

# Predicting stable gravel-bed river hydraulic geometry: A test of novel, advanced, hybrid data mining algorithms

Khabat Khosravi<sup>1</sup>, Zohreh Sheikh Khozani<sup>\*2</sup>, James R.Cooper<sup>3</sup>

1- Department of Watershed Management Engineering, Ferdowsi University of Mashhad, Mashhad, Iran. [Khabat.khosravi@gmail.com](mailto:Khabat.khosravi@gmail.com)

2- Institute of Structural Mechanics, Bauhaus Universität Weimar, 99423 Weimar, Germany. ([zohreh.khozani.sheikh@uni-weimar.de](mailto:zohreh.khozani.sheikh@uni-weimar.de))

3- Department of Geography and Planning, University of Liverpool, Liverpool, UK. [James.Cooper@liverpool.ac.uk](mailto:James.Cooper@liverpool.ac.uk)

\*Corresponding author: Zohreh Sheikh Khozani ([zohreh.khozani.sheikh@uni-weimar.de](mailto:zohreh.khozani.sheikh@uni-weimar.de))

## Abstract

Accurate prediction of stable alluvial hydraulic geometry, in which erosion and sedimentation are in equilibrium, is one of the most difficult but critical topics in the field of river engineering. Data mining algorithms have been gaining more attention in this field due to their high performance and flexibility. However, an understanding of the potential for these algorithms to provide fast, cheap, and accurate predictions of hydraulic geometry is lacking. This study provides the first quantification of this potential. Using at-a-station field data, predictions of flow depth, water-surface width and longitudinal water surface slope are made using three standalone data mining techniques - , Instance-based Learning (IBK), KStar, Locally Weighted Learning (LWL) - along with four types of novel hybrid algorithms in which the standalone models are trained with Vote, Attribute Selected Classifier (ASC), Regression by Discretization (RBD), and Cross-validation Parameter Selection (CVPS) algorithms (Vote-IBK, Vote-Kstar, Vote-LWL, ASC-IBK, ASC-Kstar, ASC-LWL, RBD-IBK, RBD-Kstar, RBD-LWL, CVPS-IBK, CVPS-Kstar, CVPS-LWL). Through a comparison of their predictive performance and a sensitivity analysis of the driving variables, the results reveal: (1) Shield stress was the most effective parameter in the prediction of all geometry dimensions; (2) hybrid models had a higher prediction power than standalone data mining models, empirical equations and traditional machine learning algorithms; (3) Vote-Kstar model had the highest performance in predicting depth and width, and ASC-Kstar in

28 estimating slope, each providing very good prediction performance.. Through these algorithms, the  
29 hydraulic geometry of any river can potentially be predicted accurately and with ease using just a few,  
30 readily available flow and channel parameters. Thus, the results reveal that these models have great  
31 potential for use in stable channel design in data poor catchments, especially in developing nations where  
32 technical modelling skills and understanding of the hydraulic and sediment processes occurring in the  
33 river system may be lacking.

34

35 **Keywords:** gravel-bed rivers, hydraulic geometry, modelling, artificial intelligence, data mining, machine  
36 learning.

37

## 38 **1. Introduction**

39 Alluvial rivers form their own geometry in plan and cross-section, adjusting according to flow and  
40 sediment transport conditions. A river in a state of equilibrium over a specified period of time is said to be  
41 in regime or stable (Singh and Zhang 2008). This state of dynamic equilibrium occurs if the sediment  
42 transport rate is approximately equal to the upstream sediment supply, meaning that channel  
43 dimensions/geometry are maintained over this time period. Channel stability analysis involves analyzing  
44 how a channel adjusts its hydraulic geometry in response to changes in water and sediment discharge  
45 using river channel adjustment approaches (Gholami et al. 2017). This geometry is specified in terms of  
46 river flow width, depth, velocity and slope, and understanding how these hydraulic parameters vary with  
47 other variables, such as discharge, shear stress and median bed grain-size, is of paramount importance in  
48 stable channel design. The change in geometry is considered either over time at one cross-section (called  
49 at-a-station hydraulic geometry), focussing on temporal variations in the river geometry, or along the river  
50 length (called downstream hydraulic geometry). To design a stable geometry, accurate prediction of  
51 channel form in relation to the temporal and spatial variation in river hydraulics and sediment transport  
52 dynamics, is therefore required.

53 Thus far, various methods have been used to develop functional relationships for predicting stable  
54 hydraulic geometry dimensions. These approaches can be broadly classified into three methods, each  
55 using the same basic assumption of steady and uniform flow to achieve channel equilibrium. First,  
56 empirical equations of the regime have been obtained from the statistical rule/regression analysis of  
57 channel geometry data from different rivers (Thomas Blench 1952; Bray 1982; Hey and Thorne 1986a;  
58 Leopold and Wolman 1957; Wolman 1954). In these equations, flow discharge, bed shear stress and bed-  
59 grain diameters have been considered as the most effective parameters to predict the geometry of stable  
60 rivers (Deshpande and Kumar 2012; Parker et al. 2007). The major drawbacks of this approach is the lack  
61 of hydraulic, theoretical basis to the equations (Eaton and Church 2007; Hey and Thorne 1986b), and  
62 consequently low generalization and limited accuracy when applied to rivers in conditions that fall  
63 outside those used in the development of the equations (Bose 1936; Stevens and Nordin 1987). Another  
64 shortcoming of this method is that the equations are most often only developed only with flow discharge  
65 and bed-grain diameter as driving variables, while other important variables such as sediment transport  
66 rate or sediment concentration are neglected.

67 Secondly, theoretical and analytical models have been developed by river engineers and  
68 geomorphologists. For example, many studies have developed models based on regime theory (T Blench  
69 1969; Hey and Thorne 1986a; Huang and Nanson 1998; Lane 1957), quantifying the critical control of  
70 bed and bank materials on river channel form either through using a 'silt factor' or by developing regime  
71 relations based on the character of these materials. However, no study has proposed a universally  
72 accepted rational theory, nor defined universal formulations for its parameters (Gleason 2015). Analytical  
73 models have been developed by solving the governing hydraulic equations, most often based on field  
74 observations (Henderson 1961). For example, Julien and Wargadalam (1995) created analytical equations  
75 for downstream hydraulic geometry as a function of flow discharge, sediment size, Shields number and  
76 streamline deviation angle. They argued these models are more accurate and reliable than empirical  
77 equations because they are based on the physics and theory of the process. Afzalimehr *et al.* (2010) tested

78 the performance of these analytical equations against empirical equations based on 85 at-a-station datasets  
79 from Iranian rivers and found contrasting results. These contrasting results were reported because the  
80 empirical equations were only tested with the datasets from which they were developed. This paper also  
81 found that the grain size and the Shields parameter need not be taken into account when evaluating the  
82 width and depth of an alluvial channel at a site.

83 Thirdly, numerical models have been developed based on the solution of flow friction equations, the law  
84 of continuity, sediment transport capacity, and in some cases, the stability of the river banks (Chang 1980;  
85 Millar 2005; White 1982). Although analytical equations provide a stronger logical framework for  
86 examining possible changes in prevailing conditions (Ferguson 1986), the prediction performance of  
87 numerical solutions can be similar to those of empirical models (Millar, 2005). Examples of numerical  
88 equations for stable hydraulic geometry prediction are provided in commonly-used software, such as  
89 HEC-RAS (Mehta et al. 2013; Shelley and Parr 2009). Although this type of model is developed based on  
90 the physics of the process, they require lots of data to provide good model performance, and calibration is  
91 difficult and time-consuming. Therefore, new ways to predict stable hydraulic geometry, that are  
92 computationally simple, flexible, reliable and require small datasets, are required.

93 Since the 1980s, several Artificial Intelligence (AI) algorithms have been developed successfully to solve  
94 hydraulic problems and are gaining more attention due to their high performance and flexibility. These  
95 algorithms utilize data with different scale and are insensitive to missing data and the length of data. One  
96 of the most commonly-used AI models in hydraulics is the Artificial Neural Network (ANN). This  
97 algorithm has been used by many researchers to estimate hydraulic parameters, such as bed shear stress,  
98 as well as inform the design of alluvial irrigation canals (Mohamed 2013; Sheikh Khozani et al. 2017;  
99 Wan Mohtar et al. 2018), rainfall-runoff modelling (Antar et al. 2006), rainfall prediction (Mislan et al.  
100 2015) and water quality assessment (Cuest Cordoba et al. 2014). ANN models can implicitly identify  
101 complicated, nonlinear connections between independent and dependent parameters and can detect all  
102 potential interactions across the predictor parameters. Given the nonlinear relationship between hydraulic

103 and sediment transport parameters, ANN models have thus been used in the prediction of channel  
104 geometry. For example, Khadangi et al. (2009) predicted three channel parameters (width, depth, and  
105 slope) using data collected from 371 rivers, and examined the prediction performance of two different  
106 ANNs structures. Their results showed good performance in the evaluation phase compared with  
107 measured values, performing better in estimating channel width than depth and slope. Mohamed (2013)  
108 applied an ANN model based on a back-propagation algorithm to estimate the wetted perimeter, hydraulic  
109 radius and water surface slope of 61 Egyptian irrigation canals. The prediction performance of  
110 Mohamed's (2013) model was compared against three empirical equations frequently used to predict  
111 hydraulic geometry. The ANN model had superior performance in all cases. Gholami *et al.* (2017)  
112 showed this was also the case for gravel-bed rivers. In another study, Tahershamsi *et al.* (2012)  
113 investigated the performance of multi-layer perceptron (MLP) and Radial Basis Function (RBF) models  
114 to forecast the width of alluvial channels. Both models had good prediction performance. However,  
115 despite these promising results, ANN models have slow coverage speed during the training procedure,  
116 and model performance can decrease significantly if the training dataset is not carefully chosen (i.e. when  
117 the testing dataset is out of range of the training dataset; Choubin et al., 2018).

118 Evolutionary models have gained a lot of attention in recent years (Ferreira 2001; Wang et al. 2016). In  
119 particular, Gene Expression Programming (GEP) is recognised as a strong and problem-independent  
120 technique for multivariate optimization (Ferreira 2002; Wu et al. 2013). Shaghaghi *et al.* (2018) applied  
121 three Non-linear Regression (NLR), GEP and, Generalized Structure of Group Method of Data Handling  
122 (GS-GMDH) models to estimate alluvial channel width, depth and slope. The Group Method of Data  
123 Handling (GMDH) model relates to the deterministic self-organizing method group, where the principle  
124 of a black box, connectionism and induction is used (Anastasakis and Mort 2001). Shaghaghi *et al.* (2018)  
125 investigated the impact of different input variable combinations and found that the most effective  
126 parameters in estimating width and depth were discharge and mean particle size, while for channel slope,  
127 the Shields parameter was the most effective. They compared the accuracy of their three models and

128 deduced that GEP and GS-GMDH had better predictive performance than the NLR model. However the  
129 weakness of the GMDH algorithm lies in its fixed configuration, using a deterministic approach to find  
130 the optimal partition of datasets and parameters (Robinson 1998). Sheikh Khozani *et al.* (2017) predicted  
131 shear stress distribution in circular channels by applying GEP and evaluating the performance of different  
132 input combinations. Their model showed better performance in estimating shear stress distribution than a  
133 Shannon entropy-based equation presented by Sterling and Knight (2002). Noori *et al.* (2016) compared  
134 ANN, Adaptive Neuro-fuzzy Inference system (ANFIS), and Support Vector Machine (SVM) models for  
135 predicting the longitudinal dispersion coefficient in rivers and reported that SVM had a higher  
136 performance followed by ANFIS and ANN. The ANFIS algorithm, however, suffers from a large number  
137 of model operators, each of which needs to be set accurately, especially the weights of membership  
138 function. Although SVM has a higher prediction power, the model can be time-consuming to train, since  
139 it is susceptible to hyper-parameter selection (Ahmad *et al.* 2018), and choosing the best kernel is  
140 problematic, reducing its wider application.

141 Consequently, a new form of AI, data mining, has been applied in the fields of hydrology and hydraulics  
142 to overcome the aforementioned weaknesses in traditional AI models. Some of these new algorithms,  
143 such as tree-based models [i.e. Random Tree (RT), Random Forest (RF), M5 Prime (M5P), Reduced  
144 Error Pruning Tree (REPT)], re-sampling algorithm (i.e. Bootstrap Aggregation, also called bagging), and  
145 other algorithms such as Random Subspace, and k nearest neighbor (IBK), were used to estimate apparent  
146 shear stress in a compound river cross section (Sheikh Khozani *et al.* 2019), suspended sediment transport  
147 (Khosravi *et al.* 2018), nitrate and strontium concentrations in groundwater (Bui *et al.* 2020). These data  
148 mining algorithms have higher predictive power than traditional AI models. For example, Hussain and  
149 Khan (2020) found RF had a 17.8 % and 33.6 % higher performance than ANN and SVM for predicting  
150 river streamflow. Further, Shamshirband *et al.* (2020) demonstrated the superiority of M5P over SVM for  
151 standardized streamflow index prediction. Also Khosravi *et al.* (2019) showed that data mining  
152 algorithms outperform standalone ANFIS algorithms in the prediction of reference evaporation, while

153 optimized ANFIS using metaheuristic algorithms performed slightly better than standalone data mining  
154 algorithms. Also some researchers have reported that hybridized algorithms improve the performance of  
155 standalone algorithms, not only for traditional AI algorithms, but also for data mining models in the  
156 prediction of water quality index and bedload transport rate (Bui et al. 2020a; Bui et al. 2020b; Khosravi  
157 et al. 2020). However, these new data mining algorithms have yet to be applied for the prediction of  
158 hydraulic geometry. Thus, a significant gap exists in understanding the potential of these data mining  
159 algorithms, and in the identification of the most flexible and accurate algorithm.

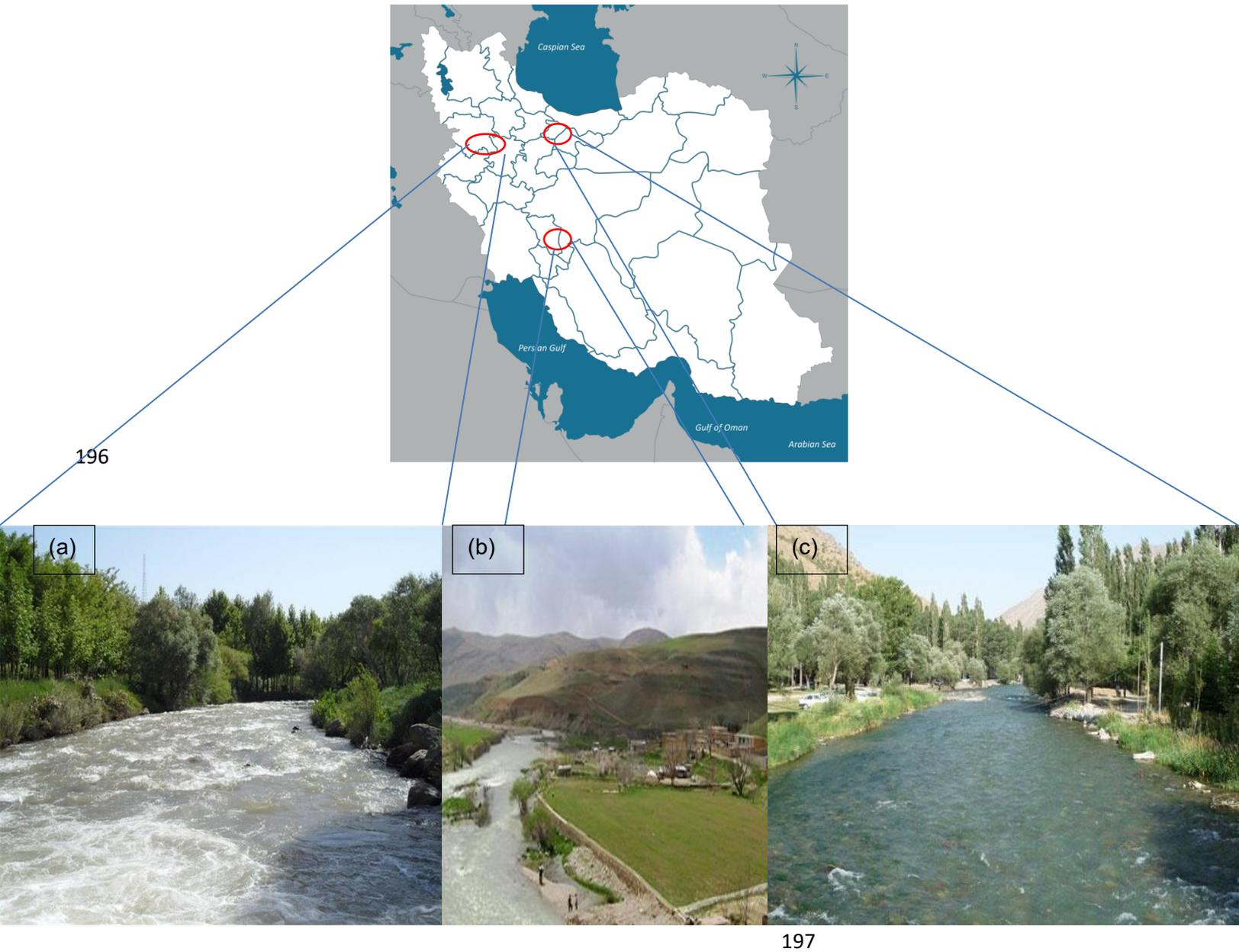
160 The present paper, therefore, aims to fill this gap in understanding by achieving the following objectives:  
161 (1) produce predictions of the three main hydraulic geometry parameters (mean flow depth, water-surface  
162 width and longitudinal water surface slope) using three standalone data mining techniques, namely  
163 Instance-based Learning (IBK), KStar, Locally Weighted Learning (LWL), along with four types of novel  
164 hybrid algorithms in which the standalone models are trained with Vote, Attribute Selected Classifier  
165 (ASC), Regression by Discretization (RBD), and Cross-validation Parameter Selection (CVPS)  
166 algorithms (Vote-IBK, Vote-Kstar, Vote-LWL, ASC-IBK, ASC-Kstar, ASC-LWL, RBD-IBK, RBD-  
167 Kstar, RBD-LWL, CVPS-IBK, CVPS-Kstar, CVPS-LWL; (2) compare the predictive power of these  
168 data-driven models; and (3) perform a sensitivity analysis of the driving variables used in each model.  
169 The performance of these algorithms is tested for the following reasons: (1) IBK can adapt to previously  
170 unseen data, storing a new instance or throwing an old instance away, making it potentially superior to  
171 other methods of machine learning. (2) The KStar algorithm uses entropic measure based on probability  
172 of transforming instance into another by randomly choosing between all possible transformations  
173 (Madhusudana et al. 2016). (3) LWL improves the overall performance of regression methods by  
174 adjusting the capacity of the models to the properties of the training data in each area of the input space  
175 (Reyes et al. 2018). (4) Vote algorithm can find the majority of a sequence of the elements by using linear  
176 time and constant space. Also, this algorithm is important for ultra-reliable system which are based on the  
177 multi-channel computation paradigm (Parhami 1994). (5) ASC model benefit from three main

178 components including base classifier, evaluator and search algorithm in its structure (Thornton et al.  
179 2013). (6) In RBD method, the estimated value is the probable value of the mean class value for each  
180 discretized interval, according to the estimated probabilities for each interval (Frank and Bouckaert 2009).  
181 (7) CVPS is a technique of selecting parameters using cross-validation sampling. To the best of our  
182 knowledge, this study is the first to apply these hybridized algorithms in any branch of geoscience. The  
183 research offers new insight into which data mining algorithms offer the potential to provide relatively  
184 cheap and fast predictions of hydraulic geometry in situations when understanding of the physical  
185 processes at play may not be well understood.

## 186 **2. Material and methods**

### 187 **2.1. Datasets**

188 The paper uses a dataset compiled by Afzalimehr *et al.* (2010) for three stable gravel-bed rivers in Iran:  
189 Karaj river in Alborz Province, Behesht-Aabad river in Charmahal-and-Bakhtiari Province and Gamasiab  
190 River in Kermanshah province (Figure 1). This dataset includes measurements of flow discharge ( $Q$ ),  
191 median sediment diameter ( $d_{50}$ ), Shields number ( $\tau^*$ ) at 85 cross-sections (Table 1), used as inputs to  
192 predict hydraulic geometry. This geometry is defined by water-surface width ( $w$ ), mean flow depth ( $h$ )  
193 and longitudinal water surface slope ( $S$ ). Flow discharge in a cross section was estimated through three to  
194 five velocity profiles, with each profile containing 13-16 velocity measurements at different heights above  
195 the sediment bed, totalling 425 profiles



198 Fig. 1. Map and illustrative photographs of the three studied rivers: (a) Gamasiab river, (b) Behesht-Abad  
199 river and (c) Karaj river.

200 and 6000 point velocities. At each cross section the top width (channel width at the water surface) was  
201 measured along with the flow depth at 0.5 m intervals across the channel. The mean flow depth at each  
202 cross section was calculated by dividing the cross-sectional area by this top width. The Wolman's walk  
203 approach (Wolman 1954) was used to measure the bed sediment size distribution. The longitudinal water

204 surface slope was determined by dividing the difference in water surface elevations between two cross  
 205 sections along the central axis of the reach. . Shields parameter was computed based on the following  
 206 equation:

$$207 \quad \tau^* = \frac{\tau}{(\rho_s - \rho)gd_{50}} \quad (1)$$

208 where  $\tau$  is the shear stress [-],  $\rho_s$  is sediment density [-],  $\rho$  is water density [ $\text{kg m}^{-3}$ ] and  $g$  is  
 209 gravitational acceleration [ $\text{m s}^{-2}$ ]. Shear stress was calculated as follows:

$$210 \quad \tau = (\rho \cdot v^*)^2 \quad (2)$$

211 where  $v^*$  is the shear velocity [ $\text{m s}^{-1}$ ] which was calculated as  $v^* = (ghS)^{0.5}$ . More information about the  
 212 data collection methodology can be found in Afzalimehr *et al.* (2010).

213 Table 1. Descriptive statistics of the training and testing dataset

	Training dataset						Testing dataset					
	max	min	mean	Std	SK	K	max	min	mean	Std	SK	K
$Q$ ( $\text{m}^3/\text{s}$ )	5.810	0.500	2.245	1.430	0.374	-0.709	5.300	0.550	2.299	1.338	0.181	-0.760
$d_{50}$ (m)	0.130	0.004	0.032	0.031	1.316	1.080	0.094	0.004	0.031	0.029	0.819	-0.612
$\tau^*$ (-)	0.814	0.000	0.121	0.170	2.059	4.393	0.481	0.001	0.105	0.139	1.704	2.124
$S$ (-)	0.028	0.0001	0.006	0.005	2.379	8.151	0.016	0.0001	0.005	0.003	1.686	3.774
$h$ (m)	0.570	0.180	0.344	0.094	0.313	-0.877	0.570	0.230	0.337	0.085	0.937	0.738
$w$ (m)	27.000	5.500	14.582	5.507	0.401	-0.749	23.000	7.000	14.072	4.748	0.217	-1.054

214 where max = maximum, min = minimum, Std = standard deviation, SK = skewness and K = kurtosis

## 215 2.2. Dataset preparation and sample size

216 The 85 datasets were split into two subgroups; 70% of the datasets were selected randomly to be used as  
 217 training data for model development and the remaining 30% was applied as testing data for model  
 218 validation. There is no agreement in the literature on this ratio. Some have used ratios of 80:20  
 219 (Zounemat-Kermani *et al.* 2019), and 75:25 (Hooshyaripor *et al.* 2014). Palani *et al.*, (2008) and Barzegar  
 220 *et al.* (2016) declared that the testing dataset should represent approximately 10 – 40% of the size of the

221 whole dataset. Also, Kisi et al. (2019) showed that by increasing the length of the training dataset from  
 222 50% to 75%, the modelling performance increased. With these considerations in mind a 70:30 ratio is the  
 223 most commonly used (Bui et al. 2018; Chen et al. 2017; Taheri et al. 2019).

### 224 2.3. Model input, calibration and sensitivity analysis

225 Flow discharge, median sediment diameter, and Shields number are the three most important and widely  
 226 used variables which affect stable river geometry (Deshpande and Kumar 2012; Gholami et al. 2017;  
 227 Parker et al. 2007; Shaghghi et al. 2018). These parameters were therefore used as an input in each  
 228 model to predict the top width, flow depth and longitudinal slope at each river cross-section.

229 There are two main steps in using AI algorithms: (i) determination of the best input variable combination;  
 230 and (ii) identifying the operator's optimum values. Each input variable has a differing impact on these  
 231 hydraulic geometry parameters. Thus different input combinations were constructed and examined to find  
 232 the most effective input combination (Table 2). These combinations were constructed by beginning with  
 233 the variable with the highest Pearson correlation coefficient (*PCC*) (a measure of linear correlation  
 234 between two sets of data) ( $\tau^*$  for  $h$  and  $S$ , and  $d_{50}$  for  $w$ ), and then exploring all other input  
 235 combinations.. The effect of each input variables on the output was examined through a sensitivity  
 236 analysis. To explore the most effective combination, the models were implemented using default models  
 237 operators. Their effectiveness was assessed using Root Mean Square Error (*RMSE*); the lower the *RMSE*,  
 238 the higher the effectiveness of the input combination.

239 Table 2. Different input combinations constructed to explore the most effective combination for model  
 240 calibration.

241

No.	Input	Output		No.	Input	Output
1	$\tau^*$	$h, S$		1	$d_{50}$	$w$
2	$\tau^*, Q$	$h, S$		2	$\tau^*, Q$	$w$

3	$\tau^*, d_{50}$	$h, S$	3	$\tau^*, d_{50}$	$w$
4	$Q, d_{50}$	$h, S$	4	$Q, d_{50}$	$w$
5	$\tau^*, Q, d_{50}$	$h, S$	5	$\tau^*, Q, d_{50}$	$w$

242

243

244 Along with data quality, length of data, and input variable choice, the calibration of model operator values  
 245 has an important impact on prediction performance. There are no optimum operator values which work  
 246 globally for model calibration. Hence, to enhance the prediction power of each algorithm, these values  
 247 were set after the determination of the best input combination. At first, default values of each operator  
 248 were considered, and then based on this result, lower and higher values were selected to find the optimum  
 249 value. The most widely used approach of trial and error was performed in Waikato Environment for  
 250 Knowledge Analysis (WEKA 3.9) software. The optimum operator values were achieved by minimizing  
 251 the Root Mean Square Error (*RMSE*) during the testing phase.

252

## 253 **2.4. Model descriptions**

### 254 **2.4.1. Instance-based Learning (IBK)**

255 Instance-based Learning, also known as K-Nearest Neighbor classification, is a lazy learning algorithm,  
 256 well known for its ability to recognise data patterns. The algorithm applies a relatively simple method to  
 257 store training data and identify new undefined data by measuring the distance between similar recorded  
 258 samples. The IBK utilises an election system to determine the class of new samples; the number of votes  
 259 defines the  $k$  value. The distance is defined after the  $k$  value is determined. The application of the IBK  
 260 algorithm involves three steps. (1) reading the  $k$  value, distance type and test data, (2) finding the  $k$   
 261 nearest neighbor to the test data, and (3) setting the maximum label class to the test data. The WEKA  
 262 Machine Learning Software (Witten et al. 2016) was utilized for running the IBK algorithm.

### 263 **2.4.2. Kstar**

264 The Kstar algorithm, first introduced by Cleary and Trigg (1995), is another type of lazy algorithm, which  
 265 uses an entropy-based distance function to transform one sample probability to another by selecting  
 266 arbitrarily all feasible transformations. The classification with Kstar is performed by summing the new  
 267 instance probabilities to all the members of a group. This classification must be achieved for the other  
 268 groups in order to eventually choose the one with the highest probability (Cleary and Trigg 1995). For  
 269 missing values, Cleary and Trigg (1995) assumed that the likelihood of transformation to these values is  
 270 the average of the likelihood of transformation to each of the defined attributes in the whole dataset. The  
 271 algorithm is defined as follows. Consider  $I$  as a set of instances and  $T$  as a set of transformations on  $I$   
 272 (Clearly *et al.*, 1995). Each instance ( $t \in T$ ) maps to another instance as  $t: I \rightarrow I$ .  $T$  has a special member  
 273  $\omega$  to map samples to themselves ( $\omega(a) = a$ ). Let  $P$  be the set of all prefix codes from  $T^*$ , which is  
 274 terminated by  $\omega$ . The  $T^*$  members define a transformation on  $I$ :

$$275 \quad \bar{t}(a) = t_n(t_{n-1}(\dots t_1(a) \dots)) \quad \text{where } \bar{t} = t_1, \dots, t_n \quad (3)$$

276 The probability function of  $T^*$  is defined as  $p$ :

$$277 \quad 0 \leq \frac{p(\bar{t}u)}{p(\bar{t})} \leq 1 \quad (4)$$

$$278 \quad \sum_u p(\bar{t}u) = p(\bar{t}) \quad (5)$$

$$279 \quad \sum_{\bar{t} \in P} p(\bar{t}) = 1 \quad (6)$$

280 Furthermore, the probability of the entire path from such an instance to  $b$  is defined as  $P^*$ :

$$281 \quad P^*(b|a) = \sum_{\bar{t} \in P: \bar{t}(a)=b} p(\bar{t}) = 1 \quad (7)$$

### 282 2.4.3. Locally Weighted Learner (LWL)

283 Locally weighted Learner is another lazy learning algorithm, The algorithm has an optimal convergence  
 284 speed and its minimum performance is higher than all possible linear regressions ( Stone 1982). The LWL

285 method is able to manage a wide range of data distribution types and can prevent boundary and cluster  
286 effects (Hastie and Loader 1993). The LWL depends on the distance function, which is used to recover  
287 the nearest neighbours of a given query example (Atkeson *et al.*, 1997). The method also depends on a  
288 smoothing parameter and weighting function. The weighting function calculates the weight of the sample  
289 neighbor query. This function should have a maximum value at a distance of zero, and as the distance  
290 increases, the performance slowly decreases., A bandwidth parameter ( $k$ ) acts as the smoothing  
291 parameter, determining the size or the range in which generalisation is accomplished. This parameter is  
292 defined as follows.

293 Let a non-linear system be defined as (Arif et al. 2001):

$$294 \quad y(k) = z(x(k), u(k)) \quad (8)$$

$$295 \quad u_d(k) = z^{-1}(x_d(k), y_d(k)) \quad (9)$$

296 in which a non-linear function is defined as  $z(\cdot)$ , the states as  $x_d(k)$ , and the output parameter as  $y_d(k)$ .

297

#### 298 **2.4.4. Vote**

299 The meta algorithm Vote was used to train the IBK, Kstar, and LWL models and produced three hybrid  
300 models, Vote-IBK, Vote-Kstar, and Vote-LWL. This algorithm combined each basic-level classifier using  
301 a vote approach. The simplest voting approach is majority voting, in which the basic-level classifier casts  
302 one vote for its predictions. The instance is categorised into the class which obtains the most votes. For  
303 the situation where class probability distributions are estimated by the basic-level classifiers, the plurality  
304 voting method is modified (Dietterich 1997), defined as follows Assume  $P_s(x)$  is the estimated class  
305 probability distribution by the basic-level classifier  $S$  on sample  $x$ . The probability distribution  
306 components restored by the basic-level classifiers are summed to reach the probability distribution class  
307 of meta-level voting classifier as:

$$P_{S(ML)}(x) = \frac{1}{|S|} \sum P_s(x) \quad (10)$$

#### 2.4.5. Attribute Selection Committees (ASC)

The Attribute Selected Classifier algorithm was applied to train the IBK, Kstar, and LWL models and produced three hybrid models, ASC-IBK, ASC-Kstar, and ASC-LWL. The Attribute Selected Classifier is an ensemble technique, generally considered as a black-box form of classifier. The structure of ensemble classifiers is such that much information can be obtained by using bi-product data (Gislason et al. 2006). making it possible to determine an attribute based on the training set before learning the predefined classification.

The advantages of applying the attribute subsets in ensemble learning are, according to Thornton et al. (2013): (1) reduction in the dimension of the data, which decreases the effect of the “curse of dimensionality”; (2) decrease in the connection between classifiers through training them on several characteristics; and (3) improvement in the classifiers output of the ensemble.

#### 2.4.6. Regression by Discretization (RBD)

The Regression by Discretization algorithm was used to train the standalones models and produce the following hybrid models: RBD-IBK, RBD-Kstar, and RBD-LWL. This algorithm is a meta classifier technique, based on conditional density prediction via the class probabilities. The output parameter is discretized in non-overlapping periods which are called “bins”. These bins can be produced of equal frequency and equal width. If a bin is defined as  $k_y$  which consist of the output value  $y$ , the whole number of output values in the training stage is  $n$ , the number of output values in bin  $m$  is  $n_m$  and  $p(k_y|X)$  is the estimated probability of specified class  $X$  forecasted from the class probability predictor. The weight, for a specified output value  $y_i$  in case  $X$ , was computed as:

$$w(y_i|X) = m \frac{p(k_{y_i}|X)}{n_{k_{y_i}}} \quad (11)$$

330 The weight  $w(y_i|X)$  can be seen as an approximation of the likelihood of a future target predicted value  
331 correlated with  $X$  being close to  $y_i$ , based on the class probability prediction model derived from discrete  
332 training data.

#### 333 **2.4.7. Cross-Validation Parameter Selection (CVPS)**

334 The Cross-Validation Parameter Selection algorithm was used to train the standalone models IBK, Kstar,  
335 and LWL and produce the following three hybrid models: CVPS-IBK, CVPS-Kstar, and CVPS-LWL.  
336 Cross-validation is one of the most widely used statistical methods for assessing predictor model  
337 performance by using an *a priori* modelling procedure ( Stone 1974). The method is based on data  
338 splitting; a portion of the data is used to fit each competing method and the remaining data is used to  
339 calculate the predictive model's performance, and the model with the best overall efficiency is chosen.  
340 Using continuous cycles, the training and validation sets are cross-overed so that each data point has a  
341 chance of being verified against all other data points. The CVPS algorithm is one of the meta-classifier  
342 techniques which was extended in WEKA environment by Garg and Khurana (2014) and is used to  
343 improve the prediction power of standalone algorithms through hybridization.

#### 344 **2.5. Model validation**

345 Five frequently used metrics for assessing model performance were applied: coefficient of determination  
346 ( $R^2$ ), Root Mean Square Error ( $RMSE$ ), Mean Absolute Error ( $MAE$ ), Nash-Sutcliffe Efficiency ( $NSE$ ) and  
347 percent bias ( $PBIAS$ ). These metrics were calculated as follows (Dawson et al. 2007; Legates and  
348 McCabe Jr 1999; Moriasi et al. 2007):

349

$$R^2 = \left( \frac{\sum_{i=1}^n (X_o - \bar{X}_o)(X_e - \bar{X}_e)}{\sqrt{\sum_{i=1}^n (X_o - \bar{X}_o)^2 \sum_{i=1}^n (X_e - \bar{X}_e)^2}} \right)^2 \quad (12)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_e - X_o)^2} \quad (13)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |X_e - X_o| \quad (14)$$

$$NSE = 1 - \frac{\sum_{i=1}^n (X_e - X_o)^2}{\sum_{i=1}^n (X_o - \bar{X}_o)^2} \quad (15)$$

$$PBIAS = \left( \frac{\sum_{i=1}^n (X_o - X_e)}{\sum_{i=1}^n X_e} \right) * 100 \quad (16)$$

350 where,  $X_o$  and  $X_e$  are observed and predicted values,  $\bar{X}_o$  and  $\bar{X}_e$  are mean observed and predicted  
 351 values, respectively, and  $n$  is the number of data points. The performance classification of the model  
 352 evaluation metrics is shown in Table 3. The *PBIAS* metric reports over- ( $PBIAS < 0$ ) or under-prediction  
 353 ( $PBIAS > 0$ ).

354

355

Table 3. Performance classification of the model evaluation metrics

Objective function	Value range	Performance classification	References
$R^2$	$0.7 < R^2 < 1$ $0.6 < R^2 < 0.7$	Very good Good Satisfactory	Moriasi et al. (2007); Ayele et

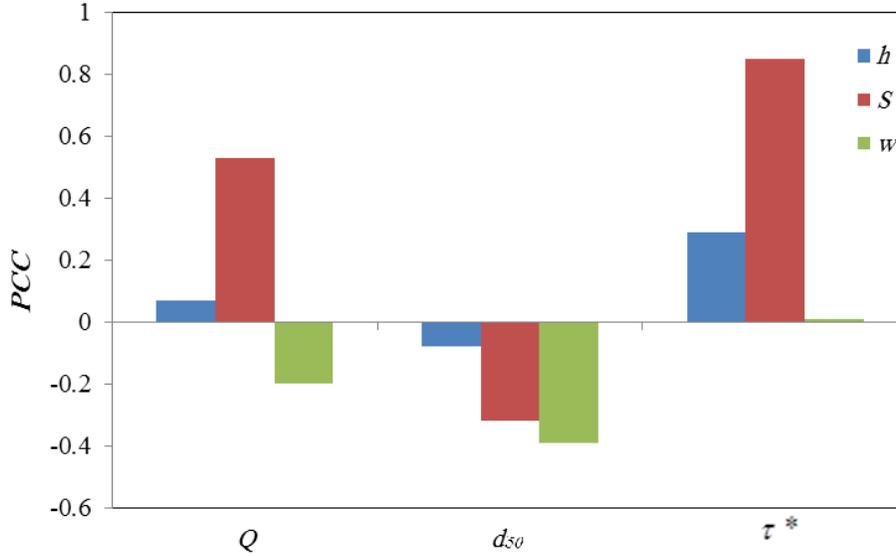
	$0.5 < R^2 < 0.6$	Unsatisfactory	al. (2017)
	$R^2 < 0.5$		
<i>RMSE</i>		The lower the <i>RMSE</i> , the better the model performance	Dawson et al. (2006)
<i>MAE</i>		The lower the <i>MAE</i> , the better the model performance	Dawson et al. (2006)
<i>NSE</i>	$0.75 < NSE \leq 1.00$	Very good	Moriassi et al. (2007); Boskidis et al. (2012)
	$0.65 < NSE \leq 0.75$	Good	
	$0.50 < NSE \leq 0.65$	Satisfactory	
	$0.4 < NSE \leq 0.50$	Acceptable	
	$NSE \leq 0.4$	Unsatisfactory	
<i>PBIAS</i>	$PBIAS < \pm 10$	Very good	Legates et al. (1999)
	$10 \leq  PBIAS  < 15$	Good	
	$15 \leq  PBIAS  < 25$	Satisfactory	
	$PBIAS \geq \pm 25$	Unsatisfactory	

356

357 For a visual assessment of the applied models, boxplots of observed and predicted values were compared  
358 (Figure A, Supplementary material). These were used to shows how well a model predicts extreme,  
359 median and quartile values.

### 360 **3. Results**

361 The *PCC* values in Figure 2 show the level of correlation between input variables and hydraulic geometry  
362 parameters. First, Shields's stress had the highest correlation with longitudinal slope ( $PCC = 0.85$ )  
363 followed by flow depth ( $PCC=0.29$ ) and width ( $PCC=0.01$ ). Second, median sediment diameter  $e$  had the  
364 highest correlation with width ( $PCC = -0.39$ ), followed by slope ( $PCC = -0.32$ ) and depth ( $PCC = 0.08$ ).  
365 Finally, discharge had the highest correlation coefficient with longitudinal slope ( $PCC=0.53$ ) followed by  
366 width ( $PCC=0.2$ ) and depth ( $PCC=0.07$ ).



367  
 368 Fig 2. Pearson correlation coefficient ( $PCC$ ) between input variables and hydraulic geometry parameters  
 369

370 **3.1. Determination of the best input variable combination**

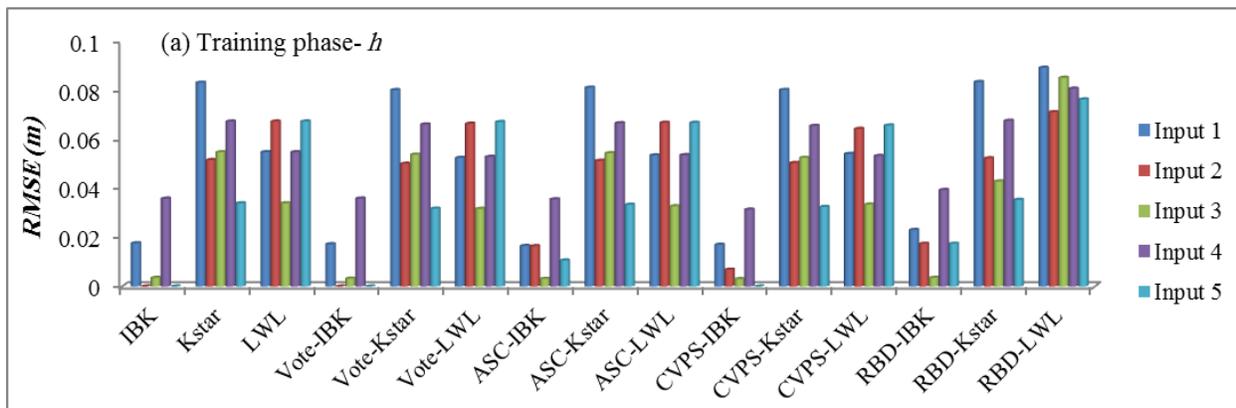
371 Figure 3 shows that, due to the different structures of each model, the optimal input variable combinations  
 372 differ between the models. Input combination No. 3 ( $d_{50}$  and  $\tau^*$ ) and No. 5 ( $Q$ ,  $d_{50}$  and  $\tau^*$ ) were most  
 373 influential in both the training and testing phase for flow depth prediction; No. 2 input combination, was  
 374 only the most effective for the RBD-LWL algorithm. This result reveals that overall  $Q$  is not a  
 375 particularly effective variable, as neither No. 2 nor No.4 input combinations could predict flow depth  
 376 accurately. This finding is in accordance with the  $PCC$  values displayed in Figure 2.

377 The best input combinations for predicting longitudinal slope were No. 2 and 3. Combination No. 4 ( $Q$ ,  
 378  $d_{50}$ ) could not predict slope accurately, revealing that Shields stress was the most effective parameter.

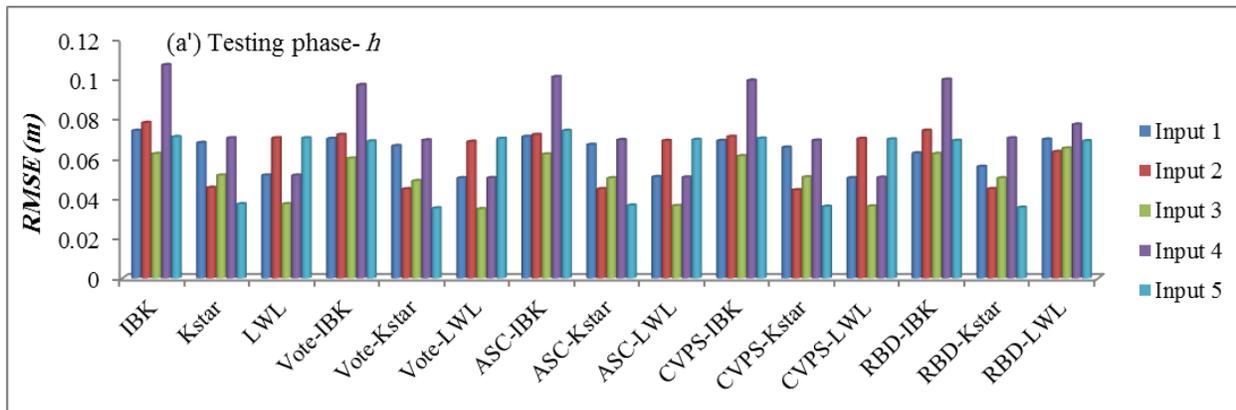
379 Contrasting results were found for predicting width. No. 2 and 3 were the optimum input combination for  
 380 just a few of the models, while No.1 and No.5 input combinations were the optimum combination in most  
 381 cases. Input No.1, which only contains  $d_{50}$ , predicted flow depth accurately in all models, reflecting its  
 382 high  $PCC$  value (Figure 2). In all models, the  $RMSE$  is larger for the testing than the training phase as

383 commonly found in AI methods because the training data are assessed on the same data that have been  
 384 learnt before, while the test dataset has data that is unknown to the algorithm and gives rise to more errors  
 385 or misclassification. Overall, the results show the single, most effective parameter is not able to predict  
 386 hydraulic geometry dimensions with the highest degree of accuracy.

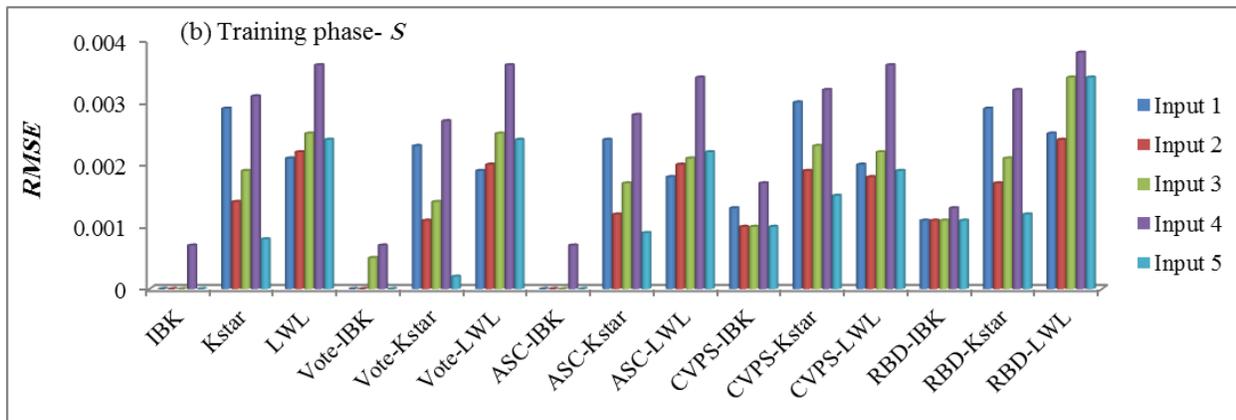
387

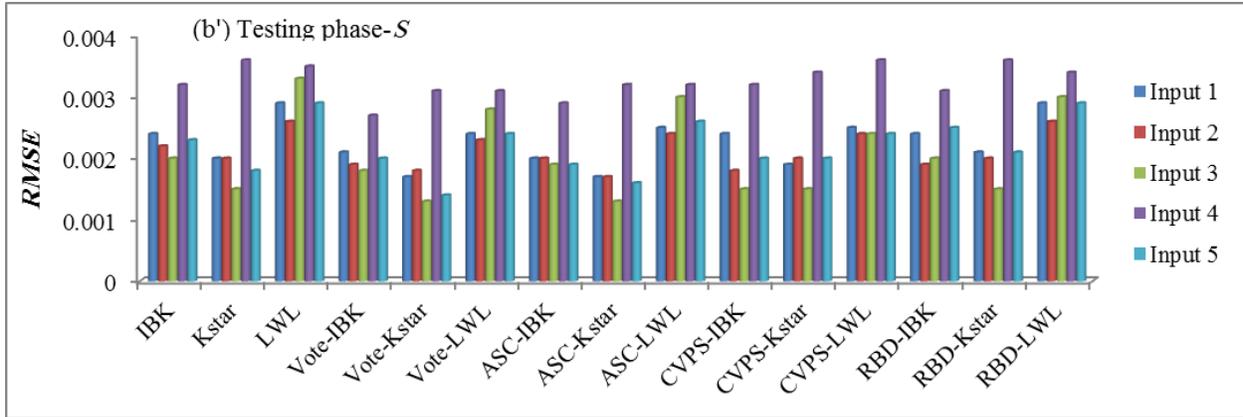


388

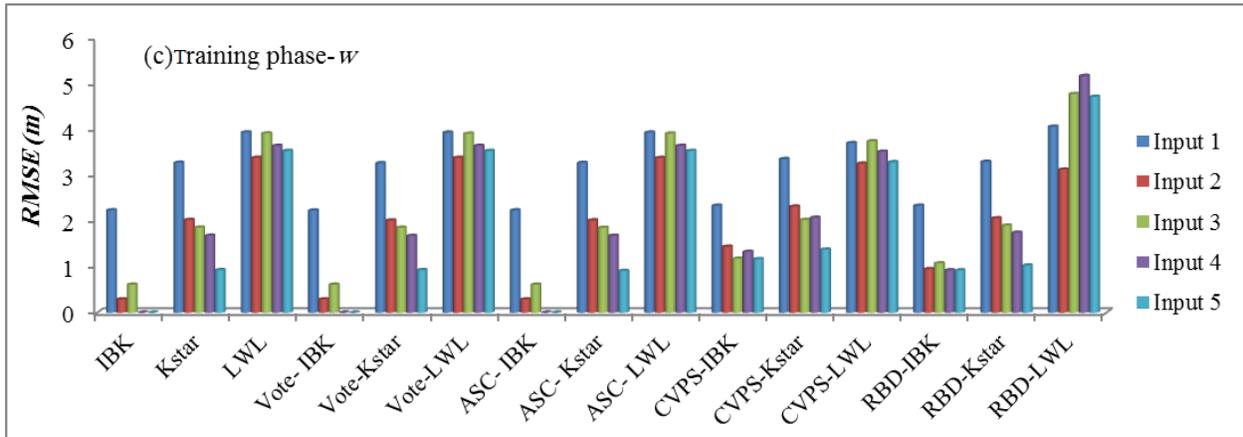


389

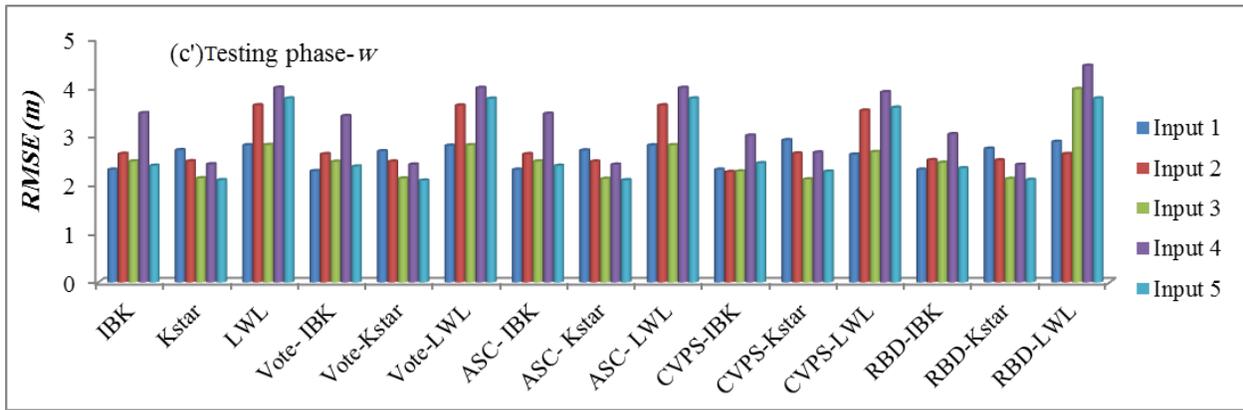




390



391



392

393 Fig 3. The change in model  $RMSE$  for different input variable combinations: (a) training phase,  $h$ ; (a')  
 394 testing phase,  $h$ ; (b) training phase,  $S$ ; (b') testing phase,  $S$ ; (c) training phase,  $w$ ; (c') testing phase,  $w$ .

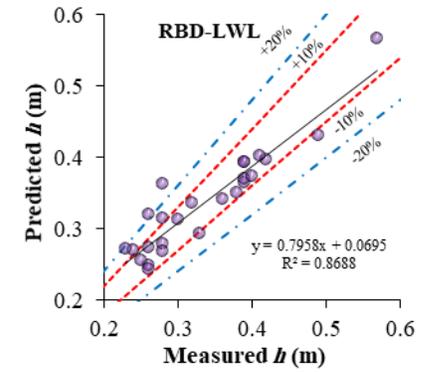
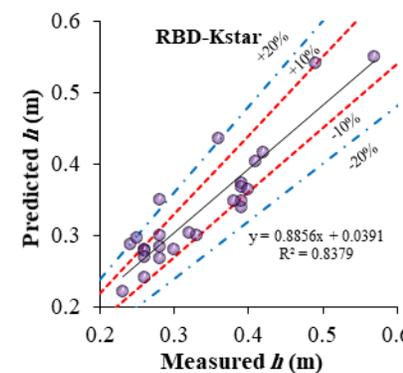
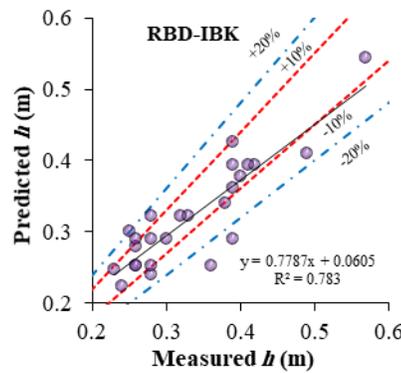
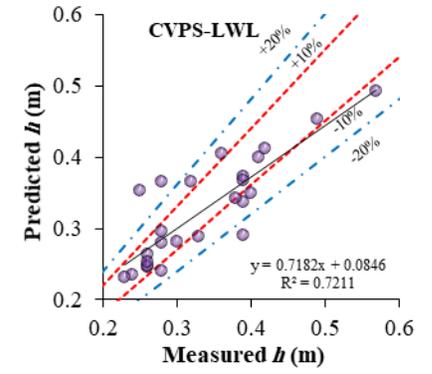
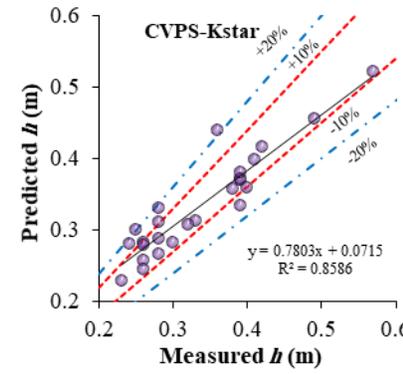
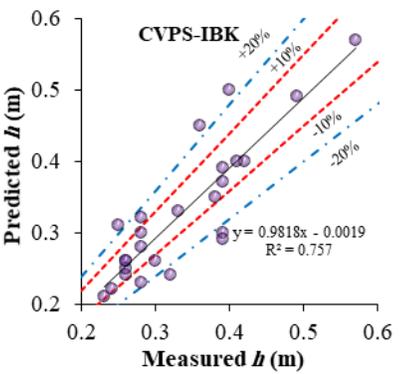
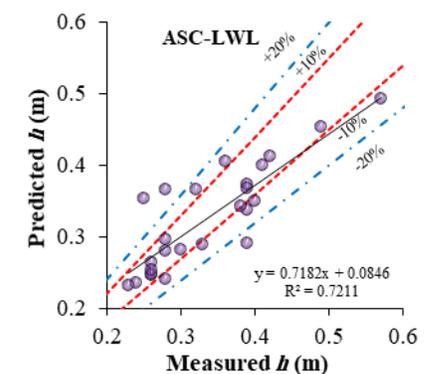
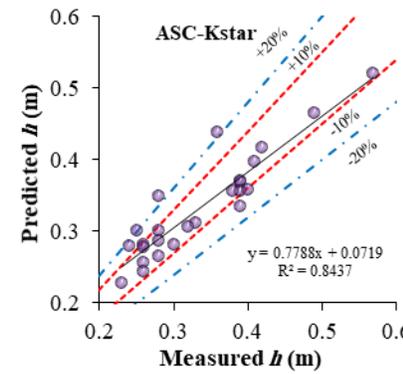
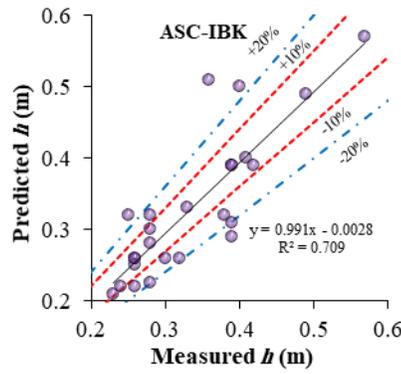
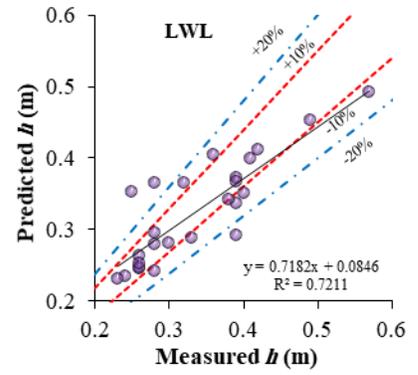
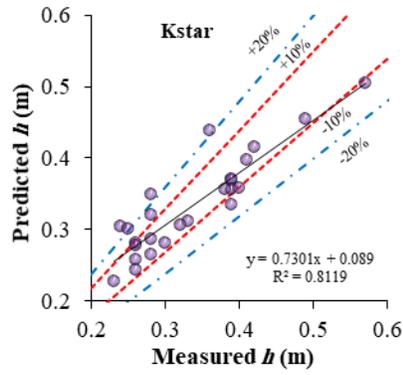
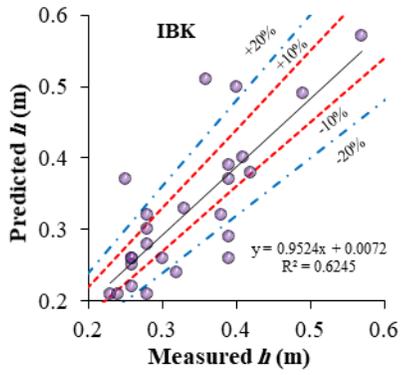
395

396 **3.2. Model performance**

397 After determination of the most effective input variables and optimised model operators, three standalone  
398 data mining models, along with 12 types of novel hybrid models were developed to predict the hydraulic  
399 geometry. The models were built by a training dataset. A comparison of the observed and predicted  
400 values from the testing dataset (Figure 4) shows that of the three standalone models, IBK had the lowest  
401 prediction power for flow depth ( $R^2 = 0.624$ ), and Kstar had the highest ( $R^2 = 0.812$ ). All hybrid  
402 algorithms performed better than the standalone models, with, the hybrid Vote-Kstar algorithm  
403 performing the best of all models ( $R^2 = 0.889$ ).

404 Kstar was also the best performing standalone model for predicting slope ( $R^2 = 0.792$ ; Figure 5) and width  
405 ( $R^2 = 0.792$ ; Figure 6), and LWL was the lowest ( $R^2 = 0.792$ ;  $R^2 = 0.754$ , respectively). Hybridization of  
406 the standalone algorithms increased the model performance for slope and width by a greater degree than  
407 for flow depth. The RBD-IBK algorithm outperformed all other algorithms in the prediction of slope ( $R^2 =$   
408  $0.913$ ), followed very closely by RBD-LWL ( $R^2 = 0.910$ ) and ASC-Kstar ( $R^2 = 0.909$ ). Whereas for width,  
409 this order was CVPS-Kstar ( $R^2 = 0.914$ ), Vote-Kstar ( $R^2 = 0.911$ ) and RBD-IBK ( $R^2 = 0.908$ ). According  
410 to the classification of performance based on the  $R^2$  metric (Ayele et al. 2017; Legates and McCabe Jr  
411 1999; Moriasi et al. 2007), all models had a ‘very good’ performance, except the IBK model for depth  
412 and slope, and the LWL model for slope, which had ‘good’ performance.

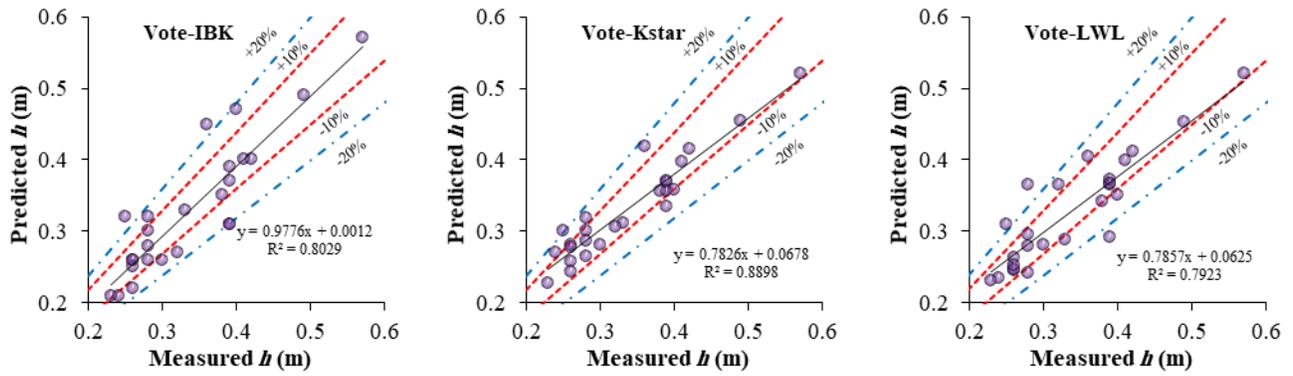
413



414

415

416



417

418

Fig 4. Scatter plot of measured versus predicted flow depth  $h$ .

419

420

421

422

423

424

425

426

427

428

429

430

431

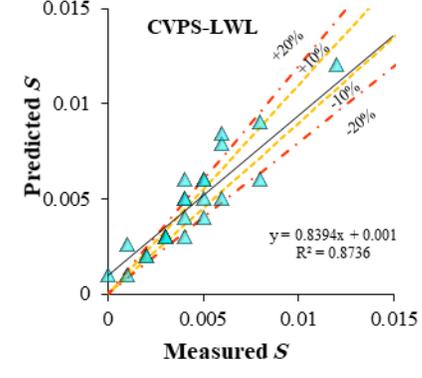
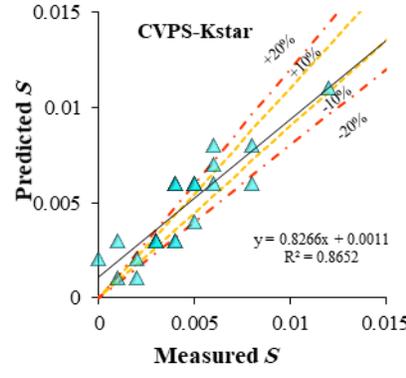
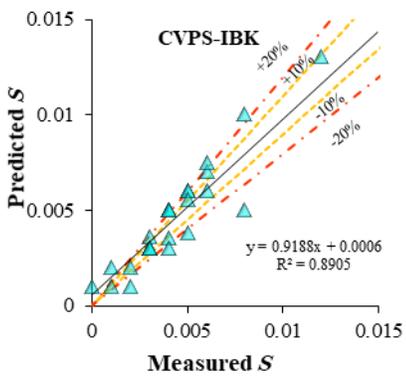
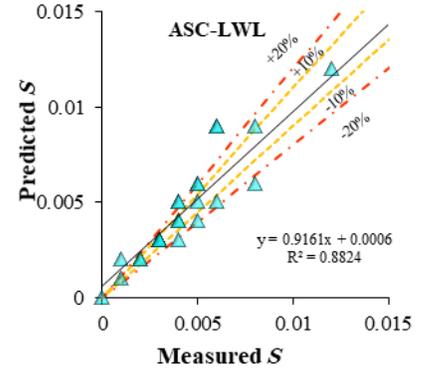
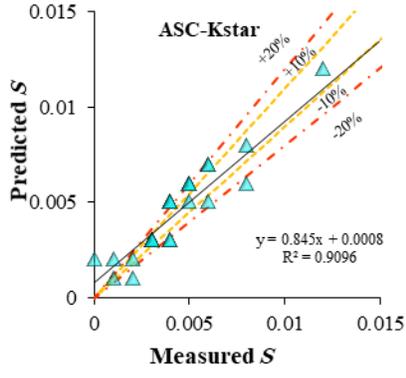
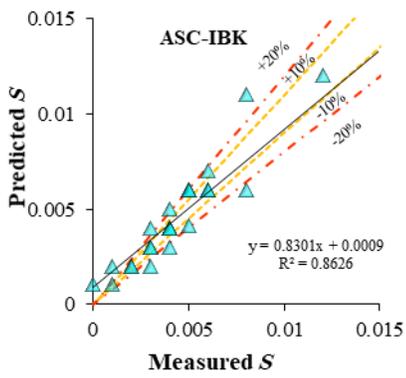
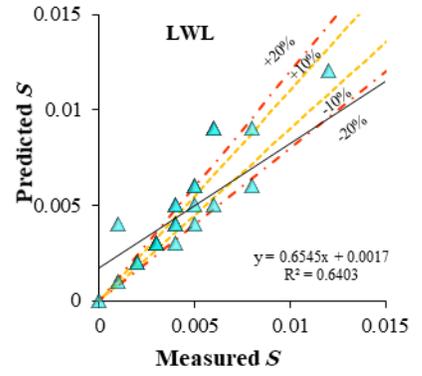
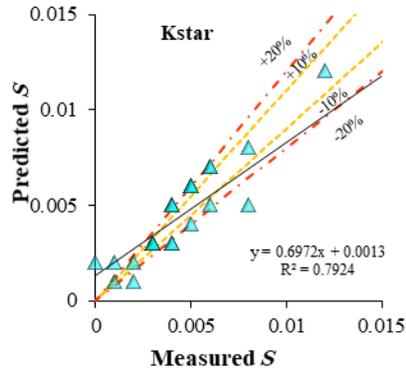
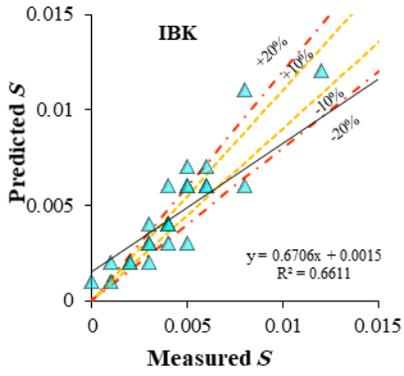
432

433

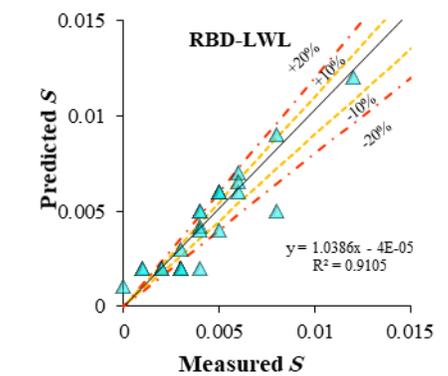
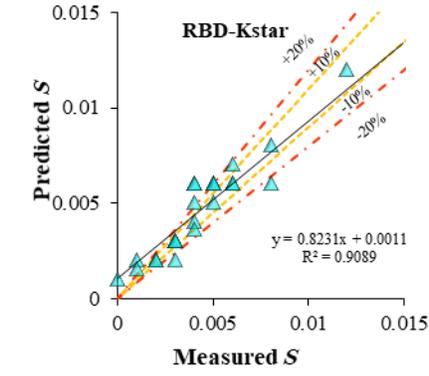
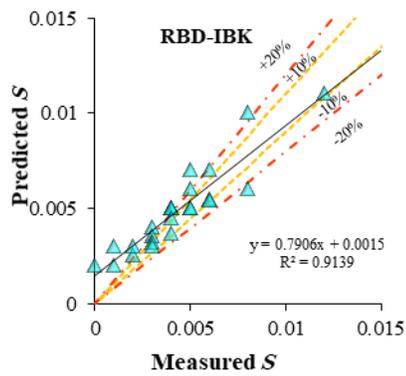
434

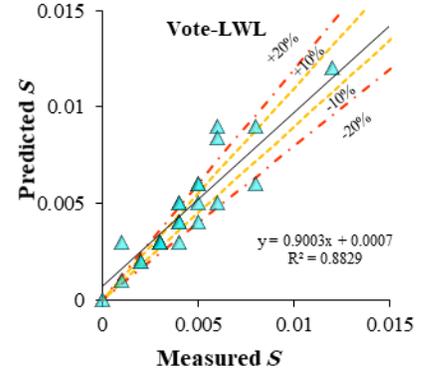
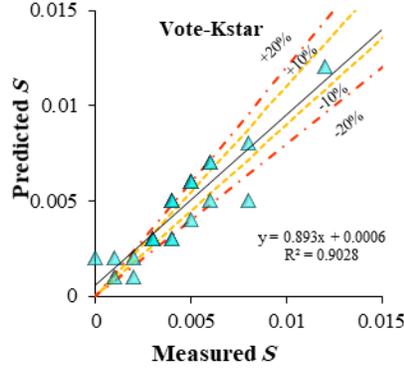
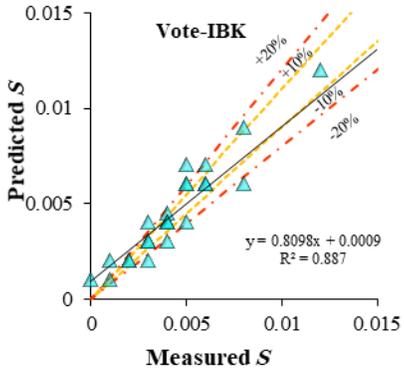
435

436



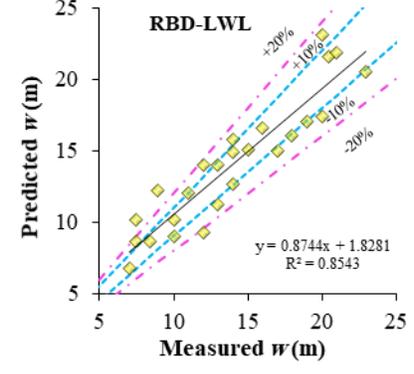
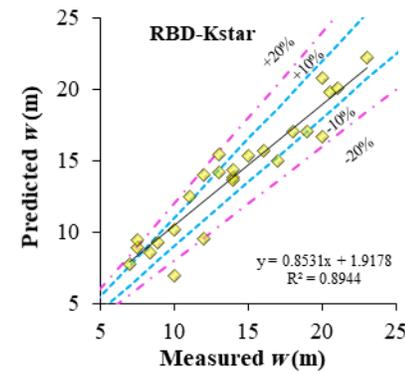
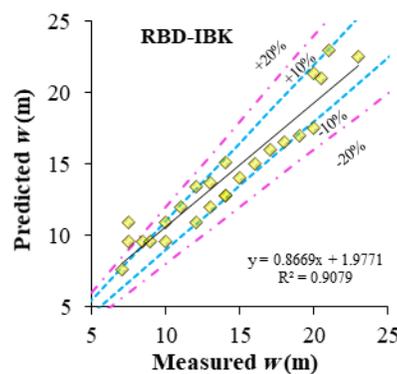
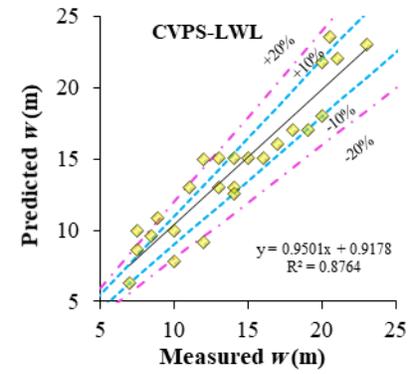
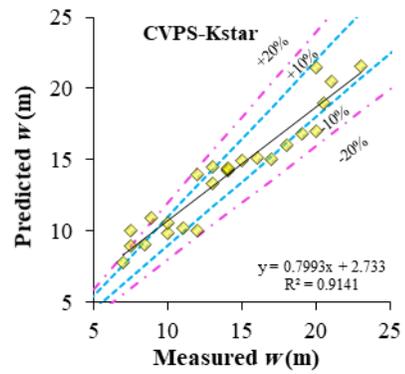
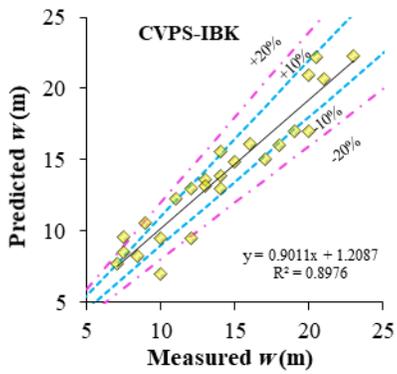
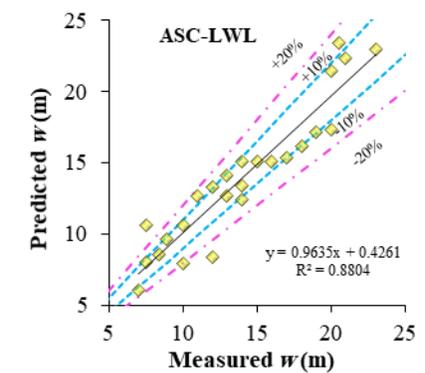
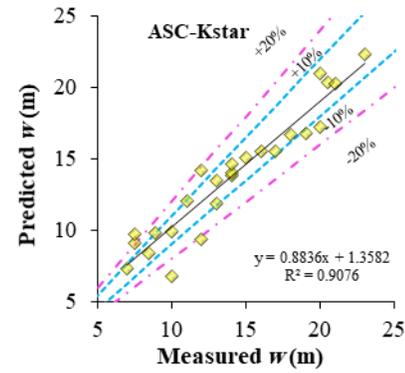
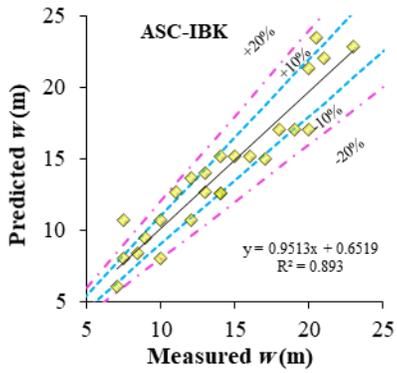
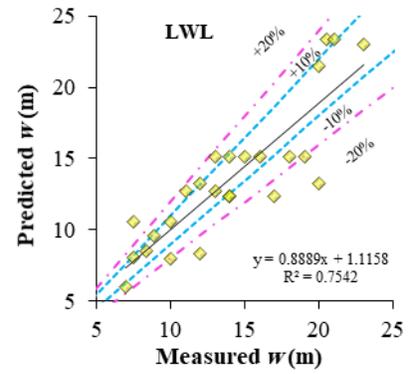
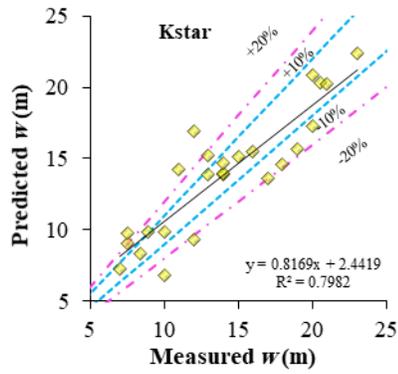
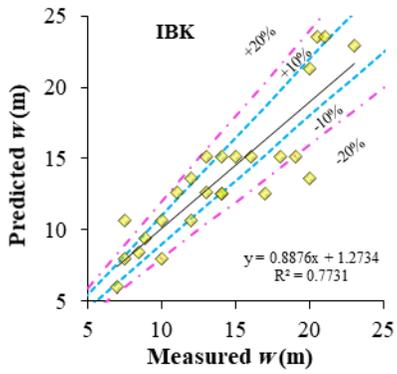
437

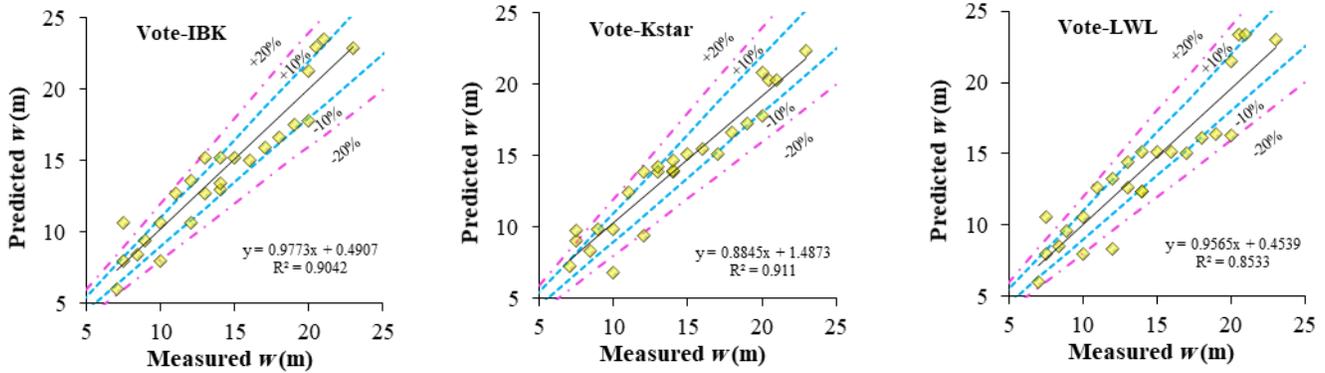




438  
 439  
 440  
 441  
 442  
 443  
 444

Fig 5. Scatter plot of measured versus predicted longitudinal slope  $S$ .





446  
447

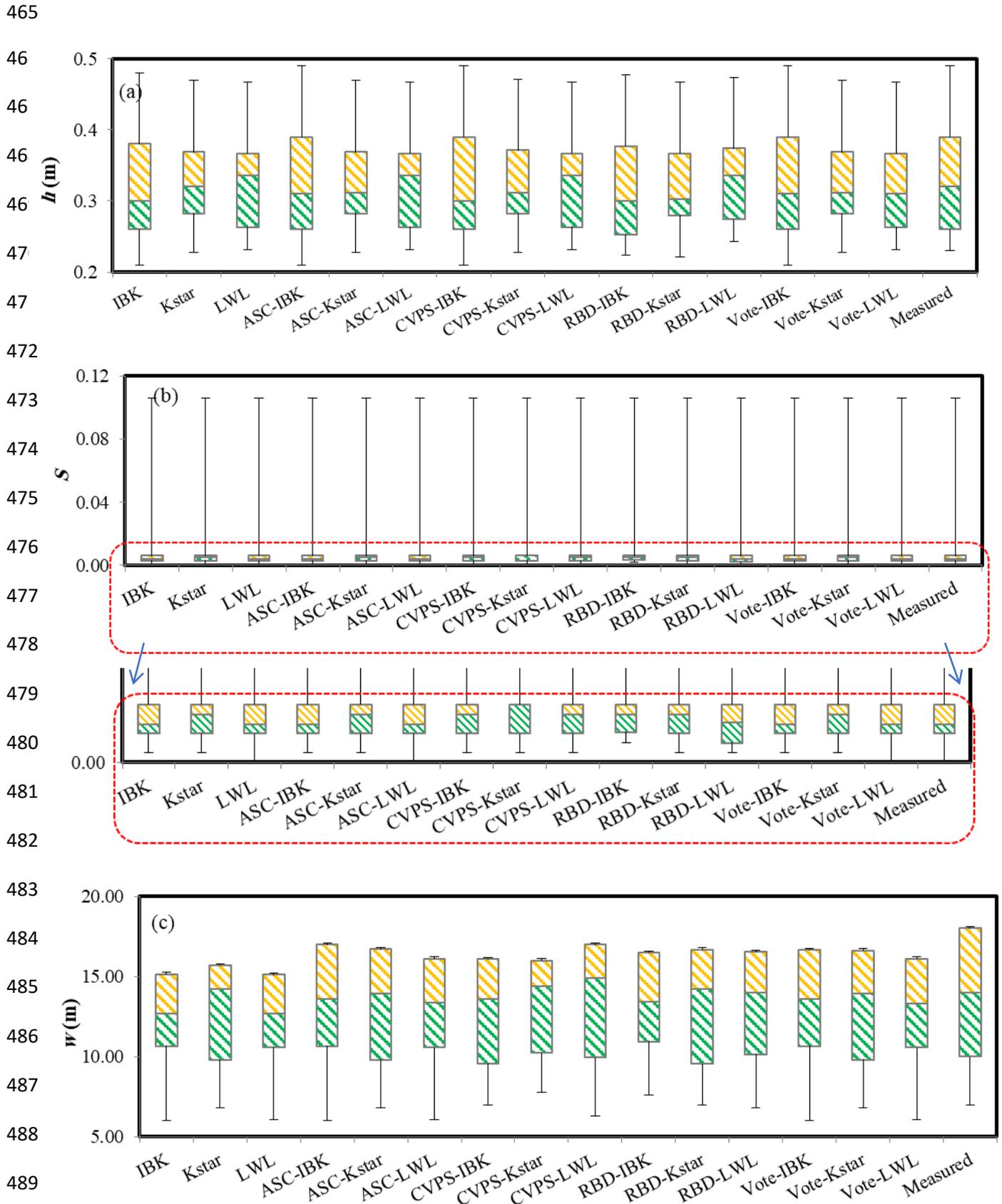
448 Fig 6. Scatter plot of measured versus predicted water-surface width  $w$ .

449  
450

451 Box plots of measured and predicted hydraulic geometry dimensions shows that the hybrid models ASC-  
 452 IBK, CVPS-IBK and Vote-IBK predicted the maximum and third quartile depth well, and the IBK  
 453 standalone algorithm was reasonably accurate in predicting the maximum. Kstar was the only model to  
 454 predict the median depth well. In terms of the first quartile, IBK, LWL, ASC-LWL, CVPS-IBK, CVPS-  
 455 LWL, Vote-IBK, and Vote-LWL were the most accurate, and the LWL, ASC-LWL, CVPS-LWL, Vote-  
 456 Kstar, Vote-LWL, CVPS-Kstar, ASC-Kstar and Kstar model were the most accurate for the minimum  
 457 channel depth.

458 All algorithms predicted the maximum and third quartile slope well (Figure 7(b)), but only RBD-LWL,  
 459 Vote-IBK and Vote-LWL were able to predict median slope accurately. All algorithms provided good  
 460 estimates of the first quartile, except the standalone algorithms and RBD-LWL. The minimum slope was  
 461 well reproduced by the LWL, ASC-LWL and Vote-LWL models.

462 In contrast, none of the algorithms were able to predict maximum and third quartile width accurately. But  
 463 Vote-Kstar, RBD-LWL and RBD-Kstar models predicted median values very well, and RBD-LWL and  
 464 CVPS-Kstar did likewise for the third quartile.



490 Fig 7. Box plot of measured and predicted hydraulic geometry: (a) flow depth, (b) longitudinal water  
491 surface slope and (c) water surface width.

492

493 Since the coefficient of determination  $R^2$  is standardised for differences between the mean and variance of  
494 measured and predicted values, this metric is sensitive to outliers and should not be used for model  
495 evaluation alone (Legates and McCabe, 1999; Shiri and Kisi, 2012). Thus other evaluation metrics were  
496 considered and are shown in Table 3. The metrics of model performance reveal that Vote-Kstar algorithm  
497 had the highest prediction power for depth ( $RMSE = 0.0292$  m,  $MAE = 0.0241$ ,  $NSE = 0.872$ ) followed by  
498 RBD-LWL ( $RMSE = 0.0304$  m,  $MAE = 0.0229$  m,  $NSE = 0.862$ ) and CVPS-Kstar ( $RMSE = 0.0317$  m,  
499  $MAE = 0.0251$  m,  $NSE = 0.850$ ) (Table 4). The best performing model (Vote-Kstar) had 49.5 %, 7.8 %  
500 and 19.2 % higher prediction capability than the IBK, Kstar and LWL standalone algorithms, based on  
501 the  $NSE$  metric. According to the  $NSE$  values, IBK model had an ‘acceptable’ performance, ASC-IBK  
502 had a ‘satisfactory’ performance, LWL, ASC-LWL, CVPS-IBK and CVPS-LWL had ‘good’ prediction  
503 power, and the rest of algorithms had ‘very good’ performance.

504 Differing results were found in the prediction of slope. The ASC-Kstar algorithm ( $RMSE = 0.001$  m,  
505  $MAE = 0.0008$  m, and  $NSE = 0.904$ ) outperformed other algorithms, followed by Vote-Kstar ( $RMSE =$   
506  $0.001$ ,  $MAE = 0.0008$  m,  $NSE = 0.902$ ), RBD-Kstar ( $RMSE = 0.0011$  m,  $MAE = 0.0007$  m,  $NSE = 0.897$ ).  
507 In terms of  $NSE$ , ASC-Kstar, as the most accurate model, had 27.0 %, 13.8 % and 29.3 % higher  
508 prediction power than the IBK, Kstar and LWL standalone models respectively. LWL had an ‘acceptable’  
509 performance, IBK had a ‘good’ performance, and other algorithms had a ‘very good’ prediction power.

510 Vote-Kstar outperformed all algorithms ( $RMSE = 1.373$  m,  $MAE = 1.059$  m,  $NSE = 0.909$ ) for the  
511 estimation of width, as also observed for depth,, followed by RBD-IBK ( $RMSE = 1.401$  m,  $MAE = 1.206$   
512 m, a  $NSE = 0.905$ ), ASC-Kstar ( $RMSE = 1.418$  m,  $MAE = 1.065$  m,  $NSE = 0.903$ ). In terms of  $NSE$ , the  
513 Vote-Kstar model had about 17.4 %, 12.4 % and 20.8 % higher performance than the standalone models.  
514 LWL had ‘good’ prediction and other algorithms had ‘very good’ performance.

515 According to the *PBIAS* metric, all developed algorithms under-estimated depth except RBD-Kstar and  
 516 RBD-LWL models, over-estimated slope except IBK and Kstar, and under-estimated width except CVPS-  
 517 LWL, RBD-IBK, RBD-LWL and Vote-IBK.

518 All model performance metrics reveal that although hybridisation enhances the prediction power of  
 519 standalone algorithms, the level of enhancement and overall performance of hybridised algorithms were  
 520 strongly dependent upon the choice of standalone algorithm. For instance, in the prediction of depth, the  
 521 use of Vote to hybridise IBK increased the *NSE* by 72 % but by just 9 % in the case of Kstar. But despite  
 522 this increase, the standalone Kstar algorithm (*NSE* = 0.80) still had a higher performance than the hybrid  
 523 Vote-IBK model (*NSE* = 0.76).

524

525 Table 4. Evaluation of model performance

Variable	Models	$R^2$	<i>RMSE</i> (m)	<i>MAE</i> (m)	<i>NSE</i>	<i>PBIAS</i> (%)	Rank based on <i>NSE</i>	Percentage lower performance than the best model according to <i>NSE</i>
<i>h</i>	IBK	0.62	0.06	0.04	0.44	2.61	13	49.46
	Kstar	0.81	0.04	0.03	0.80	0.56	6	7.98
	LWL	0.72	0.04	0.03	0.71	3.06	10	19.22
	ASC-IBK	0.70	0.05	0.03	0.59	1.72	12	32.16
	ASC-Kstar	0.84	0.03	0.03	0.84	0.78	4	4.12
	ASC-LWL	0.72	0.04	0.03	0.71	3.06	10	19.22
	CVPS-IBK	0.75	0.05	0.03	0.68	2.40	11	22.08
	CVPS-Kstar	0.85	0.03	0.03	0.85	0.75	3	2.60
	CVPS-LWL	0.72	0.04	0.03	0.71	3.06	10	19.22
	RBD-IBK	0.78	0.04	0.03	0.75	4.16	9	13.60
	RBD-Kstar	0.83	0.03	0.03	0.84	-0.17	5	4.89
	RBD-LWL	0.86	0.03	0.02	0.86	-0.20	2	1.46
	Vote-IBK	0.80	0.04	0.03	0.76	1.90	8	13.02
	Vote-Kstar	0.88	0.03	0.02	0.87	1.62	1	-----
	Vote-LWL	0.79	0.04	0.03	0.78	2.87	7	11.76
	Models	$R^2$	<i>RMSE</i>	<i>MAE</i>	<i>NSE</i>	<i>PBIAS</i> (%)	Rank based on <i>NSE</i>	Lower performance than the best model (%) according to

								NSE
S	IBK	0.66	0.0019	0.0010	0.6608	0.8333	14	26.991
	Kstar	0.79	0.0016	0.0010	0.7797	2.5000	13	13.827
	LWL	0.64	0.0020	0.0010	0.6399	-0.8333	15	29.314
	ASC-IBK	0.86	0.0012	0.0008	0.8608	-1.7500	11	4.867
	ASC-Kstar	0.91	0.0010	0.0008	0.9042	-2.0583	1	-----
	ASC-LWL	0.88	0.0012	0.0007	0.8776	-4.1667	9	2.986
	CVPS-IBK	0.89	0.0011	0.0009	0.8862	-4.0833	5	1.991
	CVPS-Kstar	0.86	0.0013	0.0010	0.8566	-5.8333	12	5.309
	CVPS-LWL	0.87	0.0012	0.0008	0.8674	-4.9167	10	4.092
	RBD-IBK	0.91	0.0012	0.0009	0.8796	-9.3750	7	2.765
	RBD-Kstar	0.90	0.0011	0.0007	0.8972	-4.2500	3	1.32
	RBD-LWL	0.91	0.0011	0.0008	0.8906	-3.0833	4	1.548
	Vote-IBK	0.88	0.0011	0.0007	0.8802	-0.4167	6	2.654
	Vote-Kstar	0.90	0.0010	0.0008	0.9021	-1.6667	2	0.221
	Vote-LWL	0.88	0.0012	0.0007	0.8785	-4.5000	8	2.876

								Lower performance than the best model (%) according to NSE
Models	$R^2$	$RMSE$ (m)	$MAE$ (m)	$NSE$	PBIAS (%)	Rank based on $NSE$		
IBK	0.77	2.27	1.71	0.75	2.19	14	17.38	
Kstar	0.79	2.06	1.50	0.80	0.96	13	12.43	
LWL	0.75	2.41	1.82	0.72	3.18	15	20.79	
ASC-IBK	0.89	1.52	1.24	0.89	0.24	8	2.20	
ASC-Kstar	0.90	1.42	1.07	0.90	1.99	3	0.66	
ASC-LWL	0.88	1.63	1.32	0.87	0.62	9	4.07	
CVPS-IBK	0.89	1.47	1.17	0.90	1.30	6	1.43	
CVPS-Kstar	0.91	1.45	1.18	0.90	0.65	4	1.10	
CVPS-LWL	0.87	1.66	1.38	0.87	-1.53	10	5.39	
RBD-IBK	0.90	1.40	1.21	0.91	-0.74	2	0.44	
RBD-Kstar	0.89	1.50	1.18	0.89	1.06	7	1.98	
RBD-LWL	0.85	1.74	1.45	0.85	-0.43	11	6.16	
Vote-IBK	0.90	1.46	1.21	0.90	-1.22	5	1.32	
Vote-Kstar	0.91	1.37	1.06	0.91	0.98	1	-----	
Vote-LWL	0.85	1.83	1.49	0.84	1.13	12	7.70	

526

527

## 528 4. Discussion

### 529 4.1 Effect of input variables on model prediction performance

530 The combination of input variables had a strong effect on model prediction power, confirming that the  
531 determination of the optimum combination is one of the most significant steps in producing an accurate  
532 data mining model. For example, the best input combination for the prediction of flow depth using the  
533 Vote-LWL model had ~51 % higher prediction performance (in terms of *NSE*) than the worst input  
534 combination. The optimum input variable combination was different from one model to another one,  
535 resulting from the different structure of each model, particularly in terms of their flexibility, computing  
536 capability and complexity. Thus a range of different input variable combinations must be considered in  
537 the optimisation of data mining models.

538 To determine this optimum input combination, this paper used a manual approach, building and testing  
539 numerous input combinations. Others have used Principal Component Analysis (*PCA*) (e.g. Barzegar *et*  
540 *al.*, 2017) or a gamma test (e.g. Ahmadi *et al.*, 2015) of the input and output data to determine just one  
541 input combination automatically. Determining the optimum combination manually can produce models  
542 with a higher prediction performance. For example, in the prediction of fluoride concentration in  
543 groundwater Khosravi *et al.* (2019a) built eight different input combinations, and Barzegar *et al.* (2017),  
544 using the same dataset, applied *PCA* to extract the best input combination. The manually derived input  
545 variable combination produced a 27.5 % higher prediction performance (in terms of *NSE*) than the one  
546 extracted by *PCA*, highlighting the need to first conduct a sensitivity analysis to establish the range of  
547 input combinations that need to be considered manually.

548 According to the findings of this paper, excluding Shields stress from the input combination in the  
549 prediction of flow depth caused a 49.2 % decrease in model prediction power, and was the most  
550 influential variable on model prediction performance, followed by *Q* (28.0 %) and *d*<sub>50</sub> (21.4 %) (Figure  
551 8). Very similar results were found for longitudinal slope; Shields stress caused a 54.0 % change in  
552 prediction power and was the most effective parameter, followed by *d*<sub>50</sub> (22.3 %) and *Q* (7.7 %), in line  
553 with previous results (Julien and Wargadalam, 1995; Afzalimehr *et al.*, 2010; Gholami *et al.*, 2017;

554 Shaghaghi *et al.*, 2018). Omitting flow discharge as an input variable increased model prediction  
555 performance, showing that Shields stress and  $d_{50}$  are only required to predict accurately slope.

556 In the prediction of top width,  $d_{50}$  caused a 15.4 % change in model prediction power, and was the most  
557 effective parameter, followed by Shields stress (13.3 %) and discharge (1.41 %). When compared to the  
558 effect on the prediction of depth and slope,  $d_{50}$  had a much lower impact on width, resulting from the low  
559 correlation between  $d_{50}$  and width in alluvial rivers. For example, width more strongly depends on the  
560 characteristics of the bank material, such as the percentage of bank vegetation growth (e.g. Hey and  
561 Thorne, 1986; Bettess *et al.*, 1988; Gholami *et al.*, 2017) than of the bed materials.

562 These results on the most effective parameters are intuitive, given the strong correlation between  
563 sediment transport rate - and thus channel form - and Shields stress, and the weaker correlations with  $Q$   
564 and  $d_{50}$  (Julien and Wargadalam 1995). For example,  $Q$  only has an in-direct influence on sediment  
565 transport through its correlation with Shields stress. In other words, two channels, one wide and one  
566 narrow, or one shallow and one steep, with the same  $d_{50}$  can experience the same  $Q$  but different Shields  
567 stress and thus sediment transport rates. However, contrasting results on the most influential factors on  
568 depth and width have been found. For example, numerous studies have found discharge to be the most  
569 important factor, followed by Shields stress and  $d_{50}$  (Afzalimehr *et al.* 2009, 2010; Bray 1982; Hey and  
570 Thorne 1986a). Abdelhaleem *et al.* (2016) showed that as well as flow discharge other controlling  
571 factors such as flow velocity must be incorporated to increase prediction accuracy. Thus, the most  
572 effective input parameter is not constant and differs from one river to another according to morphology,  
573 such as the presence of bedforms, large woody debris, vegetation and changes in channel planform.  
574 Therefore the models presented here, which are statistical in nature, apply only to the three rivers  
575 considered and river swith similar conditions, and should not be applied universally.

576  
577

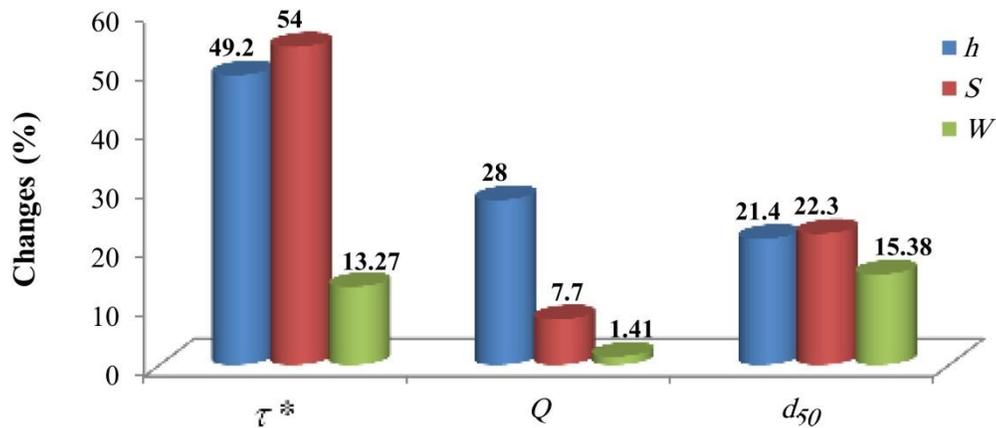


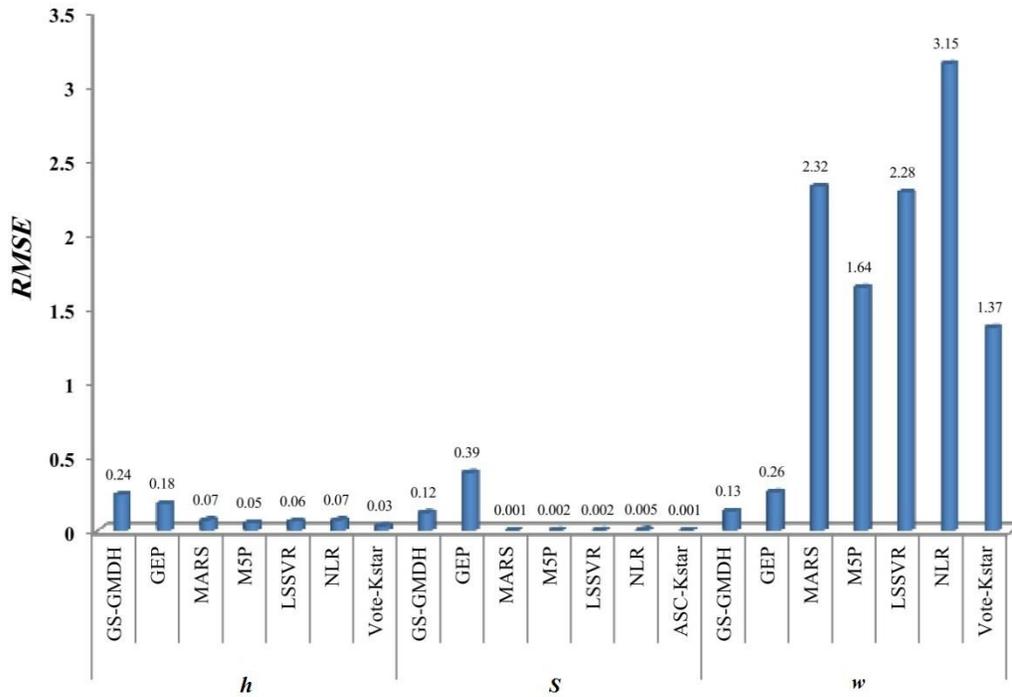
Fig 8. The percentage change in model *RMSE* with each input variable

578  
579

580

581 **4.2 Comparison in model prediction performance between empirical, traditional and advanced data**  
582 **mining models**

583 The Afzalimehr et al. (2010) dataset used in this study provides a unique opportunity to compare the  
584 performance of empirical equations, traditional machine learning algorithms with the newly developed,  
585 advanced data mining models directly: Afzalimehr et al. (2010) tested the performance of empirical and  
586 NLR models using this dataset, and Shaghaghi et al. (2018a,b) used the dataset to test traditional data  
587 mining models (a hybrid model GS-GMDH, and two standalone models, GEP, Multivariate Adaptive  
588 Regression Splines (MARS), Least Square Support Vector Regression (LSSVR) and NLR). Figure 9  
589 shows the results of this comparison in performance, revealing that, the newly developed advanced data  
590 mining models outperformed the NLR and traditional data mining models in most of the cases. For  
591 example, in modelling depth, the newly developed an ASC-Kstar model produced the lowest *RMSE* value  
592 of 0.03 m using all the variables ( $Q$ ,  $d_{50}$  and  $\tau^*$ ) as input, comparing favourably to the results for the GS-  
593 GMDH (*RMSE* = 0.24 GEP (*RMSE* = 0.18 m), MARS (*RMSE* = 0.07 m) and NLR (*RMSE* = 0.07 m)  
594 models.



595

596 Fig 9. Comparison between model performance of the present study with literature in terms of depth ( $h$ ),  
 597 slope and width variables.

598 These comparisons reveal that the new hybrid models proposed in this study are more flexible and  
 599 accurate than traditional machine learning standalone and hybrid models in most of the cases. The reasons  
 600 are three-fold. First traditional models are neuron based and need to be optimised to get high prediction  
 601 power, especially in the determination of the weights of membership function. Advanced data mining  
 602 models such as tree, lazy, and rule-based models do not have this weakness. Second, NLR models are  
 603 regression based models and due to their simple structure, are not capable of predicting complicated  
 604 phenomena accurately. Finally hybridisation improves the performance of standalone models because the  
 605 process develops a coupled model with higher flexibility, which is proven to better reproduce complex,  
 606 nonlinear processes that are at play in rivers (Khosravi *et al.* 2020).

607 **4.3 Applying advanced data mining models to forecast stable channel geometry**

608 The choice of the ‘best’ predictive model is most often a compromise between model prediction accuracy  
609 and model complexity, with the later, in data mining models, most closely related to the data input  
610 requirements. In some data mining models, the highest accuracy has been achieved using all input  
611 variables (e.g. Bui *et al.*, 2020b; Khosravi *et al.*, 2020), whilst in others, the best prediction power has  
612 been obtained with a less complex model using fewer input variables (e.g. Sheikh Khozani *et al.*, 2017b;  
613 2019). The major advantages of the data mining models developed in this paper are their simplicity, and  
614 their ease and inexpensive to build and run, unlike theoretical and numerical models, whilst providing  
615 little compromise on model performance. In other words, a number of the advanced hybrid data mining  
616 models provided very good prediction performance for depth and width based on just three input  
617 parameters, and for slope based on just two. In stable channel design, predictions of channel geometry are  
618 often constrained by the availability of channel data, making less complex models more desirable. Thus  
619 the results reveal that these models have great potential for use in stable channel design in data poor  
620 catchments, especially in developing nations where technical modelling skills and understanding of the  
621 hydraulic and sediment processes occurring in the river system may be lacking.

622 The major disadvantages of these types of model however are two-fold. First, like all statistical methods,  
623 the developed models only relate directly to the rivers being considered, and thus their application to  
624 other rivers may prove inappropriate. Future studies should apply the developed models to rivers with  
625 differing morphologies to discover whether this is the case. Second, due to their ‘black-box’ structure,  
626 they provide poor explanatory power, and thus are unable to extract understanding on the physics that  
627 determine hydraulic geometry.

628 With these considerations in mind, the use of data mining techniques may not simply lie in predicting  
629 stable channel parameters, but integrating these techniques into process-based models to help identify and  
630 optimise model parameters and mitigate uncertainty in model estimates (e.g. Vojinovic *et al.*, 2013),d  
631 help recognise patterns within observational data to unveil critical details about behavior, and possibly  
632 reveal new environmental relationships. Future studies should seek to explore this potential.

633 This study has only considered three controlling parameters. Where data is available, future studies  
634 should consider other factors in data mining models, such as flow velocity, relative roughness, suspended  
635 sediment load, and bed load transport rate, vegetation form, channel planform, channel roughness, Froude  
636 and Reynolds number, and sediment composition (e.g. Abernethy, 2000; Davidson and Hey, 2011),  
637 helping to determine the most influential parameters on stable hydraulic geometry and why they vary  
638 between rivers.

639

## 640 **5. Conclusion**

641  
642 Using at-a station field data, this paper has quantified, for the first time, the potential of advanced data  
643 mining algorithms to provide accurate predictions of stable hydraulic geometry. Predictions of mean flow  
644 depth, top-width and longitudinal slope were made using three standalone data mining techniques -  
645 Instance-based Learning (IBK), KStar, Locally Weighted Learning (LWL) - along with 12 types of novel  
646 hybrid algorithms in which the standalone models were trained with Vote, Attribute Selection  
647 Committees (ASC), Regression by Discretization (RBD), and Cross-validation Parameter Selection  
648 (CVPS) algorithms. A comparison was made of the predictive power of these data-driven models, and a  
649 sensitivity analysis of three driving variables (discharge, median bed grain diameter and Shields stress)  
650 was performed. The main findings were as follows:

- 651 1- Shield stress was the most effective variable on flow depth and slope prediction; excluding it as  
652 an input variable to models caused a 49.2 % and 54 % increase in relative error. Median sediment  
653 size had the greatest effect on width prediction power, and excluding this parameter caused a 15.4  
654 % increase in relative error. Overall, Shield stress parameter was the most effective parameter on  
655 all geometry dimensions.
- 656 2- The hybrid data mining models had a higher prediction power than standalone models, empirical  
657 equations and traditional machine learning algorithms because the hybrid models were more

658 flexible and thus could better reproduce the nonlinear interactions between input variables and  
659 hydraulic geometry. In particular, Vote-Kstar model had the highest prediction capability for  
660 depth and width prediction, and ASC-Kstar for slope,

661 3- According to Nash-Sutcliffe Efficiency values, the IBK model had an acceptable performance,  
662 ASC-IBK a satisfactory performance, LWL, ASC-LWL, CVPS-IBK and CVPS-LWL a good  
663 prediction power and the rest of the algorithms had a very good performance in flow depth  
664 prediction. In estimating slope, LWL had an acceptable performance, IBK a good performance  
665 and all other algorithms had a very good prediction accuracy. LWL had a good prediction power  
666 in width prediction and all other algorithms had a very good performance.

667 The strength of these hybrid algorithms lies in their ease to implement, use of a small number of input  
668 variables, and being inexpensive to build and run-in comparison to theoretical and numerical models,  
669 whilst providing little compromise on model performance. Together, these findings reveal that hybrid  
670 data mining models have great potential for use in stable channel design, especially in situations when  
671 understanding of the physical processes at play may not be well understood. Thus, understanding  
672 more about this potential for different river conditions and input variables represents a vital research  
673 avenue.

674

#### 675 **Authorship contribution statement**

676 **Khabat Khosravi:** Conceptualization, formal analysis, writing original draft (result and discussion),  
677 review & editing, **Zohreh Sheikh Khozani:** Data collection and writing original draft (Introduction,  
678 methodology and models description), **James R. Cooper:** writing, review & editing.

679

680 **Authors' Note:** The authors do not have any conflicts of interest or financial disclosures to report.

#### 681 ***Software availability***

#### 682 **Software Name**

683 Waikato Environment for Knowledge Analysis (WEKA) software.

684

## 685 **Availability**

686 Weka is open-source software and has been written in Java and developed at the University of Waikato,

687 New Zealand. It is free software licensed under the GNU General Public License. Software and

688 documentation (user manual and training material) are freely available at:

689 <https://www.cs.waikato.ac.nz/ml/weka/downloading.html>

690

## 691 **References**

692 Abdelhaleem, F. S., Amin, A. M., & Ibraheem, M. M. (2016). Updated regime equations for alluvial

693 Egyptian canals. *Alexandria Engineering Journal*, 55(1), 505–512.

694 <https://doi.org/10.1016/j.aej.2015.12.011>

695 Abernethy B, R. I. (2000). The effect of riparian tree roots on the mass-stability of riverbanks. *Earth*

696 *Surface Processes and Landforms*, 25(9), 921–937. <https://doi.org/10.1002/1096->

697 9837(200008)25:9<921::AID-ESP93>3.0.CO;2-7

698 Afzalimehr, H., Abdolhosseini, M., & Singh, V. P. (2010). Hydraulic geometry relations for stable

699 channel design. *Journal of Hydrologic Engineering*, 15(10), 859–864.

700 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000260](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000260)

701 Afzalimehr, H., Singh, V. P., & Abdolhosseini, M. (2009). Effect of nonuniformity of flow on hydraulic

702 geometry relations. *Journal of Hydrologic Engineering*, 14(9), 1028–1034.

703 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000095](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000095)

704 Ahmad, M. W., Reynolds, J., & Rezgui, Y. (2018). Predictive modelling for solar thermal energy

705 systems: A comparison of support vector regression, random forest, extra trees and regression trees.

706 *Journal of Cleaner Production*, 203, 810–821. <https://doi.org/10.1016/j.jclepro.2018.08.207>

707 Ahmadi, A., Han, D., Lafdani, E. K., & Moridi, A. (2015). Input selection for long-lead precipitation

708 prediction using large-scale climate variables: A case study. *Journal of Hydroinformatics*, 17(1),

709 114–129. <https://doi.org/10.2166/hydro.2014.138>

710 Anastasakis, L., & Mort, N. (2001). *The development of self-organization techniques in modelling: a*

711 *review of the group method of data handling (GMDH)*. *gmdhsoftware.com*. United Kingdom.

712 Antar, M. A., Ellassiouti, I., & Allam, M. N. (2006). Rainfall-runoff modelling using artificial neural

713 networks technique: A Blue Nile catchment case study. *Hydrological Processes*, 20(5), 1201–1216.  
714 <https://doi.org/10.1002/hyp.5932>

715 Arif, M., Ishihara, T., & Inooka, H. (2001). Incorporation of experience in iterative learning controllers  
716 using locally weighted learning. *Automatica*, 37(6), 881–888. <https://doi.org/10.1016/S0005->  
717 1098(01)00030-9

718 Atkeson, C. G., Moore, A. W., & Schaal, S. (1997). Locally Weighted Learning. *Artificial Intelligence*  
719 *Review*, 11(1–5), 11–73. [https://doi.org/10.1007/978-94-017-2053-3\\_2](https://doi.org/10.1007/978-94-017-2053-3_2)

720 Ayele, G. T., Teshale, E. Z., Yu, B., Rutherford, I. D., & Jeong, J. (2017). Streamflow and sediment yield  
721 prediction for watershed prioritization in the upper Blue Nile river basin, Ethiopia. *Water*  
722 *(Switzerland)*, 9(10), 782. <https://doi.org/10.3390/w9100782>

723 Barzegar, R., Asghari Moghaddam, A., Adamowski, J., & Fijani, E. (2017). Comparison of machine  
724 learning models for predicting fluoride contamination in groundwater. *Stochastic Environmental*  
725 *Research and Risk Assessment*, 31(10), 2705–2718. <https://doi.org/10.1007/s00477-016-1338-z>

726 Blench, T. (1969). *Mobile-bed fluviology*. Alberta Press, Edmonton, Canada. <https://doi.org/10.1016/0022->  
727 1694(70)90091-0

728 Blench, Thomas. (1952). Regime theory for self-formed sediment-bearing channels. *Transactions of the*  
729 *American Society of Civil Engineers*, 117(1), 383–400.

730 Bose, N. K. (1936). Silt movement and design of channels. In *Punjab Eng Congr*. Punjab, India.

731 Bray, D. I. (1982). Regime equations for gravel-bed rivers. In *Gravel bed rivers: Fluvial processes,*  
732 *engineering and management*, R. D. Hey, J. C. Bathurst, and C. R. Thorne, eds., Wiley (pp. 517–  
733 552). Chichester, U.K.

734 Bui, D. T., Khosravi, K., Karimi, M., Busico, G., Khozani, Z. S., Nguyen, H., et al. (2020). Enhancing  
735 nitrate and strontium concentration prediction in groundwater by using new data mining algorithm.  
736 *Science of the Total Environment*, 715, 136836. <https://doi.org/10.1016/j.scitotenv.2020.136836>

737 Bui, D. T., Khosravi, K., Li, S., Shahabi, H., Panahi, M., Singh, V. P., et al. (2018). New hybrids of  
738 ANFIS with several optimization algorithms for flood susceptibility modeling. *Water (Switzerland)*,  
739 10(9). <https://doi.org/10.3390/w10091210>

740 Bui, D. T., Khosravi, K., Tiefenbacher, J., Nguyen, H., & Kazakis, N. (2020). Improving prediction of  
741 water quality indices using novel hybrid machine-learning algorithms. *Science of the Total*  
742 *Environment*, 721. <https://doi.org/10.1016/j.scitotenv.2020.137612>

743 Chang, H. H. (1980). Stable alluvial canal design. *Journal of the Hydraulics Division, ASCE*, 106(HY5,  
744 Proc. Paper, 15420), 873–891.

745 Chen, W., Panahi, M., & Pourghasemi, H. R. (2017). Performance evaluation of GIS-based new ensemble  
746 data mining techniques of adaptive neuro-fuzzy inference system (ANFIS) with genetic algorithm

747 (GA), differential evolution (DE), and particle swarm optimization (PSO) for landslide spatial  
748 modelling. *Catena*, 157, 310–324. <https://doi.org/10.1016/j.catena.2017.05.034>

749 Choubin, B., Darabi, H., Rahmati, O., Sajedi-Hosseini, F., & Kløve, B. (2018). River suspended sediment  
750 modelling using the CART model: A comparative study of machine learning techniques. *Science of*  
751 *the Total Environment*, 615, 272–281. <https://doi.org/10.1016/j.scitotenv.2017.09.293>

752 Cleary, J. G., & Trigg, L. E. (1995). K\*: An Instance-based Learner Using an Entropic Distance Measure.  
753 In *Machine Learning Proceedings 1995* (pp. 108–114). [https://doi.org/10.1016/b978-1-55860-377-](https://doi.org/10.1016/b978-1-55860-377-6.50022-0)  
754 [6.50022-0](https://doi.org/10.1016/b978-1-55860-377-6.50022-0)

755 Cuest Cordoba, G. A., Tuhovčák, L., & Tauš, M. (2014). Using artificial neural network models to assess  
756 water quality in water distribution networks. In *Procedia Engineering* (Vol. 70, pp. 399–408).  
757 Elsevier Ltd. <https://doi.org/10.1016/j.proeng.2014.02.045>

758 Davidson, S. K., & Hey, R. D. (2011). Regime equations for natural meandering cobble- and gravel-bed  
759 rivers. *Journal of Hydraulic Engineering*, 137(9), 894–910.  
760 [https://doi.org/10.1061/\(ASCE\)HY.1943-7900.0000408](https://doi.org/10.1061/(ASCE)HY.1943-7900.0000408)

761 Dawson, C. W., Abrahart, R. J., & See, L. M. (2007). HydroTest: A web-based toolbox of evaluation  
762 metrics for the standardised assessment of hydrological forecasts. *Environmental Modelling and*  
763 *Software*, 22(7), 1034–1052. <https://doi.org/10.1016/j.envsoft.2006.06.008>

764 Deshpande, V., & Kumar, B. (2012). Review and assessment of the theories of stable alluvial channel  
765 design. *Water Resources*, 39(4), 481–487. <https://doi.org/10.1134/S0097807812040033>

766 Dietterich, T. G. (1997). Machine learning research\_ four current directions. *AI Magazine*, 18(4), 97–1.

767 Eaton, B. C., & Church, M. (2007). Predicting downstream hydraulic geometry: A test of rational regime  
768 theory. *Journal of Geophysical Research: Earth Surface*, 112(3), 3025.  
769 <https://doi.org/10.1029/2006JF000734>

770 Ferguson, R. I. (1986). Hydraulics and hydraulic geometry. *Progress in Physical Geography*, 10(1), 1–31.  
771 <https://doi.org/10.1177/030913338601000101>

772 Ferreira, C. (2001). Gene Expression Programming: a New Adaptive Algorithm for Solving Problems.  
773 *Complex System*, 13(2), 87–129.

774 Ferreira, C. (2002). Genetic representation and genetic neutrality in gene expression programming.  
775 *Advances in Complex Systems*, 05(04), 389–408. <https://doi.org/10.1142/s0219525902000626>

776 Frank, E., & Bouckaert, R. R. (2009). Conditional density estimation with class probability estimators. In  
777 *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and*  
778 *Lecture Notes in Bioinformatics)* (Vol. 5828 LNAI, pp. 65–81). Springer, Berlin, Heidelberg.  
779 [https://doi.org/10.1007/978-3-642-05224-8\\_7](https://doi.org/10.1007/978-3-642-05224-8_7)

780 Garg, T., & Khurana, S. S. (2014). Comparison of classification techniques for intrusion detection dataset

781 using WEKA. In *International Conference on Recent Advances and Innovations in Engineering,*  
782 *ICRAIE 2014.* <https://doi.org/10.1109/ICRAIE.2014.6909184>

783 Gholami, A., Bonakdari, H., Ebtehaj, I., Shaghghi, S., & Khoshbin, F. (2017). Developing an expert  
784 group method of data handling system for predicting the geometry of a stable channel with a gravel  
785 bed. *Earth Surface Processes and Landforms*, 42(10), 1460–1471. <https://doi.org/10.1002/esp.4104>

786 Gislason, P. O., Benediktsson, J. A., & Sveinsson, J. R. (2006). Random forests for land cover  
787 classification. In *Pattern Recognition Letters* (Vol. 27, pp. 294–300). North-Holland.  
788 <https://doi.org/10.1016/j.patrec.2005.08.011>

789 Gleason, C. J. (2015). Hydraulic geometry of natural rivers: A review and future directions. *Progress in*  
790 *Physical Geography*, 39(3), 337–360. <https://doi.org/10.1177/0309133314567584>

791 Hastie, T., & Loader, C. (1993). Local regression: Automatic Kernel carpentry. *Statistical Science*, 8(2),  
792 120–129. <https://doi.org/10.1214/ss/1177011002>

793 Henderson, F. M. (1961). Stability of alluvial channels. *Journal of the Hydraulics Division*, 87(6), 109–  
794 138.

795 Hey, R. D., & Thorne, C. R. (1986a). Stable channels with mobile gravel beds. *Journal of Hydraulic*  
796 *Engineering*, 112(8), 671–689. [https://doi.org/10.1061/\(ASCE\)0733-9429\(1986\)112:8\(671\)](https://doi.org/10.1061/(ASCE)0733-9429(1986)112:8(671))

797 Hey, R. D., & Thorne, C. R. (1986b). Stable Channels with Mobile Gravel Beds. *Journal of Hydraulic*  
798 *Engineering*, 112(8), 671–689. [https://doi.org/10.1061/\(asce\)0733-9429\(1986\)112:8\(671\)](https://doi.org/10.1061/(asce)0733-9429(1986)112:8(671))

799 Hooshyaripor, F., Tahershamsi, A., & Golian, S. (2014). Application of copula method and neural  
800 networks for predicting peak outflow from breached embankments. *Journal of Hydro-Environment*  
801 *Research*, 8(3), 292–303. <https://doi.org/10.1016/j.jher.2013.11.004>

802 Huang, H. Q., & Nanson, G. C. (1998). The influence of bank strength on channel geometry: an  
803 integrated analysis of some observations. *Earth Surface Processes and Landforms*, 23(10), 865–876.  
804 [https://doi.org/10.1002/\(SICI\)1096-9837\(199810\)23:10<865::AID-ESP903>3.0.CO;2-3](https://doi.org/10.1002/(SICI)1096-9837(199810)23:10<865::AID-ESP903>3.0.CO;2-3)

805 Julien, P. Y., & Wargadalam, J. (1995). Alluvial channel geometry: Theory and applications. *Journal of*  
806 *Hydraulic Engineering*, 121(4), 312–325. [https://doi.org/10.1061/\(ASCE\)0733-9429\(1995\)121:4\(312\)](https://doi.org/10.1061/(ASCE)0733-9429(1995)121:4(312))

807

808 Khadangi, E., Madvar, H. R., & Kiani, H. (2009). Application of artificial neural networks in establishing  
809 regime channel relationships. In *2009 2nd International Conference on Computer, Control and*  
810 *Communication, IC4 2009.* <https://doi.org/10.1109/IC4.2009.4909224>

811 Khosravi, K., Barzegar, R., Miraki, S., Adamowski, J., Daggupati, P., Alizadeh, M. R., et al. (2019).  
812 Stochastic Modeling of Groundwater Fluoride Contamination: Introducing Lazy Learners.  
813 *Groundwater*, gwat.12963. <https://doi.org/10.1111/gwat.12963>

814 Khosravi, K., Cooper, J. R., Daggupati, P., Thai Pham, B., & Bui, D. T. (2020). Bedload transport rate

815 prediction: Application of novel hybrid data mining techniques. *Journal of Hydrology*, 585, 124774.  
816 <https://doi.org/10.1016/j.jhydrol.2020.124774>

817 Khosravi, K., Daggupati, P., Alami, M. T., Awadh, S. M., Ghareb, M. I., Panahi, M., et al. (2019).  
818 Meteorological data mining and hybrid data-intelligence models for reference evaporation  
819 simulation: A case study in Iraq. *Computers and Electronics in Agriculture*, 167.  
820 <https://doi.org/10.1016/j.compag.2019.105041>

821 Khosravi, K., Mao, L., Kisi, O., Yaseen, Z. M., & Shahid, S. (2018). Quantifying hourly suspended  
822 sediment load using data mining models: Case study of a glacierized Andean catchment in Chile.  
823 *Journal of Hydrology*, 567, 165–179. <https://doi.org/10.1016/j.jhydrol.2018.10.015>

824 Lane, E. W. (1957). A study of the shape of channels formed by natural streams flowing in erodible  
825 material. *U.S. Army Corps of Engineers Report*, (9), 141. [https://www.worldcat.org/title/study-of-](https://www.worldcat.org/title/study-of-the-shape-of-channels-formed-by-natural-streams-flowing-in-erodible-material/oclc/6671104)  
826 [the-shape-of-channels-formed-by-natural-streams-flowing-in-erodible-material/oclc/6671104.](https://www.worldcat.org/title/study-of-the-shape-of-channels-formed-by-natural-streams-flowing-in-erodible-material/oclc/6671104)  
827 Accessed 26 June 2021

828 Legates, D. R., & McCabe Jr, G. J. (1999). Evaluating the use of “goodness-of-fit” measures in  
829 hydrologic and hydroclimatic model validation. *Water resources research*, 35(1), 233–241.

830 Leopold, L., & Wolman, M. (1957). River channel patterns: braided, meandering, and straight. *USGS*  
831 *Professional Paper*, 282-B, 51.

832 Mehta, D., Yadav, S., & Anal, S. (2013). Geomorphic channel design and analysis using HEC-RAS  
833 hydraulic design functions. *Journal of Global Analysis*, 2(4), 90–93.

834 Millar, R. G. (2005). Theoretical regime equations for mobile gravel-bed rivers with stable banks.  
835 *Geomorphology*, 64(3–4), 207–220. <https://doi.org/10.1016/j.geomorph.2004.07.001>

836 Mislán, Haviluddin, Hardwinarto, S., Sumaryono, & Aipassa, M. (2015). Rainfall Monthly Prediction  
837 Based on Artificial Neural Network: A Case Study in Tenggara Station, East Kalimantan -  
838 Indonesia. In *Procedia Computer Science* (Vol. 59, pp. 142–151). Elsevier.  
839 <https://doi.org/10.1016/j.procs.2015.07.528>

840 Mohamed, H. I. (2013). Design of alluvial Egyptian irrigation canals using artificial neural networks  
841 method. *Ain Shams Engineering Journal*, 4(2), 163–171. <https://doi.org/10.1016/j.asej.2012.08.009>

842 Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D., & Veith, T. L. (2007).  
843 Model evaluation guidelines for systematic quantification of accuracy in watershed simulations.  
844 *Transactions of the ASABE*, 50(3), 885–900. <https://doi.org/10.13031/2013.23153>

845 Noori, R., Deng, Z., Kiaghadi, A., & Kachosangi, F. T. (2016). How reliable are ANN, ANFIS, and  
846 SVM techniques for predicting longitudinal dispersion coefficient in natural rivers? *Journal of*  
847 *Hydraulic Engineering*, 142(1). [https://doi.org/10.1061/\(ASCE\)HY.1943-7900.0001062](https://doi.org/10.1061/(ASCE)HY.1943-7900.0001062)

848 Parhami, B. (1994). Voting Algorithms. *IEEE Transactions on Reliability*, 43(4), 617–629.

849 <https://doi.org/10.1109/24.370218>

850 Parker, G., Wilcock, P. R., Paola, C., Dietrich, W. E., & Pitlick, J. (2007). Physical basis for quasi-  
851 universal relations describing bankfull hydraulic geometry of single-thread gravel bed rivers.  
852 *Journal of Geophysical Research: Earth Surface*, 112(4). <https://doi.org/10.1029/2006JF000549>

853 Reyes, O., Cano, A., Fardoun, H. M., & Ventura, S. (2018). A locally weighted learning method based on  
854 a data gravitation model for multi-target regression. *International Journal of Computational*  
855 *Intelligence Systems*, 11(1), 282–295. <https://doi.org/10.2991/ijcis.11.1.22>

856 Robinson, C. (1998). *Multi-objective optimisation of polynomial models for time series prediction using*  
857 *genetic algorithms and neural networks*. University of Sheffield, UK.

858 Shaghghi, S., Bonakdari, H., Gholami, A., Kisi, O., Shiri, J., Binns, A. D., & Gharabaghi, B. (2018).  
859 Stable alluvial channel design using evolutionary neural networks. *Journal of Hydrology*, 566, 770–  
860 782. <https://doi.org/10.1016/j.jhydrol.2018.09.057>

861 Shamshirband, S., Hashemi, S., Salimi, H., Samadianfard, S., Asadi, E., Shadkani, S., et al. (2020).  
862 Predicting Standardized Streamflow index for hydrological drought using machine learning models.  
863 *Engineering Applications of Computational Fluid Mechanics*, 14(1), 339–350.  
864 <https://doi.org/10.1080/19942060.2020.1715844>

865 Sheikh Khozani, Z., Bonakdari, H., & Ebtehaj, I. (2017). An analysis of shear stress distribution in  
866 circular channels with sediment deposition based on Gene Expression Programming. *International*  
867 *Journal of Sediment Research*, 32(4), 575–584. <https://doi.org/10.1016/J.IJSRC.2017.04.004>

868 Sheikh Khozani, Z., Bonakdari, H., & Zaji, A. H. (2017). Estimating the shear stress distribution in  
869 circular channels based on the randomized neural network technique. *Applied Soft Computing*, 58,  
870 441–448. <https://doi.org/10.1016/j.asoc.2017.05.024>

871 Sheikh Khozani, Z., Khosravi, K., Pham, B. T., Kløve, B., Wan Mohtar, W. H. M., & Yaseen, Z. M.  
872 (2019). Determination of compound channel apparent shear stress: application of novel data mining  
873 models. *Journal of Hydroinformatics*, 21(5), 798–811. <https://doi.org/10.2166/hydro.2019.037>

874 Shelley, J., & Parr, A. D. (2009). Using HEC-RAS hydraulic design functions for geomorphic channel  
875 design and analysis. In *Proceedings of World Environmental and Water Resources Congress 2009 -*  
876 *World Environmental and Water Resources Congress 2009: Great Rivers* (Vol. 342, pp. 3722–  
877 3731). Reston, VA: American Society of Civil Engineers. [https://doi.org/10.1061/41036\(342\)374](https://doi.org/10.1061/41036(342)374)

878 Singh, V. P., & Zhang, L. (2008). At-a-station hydraulic geometry relations, 1: Theoretical development.  
879 *Hydrological Processes*, 22(2), 189–215. <https://doi.org/10.1002/hyp.6411>

880 Sterling, M., & Knight, D. (2002). An attempt at using the entropy approach to predict the transverse  
881 distribution of boundary shear stress in open channel flow. *Stochastic environmental research and*  
882 *risk assessment*, 16(2), 127–142.

883 Stevens, M. A., & Nordin, C. F. (1987). Critique of the regime theory for alluvial channels. *Journal of*  
884 *Hydraulic Engineering*, 113(11), 1359–1380. [https://doi.org/10.1061/\(ASCE\)0733-](https://doi.org/10.1061/(ASCE)0733-)  
885 9429(1987)113:11(1359)

886 Stone, C. (1982). Optimal global rates of convergence for nonparametric regression. *Annals of Statistics*,  
887 10, 1040–1053.

888 Stone, M. (1974). Cross-Validatory Choice and Assessment of Statistical Predictions. Journal of the  
889 Royal Statistical Societ. *Journal of the Royal Statistical Society*, 36(2), 111–147.  
890 <https://doi.org/10.2307/2984809>

891 Taheri, K., Shahabi, H., Chapi, K., Shirzadi, A., Gutiérrez, F., & Khosravi, K. (2019). Sinkhole  
892 susceptibility mapping: A comparison between Bayes-based machine learning algorithms. *Land*  
893 *Degradation and Development*, 30(7), 730–745. <https://doi.org/10.1002/ldr.3255>

894 Tahershamsi, A., Majdzade Tabatabai, M. R., & Shirkhani, R. (2012). An evaluation model of artificial  
895 neural network to predict stable width in gravel bed rivers. *International Journal of Environmental*  
896 *Science and Technology*, 9(2), 333–342. <https://doi.org/10.1007/s13762-012-0036-8>

897 Thornton, C., Hutter, F., Hoos, H. H., & Leyton-Brown, K. (2013). Auto-WEKA: Combined selection  
898 and hyperparameter optimization of classification algorithms. In *Proceedings of the ACM SIGKDD*  
899 *International Conference on Knowledge Discovery and Data Mining* (Vol. Part F1288, pp. 847–  
900 855). Association for Computing Machinery. <https://doi.org/10.1145/2487575.2487629>

901 Vojinovic, Z., Abebe, Y. A., Ranasinghe, R., Vacher, A., Martens, P., Mandl, D. J., et al. (2013). A  
902 machine learning approach for estimation of shallow water depths from optical satellite images and  
903 sonar measurements. In *Journal of Hydroinformatics* (Vol. 15, pp. 1408–1424). IWA Publishing.  
904 <https://doi.org/10.2166/hydro.2013.234>

905 Wan Mohtar, W. H. M., Afan, H., El-Shafie, A., Bong, C. H. J., & Ab. Ghani, A. (2018). Influence of bed  
906 deposit in the prediction of incipient sediment motion in sewers using artificial neural networks.  
907 *Urban Water Journal*, 15(4), 296–302. <https://doi.org/10.1080/1573062X.2018.1455880>

908 Wang, Z., Xing, H., Li, T., Yang, Y., Qu, R., & Pan, Y. (2016). A modified ant colony optimization  
909 algorithm for network coding resource minimization. *IEEE Transactions on Evolutionary*  
910 *Computation*, 20(3), 325–342. <https://doi.org/10.1109/TEVC.2015.2457437>

911 White, W. R. (1982). Analytical approach to river regime. *Journal of the Hydraulics Division*, 108(10),  
912 1179–1193.

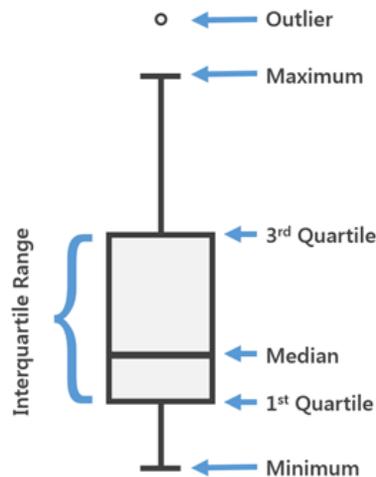
913 Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining: Practical Machine Learning Tools*  
914 *and Techniques*. Morgan Kaufmann. <https://doi.org/10.1016/c2009-0-19715-5>

915 Wolman, M. G. (1954). A method of sampling coarse river-bed material. *Eos, Transactions American*  
916 *Geophysical Union*, 35(6), 951–956. <https://doi.org/10.1029/TR035i006p00951>

917 Wu, C. H., Lin, I. S., Wei, M. L., & Cheng, T. Y. (2013). Target position estimation by genetic  
918 expression programming for mobile robots with vision sensors. *IEEE Transactions on*  
919 *Instrumentation and Measurement*, 62(12), 3218–3230. <https://doi.org/10.1109/TIM.2013.2272173>  
920 Zounemat-Kermani, M., Seo, Y., Kim, S., Ghorbani, M. A., Samadianfard, S., Naghshara, S., et al.  
921 (2019). Can decomposition approaches always enhance soft computing models? Predicting the  
922 dissolved oxygen concentration in the St. Johns River, Florida. *Applied Sciences (Switzerland)*,  
923 9(12). <https://doi.org/10.3390/app9122534>

924  
925  
926  
927  
928

### Supplementary material



929  
930

Fig A. Box-plot and its component in details