

## Predicting non-deposition sediment transport in sewer pipes using Random forest

Montes, Carlos; Kapelan, Zoran; Saldarriaga, Juan

**DOI**

[10.1016/j.watres.2020.116639](https://doi.org/10.1016/j.watres.2020.116639)

**Publication date**

2021

**Document Version**

Accepted author manuscript

**Published in**

Water Research

**Citation (APA)**

Montes, C., Kapelan, Z., & Saldarriaga, J. (2021). Predicting non-deposition sediment transport in sewer pipes using Random forest. *Water Research*, 189, 1-11. Article 116639. <https://doi.org/10.1016/j.watres.2020.116639>

**Important note**

To cite this publication, please use the final published version (if applicable). Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

1 **Predicting non-deposition sediment transport in sewer pipes using**  
2 **Random Forest**

3 Carlos Montes<sup>a\*</sup>, Zoran Kapelan<sup>b</sup> and Juan Saldarriaga<sup>c</sup>

4 *<sup>a</sup>Department of Civil and Environmental Engineering, Universidad de los Andes, Bogotá,*  
5 *Colombia; e-mail: cd.montes1256@uniandes.edu.co*

6 *<sup>b</sup>Department of Water Management, Delft University of Technology, Delft, Netherlands;*  
7 *e-mail: Z.Kapelan@tudelft.nl*

8 *<sup>c</sup>Department of Civil and Environmental Engineering, Universidad de los Andes, Bogotá,*  
9 *Colombia; e-mail: jsaldarr@uniandes.edu.co*

10 \*corresponding author; Correspondence address: Cra 1 Este No. 19A – 40 Bogota  
11 (Colombia); Tel.: +57-1-339-49-49 (ext. 1765)

# Predicting non-deposition sediment transport in sewer pipes using Random Forest

## Abstract

Sediment transport in sewers has been extensively studied in the past. This paper aims to propose a new method for predicting the self-cleansing velocity required to avoid permanent deposition of material in sewer pipes. The new Random Forest (RF) based model was implemented using experimental data collected from the literature. The accuracy of the developed model was evaluated and compared with ten promising literature models using multiple observed datasets. The results obtained demonstrate that the RF model is able to make predictions with high accuracy for the whole dataset used. These predictions clearly outperform predictions made by other models, especially for the case of non-deposition with deposited bed criterion that is used for designing large sewer pipes. The volumetric sediment concentration was identified as the most important parameter for predicting self-cleansing velocity.

Keywords: non-deposition; random forest; sediment transport; self-cleansing; sewer systems.

## 1. INTRODUCTION

Designing sediment-carrying sewer systems is a well-known field of research in hydraulic engineering. This interest is explained by the problems related to the presence of material in the systems. Due to the varying environmental conditions (i.e. loading and sediment characteristics and intermittent flow), the risk of building up a permanent sediment deposit increases during dry weather seasons. These deposits lead to problems such as reduced pipe capacity, increased roughness, and premature overflows. As an example, Ackers et al. (2001) showed that the presence of a permanent deposit at the bottom of sewer pipes increases hydraulic roughness and reduces discharge capacity by about 20%.

38           The most common criterion to avoid permanent deposit of material in sewer pipes  
39 is known as non-deposition. Several authors (Safari et al., 2018; Vongvisessomjai et al.,  
40 2010) have classified this criterion into two subgroups: 1) Non-deposition without  
41 deposited bed and 2) Non-deposition with deposited bed. Both groups are based on the  
42 presence of sediments at the bottom of the pipe. In the first case, high water velocities  
43 produce an individual and separate movement of the particles by slicing or rolling over  
44 the pipe invert, i.e. without deposited bed. In contrast, the second case is seen when lower  
45 water velocities are presented and the particles are grouped and move as a transitional  
46 deposited bed.

47           In the case of ‘without deposited bed’, traditional criteria of minimum velocities  
48 and shear stress values are commonly found in water utilities standards and industry  
49 design codes. Generally, these standards and codes suggest values ranging from 0.30 m  
50 s<sup>-1</sup> to 1.0 m s<sup>-1</sup> for minimum velocity and from 1.0 Pa to 4.0 Pa for shear stress (Montes  
51 et al., 2019; Nalluri and Ab Ghani, 1996; Vongvisessomjai et al., 2010). Several authors  
52 (Merritt and Enfinger, 2019; Nalluri and Ab Ghani, 1996) have shown how traditional  
53 threshold values lead to over-design of small diameter pipes and under-design of large  
54 diameter pipes (as a rule-of-thumb, pipes with diameter greater than 500 mm).  
55 Consequently, large sewers commonly require frequent removal of sediment deposits  
56 (Ackers et al., 2001) because of the minimum self-cleansing value adopted during the  
57 design stage. A unique design value is inadequate; hence sediment characteristics and  
58 hydraulic conditions must be included in the definition of the self-cleansing design  
59 criterion.

60           According to Safari and Aksoy (2020), existing traditional self-cleansing criteria  
61 can be up to 20% different from laboratory-scale measured values. The channel cross-  
62 section is relevant in the choice of the self-cleansing criterion. For example, rectangular

63 cross-sections require lower velocities compared to V-bottom or U-shape channels. Even  
 64 criteria based on the Shields diagram, such as the Camp criterion, seem to be inadequate  
 65 to define the self-cleansing value due to the non-inclusion of sediment concentration.

66 The above has motivated extensive experimental research (Ab Ghani, 1993; El-  
 67 Zaemey, 1991; May, 1993; May et al., 1989; Mayerle, 1988; Montes et al., 2020a, 2020b;  
 68 Ota, 1999; Perrusquía, 1991; Vongvisessomjai et al., 2010) aiming to collect data and  
 69 developing models for predicting the self-cleansing velocity as a function of sediment  
 70 characteristics and system hydraulics, based on the concept of non-deposition. These  
 71 studies have been carried out at laboratory scale under well-controlled and steady flow  
 72 conditions, using non-cohesive sediments. Different authors collected data in pipes with  
 73 different materials (e.g. concrete, acrylic or PVC, among other materials) and internal  
 74 diameters, ranging from 100 mm to 595 mm. In the end, all these studies proposed a  
 75 model for predicting the self-cleansing conditions in practice that was either developed  
 76 with their own experimental data or using the benchmark data reported in the literature.  
 77 Most models developed are regression-based and include the group of input parameters  
 78 that most affect the prediction of the self-cleansing velocity (Ackers et al., 2001; Ebtehaj  
 79 and Bonakdari, 2016a; May et al., 1996). Most of these models are in the form of:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = aC_v^b \left(\frac{d}{R} \text{ or } \frac{d}{D}\right)^c \lambda^e D_{gr}^f \left(\frac{W_b}{Y} \text{ or } \frac{y_s}{Y} \text{ or } \frac{y_s}{D}\right)^g \left(\frac{P}{B}\right)^h \quad (1)$$

80 where  $V_l$  is the self-cleansing velocity,  $d$  the mean particle diameter,  $g$  the gravity  
 81 acceleration coefficient,  $S_s$  the specific gravity of sediments,  $C_v$  the volumetric sediment  
 82 concentration,  $R$  the hydraulic radius,  $D$  the pipe diameter,  $\lambda$  the channel friction factor,  
 83  $D_{gr}$  the dimensionless grain size  $\left(= \left(\frac{(S_s-1)gd^3}{\nu^2}\right)^{\frac{1}{3}}\right)$ ,  $\nu$  the water kinematic viscosity,  $W_b$   
 84 the sediment deposited width,  $P$  the wetted perimeter,  $y_s$  the sediment deposited  
 85 thickness,  $B$  the water surface width,  $Y$  the water level and  $a, b, c, e, f, g$  and  $h$  regression

86 coefficients. Other parameters as  $V_t$  the threshold velocity required to initiate movement  
87  $(= 0.125(gd(S_s - 1))^{0.5}(Y/d)^{0.47})$  and  $S_o$  the pipe slope have also been included in  
88 regression models (May et al., 1996; Montes et al., 2020a).

89 Most of above studies for both non-deposition criteria, have developed predictive  
90 models which tend to be overfitted to their own experimental data. This problem can be  
91 seen especially in the earlier works, where no advanced techniques were used to develop  
92 regression models. For example, several authors (Montes et al., 2020b; Safari et al., 2018)  
93 have pointed out that early work of Mayerle's (1988) has developed a model that shows  
94 high accuracy prediction with its data and poor prediction when other datasets are used.  
95 In contrast, recent regression-models, which used novel techniques such as Evolutionary  
96 Polynomial Regression – Multi-Objective Genetic Algorithm (EPR-MOGA) and Least  
97 Absolute Shrinkage and Selection Operator (LASSO) have demonstrated better  
98 prediction results (Montes et al., 2020a, 2020b).

99 In order to address the above overfitting issue in regression models, new Machine  
100 Learning (ML) and Artificial Intelligence (AI) techniques have been introduced for  
101 predicting the self-cleansing velocity based on the concept of non-deposition sediment  
102 transport. Examples of models developed for the 'without deposited bed' case include  
103 using techniques such as Artificial Neural Network (ANN) (Ebtehaj and Bonakdari,  
104 2013), Support Vector Regression (SVR) coupled with the Firefly Algorithm (Ebtehaj  
105 and Bonakdari, 2016b), the Group Method of Data Handling (GMDH) (Ebtehaj and  
106 Bonakdari, 2016a), neuro-fuzzy inference system combined with the Particle Swarm  
107 Optimisation (ANFIS-PSO) (Ebtehaj et al., 2019), Decision Trees (DT), Generalised  
108 Regression Neural Network (GRNN), Multivariate Adaptive Regression Splines (MARS)  
109 (Safari, 2019) and Extreme Learning Machine (ELM) (Ebtehaj et al., 2020). For the other  
110 case of 'non-deposition with deposited bed', fewer ML/AI type models have been

111 developed. Examples include models based on Particle Swarm Optimisation (PSO)  
112 algorithm (Safari et al., 2017), Gene Expression Programming (GEP) (Roushangar and  
113 Ghasempour, 2017) and Multigene Genetic Programming (MGP) (Safari and Danandeh  
114 Mehr, 2018).

115         The above models, developed using different ML/AI techniques (for both non-  
116 deposition criteria), have improved the prediction accuracy of self-cleansing velocities  
117 and addressed the issues of model overfitting but only partially. As noted by Zendehboudi  
118 et al. (2018), these models still tend to have rather limited extrapolation capabilities  
119 meaning that once they are applied to datasets that were not used for their training they  
120 tend to underperform. Also, the ML/AI based models developed so far are largely black-  
121 box type models (e.g. ANN) meaning that, unlike white-box type regression models, they  
122 suffer from low interpretability of physical significance of model inputs (i.e. explanatory  
123 factors), and interactions with the model output.

124         The aim of this paper is to overcome above deficiencies using the Random Forest  
125 (RF) technique for predicting self-cleansing sewer velocities. RF (Breiman, 2001) is a  
126 flexible and interpretable supervised ML technique that combines the results (outputs) of  
127 multiple individual decision trees to make a prediction of interest. Due to its good  
128 characteristics and easy application, it has been a widely used for addressing many other  
129 problems in water engineering. Tyrallis et al. (2019) showed a full review of studies in  
130 which RF was successfully applied to water resources problems.

131         Using the RF technique, a new predictive self-cleansing model is developed and  
132 presented here for both non-deposition criteria (with and without deposited bed). This  
133 model aims to increase prediction accuracy whilst avoiding overfitting issues and  
134 enabling interpretability of results obtained. The new modelling technique is compared  
135 to ten literature models using multiple datasets.

## 136 2. DATA

### 137 2.1. *Non-deposition without deposited bed data*

138 Several experimental data were collected from the literature to implement the RF  
139 method. Mayerle (1988) studied the sediment transport in a 152 mm diameter pipe and in  
140 two rectangular channels of 311.5 mm and 462.3 mm bottom width ( $W$ ) using granular  
141 sands ranging from 0.50 mm to 8.74 mm. Ab Ghani (1993) collected 221 data in 154 mm,  
142 305 mm and 450 mm diameter pipes, testing sands between 0.46 mm and 8.40 mm. Ota  
143 (1999) used a 225 mm concrete pipe with a constant slope of 0.002, varying the  
144 volumetric sediment concentration between 4.2 ppm to 59.4 ppm. Vongvisessomjai et al.  
145 (2010) used two circular PVC pipes of 100 mm and 150 mm diameter to study the bedload  
146 and suspended load transport. Montes et al. (2020a) collected experimental data in a 242  
147 mm acrylic pipe using granular material with a mean particle diameter of 0.35 mm and  
148 1.51 mm. Montes et al. (2020b) carried out 107 experiments in a 595 mm PVC pipe, using  
149 sediments ranging from 0.35 mm to 2.6 mm.

### 150 2.2. *Non-deposition with deposited bed data*

151 For the non-deposition with deposited bed, El-Zaemey (1991) studied the  
152 sediment transport in a 305 mm diameter pipe, using granular particles ranging from 0.53  
153 mm to 8.40 mm. Perrusquía (1991) carried out experiments in a 225 mm diameter pipe,  
154 varying the sediment concentration from 18.7 ppm to 408.0 ppm. Ab Ghani (1993)  
155 collected the deposited bed data only in the 450 mm concrete pipe and using granular  
156 sand with a mean particle diameter of 0.72 mm. May (1993) extended their previous study  
157 (May et al., 1989) and collected experimental data with sediment thickness varying from  
158 57.6 mm to 129.6 mm. Finally, Montes et al. (2020b) carried out experiments in a 595  
159 mm PVC pipe, considering a relative sediment thickness ( $y_s/D$ ) between 0.13% and



160 1.11%. Table 1 outlines the characteristics of the data used for developing the RF  
161 algorithm.

162 **[Table 1 near here]**

163 As shown in Table 1, a total of 664 and 454 data are available for the development  
164 of models without deposited bed and the deposited bed criteria, respectively.

### 165 **3. MEHODOLOGY**

#### 166 **3.1. Random Forest Model**

167 Random Forest model developed here predicts the particle Froude number ( $F_r^*$ ) as a  
168 function of several well-known dimensionless explanatory factors (Kargar et al., 2019;  
169 Vongvisessomjai et al., 2010):

$$F_r^* = \frac{V_l}{\sqrt{gd(S_s - 1)}} = f\left(C_v, D_{gr}, \frac{d}{R}, \lambda, \frac{y_s}{D}\right) \quad (2)$$

170 Random forest (RF) is a bagging algorithm for regression and classification  
171 problem proposed by Breiman (2001). This is a low-variance method, which randomly  
172 split the training data and the input variables predictors to build a set of  $b$  decision trees  
173 ( $B_t$ ). The results of all decision trees generated from bootstrapped training samples  
174 ( $T_b(x; \theta_b)$ ) are then averaged, i.e. the final result ( $\hat{y}(x)$ ) is the average of the output of  
175 all decision trees (as shown in Eq. (3)). This procedure ensures the reduction of the model  
176 variance and consequently, the reduction of the risk of overfitting. A simplified  
177 conceptual diagram of the RF method is shown in Figure 1.

$$\hat{y}(x) = \frac{1}{B_t} \sum_{b=1}^{B_t} T(x; \theta_b) \quad (3)$$

178 **[Figure 1 near here]**

179 In this paper, the R package ‘RandomForest’ (Liaw and Wiener, 2002) was used  
180 for constructing both non-deposition, without deposited bed and deposited bed, self-  
181 cleansing models. The number of predictors considered at each split ( $mtry$ ) and the  
182 number of trees in the forest ( $B_t$ ) are the parameters that define the structure of the RF  
183 regression model. The  $mtry$  parameter is estimated by using the `rfcv()` function, which  
184 shows the cross-validation performance for each number of predictors. In addition, the  
185 optimal number of trees is defined as the value that minimises the Mean Square Error  
186 (MSE) value of the training data. These parameters are estimated and the results are  
187 shown in Figure 2. According to this figure, the optimal number of features (i.e. the  
188 random predictors used in each tree) are three and four non-dimensional parameters for  
189 the cases of without deposited bed and with deposited bed, respectively. Similarly, the  
190 optimal number of trees is 471 for without deposited bed and 229 for with deposited bed.

191 **[Figure 2 near here]**

192 Cross-validation is carried out during the training stage using out-of-bag (OOB)  
193 samples. As mentioned above, the method randomly bootstraps the training sample, that  
194 is, some of the training data are left out to build each decision tree. Only two out of three  
195 parts of the total training data are used to build the tree (Breiman, 2001). Based on this,  
196 data not included in the bootstrapped sample (OOB data) are predicted, and the prediction  
197 error is averaged over the trees that do not include these data (OOB Error).

### 198 **3.1.1. *Splitting of training and testing data***

199 The whole benchmarking data collected from the literature are used for both training and  
200 testing stages of the RF model. Usually, 75% of the data is used during the training stage  
201 of the model and the other 25% to validate the results. According to Safari (2020), the  
202 range of variation in the training data has direct implications for model performance (i.e.  
203 accuracy). As a result, the model can show overfitting issues and poor extrapolation



228 Using the above considerations, the RF model is implemented with the optimal  
229 parameters defined in Figure 2 and using the ranges of variation of the training data  
230 outlined in Table 2. The full data collected from the literature are shown in the  
231 Supplementary material. Table S1 and Table S2 show the data for non-deposition without  
232 and with deposited bed, respectively, and the corresponding RF particle Froude number  
233 predictions. The implemented code for the RF method is shown in Figure 4. An example  
234 of one of the 471 decision trees generated by the RF model, for the non-deposition without  
235 deposited bed, is shown in Figure S1, in the Supplementary material.

236 **[Figure 4 near here]**

### 237 3.1.2. *Measure of feature importance*

238 Note that in this paper, a decrease in model accuracy when the  $j$ th variable is  
239 permuted (i.e. the percentage of the increase in the MSE,  $\%IncMSE$ ) is considered as a  
240 measure of the importance of a model input variable. This index shows the strength of  
241 each explanatory variable based on the reduction of the MSE. The step-by-step to  
242 calculate the  $\%IncMSE$  is shown as follows (Hastie et al., 2009):

- 243 (1) Calculate the MSE of the OOB-sample data in each tree of the forest ( $MSE_b$ ).
- 244 (2) Randomly permute the value of the  $j$ th explanatory variable and calculate the MSE  
245 ( $MSE_j$ ).
- 246 (3) Finally, calculate  $\%IncMSE$  for each explanatory variable as:

$$\%IncMSE = 100 \cdot \frac{MSE_j - MSE_b}{MSE_b} \quad (4)$$

247 As a result, the more the  $\%IncMSE$  increases for a variable, the more important  
248 it is.

## 249 3.2. Performance Assessment

### 250 3.2.1. Models used for comparing the RF results

251 In order to evaluate the RF model performance, it is compared to several literature  
252 models. The models selected for comparison are the replicable white-box models with  
253 high prediction accuracy reported in the literature and two black-box models where the  
254 implementing code is provided in the original papers. Other black-box models cannot be  
255 evaluated due to the limited replicability shown by these models (e.g. ANN). Based on  
256 this, in the case of non-deposition without deposited bed, seven models selected are the  
257 EPR-MOGA model (Montes et al., 2020a), the GEP model (Kargar et al., 2019), the  
258 MARS model (Safari, 2019), the May et al. (1996) model, the Safari and Aksoy (2020)  
259 model, the ANFIS-PSO model (Ebtehaj et al., 2019) and the ELM model (Ebtehaj et al.,  
260 2020). In the case of non-deposition with deposited bed, three models used for  
261 comparison are the PSO model (Safari and Shirzad, 2019), the LASSO model (Montes et  
262 al., 2020b) and the MGP model (Safari and Danandeh Mehr, 2018). The EPR-MOGA,  
263 LASSO, May et al. (1996) and Safari and Aksoy (2020) are the regression type models  
264 whilst GEP, MARS, ANFIS-PSO, ELM, PSO and MGP models make use of ML/AI  
265 techniques.

266 The equations used by above ten models are as follows:

267 EPR-MOGA:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = 5.6C_v^{0.16} \left(\frac{d}{R}\right)^{-0.58} S_o^{0.14} D_{gr}^{0.02} \quad (5)$$

268 GEP:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = \frac{3.05C_v^{0.16}}{\operatorname{atan}\left(\operatorname{atan}\left(\sqrt{\frac{d}{R}}\right)\right)} + \operatorname{atan}(3.41 - \ln(D_{gr})) \\ + \operatorname{atan}\left(\tan\left(\left(8.37 - 7.99\lambda + \frac{d}{R}\lambda\right)^2\right)\right) + \ln\left(\left(\left(\frac{d}{R}\right)^3\right)^{2\lambda}\right) \quad (6)$$

269 MARS:

$$\begin{aligned} \frac{V_l}{\sqrt{gd(S_s - 1)}} = & 7.26 - 1.75 \cdot \max(0, d/R - 0.12) + 2 \\ & \cdot \max(0, 0.12 - d/R) + 15.89 \cdot \max(0, C_v - 0.44) - 16.42 \\ & \cdot \max(0, 0.44 - C_v) + 0.47 \cdot \max(0, D_{gr} - 0.29) - 7.25 \\ & \cdot \max(0, \lambda - 0.3) - 16.03 \cdot \max(0, C_v - 0.01) + 3.7 \\ & \cdot \max(0, D_{gr} - 0.12) - 4.33 \cdot \max(0, D_{gr} - 0.08) + 0.43 \\ & \cdot \max(0, \lambda - 0.59) + 6.75 \cdot \max(0, \lambda - 0.28) + 1.67 \\ & \cdot \max(0, d/R - 0.07) \end{aligned} \quad (7)$$

270 May et al. (1996):

$$C_v = 0.0303 \left( \frac{D^2}{A} \right) \left( \frac{d}{D} \right)^{0.6} \left( 1 - \frac{V_t}{V_l} \right)^4 \left( \frac{V_l^2}{gd(S_s - 1)} \right)^{1.5} \quad (8)$$

271 Safari and Aksoy (2020):

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = 4.83 C_v^{0.09} \left( \frac{d}{R} \right)^{-0.32} D_{gr}^{-0.14} \left( \frac{P}{B} \right)^{0.20} \quad (9)$$

272 ANFIS-PSO:

273 No equation. The Matlab code can be found in Ebtehaj et al. (2019).

274 ELM:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = \left[ \frac{1}{(1 + \exp(-InW \cdot InV + BHI))} \right]^T \cdot OutW \quad (10)$$

275 where  $InW$  and  $OutW$  are the input and output weights,  $BHI$  the bias of the hidden

276 neurons and  $InV$  the input variables (i.e.  $C_v$ ,  $d/R$ ,  $D^2/A$ ,  $R/D$ ,  $D_{gr}$ ,  $d/D$  and  $\lambda$ ). Full

277 details of the values chosen for each parameter are shown in Ebtehaj et al. (2020).

278 PSO:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = 3.66 C_v^{0.16} \left( \frac{d}{R} \right)^{-0.40} \left( \frac{y_s}{Y} \right)^{-0.10} \quad (11)$$

279 LASSO:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = 5.83 C_v^{0.144} \left( \frac{d}{R} \right)^{-0.305} \lambda^{-0.059} D_{gr}^{-0.169} \left( \frac{y_s}{D} \right)^{-0.104} \quad (12)$$

280 MGP:

$$\frac{V_l}{\sqrt{gd(S_s - 1)}} = 1.96 - 0.61\lambda - 0.51C_v + 1.18D_{gr}^{0.50}\lambda^{1.50} + 0.61\left(2C_v + \frac{d}{R}\right)^{0.50} - 2.45\left(\frac{d}{R}\right)^{1/8} \quad (13)$$

### 281 3.2.2. Performance Indices

282 The RF model performance is evaluated and compared to above ten models using  
 283 three performance indicators. These are the Coefficient of Determination ( $R^2$ ), the Root  
 284 Mean Square Error ( $RMSE$ ) and the Mean Absolute Percentage Error ( $MAPE$ ), defined  
 285 as follows:

$$R^2 = 1 - \frac{\sum_{i=1}^n (F_{r_{OBS}}^* - F_{r_{MOD}})^2}{\sum_{i=1}^n (F_{r_{OBS}}^* - \overline{F_{r_{OBS}}^*})^2} \quad (14)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (F_{r_{OBS}}^* - F_{r_{MOD}})^2} \quad (15)$$

$$MAPE = \frac{100}{n} \sum_{i=1}^n \left| \frac{F_{r_{OBS}}^* - F_{r_{MOD}}}{F_{r_{OBS}}^*} \right| \quad (16)$$

286 where  $F_{r_{OBS}}^*$  is the particle Froude number observed data,  $F_{r_{MOD}}$  the particle Froude  
 287 number estimated by RF algorithm (or other predictive model),  $n$  the number of data and  
 288  $\overline{F_{r_{OBS}}^*}$  the mean of observed particle Froude number data.

289 The Coefficient of Determination measures the percentage of the model variance  
 290 that can be explained. This coefficient varies between 0 and 1, with a value of 1 denoting  
 291 a perfect match between observed and modelled data. The Root Mean Square Error  
 292 measures the standard deviation of the residuals. Note that a value close to 0 indicates  
 293 high model prediction accuracy. Finally, the Mean Absolute Percentage Error assesses  
 294 the model prediction accuracy (i.e. bias) as a percentage of the observed value. Value of  
 295 0 indicates the perfect model where there are no differences between predictions and

296 observations.

#### 297 4. RESULTS

298 The results obtained by using the methodology shown in the previous section are  
299 presented in Table 3 and Table 4, for without deposited bed and deposited bed criteria,  
300 respectively. Graphically, these results are shown in Figure 5 and Figure 6. As shown in  
301 these tables, for the MARS, ANFIS-PSO, ELM and MGP models, the outliers of the  
302 particle Froude number (i.e.  $F_r^* < 0.00$  and  $F_r^* > 20.00$ ) were removed. This is because  
303 these models can produce extreme values (e.g.  $F_r^* = -58.67$  or  $F_r^* = 163.59$ , among  
304 others) that misrepresent the model comparison when evaluating the performance indices.

#### 305 [Table 3 near here]

306 As it can be seen from Table 3, Random Forest model shows a better  
307 generalisation capacity than other models shown, as demonstrated in high prediction  
308 accuracy observed for all available datasets ( $0.88 > R^2 > 0.98$ ,  $0.24 > RMSE > 0.73$  and  
309  $4.36\% > MAPE > 11.09\%$ ). The following observations can be made from the  
310 performance of the other models evaluated:

- 311 • EPR-MOGA, similarly to RF, shows good results but has inferior accuracy in  
312 large sewer pipes ( $R^2 = 0.86$ ,  $RMSE = 1.03$  and  $MAPE = 11.31\%$ ). In addition,  
313 EPR-MOGA model shows limitations for predicting the particle Froude number  
314 in non-circular sections (as shown in the Mayerle (1988) rectangular data). This  
315 equation shows good extrapolation capabilities because of the inclusion of the  
316 pipe slope as input feature for the self-cleansing prediction.
- 317 • GEP shows acceptable results ( $0.79 > R^2 > 0.87$ ,  $0.66 > RMSE > 0.89$  and  $11.45\%$   
318  $> MAPE > 22.33\%$ ) for the datasets used for its development in circular channels  
319 (Ab Ghani, 1993; Mayerle, 1988; Vongvisessomjai et al., 2010) and poor



320 performance for other datasets ( $0.00 > R^2 > 0.76$ ,  $1.00 > RMSE > 1.95$  and  $14.35\%$   
321  $> MAPE > 37.92\%$ ). This model presents good performance for large sewer pipes.  
322 In contrast, for non-circular channels the model quickly loss accuracy.

- 323 • According to Safari (2019), MARS model was developed by using the  
324 experimental data collected by Mayerle (1988) (in both circular and rectangular  
325 channels), May (1993), Ab Ghani (1993) and Vongvisessomjai et al. (2010). As a  
326 result, this model shows acceptable performance for these datasets ( $0.49 > R^2 >$   
327  $0.87$ ,  $0.81 > RMSE > 1.15$  and  $13.63\% > MAPE > 28.08\%$ ) but poor performance  
328 for the remaining datasets ( $R^2 = 0.00$ ,  $1.48 > RMSE > 2.88$  and  $29.14\% > MAPE$   
329  $> 51.28\%$ ). Based on the above, and compared to the RF model, limited  
330 extrapolation capabilities are identified for the MARS model.
- 331 • May et al. (1996) is the best regression-based equation reported in the literature  
332 (Ackers et al., 2001; Ebtahaj et al., 2014), as it was developed using several  
333 experimental datasets. This is the equation proposed by the Construction Industry  
334 Research and Information Association (CIRIA) for designing self-cleansing  
335 sewer pipes transporting coarser granular material as bedload (Ackers et al.,  
336 2001). This model shows good performance for pipe diameters less than 500 mm  
337 ( $0.83 > R^2 > 0.99$ ,  $0.13 > RMSE > 0.82$  and  $2.38\% > MAPE > 11.61\%$ ). In  
338 contrast, limited extrapolation for large sewer pipes is identified as the low  
339 performance indices values obtained ( $R^2 = 0.00$ ,  $RMSE = 4.88$  and  $MAPE =$   
340  $48.97\%$ ). This equation shows better performance than the RF model when  
341 compared to data from Vongvisessomjai et al. (2010), but lower accuracy when  
342 applied to the rest of the datasets.
- 343 • Safari and Aksoy (2020) model is a competitive equation for predicting the self-  
344 cleansing velocity in both circular and non-circular channels. This model shows

345 similar but inferior performance to EPR-MOGA model in small sewer pipes ( $0.67$   
346  $> R^2 > 0.97$ ,  $0.25 > RMSE > 1.12$  and  $7.90\% > MAPE > 15.60\%$ ), but in large  
347 sewers the accuracy is quickly lost ( $R^2 = 0.34$ ,  $RMSE = 2.26$  and  $MAPE =$   
348  $23.46\%$ ). In contrast, this model outperforms the results, compared to other  
349 regression models (EPR-MOGA, GEP and MARS) and ML/AI models (ANFIS-  
350 PSO and ELM), in non-circular channels ( $R^2 = 0.87$ ,  $RMSE = 0.66$  and  $MAPE =$   
351  $13.41\%$ ), which is a competitive performance compared to the RF model ( $R^2 =$   
352  $0.89$ ,  $RMSE = 0.61$  and  $MAPE = 10.05\%$ ). This is because of the inclusion of the  
353  $P/B$  relation as explanatory variable for predicting the particle Froude number.  
354 This model is competitive and shows good generalisation of the problem for  
355 designing sewers under the non-deposition without deposited bed criterion.

356 • According to Ebtehaj et al. (2019), ANFIS-PSO model was developed by using  
357 the experimental data collected by Ab Ghani (1993), Ota (1999) and  
358 Vongvisessomjai et al. (2010). As a result, this model shows good performance  
359 for these datasets ( $0.88 > R^2 > 0.97$ ,  $0.22 > RMSE > 0.74$  and  $3.62\% > MAPE >$   
360  $10.34\%$ ). In large sewers and non-circular channels, the model losses accuracy  
361 ( $R^2 = 0.00$ ,  $2.74 > RMSE > 3.01$  and  $30.56\% > MAPE > 45.28\%$ ). This model  
362 produces some extreme values when the particle Froude number is calculated,  
363 especially in the Montes et al. (2020b) dataset. The RF model generates better  
364 results compared to this model.

365 • ELM was trained with the same dataset used for the ANFIS-PSO model. Not  
366 satisfactory results are obtained when this model is applied on the dataset  
367 considered in this study ( $0.00 > R^2 > 0.55$ ,  $0.90 > RMSE > 3.1$  and  $19.54\% >$   
368  $MAPE > 39.30\%$ ). Same comments, as mentioned above for the ANFIS-PSO  
369 model, can be shown here.

370

[Figure 5 near here]

371

[Table 4 near here]

372

373

374

375

376

According to the results shown in Table 4 (deposited bed criterion), RF model outperforms the other models for the entire considered dataset. This model shows good accuracy levels ( $0.84 > R^2 > 0.98$ ,  $0.32 > RMSE > 0.81$  and  $4.70\% > MAPE > 12.10\%$ ) for all the range of variation of the hydraulics and sediment characteristics. Comments related to the other models studied are as follows:

377

378

379

380

381

382

- PSO model was developed by using the experimental data collected by El-Zaemey (1991), Perrusquía (1991), May (1993) and Ab Ghani (1993). As a result, this model shows good performance for these datasets ( $0.56 > R^2 > 0.78$ ,  $0.49 > RMSE > 1.32$  and  $10.15\% > MAPE > 16.26\%$ ). However, when the model is compared to the data collected in the large sewer pipe, the accuracy quickly decreases ( $R^2 = 0.00$ ,  $RMSE = 3.06$  and  $MAPE = 21.05\%$ ).

383

384

385

386

- LASSO model reports good accuracy levels for all the datasets considered ( $0.62 > R^2 > 0.83$ ,  $0.50 > RMSE > 1.56$  and  $10.36\% > MAPE > 14.26\%$ ). However, the accuracy is still inferior compared to the RF model. This model shows good extrapolation capabilities and generalisation of the problem.

387

388

389

390

391

392

393

- MGP was developed by using the same experimental datasets of the PSO model. This model shows less accuracy compared to the PSO model ( $0.00 > R^2 > 0.54$ ,  $1.08 > RMSE > 5.54$  and  $13.07\% > MAPE > 58.79\%$ ). In large sewer pipes, the model shows poor performance. In contrast to other models, the MGP was developed by using normalised values. Based on this, the range of variation used for training the model can potentially affect the final form/structure of the final expression shown by the MGP.

394 **[Figure 6 near here]**

395 RF accuracy shown in the Montes et al. (2020b) data is especially important due  
396 to the relative sediment thickness ( $y_s/D$ ) used at laboratory scale in that study. As Table  
397 1 shows, the sediment thickness used at laboratory scale ranging from 0.8 mm (for Montes  
398 et al. (2020b) data) to 129.6 mm (for May (1993) data), i.e. the variation of  $y_s/D$  is from  
399 1.1% to 20.0% of the pipe diameter. Values of  $y_s/D = 20%$  is an unrealistic consideration  
400 since the optimal sediment thickness design has been defined as 1% of the pipe diameter  
401 (May et al., 1989; Safari and Shirzad, 2019). Data collected by Montes et al. (2020b)  
402 seem to be the closer representation of the real conditions found in sewer systems. Based  
403 on this, RF is the model that best predicts the self-cleansing velocity for data close to real  
404 conditions.

405 **4.1. Variable importance**

406 RF model input variable importance is presented in Figure 7. As shown in this figure, for  
407 both non-deposition criteria the most important variable is the volumetric sediment  
408 concentration, followed by the dimensionless grain size and the relative grain size . This  
409 result is consistent with previous findings reported in the literature (Ackers et al., 2001;  
410 Ebtehaj et al., 2020). Less important parameters for predicting the particle Froude number  
411 and thus the self-cleansing velocity, are the relative sediment thickness and the channel  
412 friction factor, for the deposited bed criterion.

413 Parameter importance shown by EPR-MOGA, Safari and Aksoy (2020), PSO and  
414 LASSO is quite different. In these techniques, the most important parameter is the relative  
415 grain size due to the highest values of the regression coefficients  $\left(\left(\frac{d}{R}\right)^{-c}\right)$ ;  $0.305 < c <$   
416  $0.58$ ), as shown in Eq. (5), Eq. (9), Eq. (11) and Eq. (12). The parameter importance for  
417 the GEP, MARS and MGP model is less intuitive because of the form of the equations,

418 as shown in Eq. (6), Eq. (7) and Eq. (13), which include logarithmic and inverse tangent  
419 functions for calculating the particle Froude number. Less comparable are the results  
420 shown by ANFIS-PSO and ELM since no practical equation is provided.

421 **[Figure 7 near here]**

422 Based on the above results shown in Figure 7, a good estimate of the volumetric  
423 sediment concentration seems to be essential for increasing the accuracy of the calculation  
424 of the particle Froude number and consequently the minimum self-cleansing velocity for  
425 both non-deposition criteria. In addition, hydraulic characteristics of the pipe (defined by  
426 the hydraulic radius) and the sediment characteristics (i.e. particle diameter and specific  
427 gravity) are proportionally important for model performance.

## 428 **5. DISCUSSION**

429 The prediction of self-cleansing conditions in sewers remains a challenge despite multiple  
430 models and equations developed and reported in the literature. Existing regression-based  
431 equations and AI/ML models show limited generalisation capabilities and overfitting  
432 problems. In this paper, a new approach for addressing these issues is proposed by using  
433 the Random Forest method.

434 Due to the nature of the RF method, where the model variance is reduced by  
435 averaging the results from an ensemble of decision trees, the risk of overfitting is low. By  
436 using a reduced number of input features for constructing each decision tree in the forest,  
437 the correlation between base trees is avoided. This is an improvement of the method  
438 compared to a single decision tree, which can be overtrained (i.e. the tree learns the noise  
439 from the training data) and thus shows poor performance in the testing dataset.

440 RF model showed good generalisation capabilities when the whole dataset is  
441 divided into 75% for the training stage and 25% for the testing stage. For this percentage  
442 of split data, the testing error presented a low variance. In contrast, by increasing the

443 number of data used in the training stage (e.g. 95% of the whole data) the testing error  
444 showed high variance, which is an indicator of an over-trained model with limited  
445 extrapolation capabilities (as shown in Figure 3b). Therefore, choosing the right  
446 percentage split is critical to avoid model overfitting.

447 Variable importance analysis showed that the volumetric sediment concentration  
448 is the most relevant feature for predicting the self-cleansing velocity in practice for both  
449 non-deposition criteria, followed by the dimensionless grain size. The self-cleansing  
450 prediction is no conditioned by the channel material, as the low variable importance  
451 shown by the channel friction factor.

452 RF results are compared to existing models reported in the literature and showed  
453 better performance for the whole dataset for both non-deposition without and with  
454 deposited bed criteria. This is explained by several factors, such as:

- 455 • RF is able to better capture the non-linearity in the data compared to linear  
456 regression models (i.e. regression-based models proposed by May et al. (1996)  
457 and Safari and Aksory (2020)). The RF model also better captures complex  
458 interactions between features. This is because of RF model's ability to capture  
459 effectively non-linear patterns in data.
- 460 • RF showed a good bias-variance trade-off (i.e. low bias and low variance) for both  
461 non-deposition criteria. In contrast, existing non-regression models reported in the  
462 literature (i.e. MARS, ANFIS-PSO and ELM), and compared to the RF model in  
463 this paper, in some cases presented low bias and high variance (i.e. overfitting)  
464 for the non-deposition without deposited bed criterion, as shown in Figure 5. For  
465 the non-deposition with deposited bed criterion, the existing models (i.e. PSO,  
466 LASSO and MGP) showed high bias, since these models systematically

467 underestimate the particle Froude number in the testing dataset (as shown in  
468 Figure 6).

469 • The range of variation used for training and testing the RF model is much larger  
470 than the dataset used in the literature for developing the existing predictive  
471 models. For example, the ANFIS-PSO and ELM were trained and testing with the  
472 Ab Ghani (1993), Ota (1999) and Vongvisessomjai et al. (2010) data (i.e. 290 data  
473 approx.). Given this, the RF model developed here is able to predict the particle  
474 Froude number for a larger range of variation of the input conditions. An example  
475 of this is shown in Figure 6 where the existing models reported for the non-  
476 deposition with deposited bed criterion underestimate the particle Froude number  
477 for values above 9.0 ( $F_r^* > 9.0$ ).

478 Despite the RF presented in this study outperforms the existing models reported  
479 in the literature, further tests with data collected in real sewers should be conducted. The  
480 cohesive effects of the deposited material must be included for future developments.  
481 Finally, further evaluation of the performance of the model in trapezoidal, ovoid, or U-  
482 shape channels should be carried out to check the applicability of the model under these  
483 channel characteristics.

## 484 **6. CONCLUSIONS**

485 Random Forest based model was developed for predicting the self-cleansing velocity  
486 under the concept of non-deposition. This model was implemented using the experimental  
487 benchmark data reported in the literature. The RF model was compared to the following  
488 ten literature models: EPR-MOGA, MARS, MGP, ANFIS-PSO, ELM, LASSO, GEP and  
489 PSO, and two regression-based equations proposed by May et al. (1996) and Safari and  
490 Aksoy (2020).

491 The following conclusions are made based on the results obtained:

492 (1) Random Forest model is able to predict the particle Froude number (i.e. minimum  
493 self-cleansing velocity) for the non-deposition self-cleansing design criteria with  
494 high accuracy on validation (i.e. unseen) data. This is due to the ability of RF to  
495 better generalise the analysed data, i.e. the ability to avoid model overfitting.

496 (2) RF model prediction accuracy is consistently superior to ten other literature  
497 models considered here. This is likely due to the reason mentioned above but also  
498 the capability to better capture the complex interactions between input variables  
499 when compared to other models considered in this paper. This is especially  
500 relevant for the non-deposition with deposited bed case where the accuracy of RF  
501 model predictions is substantially higher than in other models (i.e. LASSO, MGP  
502 and PSO models).

503 (3) The volumetric sediment concentration is the most important input variable for  
504 predicting the self-cleansing velocity in sewer pipes. A good characterisation of  
505 this parameter seems to be essential for improving the design of new self-  
506 cleansing sewers.

507 Based on the above, RF can be used for predicting self-cleansing velocity with  
508 high accuracy, especially for large sewer pipes with the presence of deposited bed. This  
509 technique can be used for designing self-cleansing sewer systems.

510 Further testing of the RF and other self-cleansing models in real sewer systems is  
511 required to further validate these models in those circumstances and ensure their  
512 applicability in engineering practice.

## 513 **7. SUPPLEMENTARY MATERIAL**

514 Data used for training and testing the Random Forest method is shown in Table S1 and



515 Table S2 for non-deposition without and with deposited bed, respectively. In addition, an  
516 example of one of the decision trees considered by the RF method is shown in Figure S1.

## 517 **FUNDING**

518 This research did not receive any specific grant from funding agencies in the public,  
519 commercial, or not-for-profit sectors.

## 520 **REFERENCES**

- 521 Ab Ghani, A., 1993. Sediment Transport in Sewers. PhD thesis, University of Newcastle  
522 upon Tyne, Newcastle upon Tyne, UK.
- 523 Ackers, J., Butler, D., Leggett, D., May, R., 2001. Designing Sewers to Control Sediment  
524 Problems, in: Urban Drainage Modeling. ASCE, Orlando, FL, pp. 818–823.  
525 [https://doi.org/10.1061/40583\(275\)77](https://doi.org/10.1061/40583(275)77)
- 526 Breiman, L., 2001. Random Forests. Mach. Learn. 45, 5–32.  
527 <https://doi.org/10.1023/A:1010933404324>
- 528 Ebtehaj, I., Bonakdari, H., 2016a. Bed Load Sediment Transport in Sewers at Limit of  
529 Deposition. Sci. Iran. 23 (3), 907–917. <https://doi.org/10.24200/sci.2016.2169>
- 530 Ebtehaj, I., Bonakdari, H., 2016b. A support vector regression-firefly algorithm-based  
531 model for limiting velocity prediction in sewer pipes. Water Sci. Technol. 73 (9),  
532 2244–2250. <https://doi.org/10.2166/wst.2016.064>
- 533 Ebtehaj, I., Bonakdari, H., 2013. Evaluation of sediment transport in sewer using artificial  
534 neural network. Eng. Appl. Comput. Fluid Mech. 7 (3), 382–392.  
535 <https://doi.org/10.1080/19942060.2013.111015479>
- 536 Ebtehaj, I., Bonakdari, H., Es-Haghi, M., 2019. Design of a Hybrid ANFIS–PSO Model  
537 to Estimate Sediment Transport in Open Channels. Iran. J. Sci. Technol. Trans.  
538 44 (4), 851–857. <https://doi.org/10.1007/s40996-018-0218-9>
- 539 Ebtehaj, I., Bonakdari, H., Safari, M., Gharabaghi, B., Zaji, A., Riahi Madavar, H., Sheikh  
540 Khozani, Z., Es-haghi, M., Shishegaran, A., Danandeh Mehr, A., 2020.  
541 Combination of sensitivity and uncertainty analyses for sediment transport  
542 modeling in sewer pipes. Int. J. Sediment Res. 35 (2), 157–170.  
543 <https://doi.org/10.1016/j.ijsrc.2019.08.005>

544 Ebtehaj, I., Bonakdari, H., Sharifi, A., 2014. Design criteria for sediment transport in  
545 sewers based on self-cleansing concept. *J. Zhejiang Univ. Sci. A* 15 (11), 914-  
546 924. <https://doi.org/10.1631/jzus.a1300135>

547 El-Zaemey A., 1991. Sediment Transport over Deposited Beds in Sewers. PhD thesis,  
548 University of Newcastle upon Tyne, Newcastle upon Tyne, UK.

549 Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical Learning: Data*  
550 *Mining, Inference, and Prediction*. Springer, New York, USA.  
551 <https://doi.org/10.1007/978-0-387-84858-7>

552 Kargar, K., Safari, M., Mohammadi, M., Samadianfard, S., 2019. Sediment transport  
553 modeling in open channels using neuro-fuzzy and gene expression programming  
554 techniques. *Water Sci. Technol.* 79 (12), 2318–2327.  
555 <https://doi.org/10.2166/wst.2019.229>

556 Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest. *R News* 2  
557 (3), 18–22.

558 May R., 1993. Sediment Transport in Pipes and Sewers with Deposited Beds. Report SR  
559 320, HR Wallingford, Oxfordshire, UK.

560 May R., Ackers, J., Butler, D., John, S., 1996. Development of design methodology for  
561 self-cleansing sewers. *Water Sci. Technol.* 33 (9), 195–205.  
562 [https://doi.org/10.1016/0273-1223\(96\)00387-3](https://doi.org/10.1016/0273-1223(96)00387-3)

563 May R., Brown P., Hare G., Jones K., 1989. Self-Cleansing Conditions for Sewers  
564 Carrying Sediment. Report SR 221, HR Wallingford, Oxfordshire, UK.

565 Mayerle R., 1988. Sediment Transport in Rigid Boundary Channels. PhD thesis,  
566 University of Newcastle upon Tyne, Newcastle upon Tyne, UK.

567 Merritt, L., Enfinger, K., 2019. Tractive Force: A Key to Solids Transport in Gravity  
568 Flow Drainage Pipes, in: *Pipelines 2019*. ASCE, Nashville, TN, pp. 349–358.

569 Montes, C., Berardi, L., Kapelan, Z., Saldarriaga, J., 2020a. Predicting bedload sediment  
570 transport of non-cohesive material in sewer pipes using evolutionary polynomial  
571 regression – multi-objective genetic algorithm strategy. *Urban Water J.* 17 (2),  
572 154–162. <https://doi.org/10.1080/1573062X.2020.1748210>

573 Montes, C., Kapelan, Z., Saldarriaga, J., 2019. Impact of Self-Cleansing Criteria Choice  
574 on the Optimal Design of Sewer Networks in South America. *Water* 11 (6), 1148.  
575 <https://doi.org/10.3390/w11061148>

576 Montes, C., Vanegas, S., Kapelan, Z., Berardi, L., Saldarriaga, J., 2020b. Non-deposition  
577 self-cleansing models for large sewer pipes. *Water Sci. Technol.* 81 (3), 606-621.  
578 <https://doi.org/10.2166/wst.2020.154>

579 Nalluri, C., Ab Ghani, A., 1996. Design options for self-cleansing storm sewers. *Water*  
580 *Sci. Technol.* 33 (9), 215–220. [https://doi.org/10.1016/0273-1223\(96\)00389-7](https://doi.org/10.1016/0273-1223(96)00389-7)

581 Ota J., 1999. Effect of Particle Size and Gradation on Sediment Transport in Storm  
582 Sewers. PhD thesis, University of Newcastle upon Tyne, Newcastle upon Tyne,  
583 UK.

584 Perrusquía, G., 1991. Bedload Transport in Storm Sewers: Stream Traction in Pipe  
585 Channels. PhD thesis, Chalmers University of Technology, Gothenburg, Sweden.

586 Roushangar, K., Ghasempour, R., 2017. Estimation of bedload discharge in sewer pipes  
587 with different boundary conditions using an evolutionary algorithm. *Int. J.*  
588 *Sediment Res.* 32 (4), 564–574. <https://doi.org/10.1016/j.ijsrc.2017.05.007>

589 Safari, M., 2019. Decision tree (DT), generalized regression neural network (GR) and  
590 multivariate adaptive regression splines (MARS) models for sediment transport  
591 in sewer pipes. *Water Sci. Technol.* 79 (6), 1113-1122.  
592 <https://doi.org/10.2166/wst.2019.106>

593 Safari, M., Danandeh Mehr, A., 2018. Multigene genetic programming for sediment  
594 transport modeling in sewers for conditions of non-deposition with a bed deposit.  
595 *Int. J. Sediment Res.* 33 (3), 262-270. <https://doi.org/10.1016/j.ijsrc.2018.04.007>

596 Safari, M., Mohammadi, M., Ab Ghani, A., 2018. Experimental Studies of Self-Cleansing  
597 Drainage System Design: A Review. *J. Pipeline Syst. Eng. Pract.* 9 (4), 04018017.  
598 [https://doi.org/10.1061/\(ASCE\)PS.1949-1204.0000335](https://doi.org/10.1061/(ASCE)PS.1949-1204.0000335)

599 Safari, M., Shirzad, A., 2019. Self-cleansing design of sewers: Definition of the optimum  
600 deposited bed thickness. *Water Environ. Res.* 91 (5), 407–416.  
601 <https://doi.org/10.1002/wer.1037>

602 Safari, M., Shirzad, A., Mohammadi, M., 2017. Sediment transport modeling in deposited  
603 bed sewers: Unified form of May's equations using the particle swarm  
604 optimization algorithm. *Water Sci. Technol.* 76 (4), 992–1000.  
605 <https://doi.org/10.2166/wst.2017.267>

606 Safari, M., 2020. Hybridization of multivariate adaptive regression splines and random  
607 forest models with an empirical equation for sediment deposition prediction in

608 open channel flow. *J. Hydrol.* 590 (November 2020), 125392.  
609 <https://doi.org/10.1016/j.jhydrol.2020.125392>

610 Safari, M., Aksoy, H., 2020. Experimental analysis for self-cleansing open channel  
611 design. *J. Hydraul. Res.* 1-12. <https://doi.org/10.1080/00221686.2020.1780501>

612 Tyrallis, H., Papacharalampous, G., & Langousis, A. 2019 A Brief Review of Random  
613 Forests for Water Scientists and Practitioners and Their Recent History in Water  
614 Resources. *Water*, 11(5), 910. <https://doi.org/10.3390/w11050910>

615 Vongvisessomjai, N., Tingsanchali, T., & Babel, M. 2010 Non-deposition design criteria  
616 for sewers with part-full flow. *Urban Water Journal*, 7(1), 61–77.  
617 <https://doi.org/10.1080/15730620903242824>

618 Zendehboudi, S., Rezaei, N., & Lohi, A. 2018 Applications of hybrid models in chemical,  
619 petroleum, and energy systems: A systematic review. *Applied Energy*, 228(2018),  
620 2539–2566. <https://doi.org/10.1016/j.apenergy.2018.06.051>

621 Table 1. Data used for implementing data mining and regression models.

Reference	Non-deposition criterion	No. of runs	Pipe diameter or bottom width (mm)	Flow Velocity (m/s)	Pipe slope (%)	Sediment Concentration (ppm)	Sediment thickness bed (mm)
Mayerle (1988) circular channel	Without deposited bed	106	152	0.37 - 1.10	0.13 - 0.56	20.0 - 1275.0	-
Mayerle (1988) rectangular channel	Without deposited bed	105	311.5 and 462.3	0.41 - 1.04	0.09 - 0.64	14.0 - 1568.0	-
Ab Ghani (1993)	Without deposited bed	221	154, 305 and 405	0.24 - 1.22	0.04 - 2.56	0.8 - 1450.0	-
Ota (1999)	Without deposited bed	36	305	0.39 - 0.74	0.2	4.2 - 59.4	-
Vongvisessomjai et al. (2010)	Without deposited bed	45	100 and 150	0.24 - 0.63	0.20 - 0.60	4.0 - 90.0	-
Montes et al. (2020a)	Without deposited bed	44	242	0.24 - 1.05	0.20 - 0.80	0.3 - 875.7	-
Montes et al. (2020b)	Without deposited bed	107	595	0.41 - 1.41	0.04 - 3.43	1.3 - 19957.0	-
El-Zaemey (1991)	With deposited bed	290	305	0.39 - 0.96	0.05 - 0.44	7.0 - 917.0	47.0 - 120.0
Perrusquía (1991)	With deposited bed	38	225	0.29 - 0.67	0.20 - 0.60	18.7 - 408.0	45.0 - 90.0
Ab Ghani (1993)	With deposited bed	26	450	0.49 - 1.33	0.07 - 0.47	21.0 - 1259.0	52.0 - 108.0
May (1993)	With deposited bed	46	450	0.39 - 1.14	0.07 - 0.97	3.5 - 823.0	57.6 - 129.6
Montes et al. (2020b)	With deposited bed	54	595	0.73 - 1.53	0.46 - 5.42	389.0 - 10275.0	0.8 - 6.6

622

623 Table 2. Variation of the data for training and testing the RF model.

Non-deposition criterion	Stage	No. of runs	Channel geometry (mm)	Flow Velocity (m/s)	Pipe slope (%)	Sediment Concentration (ppm)	Sediment thickness bed (mm)
Without deposited bed	Training	498	$D = 100.0 - 595.0$ $W = 311.5 - 462.3$	0.237 - 1.41	0.04 - 3.43	0.53 - 19957	-
	Testing	166	$D = 100.0 - 595.0$ $W = 311.5 - 462.3$	0.237 - 1.24	0.04 - 2.74	1.00 - 13840	-
With deposited bed	Training	340	$D = 225 - 595$	0.294 - 1.53	0.05 - 5.42	3.50 - 10274	0.78 - 129.6
	Testing	114	$D = 225 - 595$	0.319 - 1.28	0.05 - 2.58	17.00 - 9101	1.78 - 120.0

624

625 Table 3. Accuracy of self-cleansing models for without deposited bed criterion using  
626 performance indices for training and testing dataset. Bolded values show best  
627 performance model.

Dataset	Performance Index	Model							
		RF	EPR-MOGA	GEP	MARS	May et al. (1996) <sup>1</sup>	Safari and Aksoy (2020)	ANFIS-PSO	ELM
Training	$R^2$	<b>0.98</b>	0.90	0.75	0.00	0.27	0.74	0.51*	0.30*
	RMSE	<b>0.33</b>	0.76	1.22	2.55	2.17	1.25	1.69*	1.95*
	MAPE (%)	<b>4.88</b>	11.54	23.52	34.16	17.49	17.21	19.32*	29.76*
Testing	$R^2$	<b>0.91</b>	0.86	0.69	0.00	0.09	0.74	0.40*	0.32*
	RMSE	<b>0.73</b>	0.88	1.33	2.55	2.27	1.21	1.84*	1.92*
	MAPE (%)	<b>11.09</b>	12.35	26.43	36.57	19.15	17.24	20.95*	29.82*
Mayerle (1988) circular	$R^2$	<b>0.96</b>	0.89	0.87	0.87	0.87	0.75	0.80*	0.42
	RMSE	<b>0.45</b>	0.75	0.81	0.81	0.82	1.12	1.00*	1.71
	MAPE (%)	<b>5.62</b>	8.90	14.77	14.03	11.49	14.91	17.92*	26.75
Mayerle (1988) rectangular	$R^2$	<b>0.93</b>	0.38	0.30	0.81	-	0.87	0.00	0.47
	RMSE	<b>0.49</b>	1.44	1.54	0.81	-	0.66	2.74	1.33
	MAPE (%)	<b>8.49</b>	28.97	33.00	15.51	-	13.14	45.28	20.75
Ab Ghani (1993)	$R^2$	<b>0.97</b>	0.96	0.83	0.72	0.90	0.81	0.88	0.38
	RMSE	<b>0.36</b>	0.43	0.89	1.15	0.67	0.94	0.74	1.69
	MAPE (%)	<b>5.94</b>	9.35	22.33	28.08	10.32	15.60	10.34	23.96
Ota (1999)	$R^2$	0.97	<b>0.98</b>	0.44	0.00	0.96	0.97	0.97	0.55
	RMSE	0.24	<b>0.20</b>	1.00	1.48	0.27	0.25	0.22	0.90
	MAPE (%)	<b>5.55</b>	6.90	37.92	51.28	7.78	7.90	6.46	19.54
Vongvisessomjai et al. (2010)	$R^2$	0.88	0.95	0.79	0.49	<b>0.99</b>	0.71	0.97	0.00
	RMSE	0.49	0.33	0.66	1.03	<b>0.13</b>	0.78	0.24	1.59
	MAPE (%)	6.56	5.78	11.45	13.63	<b>2.38</b>	13.34	3.62	28.50
Montes et al. (2020a)	$R^2$	0.96	<b>0.98</b>	0.00	0.00	0.83	0.67	0.77*	0.00
	RMSE	0.31	<b>0.25</b>	1.64	2.37	0.67	0.94	0.75*	1.85
	MAPE (%)	<b>4.36</b>	4.94	28.15	49.73	11.61	15.39	12.39*	33.96
Montes et al. (2020b)	$R^2$	<b>0.94</b>	0.86	0.76	0.00*	0.00	0.34	0.00*	0.00*
	RMSE	<b>0.70</b>	1.03	1.37	2.88*	4.88	2.26	3.01*	3.10*
	MAPE (%)	<b>7.33</b>	11.31	14.35	29.14*	48.97	23.44	30.56*	39.30*

628 <sup>1</sup> Model not valid for non-circular channels

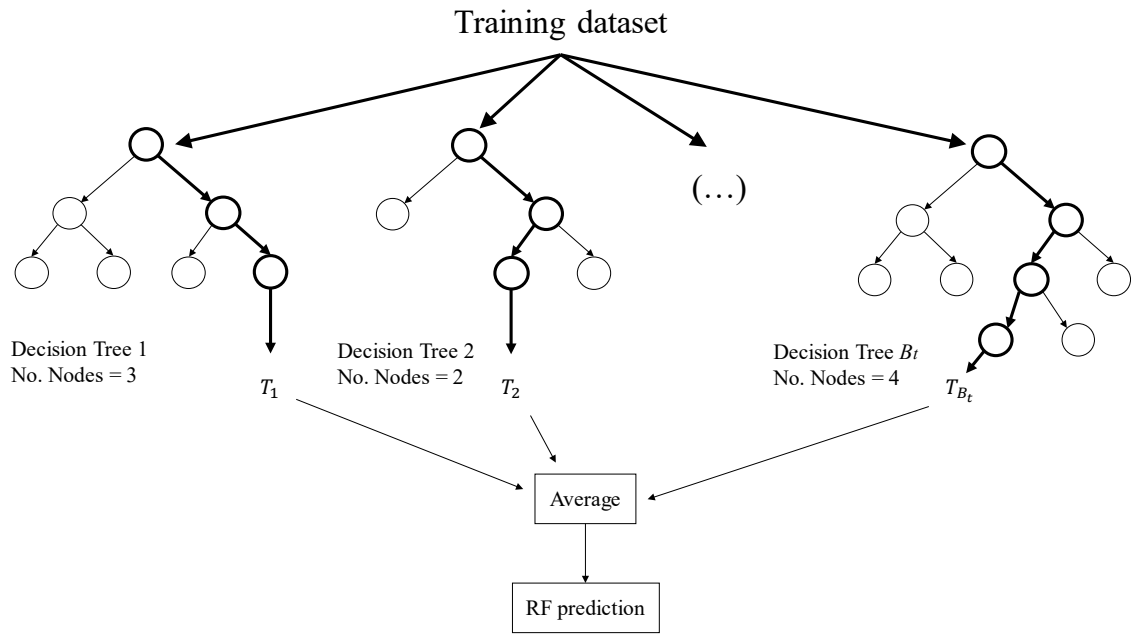
629 \* Outliers removed

630 Table 4. Accuracy of self-cleansing models for deposited bed criterion using performance  
 631 indices for training and testing dataset. Bolded values show best performance model.

Dataset	Performance Index	Model			
		RF	PSO	LASSO	MGP
Training	$R^2$	<b>0.98</b>	0.75	0.82	0.51*
	RMSE	<b>0.32</b>	1.30	1.13	1.69*
	MAPE (%)	<b>4.70</b>	14.36	13.07	28.78*
Testing	$R^2$	<b>0.91</b>	0.70	0.83	0.29*
	RMSE	<b>0.80</b>	1.47	1.10	2.19*
	MAPE (%)	<b>12.10</b>	15.94	12.59	31.36*
El-Zaemey (1991)	$R^2$	<b>0.94</b>	0.78	0.83	0.54
	RMSE	<b>0.38</b>	0.76	0.66	1.08
	MAPE (%)	<b>6.49</b>	14.28	11.97	30.19
Perrusquía (1991)	$R^2$	<b>0.84</b>	0.65	0.62	0.00
	RMSE	<b>0.33</b>	0.49	0.50	1.29
	MAPE (%)	<b>7.07</b>	10.15	12.05	30.58
Ab Ghani (1993)	$R^2$	<b>0.91</b>	0.56	0.74	0.51
	RMSE	<b>0.60</b>	1.32	1.01	1.40
	MAPE (%)	<b>6.13</b>	16.26	11.19	13.07
May (1993)	$R^2$	<b>0.90</b>	0.63	0.64	0.54
	RMSE	<b>0.62</b>	1.18	1.16	1.31
	MAPE (%)	<b>6.50</b>	13.47	14.26	14.21
Montes et al. (2020a)	$R^2$	<b>0.93</b>	0.00	0.73	0.00*
	RMSE	<b>0.81</b>	3.06	1.56	5.54*
	MAPE (%)	<b>6.84</b>	21.05	10.36	58.79*

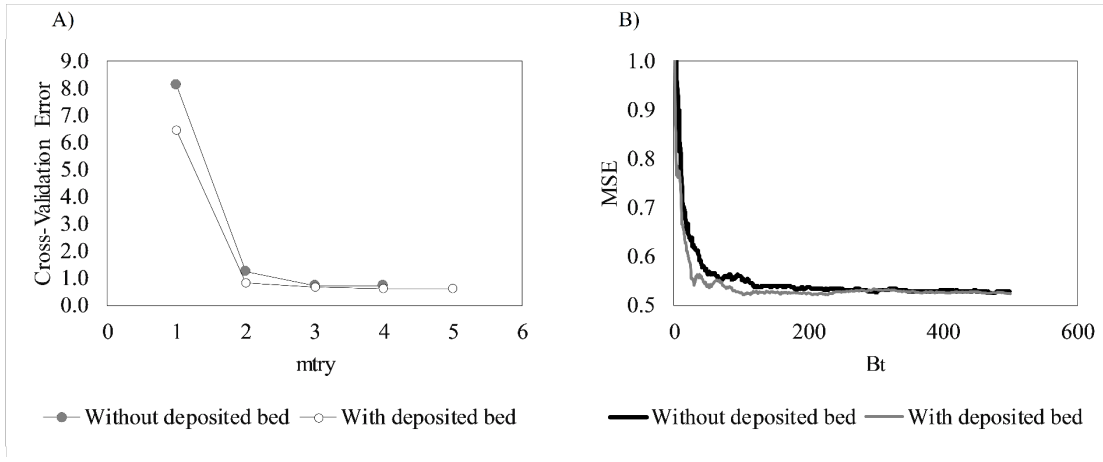
632 \* Outliers removed





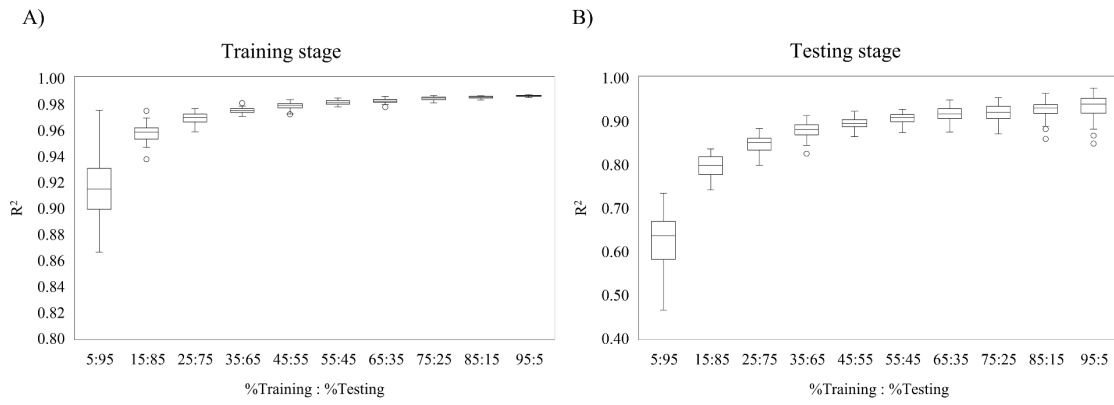
633

634 Figure 1. Simplified conceptual diagram of the RF method.



635

636 Figure 2. Selection of the optimal Random Forest parameters.



637

638 Figure 3. Variation of the training and testing error using different combination of  
 639 percentages between the training and testing dataset. A) Training stage and B) Testing  
 640 stage.

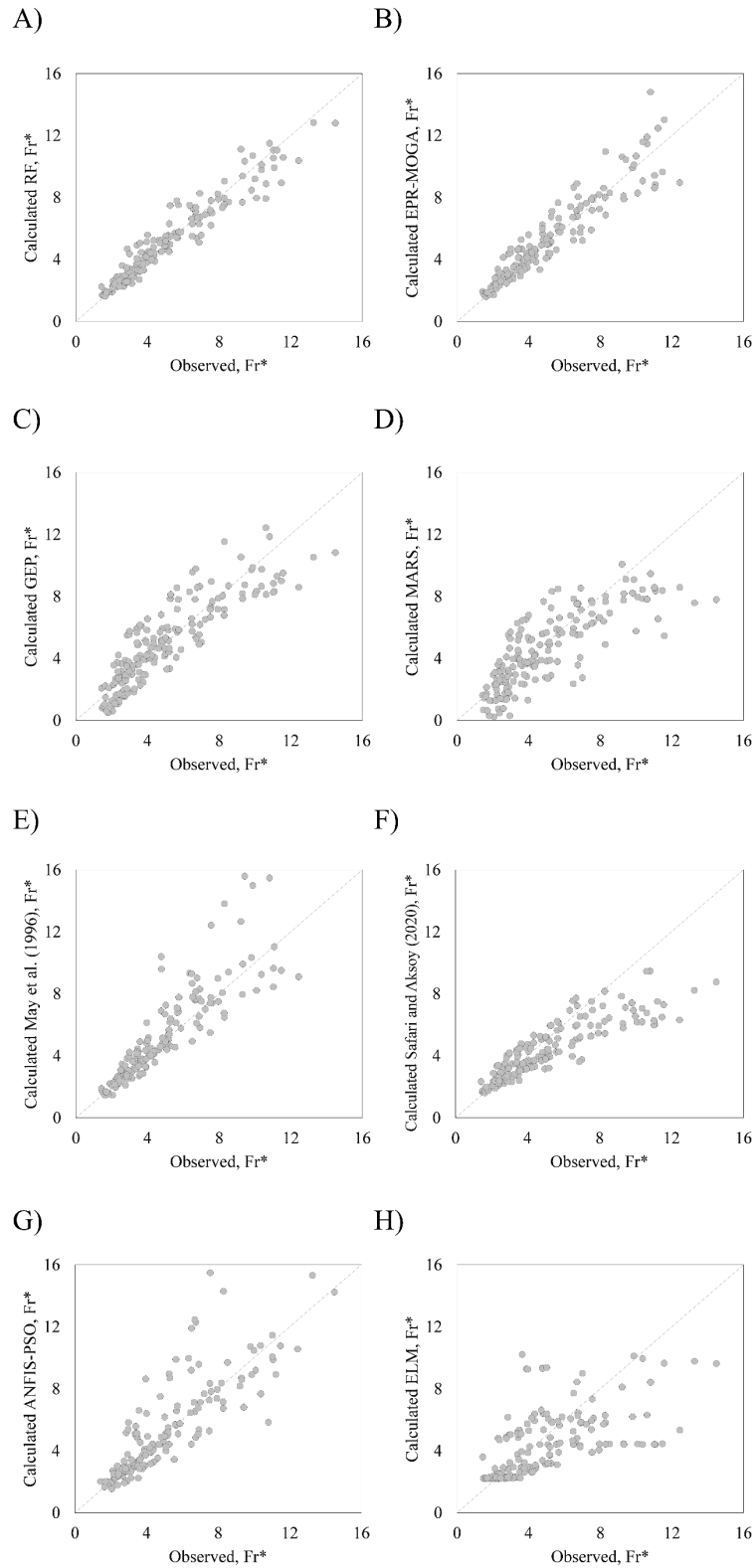
```

1 #####Random Forest model
2
3 library(randomForest)
4
5 #load data shown in Table S1 (without deposited bed)
6 #or Table S2 (with deposited bed)
7 #Please remove the "Fr* prediction RF" column
8
9 data=read.csv("Fulldata.csv",header=TRUE,sep=",")
10
11 set.seed(4260)
12
13 train=sample(1:nrow(data),nrow(data)*0.75)
14 test=data[-train,]
15 train=data[train,]
16
17 #Run Random Forest method
18 #Use ntree = 471 and mtry = 3 for without deposited bed
19 #Use ntree = 229 and mtry = 4 for with deposited bed
20
21 rf=randomForest(Frp~.,data=train,
22                 localImp=TRUE,importance=TRUE,
23                 mtry=3,ntree=471)
24
25 #RF Prediction in training and testing dataset
26
27 rf.pred.train=predict(rf,newdata=train)
28 rf.pred.test=predict(rf,newdata=test)
29
30 #Use function predict() to calculate the particle
31 #Froude number using other datasets. Use the same
32 #data frame headers.

```

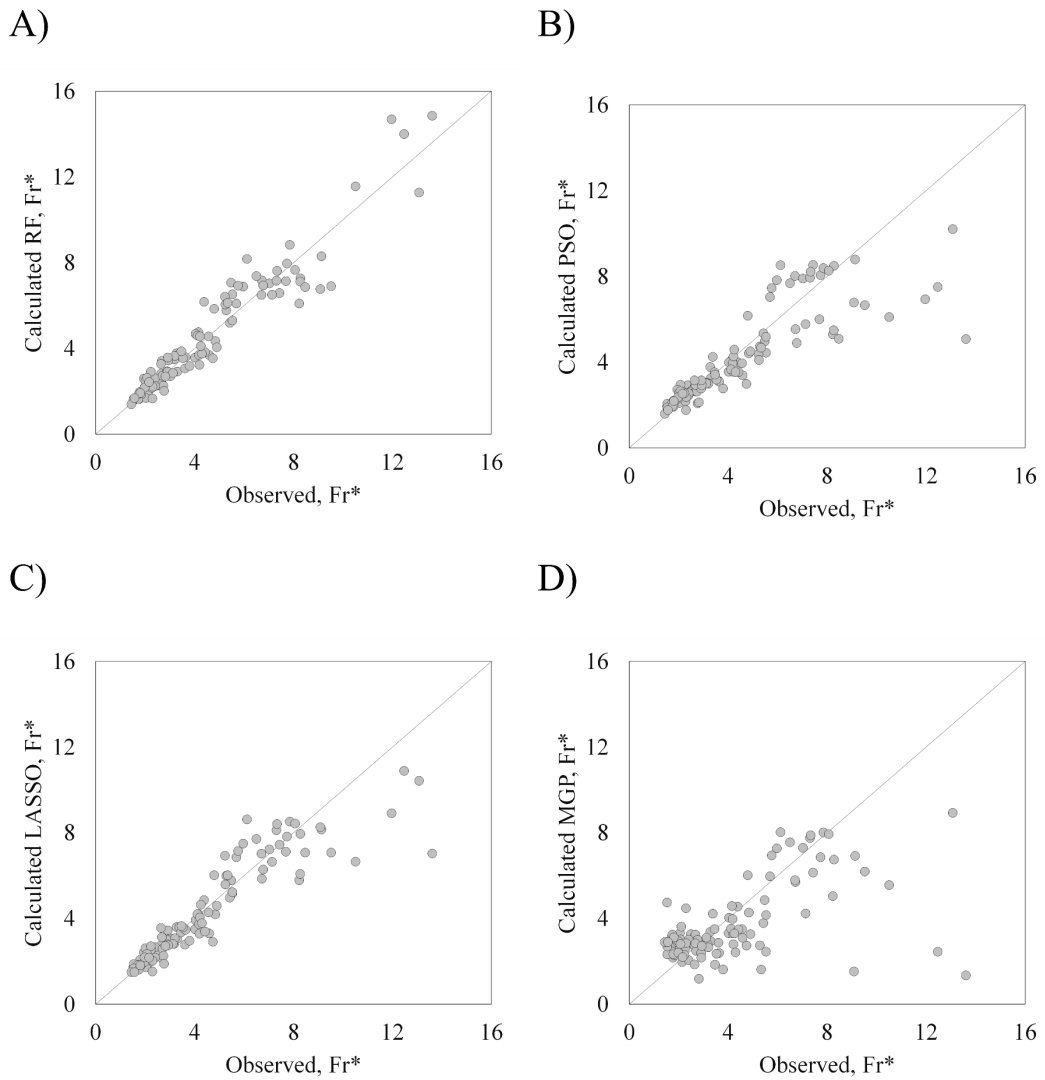
641

642 Figure 4. Random Forest code to calculate the particle Froude number in sewer pipes.



643

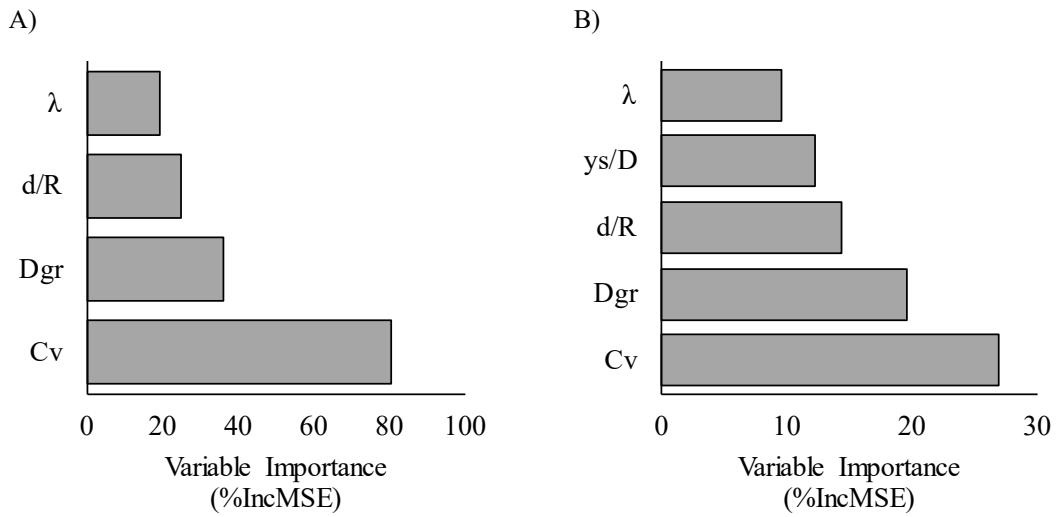
644 Figure 5. Performance of the models applied in the non-deposition without deposited bed  
 645 testing dataset.



646

647 Figure 6. Performance of the models applied in the non-deposition with deposited bed

648 testing dataset.



649

650 Figure 7. Variable importance estimated by RF model: A) without deposited bed; B) with  
 651 deposited bed.