chen, Lu, Luo, Pokhrel, Deb, Huang, & Ran, 2018), wind engineering and industrial aerodynamics (Ma, Xu, & Chen, 2019), applied geophysics (Noori, Hassani, Javaherian, Amindavar, & Torabi, 2019), applied demography (Wu & Wang, 2018), psychology (Schulz, Speekenbrink, & Krause, 2018), mechanical engineering (Kong, Chen, & Li, 2018), environmental engineering (Liu, Yang, Huang, Wang, & Yoo, 2018), tracking and positioning (Ko, Klein, Fox, & Haehnelt, 2007a), deformation observation (Rogers & Girolami, 2016), system identification and control (Ko, Klein, Fox, & Haehnelt, 2007b) and so on. However, it has not been used for predicting micro-irrigation sand filter performance.

The main objective of the this study was to predict the outlet dissolved oxygen ($DO_o$) in sand media filters operating with reclaimed effluents by using Gaussian processes (GPs) in combination with the DE parameter optimisation technique.

102 The structure of this paper is organised as follows: section 2 introduces the experimental

103 setup and variables involved in this study as well as the GPR method; section 3

104 describes the results obtained with this model by comparing the GPR results with the

105 experimental measurements, including the importance of the input variables and

106 validating the efficacy of the proposed approach; and finally, section 4 concludes this

107 study with a list of main findings.

108

109 **2. Materials and methods**

110 *2.1. Experimental setup*

111 The experimental setup was composed of 3 media filters fed with the reclaimed effluent

112 from the wastewater treatment plant of Celrà (Girona, Spain). Each filter had a different

113 underdrain design: inserted domes (model FA-F2-188, Regaber, Parets del Vallès,

114 Spain), arm collector (model FA1M, Lama, Gelves, Spain) and porous media (prototype

115 designed by Bové et al. (2017) (see Fig. 1).

116

117 Silica sand CA-07MS (Sibelco Minerales SA, Bilbao, Spain) with an effective diameter

118 (*De*, size opening which will pass 10% of the sand) of 0.48 mm and a coefficient of

119 uniformity (ratio of the sizes opening which will pass 60% and 10% of the sand

120 through, respectively) of 1.73 was used as filtration media in the three filters. Media

121 heights of 200 and 300 mm, were tested for each filter.

122

123 Each of the three filters operated alone for 8 h per day each. Nominal filtration

124 velocities 30 and 60 m h$^{-1}$ were tested in each filter. Each combination of media height

125 and filtration velocity was tested during 250 h. The filters were automatically

126   backwashed when the pressure loss across them reached 50 kPa for more than 1 min.

127   The backwashing was carried out for 3 min with previously filtered effluent that was

128   chlorinated for achieving 4 ppm target chlorine concentration.

129

130   Filtered and backwashed effluent volumes, pressures across the filter and some effluent

131   quality parameters before (pH, temperature, electrical conductivity, DO and turbidity)

132   and after (only DO and turbidity) being filtered were measured and recorded every

133   minute in a supervisory control and data acquisition system (SCADA) fully described

134   by Solé-Torres et al. (2019a). Once the experiment started, the performance of the

135   effluent quality sensors was assessed periodically by comparing its measurements with

136   results obtained by manual sampling and, if necessary, they were calibrated following

137   manufacturer recommendations.

138

139   **Fig. 1 -** Picture of the experimental set-up with the three filter designs: (a) red: arm

140   collector; (b) blue: inserted domes; and (c) green: porous media prototype.

141

142   *2.2. Variables involved in the model and materials tested*

143   The main objective of this study was to compute the outlet dissolved oxygen as a

144   function of different experimentally measured parameters that the GPR–based model

145   needs as input. The output variable was the outlet dissolved oxygen ($DO_o$), which is an

146   indicator of the quality of the filtered effluent and it is directly related to the organic

147   load and the hypoxic risk of irrigation water.

148

149    The new predictive model used eight different operating variables commonly used for

150    characterising sand media filter performance as input variables (see Table 1) (Puig-

151    Bargués et al., 2012). After removing samples with missing data from the initial 637

152    samples, 547 satisfactory samples were obtained.

153

154    **Table 1 -** Set of operation physical input variables used in this study and their names

155    along with their mean and standard deviation.

156

157    The operating input variables are as follows:

158    • Filter: three filter designs (porous, dome and arm collector underdrains) as

159        described in section 2.1. This is a categorical variable.

160    • Height of the filter bed (mm): an operating variable for sand filters. Two

161        different filter bed heights of 200 and 300 mm were tested for each filter.

162    • Filtration velocity (m h$^{-1}$): a operating variable related to filter operation. Two

163        filtration velocities (30 and 60 m h$^{-1}$) were tested for each filter since these

164        follow within the common range of velocities suggested by the manufacturers.

165    • Electrical conductivity ($\mu$S cm$^{-1}$): a general measure of water quality related to

166        salinity, which is a constraint in microirrigation (Tal, 2016).

167    • Dissolved oxygen (mg l$^{-1}$): a variable related to the ability of water to support

168        aerobic processes. This is a common parameter used for both controlling the

169        biological treatment in wastewater plants and measuring irrigation water quality.

170    • pH: a measure of water acidity or alkalinity.

171  • Water temperature (ºC): temperature of the effluent at the filter inlet.

172  • Input turbidity (FNU): a key parameter for water quality that measures water
173    clarity, which depends on suspended solid load.

174  • Filtered volume ($m^3$): a measure of the volume of effluent filtered in each
175    filtration cycle.

176

177  *2.3. Gaussian process regression (GPR)*

178  GPRs are Bayesian state-of-the-art tools for discriminative machine learning (i.e.,
179  regression, classification, and dimensionality reduction). GPs assume that a GP prior
180  governs the possible unobserved latent functions and the marginal likelihood of the
181  latent function. Thus, a priori observations shape this to produce posteriori probabilistic
182  estimates. Consequently, the joint distribution of training and test data is a
183  multidimensional GP, and the predicted distribution is estimated by conditioning based
184  on training data (Camps−Valls, 2016; Witten, Frank, Hall, & Pal, 2016).

185

186  To fix ideas, a Gaussian distribution is a probability distribution that explains the
187  random variables including vectors and scalars. On the one hand, this kind of
188  distribution is stated exactly through its mean and covariance: $x \sim N\left(\mu, \sigma^2\right)$. On the
189  other hand, a GP can be seen as a generalisation of the Gaussian probability distribution
190  and it applies over functions. From the functional space point of view, a GP is an
191  ensemble of random variables, that is to say, any finite number having a joint Gaussian
192  distribution.

193

9

194    *2.3.1. The fundamentals of GPR*

195    Let us assume that $D = \left\{ (\mathbf{x}_i, y_i) / i = 1, 2, ..., N \right\}$ depicts the training dataset of the

196    Gaussian approach and the feature vectors $\mathbf{x}_i \in \mathfrak{R}^n$ comprises the extracted features or

197    the merged features and the pertinent segregation parameters. The observed target

198    values $y_i$ reproduce the outlet dissolved oxygen measured in a filtration process,

199    respectively. $X = \left\{ \mathbf{x}_i \right\}_{i=1}^{N}$ depicts the input matrix of training dataset,

200    $\mathbf{y} = \left\{ y_i \right\}_{i=1}^{N}$ symbolises the output vector. A GP $f(\mathbf{x})$ defines a priori over functions,

201    which can be converted into a posteriori over functions once some data is obtained. A

202    GP can be fully stated exactly by using its mean function $m(\mathbf{x})$ and covariance function

203    $k(\mathbf{x}, \mathbf{x}')$. In this way, the Gaussian process is indicated as (Rasmussen & Williams,

204    2006; Marsland, 2014; Witten, Frank, Hall, & Pal, 2016):

$$f(\mathbf{x}) \sim GP\left(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')\right) \tag{1}$$

205    so that

$$m(\mathbf{x}) = E\left[ f(\mathbf{x}) \right] \tag{2}$$
$$k(\mathbf{x}, \mathbf{x}') = E\left[ \left( f(\mathbf{x}) - m(\mathbf{x}) \right) \left( f(\mathbf{x}') - m(\mathbf{x}') \right)^T \right]$$

206    The mean function $m(\mathbf{x})$ depicts the anticipated value of the function $f(\mathbf{x})$ at the input

207    point $\mathbf{x}$. The covariance function $k(\mathbf{x}, \mathbf{x}')$ can be taken into account as a measurement

208    of the confidence level for $m(\mathbf{x})$, and it is required that $k(\cdot, \cdot)$ be a positive definite

209    kernel. In general, the mean function is set to be zero for notation simplicity, but this is

210    also reasonable if there is no a priori knowledge about the mean variable, as is the case

211    in this study.

212　The choice of the covariance function is critical for the GP. It describes the assumptions

213　about the latent regression model and, therefore, is also referred to as the prior

214　(Schneider & Ertel, 2010). In this research, the affine mean function and squared-

215　exponential (SE) covariance function are expressed as follows (Shi & Choi, 2011;

216　Witten, Frank, Hall, & Pal, 2016; Kuhn & Johnson, 2018):

$$k_{\mathrm{SE}}\left(\mathbf{x}, \mathbf{x}'\right) = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right) \tag{3}$$

217　being $l$ the characteristic length-scale and $\sigma_f^2$ the signal variance. The parameter

218　selection of the SE covariance function has a direct effect on the performance of the GP.

219　Here, $l$ controls the horizontal scale over which the function changes, and $\sigma_f^2$ controls

220　the vertical scale of the function.

221

222　The function values $f(\mathbf{x})$ are not achievable in most applications. In practice, only the

223　noisy observations are available and they are given by:

$$\mathbf{y} = f(\mathbf{x}) + \varepsilon \tag{4}$$

224　so that $\varepsilon$ is the additive white noise. Besides, suppose that Gaussian noise is

225　independent and identically distributed such that $\varepsilon \sim N\left(0, \sigma_n^2\right)$, where $\sigma_n$ is the

226　standard deviation of this noise. Any finite number of the observed values can also

227　constitute an individual Gaussian process as given by (Witten, Frank, Hall, & Pal, 2016;

228　Vidales, 2019):

$$\mathbf{y} \sim GP\left(m(\mathbf{x}), k\left(\mathbf{x}, \mathbf{x}'\right) + \sigma_n^2 \delta_{ij}\right) = GP\left(0, k\left(\mathbf{x}, \mathbf{x}'\right) + \sigma_n^2 \delta_{ij}\right) \tag{5}$$

229　where $\delta_{ij}$ is the Kronecker delta function described as:

230
$$\delta_{ij} = \begin{cases} 1 & \text{if} & i = j \\ 0 & \text{otherwise} \end{cases}$$

231 The purpose of the GPR model is to foretell the function value $\bar{f}^*$ and its variance

232 $\text{cov}(f^*)$ given the new test point $\mathbf{x}^*$. In this sense, $X^*$ depicts the input matrix of test

233 dataset and $N^*$ the size of test dataset. Taking into account the definition of GP, the

234 observed values and the function values at new test points obey a joint Gaussian

235 previous distribution which can be expressed as:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}^* \end{bmatrix} \sim N \left( 0, \begin{bmatrix} K(X,X) + \sigma_n^2 I & K(X,X^*) \\ K(X^*,X) & K(X^*,X^*) \end{bmatrix} \right) \tag{6}$$

236 where:

237 • $K(X,X)$: is the covariance matrix of training dataset;

238 • $K(X^*,X^*)$: is the covariance matrix of test dataset;

239 • $K(X,X^*)$: depicts the covariance matrix obtained from the training and test

240 dataset. Furthermore $K(X^*,X) = K(X,X^*)^T$.

241 Since $\mathbf{y}$ and $\mathbf{f}^*$ are jointly distributed, it is possible to condition the prior on the

242 observations (6) and determine how likely are predictions for $\mathbf{f}^*$. This can be expressed

243 as:

$$\mathbf{f}^* | X^*, X, \mathbf{y} \sim N \left( \bar{\mathbf{f}}^*, \text{cov}(\mathbf{f}^*) \right) \tag{7}$$

244 where

$$\bar{\mathbf{f}}^* = E\left[ \mathbf{f}^* | X^*, X, \mathbf{y} \right] = K(X^*,X)\left[ K(X,X) + \sigma_n^2 I \right]^{-1} \mathbf{y} \tag{8}$$

$$\text{cov}(\mathbf{f}^*) = K(X^*, X^*) - K(X^*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X^*) \tag{9}$$

245    The subsequent distribution can be used for the forecast of new test input points.

246    Indeed, $\bar{\mathbf{f}}^*$ is the predicted output value of the GPR model for test point. Additionally,

247    confidence interval (CI) of the predicted output value can be calculated through the

248    variance $\text{cov}(\mathbf{f}^*)$. For instance, the 95% CI can be determined by

249    $\left[\bar{\mathbf{f}}^* - 2 \times \sqrt{\text{cov}(\mathbf{f}^*)}, \bar{\mathbf{f}}^* + 2 \times \sqrt{\text{cov}(\mathbf{f}^*)}\right]$. As a consequence, the GPR model not only

250    supplies the predicted values but also furnishes the confidence level of the predicted

251    results.

252

253    Finally, the GPR model is a nonparametric model since the predicted outputs rely only

254    on the inputs and the observed values $\mathbf{y}$. In this way, parameters $\Theta = \{l, \sigma_f, \sigma_n\}$ are

255    termed the hyperparameters of the GPR model.

256

257    *2.3.2. Hyperparameter estimation*

258    In order to carry out this study, the dataset was divided into a training set with 80% of

259    the data, and a testing set with the remainder 20% of the data. A model was constructed

260    and optimised with the training data. It was then tested with the test dataset and the

261    optimisation of the parameters was performed with the help of the differential evolution

262    (DE) technique.

263

264    The predictive performance of GPR model depends exclusively on the suitability of the

265    chosen kernel. To estimate the kernel hyperparameters, an exhaustive search over a

266    discrete grid of values can be used, but this can be quite slow. The most usual method

267    considers an empirical Bayes approach that maximises the marginal likelihood. That is,

268    the optimal hyperparameters are achieved by maximising the log marginal likelihood.

269    The marginal likelihood $P(\mathbf{y}|X)$ is obtained, using Bayes' rule, as:

$$P(\mathbf{y}|X) = \int P(\mathbf{y}|f,X)P(f|X)df \qquad (10)$$

270    The term marginal likelihood refers to the marginalisation over the function values $\mathbf{f}$.

271    Since $\mathbf{y} \sim \mathcal{N}[0, K(X,X)]$, the log marginal likelihood can be written as:

$$\log p(\mathbf{y}\widehat{\mathbf{u}}X) = -\frac{1}{2}\mathbf{y}K_y^{-1}\mathbf{y} - \frac{1}{2}\log\widehat{\mathbf{u}}K_y\widehat{\mathbf{u}} - \frac{N}{2}\log(2\pi) \qquad (11)$$

272    where $K_y = K + \sigma_n^2 I, K = K(X,X)$ and $\widehat{\mathbf{u}}$ is the determinant. In this expression, the

273    first term is a data-fit term, the second term (always positive), and subtracted from it, is

274    a model complexity penalty, and the last term is simply a normalisation constant. This

275    expression therefore shows that the criterion of maximum marginal likelihood avoids

276    the problem of over-fitting because if two models are explaining the observed data with

277    the simplest one being chosen (Murphy, 2012; Witten, Frank, Hall, & Pal, 2016).

278

279    Following      parameter      initialisation,      the      optimal      hyperparameters

280    $\Theta' = \arg\max_{\Theta} \log p(\mathbf{y}|X,\Theta)$ can be calculated using any standard evolutionary

281    optimiser. In this study, the metaheuristic optimisation algorithm, denominated the DE

282    algorithm (Storn & Price, 1997; Price, Storn, & Lampinen, 2005; Feoktistov, 2006;

283    Simon, 2013), was used. The process is shown in Fig. 2.

284

285    **Fig. 2** – GPR Model selection using the DE optimisation technique.

286    *2.4. The goodness–of–fit of this approach*

287    Eight predicting variables were used (see section 2.2) to construct the new GPR–based

288    model. The output predicted variable was the outlet dissolved oxygen. To predict the

289    outlet dissolved oxygen from other input operating parameters, it is necessary to choose

290    the model that best fits the experimental data. To determine the goodness–of–fit, the

291    criterion considered here was the coefficient of determination $R^2$ (Picard & Cook, 1984;

292    Freedman, Pisani, & Purves, 2007). A dataset takes values $t_i$, each of which has an

293    associated modelled value $y_i$. The former are usually termed the observed values and

294    the latter often referred to as the predicted values. The dataset variability is measured

295    through different sums of squares as follows (Freedman, Pisani, & Purves, 2007):

296    - $SS_{tot} = \sum_{i=1}^{n}(t_i - \bar{t})^2$ : the total sum of squares, proportional to the sample variance.

297    - $SS_{reg} = \sum_{i=1}^{n}(y_i - \bar{t})^2$ : the regression sum of squares, also termed the explained

298       sum of squares.

299    - $SS_{err} = \sum_{i=1}^{n}(t_i - y_i)^2$ : the residual sum of squares.

300    Note that in the previous sums, $\bar{t}$ is the mean of the $n$ observed data:

$$\bar{t} = \frac{1}{n}\sum_{i=1}^{n} t_i \tag{12}$$

301    Taking into account the above sums, the coefficient of determination is defined via:

$$R^2 \equiv 1 - \frac{SS_{err}}{SS_{tot}} \tag{13}$$

15

302     Thus, a coefficient of determination value of 1.0 indicates that the regression curve fits

303     the data perfectly.

304

305     The value of $R^2$ was calculated using the optimised model with the testing dataset. The

306     module Gpy from the Gaussian process framework found in Python (Gpy, 2014;

307     Martin, 2018), along with the DE technique (Storn & Price, 1997; Price, Storn, &

308     Lampinen, 2005; Simon, 2013) were used to construct the final regression model.

309

310     It is well known that the GPR technique depends strongly on the following

311     hyperparameters (Friedman & Roosen, 1995; Aggarwal, 2015; Larose, 2015; Witten,

312     Frank, Hall, & Pal, 2016; Tan, Steinbach, Karpatne, & Kumar, 2018):

313     • Variance ($\sigma_f^2$):the signal variance that controls the vertical scale of the kernel

314       function.

315     • Lengthscale ($\ell$):the characteristic length-scale that controls the horizontal scale

316       over which the kernel function changes.

317     • Gaussian noise variance ($\sigma_n^2$): if $\varepsilon$ is the additive white noise and the Gaussian

318       noise is independent and identically distributed such that $\varepsilon \sim N\left(0, \sigma_n^2\right)$, then $\sigma_n^2$

319       is the variance of this noise.

320     1. A novel GPR–based model was constructed selecting as the dependent variable

321       the outlet dissolved oxygen from the other eight remaining variables which were

322       designated as input variables in the granular filters (Tien, 2012; Bové, Arbat,

323       Duran–Ros, Pujol, Velayos, Ramírez de Cartagena, & Puig–Bargués, 2015) and

324       studying their effect in order to optimise calculation by analysing $R^2$.

325

326    As previously mentioned, this GPR technique is greatly dependent on the

327    hyperparameters: variance ($\sigma^2$); lengthscale ($\ell$) and the Gaussian noise variance ($\sigma_n^2$).

328    The traditional way of performing hyperparameter optimisation has been *grid search*, or

329    a *parameter sweep*, which is simply an exhaustive searching through a manually

330    specified subset of the hyperparameter space of a learning algorithm. In this study, the

331    metaheuristic optimisation algorithm, the DE algorithm (Storn & Price, 1997; Price,

332    Storn, & Lampinen, 2005; Feoktistov, 2006; Simon, 2013) was used for

333    multidimensional real-valued functions but it did not use the gradient of the problem

334    being optimised, thus the DE did not require the optimisation problem to be

335    differentiable, as is required by classic optimisation methods such as the gradient

336    descent and quasi-Newton methods. Like other algorithms in this evolutionary category,

337    the DE maintains a population of candidate solutions, which are recombined and

338    mutated to produce new individuals which are chosen according to the value of their

339    performance function (Storn & Price, 1997). What characterises DE is the use of test

340    vectors, which compete with individuals in the current population in order to survive.

341

342    Additionally, the importance of the variables was studied. As categorical variables are

343    present, the chosen method depends on removing a variable, evaluating the new model

344    performance and comparing it with the performance of the full model. The greater the

345    decrease in the goodness-of-fit parameter, the greater the importance of the removed

346    independent variable.

347    **3. Results and discussion**

348    As stated earlier, the outlet dissolved oxygen was used as output dependent variable of

349    the proposed GPR–based model. The prediction performed from the independent

350    variables (Tien, 2012) was satisfactory.

351    Table 2 shows the optimal hyperparameters of the best fitted GPR–based model found

352    with the DE technique. The objective function value, in this case the marginal

353    likelihood was optimised to a value of 239 using the DE technique using the training

354    set.

355

356    **Table 2 -** Optimal hyperparameters of the best fitted GPR–based model found with the

357    DE technique: variance $\sigma_f^2$ and lengthscale $\ell$ for the RBF kernel, the Gaussian noise

358    variance $\sigma_n^2$ for the optimised models for the training set.

359

360    Taking into account the results achieved, the GPR technique in combination with the

361    DE meta-heuristic optimisation method was able to build models with a high

362    performance for estimating the outlet dissolved oxygen in micro-irrigation sand filters

363    fed with effluents using the test set. Indeed, the coefficient of determination ($R^2$) of the

364    fitted GPR model was of 0.9023 with a correlation coefficient of 0.9499 for the outlet

365    dissolved oxygen.

366

367    A graphical representation of the terms that formed the best fitted GPR–based model for

368    the outlet dissolved oxygen ($DO_o$) is shown below in Figs. 3 and 4. The first order

369    terms, that is, the variations of the dependent variable when all the variables but one are

370    constant (its median value) is shown in Fig 3. The graphs suggest that the variable $DO_i$

371    is the main influence for the variations in $DO_i$, while other variables as pH and

372    temperature do not significantly affect this variable as these curves are almost constant.

373    The same effect can be shown in the surfaces that represent the second order

374    relationships, that is, leaving all the independent variables constant but two. Again, it

375    can be seen that the main influence in rapid change of output variable was due to the

376    $DO_i$.

377

378    **Fig. 3 -** First-order terms for some of the independent variables for the dependent

379    variable output dissolved oxygen ($DO_o$).

380

381    **Fig. 4 -** Second-order terms of some of the independent variables for the dependent

382    variable output dissolved oxygen ($DO_o$).

383

384    The significance rankings for the input variables predicting the outlet dissolved oxygen

385    (output variable) in this complex nonlinear study are shown in Table 3 and Fig. 5. As

386    there are some categorical variables such as the filter type involved, the method where

387    discarding one independent variable from the model at a time and taking into account

388    the decrease in goodness-of-fit, in this case, the marginal likelihoods, is shown in Table

389    3. The result is, that for the GPR model, the most significant variable in $DO_o$ prediction

390    is the $DO_i$, followed by (in order) the type of filter, water temperature, height of the

391    filter bed, pH, velocity, turbidity, and electrical conductivity.

392

393 **Table 3 -** Log marginal likelihood variation value between the full model and the model

394 without the variable for the outlet dissolved oxygen ($DO_o$) model.

395

396 **Fig. 5 -** Relative relevance of the variables in the GPR model for the outlet dissolved

397 oxygen ($DO_o$).

398

399 As it could be anticipated, $DO_o$ was highly dependent on $DO_i$ since organic pollutants

400 are retained across filter media and chlorination of filter backwashing water reduced

401 microorganisms level, and therefore less oxygen is consumed and dissolved oxygen

402 could increase. However, DO removal depended also on media particle size (Elbana,

403 Ramírez de Cartagena, & Puig-Bargués, 2012) and on the interaction between filter type

404 and filtration velocity, considering input inlet DO as a co-variable (Solé–Torres, Puig–

405 Bargués, Duran–Ros, Arbat, Pujol, & Ramírez de Cartagena, 2019b). The filter type had

406 also a contribution on the results since different underdrain designs affect backwashing

407 performance and frequency (Burt, 2010), which is directly related to DO removal

408 (Enciso-Medina, Multer, & Lamm, 2011; Elbana, Ramírez de Cartagena, & Puig-

409 Bargués, 2012). The third parameter is temperature, but this is also logical since DO

410 values are temperature dependent.

411

412 The importance of $DO_i$ for estimating $DO_o$ has been previously observed by Martí et al.

413 (2013) and García–Nieto et al. (2016), working with different types of models. Martí et

414 al. (2013) observed that pH, EC and pressure loss, but not temperature, García–Nieto et

415 al. (2016) found that inlet turbidity and pressure loss were also considered as influential

416    parameters for predicting $DO_o$. Thus, the results highlight the importance of correctly

417    assessing the performance of each prediction model.

418

419    In conclusion, this research was able to estimate the outlet dissolved oxygen (output

420    variable) in agreement with the actual experimental values observed using the GPR–

421    based model with accuracy as well as success. Indeed, Fig. 6 shows the comparison

422    among the $DO_o$ values observed and those predicted by using the GPR model with the

423    testing set. The values predicted by the model using the samples of the testing dataset

424    show a very good agreement with the observed values. As it can be seen, predicted

425    values are very close to the observed values or within the 95% confidence interval. This

426    is to be expected since the coefficient of determination was equal to 0.90. Therefore, in

427    order to achieve the best effective approach in this regression problem it is mandatory

428    the use of a GPR model with a DE optimisation technique.

429

430    **Fig. 6 -** Observed and predicted $DO_o$ values, taking into account the confidence interval,

431    by using the GPR–based model with the testing set ( $R^2 = 0.9023$ ).

432

433    **4. Conclusions**

434    Taking into account the experimental observations and numerical predictions, the main

435    findings of this study can be summarised as follows:

436    • Firstly, the development of novel data-driven diagnostic techniques is very

437      useful to predict the $DO_o$ from the experimental measurements. In this sense, the

438 new GPR–based method used here is useful to evaluate the outlet dissolved
439 oxygen in sand media filters used in microirrigation systems.

440 • Secondly, the assumption that the outlet dissolved oxygen diagnosis can be
441 accurately modelled by using a hybrid GPR–based model in granular filters was
442 confirmed.

443 • Thirdly, a reasonable coefficient of determination (0.9023) was obtained when
444 this GPR–based model was applied to the experimental dataset corresponding to
445 the $DO_o$.

446 • Fourthly, the significance order of the input variables involved in the prediction
447 of the outlet dissolved oxygen in sand media filters was set. This is one of the
448 main findings in this work. Specifically, input variable dissolved oxygen ($DO_i$)
449 could be considered the most influential parameter in the prediction of the $DO_o$.
450 In this regard, it is also important to highlight the influential role of the type of
451 filter in the dependent variable outlet dissolved oxygen.

452 • Finally, the influence of the hyperparameters setting of the GPR approach on the
453 $DO_o$ regression performance was set up.

454 In summary, this methodology could be applied to other filtration processes with similar
455 or distinct filter media types with success, but it is always necessary to take into account
456 the characteristics of each filter and experiment. Consequently, an effective GPR–based
457 model is a good practical solution to the problem of the determining $DO_o$ in the sand
458 media filters usually used in microirrigation systems.

459

460

467

468 **References**

469 Aggarwal, C.C. (2015). *Data mining: the textbook*. New York, USA: Springer.

470 Ait-Mouheb, N., Bahri, A., Ben Thayer, B., Benyahia, B., Bourrié, G., Cherki, B. et al.

471    (2018). The reuse of reclaimed water for irrigation around the Mediterranean Rim: a

472    step towards a more virtuous cycle? *Regional Environmental Change*, 18, 693–705.

473 Bhattarai, S.P., Midmore, D.J. & Pendergast, L. (2008). Yield, water-use efficiencies

474    and root distribution of soybean, chickpea and pumpkin under different subsurface

475    drip irrigation depths and oxygenation treatments in vertisols. *Irrigation Science*,

476    26(5), 439–450.

477 Bové, J., Arbat, G., Duran–Ros, M., Pujol, T., Velayos, J., Ramírez de Cartagena, F., &

478    Puig–Bargués, J. (2015). Pressure drop across sand and recycled glass media used in

479    micro irrigation filters. *Biosystems Engineering*, 137, 55–63.

480 Bové, J., Puig–Bargués, J., Arbat, G., Duran–Ros, M., Pujol, T., Pujol, J., & Ramírez de

481    Cartagena, F. (2017). Development of a new underdrain for improving the

482    efficiency of microirrigation sand media filters. *Agricultural Water Management*,

483    179, 296–305.

23

484    Camps–Valls, G., Verrelst, J., Munoz–Mari, J., Laparra, V., Mateo–Jimenez, F., &

485        Gomez–Dans, J. (2016). A survey on Gaussian processes for earth-observation data

486        analysis: a comprehensive investigation. *IEEE Geoscience and Remote Sensing*

487        *Magazine*, 4(2), 58–78.

488    Canales-Ide, F., Zubelzu, S., & Rodríguez-Sinobas, L. (2019). Irrigation systems in

489        smart cities coping with water scarcity: The case of Valdebebas, Madrid (Spain).

490        Journal of Environmental Management, 247, 187-195.

491    Chakraborty, U.K. (2008). *Advances in differential evolution*. Berlin: Springer.

492    Chen, Y., Lu, D., Luo, L., Pokhrel, Y., Deb, K., Huang, J., & Ran, Y. (2018). Detecting

493        irrigation extent, frequency, and timing in a heterogeneous arid agricultural region

494        using MODIS time series, Landsat imagery, and ancillary data. *Remote Sensing of*

495        *Environment*, 204, 197–211.

496    Ebden, M. (2015). Gaussian processes: a quick introduction.

497        https://arxiv.org/pdf/1505.02965.pdf.

498    Elbana, M., Ramírez de Cartagena, F., & Puig-Bargués, J. (2012). Effectiveness of sand

499        media filters for removing turbidity and recovering dissolved oxygen from a

500        reclaimed effluent used for micro-irrigation. *Agricultural Water Management*, 111,

501        27–33.

502    Enciso-Medina, J., Multer, W.L. & Lamm, F.R. (2011). Management, maintenance, and

503        water quality effects on the long-term performance of subsurface drip irrigation

504        systems. *Applied Engineering in Agriculture*, 27 (6), 969–978.

505    Feoktistov, V. (2006). *Differential evolution: in search of solutions*. New York:

506        Springer.

Freedman, D., Pisani, R., & Purves, R. (2007). *Statistics*. New York: W.W. Norton & Company, New York.

García–Nieto, P.J., García–Gonzalo, E., Arbat, G., Duran–Ros, M., Ramírez de Cartagena, F., & Puig–Bargués, J. (2016). A new predictive model for the filtered volume and outlet parameters in micro-irrigation sand filters fed with effluents using the hybrid PSO–SVM–based approach. *Computers and Electronics in Agriculture*, 125, 74–80.

García–Nieto, P.J., García–Gonzalo, E., Arbat, G., Duran–Ros, M., Ramírez de Cartagena, F., & Puig–Bargués, J. (2018). Pressure drop modelling in sand filters in micro-irrigation using gradient boosted regression trees. *Biosystems Engineering*, 171, 41–51.

García–Nieto, P.J., García–Gonzalo, E., Bové, J., Arbat, G., Duran–Ros, M., & Puig–Bargués, J. (2017). Modeling pressure drop produced by different filtering media in microirrigation sand filters using the hybrid ABC–MARS–based approach, MLP neural network and M5 model tree. *Computers and Electronics in Agriculture,* 139, 65–74.

GPy, 2014. A Gaussian process framework in python. http://github.com/SheffieldML/GPy.

Ko, J., Klein, D.J., Fox, D., & Haehnelt, D. (2007a). GP-UKF: Unscented Kalman filters with Gaussian process prediction and observation models. In 2007 IEEE/RSJ International Conference on Intelligent Robots and Systems (pp. 1901–1907). San Diego, CA, USA: IEEE.

Ko, J., Klein, D.J., Fox, D., & Haehnelt, D. (2007b). Gaussian processes and reinforcement learning for identification and control of an autonomous blimp. In

531       Proceedings 2007 IEEE International Conference on Robotics and Automation (pp.

532       742–747). Roma, Italy: IEEE.

533    Kong, D., Chen, Y., & Li, N. (2018). Gaussian process regression for tool wear

534       prediction. *Mechanical Systems and Signal Processing*, 104, 556–574.

535    Kuhn, M., & Johnson, K. (2018). *Applied predictive modeling*. New York, USA:

536       Springer.

537    Larose, D.T. (2015). Data mining and predictive analytics. New York, USA: Wiley.

538    Liu, H., Yang, C., Huang, M., Wang, D., & Yoo, C. (2018). Modeling of subway indoor

539       air quality using Gaussian process regression. *Journal of Hazardous Materials*, 359,

540       266–273.

541    Ma, X., Xu, F., & Chen, B. (2019). Interpolation of wind pressures using Gaussian

542       process regression. *Journal of Wind Engineering & Industrial Aerodynamics*, 188,

543       30–42.

544    Maestre-Valero, J.F., & Martínez-Álvarez, V. (2010). Effects of drip irrigation systems

545       on the recovery of dissolved oxygen from hypoxic water. *Agricultural Water*

546       *Management*, 97, 1806–1812.

547    Marsland, S. (2014). *Machine learning: an algorithmic perspective*. Boca Raton, FL,

548       USA: Chapman and Hall/CRC Press.

549    Martí, P., Shiri, J., Duran–Ros, M., Arbat, G., Ramírez de Cartagena, F., & Puig–

550       Bargués, J. (2013). Artificial neural networks vs. Gene Expression Programming for

551       estimating outlet dissolved oxygen in micro-irrigation sand filters fed with effluents.

552       *Computers and Electronics in Agriculture*, 99, 176–185.

553    Martin, O. (2018). *Bayesian analysis with python*. Birmingham, UK: Packt Publishing.

554     Murphy, K.P. (2012). *Machine learning: a probabilistic perspective*. Cambridge, MA,

555         USA: The MIT Press.

556     Nakayama, F.S., Boman, B.J., & Pitts, D.J. (2007). Maintenance. In: Lamm, F.R.,

557         Ayars, J.E. & Nakayama, F.S. (Eds.), Microirrigation for Crop Production. Design,

558         Operation, and Management (pp. 389–430). Amsterdam, Netherlands: Elsevier.

559     Noori, M., Hassani, H., Javaherian, A., Amindavar, H., & Torabi, S. (2019). Automatic

560         fault detection in seismic data using Gaussian process regression. *Journal of Applied*

561         *Geophysics*, 163, 117–131.

562     Paananen, T., Piironen, J., Andersen, M.R., & Vehtari, A. (2019). Variable selection for

563         Gaussian processes via sensitivity analysis of the posterior predictive distribution. In

564         Proceedings of the 22nd International Conference on Artificial Intelligence and

565         Statistics (AISTATS), Proceedings of Machine Learning Research (PMLR) (pp.

566         1743–1752). Naha, Okinawa, Japan: arXiv:1712.08048 [stat.ME], Cornell

567         University, USA.

568     Picard, R., & Cook, D. (1984). Cross-validation of regression models. *Journal of the*

569         *American Statistical Association*, 79(387), 575–583.

570     Piironen, J., & Vehtari, A. (2016). Projection predictive model selection for Gaussian

571         processes. In 2016 IEEE 26th International Workshop on Machine Learning for

572         Signal Processing (MLSP) (pp. 1–6). Vietri sul Mare, Italy: IEEE.

573     Price, K., Storn, R.M., & Lampinen, J.A. (2005). *Differential evolution: A practical*

574         *approach to global optimization*. Berlin: Springer.

575     Puig–Bargués, J., Duran–Ros, M., Arbat, G., Barragán, J., & Ramírez de Cartagena, F.

576         (2012). Prediction by neural networks of filtered volume and outlet parameters in

577    micro-irrigation sand filters using effluents. *Biosystems Engineering*, 111(1), 126–

578    132.

579    Rasmussen, C.E. (2003). *Gaussian processes in machine learning: summer school on*

580    *machine learning*. Berlin, Germany: Springer.

581    Rasmussen, C.E., & Williams, C.K.I. (2006). *Gaussian processes for machine learning*.

582    Cambridge, MA, USA: The MIT Press.

583    Rogers, S., & Girolami, M. (2016). *A first course in machine learning*. Boca Raton, FL,

584    USA: Chapman and Hall/CRC.

585    Schneider, M., & Ertel, W. (2010). Robot learning by demonstration with local

586    Gaussian process regression. In: The 2010 IEEE/RSJ International Conference on

587    Intelligent Robots and Systems (pp. 255–260). Taipei, Taiwan: IEEE.

588    Schulz, E., Speekenbrink, M., & Krause, A. (2018). A tutorial on Gaussian process

589    regression: Modelling, exploring, and exploiting functions. *Journal of Mathematical*

590    *Psychology*, 85, 1–16.

591    Seeger, M. (2000). Bayesian model selection for support vector machines, Gaussian

592    processes and other kernel classifiers. In NIPS'99 Proceedings of the 12th

593    International Conference on Neural Information Processing Systems (vol. 12, pp.

594    603–609). Cambridge, MA, USA: The MIT Press.

595    Shi, J.Q., & Choi, T. (2011). *Gaussian process regression analysis for functional data*.

596    Boca Raton, FL, USA: Chapman and Hall/CRC Press.

597    Simon, D. (2013). *Evolutionary optimization algorithms*. New York: Wiley.

598  Solé–Torres, C., Duran–Ros, M., Arbat, G., Pujol, J., Ramírez de Cartagena F., & Puig–

599  Bargués, J. (2019a). Assessment of field water uniformity distribution in a

600  microirrigation system using a SCADA system. *Water*, 11(7), 1346–1359.

601  Solé–Torres, C., Puig–Bargués, J., Duran–Ros, M., Arbat, G., Pujol, J., & Ramírez de

602  Cartagena, F. (2019b). Effect of underdrain design, media height and filtration

603  velocity on the performance of microirrigation sand filters using reclaimed effluents.

604  *Biosystems Engineering*, 187, 292–304.

605  Storn, R., & Price, K. (1997). Differential evolution - a simple and efficient heuristic for

606  global optimization over continuous spaces. *Journal of Global Optimization*, 11,

607  341–359.

608  Tan, P.–N., Steinbach, M., Karpatne, A., Kumar, V. (2018). *Introduction to data

609  mining*. Oxford, UK: Pearson.

610  Tien, C. (2012). *Principles of filtration*. Kidlington, Oxford, UK: Elsevier.

611  Trooien, T.P., & Hills, D.J. (2007). Application of biological effluent. In Lamm, F.R.,

612  Ayars, J.E., & Nakayama, F.S. (Eds.), Microirrigation for Crop Production. Design,

613  Operation and Management (pp. 329–356). Amsterdam: Elsevier.

614  Vidales, A. (2019). *Machine learning with MATLAB: Gaussian process regression,

615  analysis of variance and Bayesian optimization*. Independently published.

616  Witten, I.H., Frank, E., Hall, M.A., Pal, C.J. (2016). *Data mining: practical machine

617  learning tools and techniques*. Cambridge, MA, USA: Morgan Kaufmann.

618  Wu, R., & Wang, B. (2018). Gaussian process regression method for forecasting of

619  mortality rates. *Neurocomputing*, 316, 232–239.

620     Zhou, Y., Zhou, B., Xu, F., Muhammad, T., Li, Y. (2019). Appropriate dissolved

621        oxygen concentration and application stage of micro-nano bubble water oxygenation

622        in greenhouse crop plantation. *Agricultural Water Management*, 223, 105713.

623

**Fig. 1 -** Picture of the experimental set-up with the three filter designs: (a) red: arm collector; (b) blue: inserted domes; and (c) green: a porous media prototype.
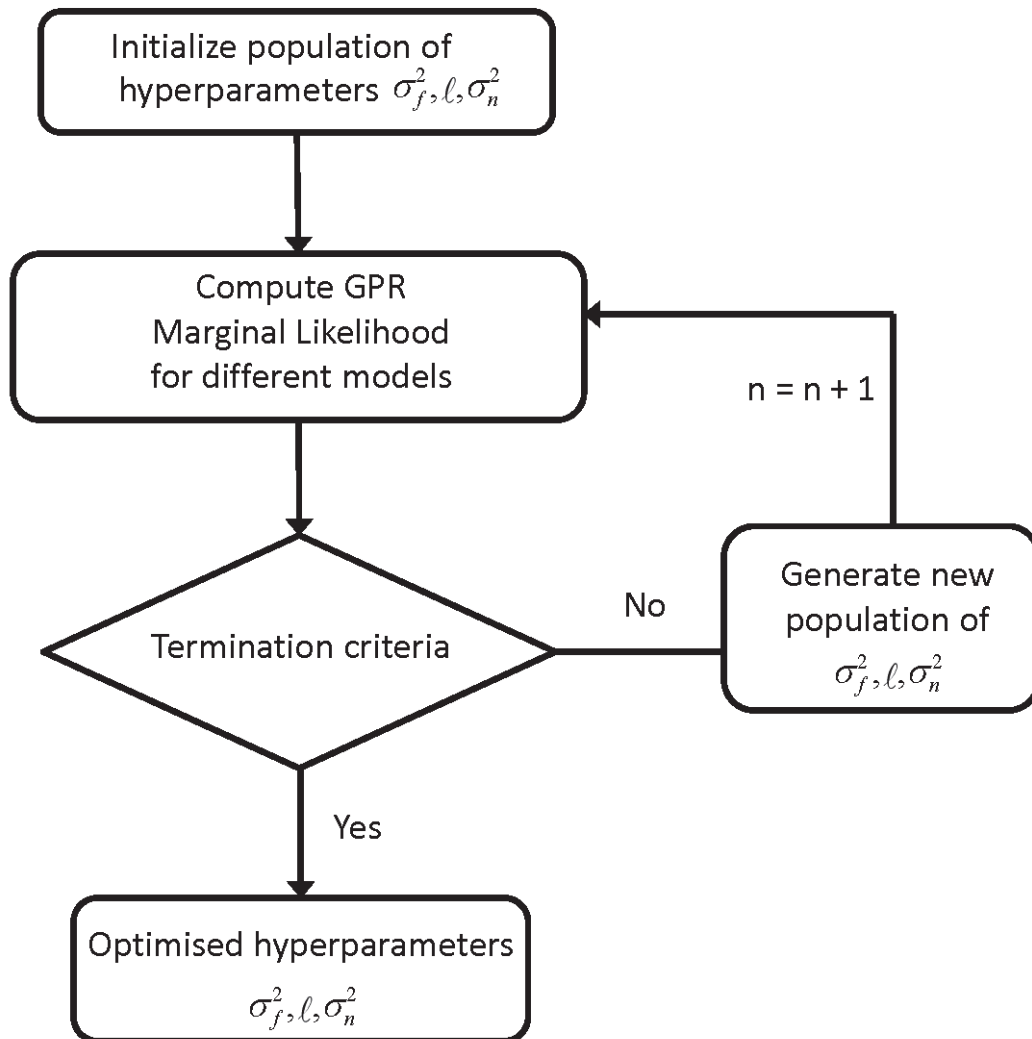
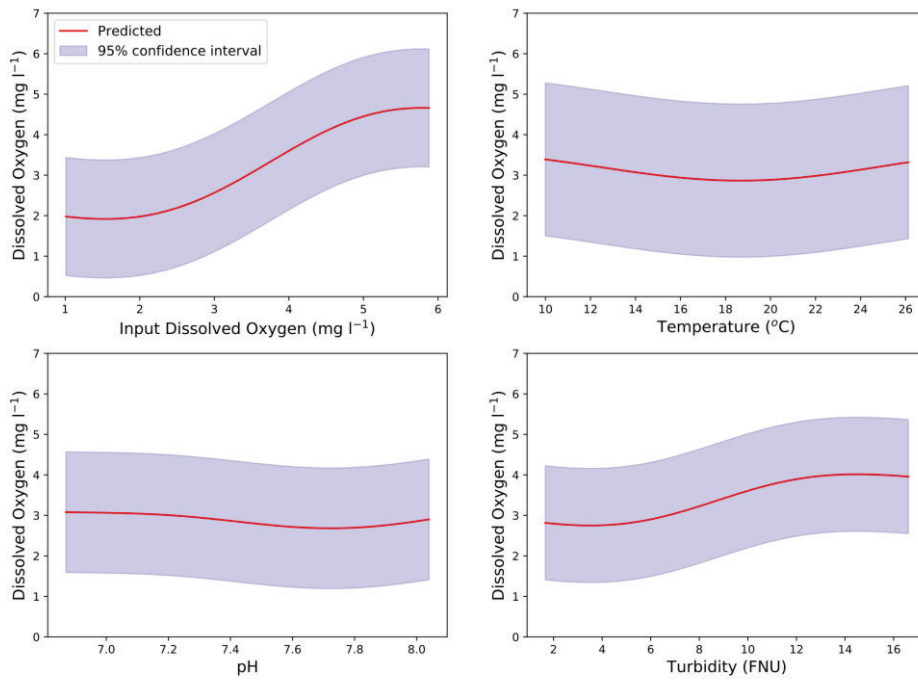**Fig. 2** – GPR Model selection using the DE optimisation technique.

**Fig. 3 -** First-order terms for some of the independent variables for the dependent variable output dissolved oxygen (*DO$_o$*).
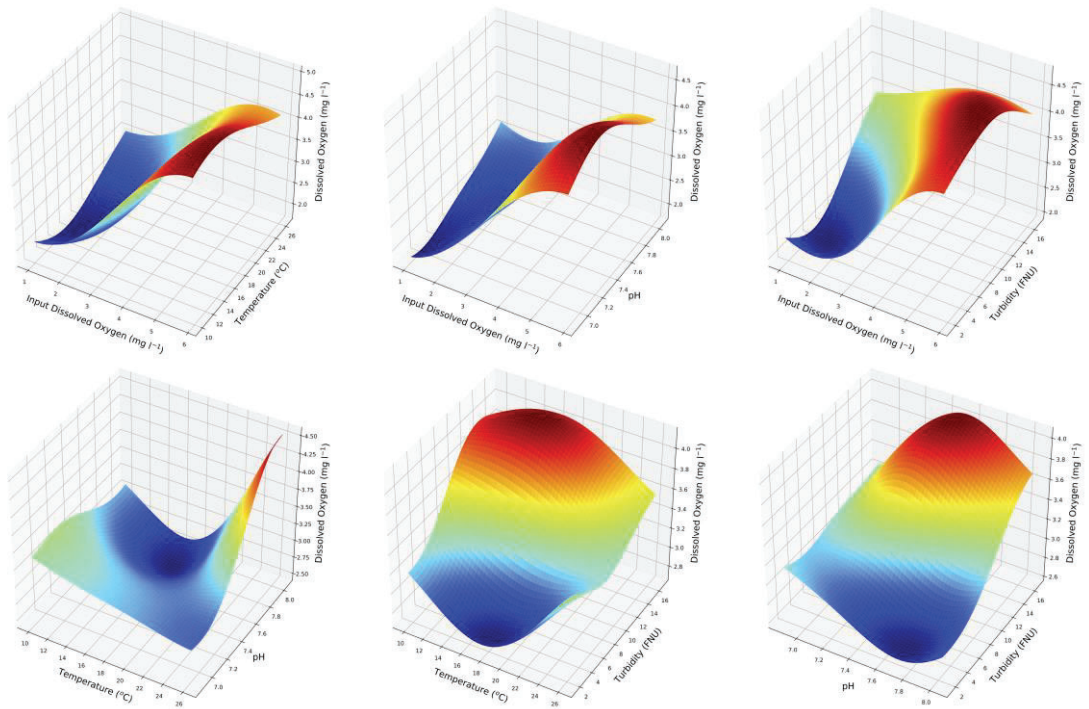
**Fig. 4 -** Second-order terms of some of the independent variables for the dependent variable output dissolved oxygen ($DO_o$).
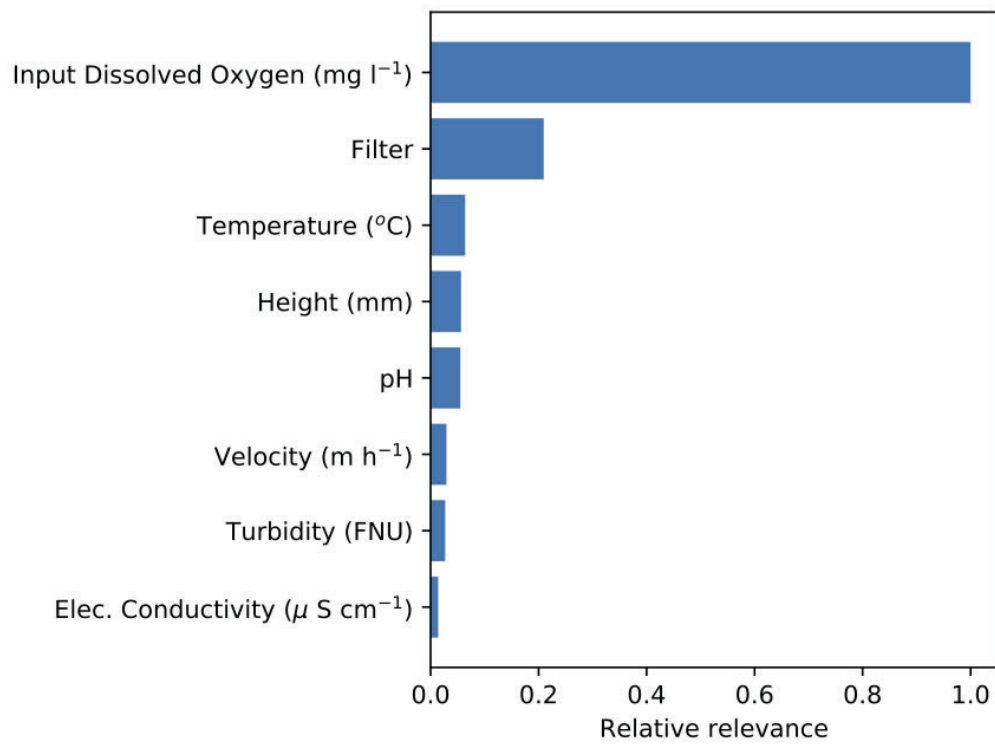
**Fig. 5 -** Relative relevance of the variables in the GPR model for the outlet dissolved oxygen ($DO_o$).
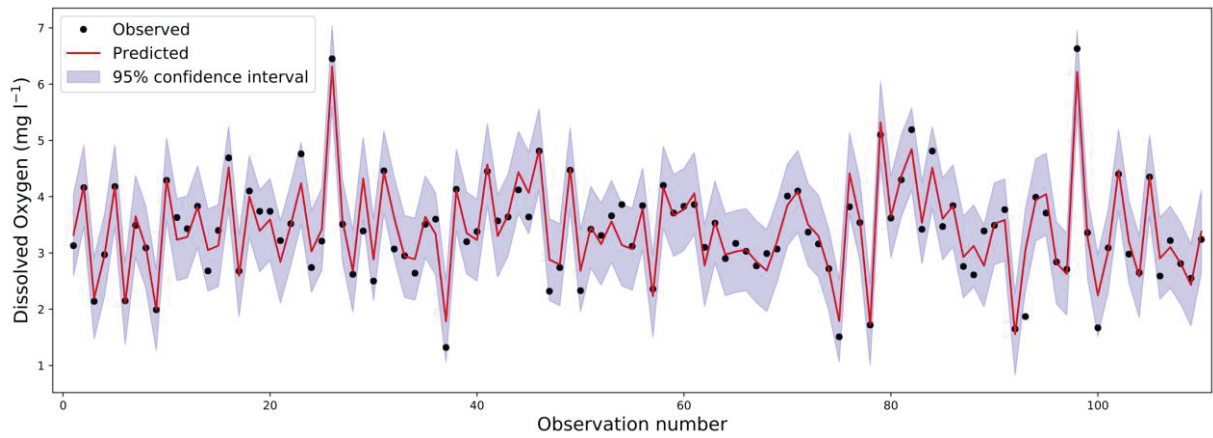
**Fig. 6 -** Observed and predicted $DO_o$ values, taking into account the confidence interval,

by using the GPR–based model with the testing set ( $R^2 = 0.9023$ ).

**Table 1 -** Set of operation physical input variables used in this study and their names along with their means and standard deviations.

| Input variables | Name of the variable | Mean | Standard deviation |
|---|---|---|---|
| Filter media type | Filter | -- | -- |
| Height of the filter bed (mm) | $H$ | 256.31 | 49.601 |
| Filtration velocity (m h$^{-1}$) | $v$ | 49.909 | 14.174 |
| Electrical conductivity ($\mu$ S cm$^{-1}$) | $CE_i$ | 2575.6 | 497.68 |
| Input dissolved oxygen (mg l$^{-1}$) | $DO_i$ | 3.3529 | 0.9860 |
| pH | $pH_i$ | 7.3526 | 0.2229 |
| Input turbidity (FNU) | $Turb_i$ | 6.1029 | 2.5898 |
| Water temperature (ºC) | $T_i$ | 20.002 | 3.3486 |

**Table 2 -** Optimal hyperparameters of the best fitted GPR–based model found with the DE technique: variance $\sigma_f^2$ and length-scale $\ell$ for the RBF kernel, the Gaussian noise variance $\sigma_n^2$, and the corresponding objective function value for the optimized models for the training set.

| Output variable | $\sigma_f^2$ | $\ell$ | $\sigma_n^2$ | Objective function value |
|---|---|---|---|---|
| $DO_o$ | 1.57 | 1.97 | 0.0636 | 239 |

**Table 3 -** Log marginal likelihood variation value between the full model and the model without the variable for the $DO_o$ model.

| Variable | Likelihood variation |
|---|---|
| Input dissolved Oxygen (mg l$^{-1}$) | 589.62 |
| Filter | 123.51 |
| Water temperature (ºC) | 37.77 |
| Height (mm) | 332.1 |
| pH | 32.45 |
| Velocity (m h$^{-1}$) | 17.31 |
| Input turbidity (FNU) | 15.96 |
| Electrical Conductivity (µS cm$^{-1}$) | 8.25 |