

Validation of a Low Complexity Machine Learning Discharge Predictive Model

Huma Zia¹, Nick Harris¹, Geoff Merrett¹, Mark Rivers²

¹*Electronics and Computer Science, University of Southampton, Southampton, United Kingdom*

²*Institute of Agriculture, University of Western Australia, Australia*

¹{hz2g11, nrh, gvm} @ ecs.soton.ac.uk, ²mark.rivers@uwa.edu.au

ABSTRACT

This paper reports on the validation of a simplified discharge prediction model that is suitable for implementation on a resourced constrained system such as a wireless sensor network, which will allow their operation to become more proactive rather than reactive. The data-driven model, utilising an M5 decision tree modelling technique, is validated using a 12-month training data set derived from published measured data. Daily runoff and drainage is predicted, and the results are compared with existing data-driven models developed in this domain. Results for the model give an R^2 of 0.82 and Root Relative Mean Square Error (RRMSE) of 35.9%. 80% of the residuals for the predicted test values fall within a ± 2 mm discharge depth/day error range. The main significance is that the proposed model gives comparable results with fewer samples and simpler parameters when compared to previous published research, which offers the potential for implementation in resource constrained monitoring and control systems.

Keyword: Wireless Sensor Networks, discharge prediction, environmental modelling, machine learning, M5 trees

1. Introduction

Over recent years, wireless sensor networks (WSNs), with their attractions of low cost and real time data availability, have received considerable attention in automating agricultural processes for economic benefits, e.g. in precision irrigation, pest control, and animal farming. However, a research gap still exists for mechanizing reutilization of resources (water and nutrients) among farms in order to additionally maximise environmental benefits. There is huge potential for leveraging existing networked agricultural activities into an integrated mechanism by sharing information about discharges (Zia et al., 2013). To illustrate this consider that the most commonly used irrigation method, surface irrigation, results in 40 to 60% of water losses in the form of runoff (Eisenhauer, 2011, Tindula et al., 2013). This runoff can transport up to 30-50% of applied nutrients to stream water and rivers (Liu et al., 2003) In the light of these figures, the motivation for this work is to develop a system that can potentially reduce water consumption and reduce outflows from farms, by predicting and

31 monitoring discharge from local areas. This will enable the development of systems that can then proactively
32 control irrigation strategies and also implement drainage reuse. This will also lead to improved water quality as
33 it will allow nutrients to be kept in the place where they can be useful where previously they would have been
34 discharged with no control into the local environment, eventually ending up in the streams and rivers. While
35 drainage reuse has been advocated and adopted in farming (Adelman, 2000, Willardson et al., 1997, Harper,
36 2012), various resource constraints and farmer's concerns regarding real time availability of information on
37 volumes, timings, and quality of discharges that will be delivered to the farms (Carr et al., 2011, Oster and
38 Grattan, 2002), currently restricts wide adoption of this mechanism in agriculture.

39 To address some of these issues, we have previously proposed a framework for water quality monitoring
40 control and management (WQMCM) using collaborative WSNs in a catchment to investigate and enable such a
41 mechanism (Zia et al., 2014a). The basic system architecture comprises various modules, one of which is a
42 discharge prediction module (Q -predictive model). The validation of this model using field data from an
43 instrumented catchment is the subject of this paper. Although previous work on the Q -predictive model has
44 shown that it works well with simulated data (Zia et al., 2014b), this paper extends this by reporting on the
45 validation of the model with field data from an instrumented catchment, and comparing its performance with
46 other published models.

47 To date, numerous physically-based hydrological models have been developed for the prediction of
48 discharges, either measured as surface runoff, groundwater leaching or stream-flow. Although these models are
49 popular in academic research and are very useful in evaluating different scenarios, their dependence on
50 acquiring numerous parameters, the need for calibrating models to individual areas, and the tremendous
51 computational burden involved in running the models makes wide-spread application complicated and difficult
52 (Basha et al., 2008, Galelli and Castelletti, 2013). In contrast, data-driven models have good prediction
53 capability and require fewer parameters, which is consistent with the requirement for a reduction in the
54 computational burden of decision making (Castelletti et al., 2010). Thus data-driven modelling, using machine
55 learning algorithms, has been widely used in hydrological modelling (Wilby et al., 2003, Rasouli et al., 2012,
56 Solomatine and Ostfeld, 2008) with artificial neural networks (ANN) being a popular choice (Dawson and
57 Wilby, 1998, Minns and Hall, 1996, Wilby et al., 2003). Recently, decision tree modelling has been investigated
58 (Galelli and Castelletti, 2013, Villa-Vialaneix et al., 2012, Fortin et al., 2014, Piñeros Garcet et al., 2006,
59 Kuzmanovski, 2012) and an interesting example of this class are M5 model trees (Quinlan, 1992). The
60 advantage of M5 model trees over ANNs are that they are faster to train and have guaranteed convergence

61 (Solomatine and Dulal, 2003). However, there are two limitations in the existing work; either the existing
62 models use simpler parameters but years of historical data with thousands of training samples to learn the
63 heterogeneity of large areas (>1000ha) (Galelli and Castelletti, 2013, Solomatine and Xue, 2004a), or they use
64 more complex models with a significant number of parameters (Bhattacharya et al., 2005, Kuzmanovski, 2012).
65 Additionally none of these approaches have been specifically targeted at sensor network applications, and the
66 data used was obtained through traditional sampling methods in gauged catchments.

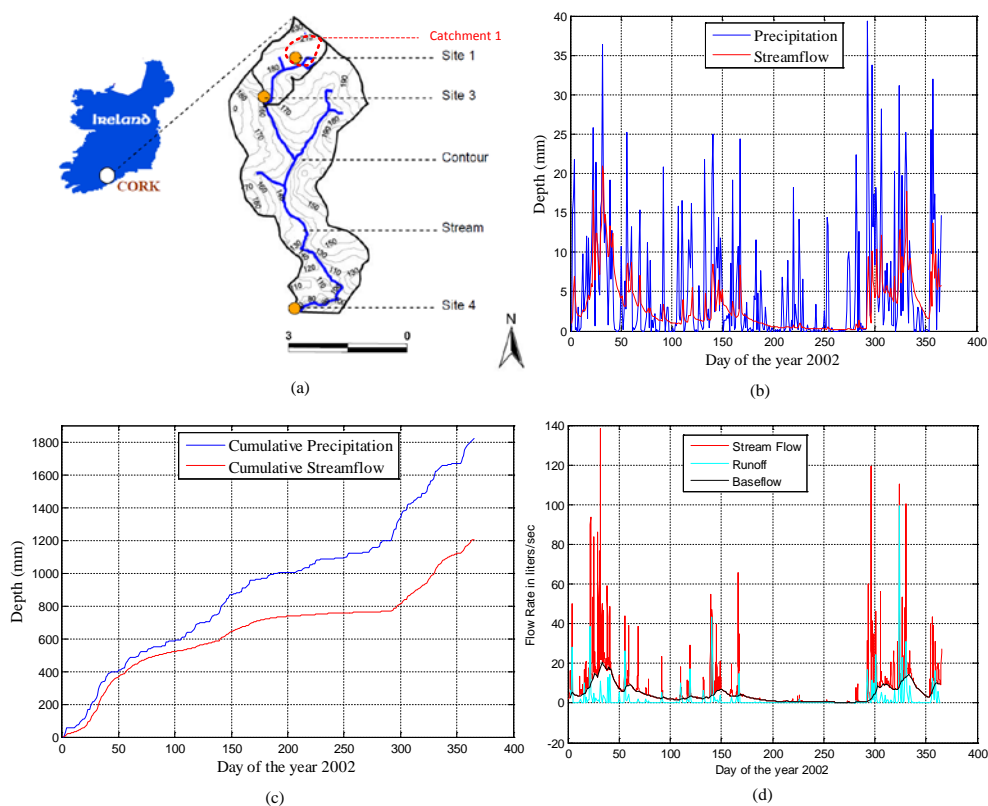
67 This highlights one of the main issues in that the historical data sets needed to develop these predictive
68 models do not exist for every farm, and even for most catchments. In addition, the strengths of a WSN
69 deployment (fine spatial and temporal measurements of dynamic parameters) requires a simplified underlying
70 physical model, and a simple machine learning model based on fewer and, ideally, real-time field parameters
71 acquired autonomously and shareable across neighbouring farms. Thus there is a requirement for a discharge
72 predictive model, which takes into account field conditions (soil moisture, vegetation cover) of the farms and
73 the drainage networks, and which could be generated with adequate performance using fewer training samples.
74 Such a model, once implemented in the network, can adaptively learn and further improve its accuracy over the
75 course of time.

76 In this paper, we recap the model simplification for the predictive model for completeness, as already
77 proposed by Zia *et.al.* (Zia et al., 2013), which is based on (but not restricted to) the popular National Resource
78 Conservation Method (NRCS curve number model). Furthermore, we explore the applicability of M5 decision
79 trees, for discharge modelling based on the proposed parameters. A year-long dataset (200 event samples)
80 consisting of daily values for precipitation, field conditions (soil moisture, vegetation cover) and discharges,
81 obtained from a grassland catchment in Ireland is used for training and testing the model. Specifically, an
82 assessment procedure with the following steps is used (i) evaluation of optimized input parameter combinations
83 with optimal performance; (ii) random sampling of the observational dataset to ensure a robust evaluation of the
84 model performance, and the use of 10-fold cross validation to avoid over fitting of the model; (iii) assessment of
85 the model performance against selected criteria; (iv) uncertainty analysis on the model residuals; and (v)
86 comparative assessment of the prediction accuracy against other similar research developed using M5 decision
87 trees.

88 2. Experimental Method

89 2.1 Specification of Catchment Data

90 The University of Cork carried out a study on the Dripsey catchment in the south of Ireland. The one-year study
91 (2002) was aimed at understanding the underlying processes of nutrient loss from soil to water bodies within the
92 catchment (Lewis, 2003) and thus fits the requirement for validating the Q -predictive model. This catchment
93 consists of smaller nested sub-catchments. Figure 1 (a) shows the location of various data collection points in
94 the stream network such as site1, site3 and site4, which collect water drained from their associated sub-
95 catchments. For the development of the Q -predictive model, data available for site1 of the stream network is
96 used. The sub-catchment which drains into this stream location is identified as ‘catchment 1’ (as shown in
97 Figure 1 (a)) consisting of 17 ha of farmland. Precipitation (mm) and stream flow (mm) data, collected every 30
98 minutes for the year 2002 is used. The data is publically available for research and education purposes via the
99 Environmental Protection Agency (EPA) website (Keily, 2003). The remainder of the data regarding field
100 conditions is extracted from catchment descriptors available in the associated documentation (Lewis, 2003).



101

102 Figure 1: (a) Location and map of the Dripsey catchment (Khandokar, 2003); (b) precipitation and stream flow (mm) at

103 site1; (c) cumulative precipitation and stream flow depth for site 1; (d) rate of stream flow, runoff, & base flow for site 1

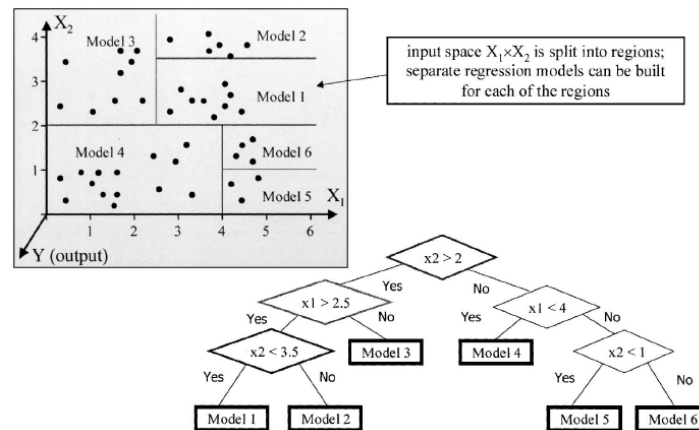
104 For catchment 1, the cumulative rainfall for the year 2002 was 1812mm. The cumulative stream flow depth
105 measured, at site1, was 1206mm of the rainfall (as shown in Figure 1 (c)). Stream flow here consists of water
106 passing this point that originated as any surface runoff, sub-surface drainage or deeper groundwater
107 contributions by catchment 1 (Khandokar, 2003). The monthly rainfall value ranges from less than 50 mm in the
108 summer months to more than 250 mm in the winter months. The mean monthly temperature is 5°C in the winter
109 and 15°C in the summer. The concentrations of total oxidised nitrate losses range from 0.5 to 6.5 mg^l⁻¹. Land
110 cover in the sub-catchment is dominated by agricultural grassland of high quality pasture and meadows. The
111 growing season in Ireland is weather-dependant but generally summer-dominant, starting in early March and
112 finishing in October. Grass is also cut as silage once or twice a year, typically at the end of May and at the end
113 of July.

114 **2.2 Modelling Technique – M5 decision Tree**

115 With WSNs, it is now possible to obtain real time field data, which presents an opportunity for the development
116 of simpler and more accurate data-driven models. These methods are based on the analysis of the data (of some
117 simplified parameters) which characterises the system under study, thereby building models of physical
118 processes. These models can complement or replace the knowledge-driven models describing behaviour of
119 physical systems, and therefore can yield low computational complexity, making them well-suited for
120 implementation on a resource constrained network.

121 As discussed in the introduction, decision tree modelling, specifically, is receiving increasing attention in
122 the hydrological literature, in comparison to other learning models. Decision tree modelling is a method of
123 approximating a target variable (output), with discrete values, from a given data set and represents the learned
124 function in form of a decision tree (Mitchell, 1999), where each leaf contains the target values. Decision trees
125 have been shown to perform well when compared to other model types (Galelli and Castelletti, 2013, Zhao and
126 Zhang, 2008) but they do have one disadvantage. In decision trees, the predicted output is composed of discrete
127 values and is reconstructed as a piecewise constant function. To ensure good prediction accuracy, the number of
128 output classes (tree leaves) should be high; however, this increases the risk of over-fitting the observed data
129 (Breiman, 1996). This can be resolved by replacing averaging in the tree leaves by fitting a linear regression
130 function to the data and obtaining a continuous representation of the output (Galelli and Castelletti, 2013). This
131 is known as M5 tree modelling, and was first introduced by Quinlan (Quinlan, 1992) and applied to hydrological
132 modelling by Solomatine (Solomatine and Dulal, 2003, Solomatine and Xue, 2004b). The M5 tree is a
133 piecewise linear model, so it takes an intermediate position between the linear models and truly nonlinear

134 models such as ANNs. Model trees have higher predictive accuracy and are able to make better predictions for
135 values outside the training data range, when compared with regression trees (Kuzmanovski, 2012).



136

137 Figure 2: A generic M5 model tree, Models1-6 are linear regression models (Solomatine and Xue, 2004b)

138 The M5 model tree is a numerical prediction algorithm, and its splitting criterion is based on the standard
139 deviation of the values in the subset of the training data that reaches a particular node. The construction of a
140 model tree is similar to that of a decision tree. Figure 2 illustrates how the splitting of space is done in a generic
141 M5 tree. Firstly, an initial tree is built and then is pruned (reduced) to overcome the over-fitting problem (that is
142 a problem when a model fits the training data set very accurately, but fails on the test set). Finally, a smoothing
143 process is employed to compensate for sharp discontinuities between adjacent linear models at the leaves of the
144 pruned tree (this operation is not needed in building a decision tree) (Solomatine and Xue, 2004b). To implement
145 the M5 model trees in this paper, MatLab toolbox M5PrimeLab (Jekabsons, 2010) is used.

146 2.3 Model Evaluation Criteria

147 The quantitative assessment of M5 tree modelling for prediction of Q based on the proposed model parameters
148 is performed using a four-step procedure (as suggested in (Galelli and Castelletti, 2013)):

149 2.3.1 Random Sampling & Cross-Validation

150 To ensure a robust evaluation of the model performance, the data set was randomly partitioned into two groups:
151 75% of the observations were used for training the model while the remaining 25% are used for validation.
152 When the available training data set is small, in order to overcome the problem of over-fitting (meaning the
153 model fits the training data but not unseen test data) and reduce the sensitivity of the model to the selected
154 training set, a cross-validation technique allows reliable model validation (Kohavi, 1995). Data is partitioned
155 into subsets; one subset is used for training the model while the other is used for testing. Multiple rounds (folds)

156 of training and testing are performed using different partitions, and the validation results are averaged over the
157 number of rounds. 10-fold is the most commonly used cross validation, which is used in this paper, where data
158 is partitioned into 10 subsets.

159 **2.3.2 Performance Measuring Parameters**

160 To determine the performance (predictive accuracy) of the learned models, a multi-assessment criterion is used
161 (Galelli and Castelletti, 2013). The performance parameters are:

- 162 • *Root mean square error (RMSE)*
- 163 • *Mean absolute error (MAE)*
- 164 • *Coefficient of determination (R^2)*
- 165 • *Root relative mean square error (RRMSE)*

166 RMSE and MAE can vary between 0 and infinity, with lower numbers indicating higher accuracy. R^2 ranges
167 between 1 and 0. It is equal to 1 if the predictions are perfect, i.e. a linear relationship exists between the
168 predicted and measured values represented by a straight line. RRMSE is the ratio of the variance of the residuals
169 to the variance of the target values themselves. Values of RRMSE can range between 0 and 1, where 0 means
170 perfect forecasting. For comparing the performance of models developed using different datasets, R^2 and
171 RRMSE are generally used. The definition of a good value for R^2 and RRMSE depends on the requirements of
172 any specific application. For example, models developed for medical sciences generally need higher accuracy,
173 whereas others might not. In this study, we compare these performance measures with those achieved by
174 previous work in this area.

175 **2.3.3 Uncertainty Analysis**

176 As in any prediction there is a potential error which needs to be accounted for (Computing, 2004). The
177 uncertainty of the predicted variable is investigated by the quantification of the residuals. Residuals represent
178 the deviation of the predicted response from the observed or measured response obtained by subtracting the two.
179 Since residuals are error, therefore, they are expected to be independently distributed. Ideally, the overall pattern
180 of the residuals should be similar to the bell-shaped pattern observed when plotting a histogram of normally
181 distributed values (Natrella, 2010). Through the analysis, we;

- 182 1. Find if residuals show any trend along the days of the year
- 183 2. Determine confidence intervals for residual errors

184 **2.3.4 Comparative Assessment of Q -predictive model with other similar models**

185 In order to evaluate the accuracy of the proposed Q -predictive model, we compare its performance with recent
186 and relevant research efforts in the same area. Hydrologic models developed using M5 trees by *Solomatine et al.*
187 (*Solomatine and Xue, 2004b*), *Kuzmanovski* (*Kuzmanovski, 2012*), *Corzo et al.* (*Corzo et al., 2007*) and *Galelli et al.*
188 (*Galelli and Castelletti, 2013*), have been used for comparison. These works used various input parameters and
189 thousands of training samples for developing the models. They were developed for daily or hourly predictions of
190 discharges measured either as flow rates in a stream (for large catchments) or drainage volumes in field drainage
191 collectors (for crop lands). For the first case, input parameters included precipitation, temperature and flow
192 values for the previous several days. Whereas for the latter, parameters related to field condition were also
193 included in the model. Most of these models did not use cross validation for performance evaluation.
194 Furthermore, not all the performance metrics were used to show the performance accuracy of the models.

195 **3. Development of Q -Predictive Model**

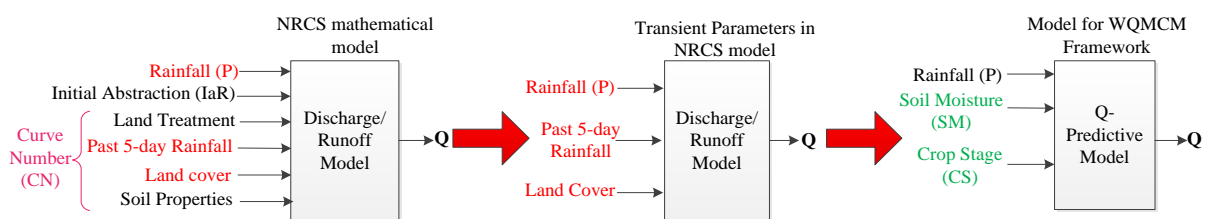
196 **3.1 Model Simplification**

197 A simplified Q -predictive model was proposed by *Zia et al.* (*Zia et al., 2014b*) based on the simplification of the
198 most popular and simplest physically-based hydrological model: the NRCS method. This simplification is
199 briefly explained here. For the prediction of discharges, the NRCS model uses rainfall, initial abstraction ratio
200 (I_a) and curve number (CN) as input parameter where I_a is the volume of rainfall either retained in surface
201 depressions or lost through evaporation or infiltration. CN is a coefficient reducing the total precipitation to
202 runoff potential after surface absorption, and is computed considering the type of land use, land treatment,
203 hydrological condition, hydrological soil group, and last 5-day rainfall (as a proxy for antecedent soil moisture
204 condition). Although, the use of last 5-day rainfall to represent soil moisture conditions has been questioned
205 (*Vellidis et al., 2011*) (*Zia et al., 2014b*), at the time the NRCS model was developed, proxy parameters and
206 manual measurements were used to represent land conditions due to the absence of inexpensive sensing
207 measures.

208 This simplification for the Q -predictive model is based on two steps: In the first stage, the transient parameters
209 from the NRCS model parameters are selected because learning models are trained only on transient values.
210 These parameters include rainfall value expected and land cover on the day of prediction as well as last 5-day
211 rainfall value. In the second stage, the transient parameters are analysed for likely improvements made possible

212 by using WSNs. For example, methods such as field imaging and signal attenuation methods have been used
 213 determine the plant biomass autonomously (Vellidis et al., 2011). This can be used to generate a crop stage (CS)
 214 parameter that can replace the land cover aspects of the NRCS model. Similarly, various applications have used
 215 sensors to monitor soil moisture conditions of the field for precision irrigation (Vellidis et al., 2008, Zia et al.,
 216 2013). Therefore, it has been proposed to use actual soil moisture values instead of the 5-day rainfall index. The
 217 model simplification is shown in Figure 3. The proposed model was developed by training it on simulated data
 218 and using M5 decision trees for learning. 10-fold cross validated results gave 90% accuracy when trained on
 219 100 samples.

220 In order to validate the Q -predictive model further, it needs to be trained using real data acquired from a
 221 catchment. Although the increasing adoption of WSNs in agriculture means that it is now possible to extract real
 222 field conditions for some parameters, free and wide access to such long-term data required for model
 223 development is still not available. Therefore, for model validation, long-term soil moisture data was not
 224 available. Hence, we use the last-5-day rainfall value as a proxy. Nevertheless, this is the best available at
 225 present, and so offers a worst case performance baseline. When real, high frequency soil moisture readings
 226 become available to the model, performance should improve. The final mode includes precipitation (P), last-5-
 227 day rainfall ($L5PPT$) and crop stage (CS) data. This is still simpler than the NRCS model and also a comparable
 228 work (Kuzmanovski, 2012) for drainage discharge prediction which used M5 trees (This used 10 parameters).
 229 Later sections will demonstrate that the proposed model still gives comparable performance with the existing
 230 models developed for discharge prediction with the advantage of using far fewer training samples.

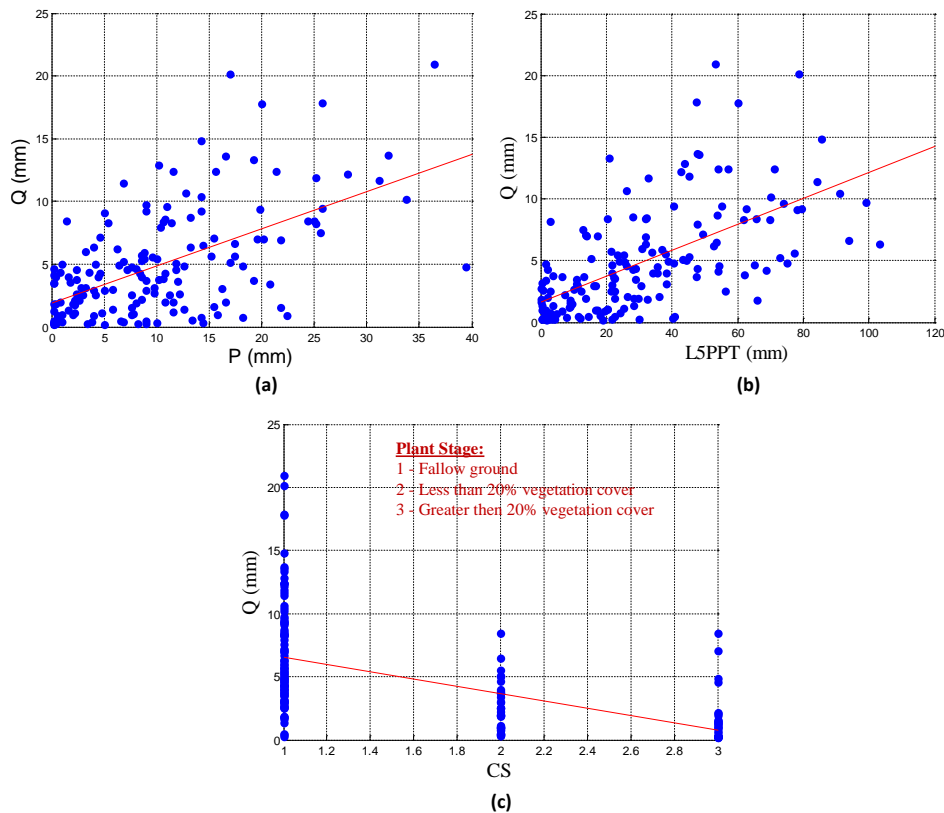


232 Figure 3: Simplification of model parameters for the Q -predictive model

233 3.2 Data Pre-processing & Sensitivity Analysis

234 The set of observations required for training the Q -predictive model was developed after implementing some
 235 pre-processing on the available dataset from the Dripsey catchment in Ireland (Keily). This is the only dataset
 236 that could be found with high temporal resolution data for an entire year. From the available dataset we used
 237 data related to 30 minute precipitation (mm) and stream flow (lsec^{-1}) for the year 2002. The remaining

238 parameters required for the Q -predictive model were either obtained using a proxy value or were extracted from
239 the information available in the documentation for this study (Lewis, 2003).



240

241 Figure 4: Sensitivity analysis of Q output with the input parameters, (a) P , (b) $LSPPT$, and (c) CS

242 Since the model is aimed at facilitating management decisions regarding drainage and nutrients reuse, we
243 convert the 30-minute values of precipitation (P) into daily depths in mm. For each of the daily P values, the
244 $LSPPT$ value is computed by aggregating the depths of rainfall received in the last 5 days. For obtaining CS
245 data, information regarding growing stages of grass in catchment 1 was assessed to obtain estimates for crop
246 coverage throughout the year. According to crop coverage values, crop levels are assigned such that fallow land
247 is referred to as stage 1, coverage less than 20% is termed as stage 2, and coverage greater than 20% is assigned
248 stage 3.

249 For measures of output of the model, discharges from the catchment, runoff and stream flow rates were all
250 available. We select stream flow as an output for our model for various reasons;

251 (i) runoff value is available as a single measurement per day, which does not provide complete
252 information for the daily runoff depths,

253 (ii) due to the presence of high base flows, as shown in Figure 1 (c), stream flow better represents the
254 discharges from the grass land

255 (iii) Pearson correlation coefficient (r) (which explains the strength and direction of the linear
256 relationship between parameters) gives better values for stream flow as compared to runoff with
257 the input parameters.

258 For stream flow, the values of ' r ' for P , $L5PPT$ and CS are 0.53, 0.32, and -0.38 respectively, while for
259 runoff, the values of ' r ' with P , $L5PPT$ and CS are 0.52, 0.11, and -0.16 respectively. Positive ' r ' values for P
260 and $L5PPT$ indicate that higher rainfall produces higher discharges. On the other hand, CS has a negative
261 correlation, which indicates that higher vegetation cover absorbs more available water and also inhibits
262 discharge flows.

263 Daily stream flow depth (mm), (Q), is computed from 30 minute stream flow rates, as the model need
264 daily depths for the decision making process. The values for ' r ' show even stronger correlation of Q with the
265 input parameters (0.57 for P , 0.61 for $L5PPT$, and -0.54 for CS). Figure 4 illustrate the sensitivity analysis of the
266 input model parameters with Q . All the three plots indicate visible correlation with all input parameters.

267 4. Results and Discussion

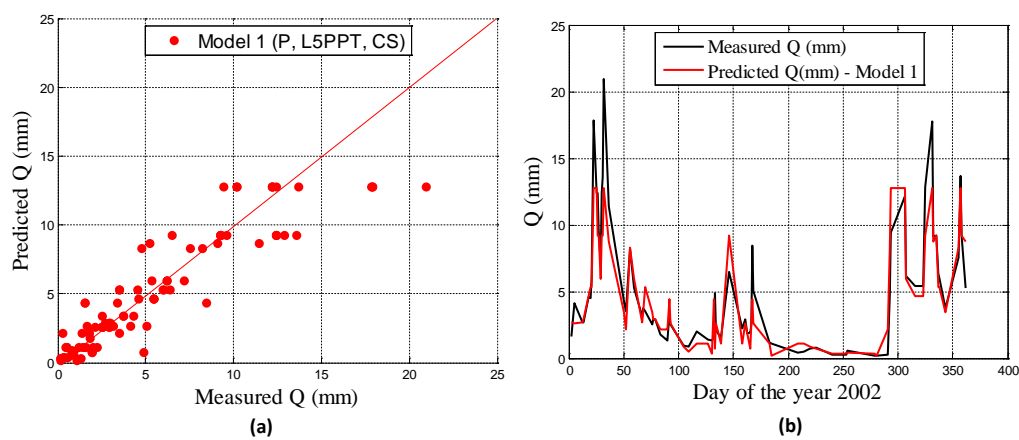
268 In order to generate the Q -predictive model, training data of 200 event instances (75% of the total event dataset)
269 are sampled randomly from the dataset which was generated as a result of the pre-processing discussed in the
270 previous section. M5 decision trees were then used to generate the tree. The average value for Q in the training
271 data is 3.80 mm, 25th percentile is 1.01, 75th percentile is 5.29, and 90th percentile is 9.22. The standard
272 deviation of Q is 3.8.

273 Table 1: Performance measure for the Q -predictive model

	<i>Performance Measure for Q-predictive Model</i>			
	R^2	<i>MAE</i>	<i>RMSE</i>	<i>RRMSE (%)</i>
Without cross validation	0.87	0.81	1.34	35.9
10-fold cross validated	0.64	1.30	1.95	55.9

274
275 Performance parameters for the generated M5 decision tree indicate a high value for R^2 of 0.823 as listed in
276 Table 1. In addition to this, the RMSE value for the model is 1.335. An accepted adequate value for RMSE in
277 hydrological modelling is normally half of the standard deviation of the training data (Singh et al., 2005), which
278 for this data is 1.9. The obtained RMSE of the trained model is well within this limit. Furthermore, the RRMSE

279 is 35.9%. Moreover, after 10-fold cross validation of the Q -predictive model, the results give R^2 as 0.63 and
280 RRMSE as 55%. These results are comparable and in some cases even better than the performance of the
281 existing models, which will be discussed in detail in the later section. However, when compared with our
282 previous work on this model in Zia *et.al.*(Zia et al., 2014b), which was developed using simulated data (from
283 NRCs simulator), R^2 was 0.99 and RRMSE was 7.5%. This is possibly because simulated data, though
284 randomly sampled, is generated on the basis of an underlying mathematical model with obvious relationships
285 between input and output, which is picked up by a learning model. Hence, models trained on simulated data,
286 even if only few parameters of the actual mathematical model are selected, have higher predictive performance.



287
288 Figure 5: (a) Scatter plot of the predicted Q using the proposed model against measured Q , (b) Curve plot
289 for predicted and measured Q plotted against days of the year

290 Furthermore, test data was used to predict Q values using the generated model. The predicted values are
291 plotted against the known measured values and against days of the year in Figure 5. The model illustrates a good
292 fit with R^2 of 0.868 as shown in Figure 5 (a). The predicted values for test data and measured values are also
293 plotted against days of the year in Figure 5 (b) to represent the difference between the two curves. As is evident
294 both curves almost overlap 50-60% of the time (covered more in section 4.2); however in the first 50 days of the
295 year, the model seems to under-predict the stream values by about 25%. Also, it is evident in Figure 5 (a) that
296 the prediction results flatten out for higher Q values. The reason is possibly because of the sparse training
297 samples available above that threshold. Further work on verifying this with either more training samples or with
298 classifying events with respect to their intensity (lower, medium and higher flows) would be able to explain
299 these variations.

300 In the later section, we evaluate if subsets of the proposed model, i.e. any further simplification to the
301 model parameters, impacts the model performance in any way. This would validate that using field condition

302 parameters in the proposed model is a better approach for prediction of outflows from smaller field plots,
 303 compared to relying only on the climatic conditions (Solomatine and Xue, 2004b, Corzo et al., 2007, Galelli and
 304 Castelletti, 2013).

305 **4.1 Further Model Simplification – Viable or not?**

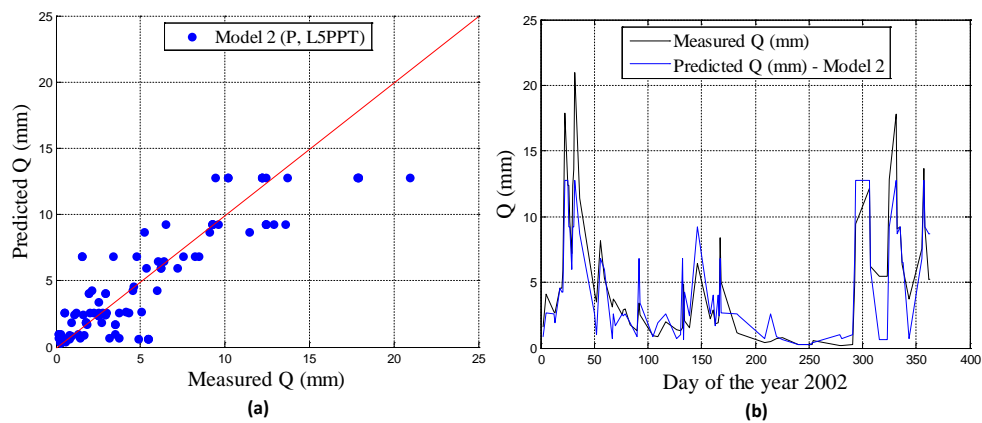
306 In this section two models are developed using ‘*P* plus *LSPPT*’ and only ‘*P*’ as input parameters. These models
 307 are called ‘model 2’ and ‘model 3’ respectively. The full model previously discussed is referred to as ‘model 1’.
 308 The performance parameters of the generated models are compared with the performance of the proposed model
 309 (as shown in Table 2). It is clear from the table that, model 1 show the best performance in comparison to the
 310 other models, although it is quite close to the performance of model 2. For instance, model 1 has 35.9% RRMSE
 311 value while model 2 has a 39.2% RRMSE value, although the R^2 values are almost similar. This shows that
 312 *LSPPT* (not included in model 2) is very significant for predicting outflows. Also, it is consistent with the
 313 correlation values of *LSPPT* with *Q* (0.609) as discussed in section 3.2. However, by using only *P* as an input in
 314 model 3, a very weak model is generated with 70% RRMSE value and 0.5 for R^2 . This means that the predicted
 315 results using model 3 will have only 30% accuracy.

316 Table 2: Comparison of *Q*-predictive models developed using different combination of input parameters

<i>Model No.</i>	<i>Features for Q-</i> <i>predictive model</i>	<i>Performance Metrics</i>				<i>10-fold cross validated</i> <i>Performance Metrics</i>			
		R^2	<i>MAE</i>	<i>RMSE</i>	<i>RRMSE</i> (%)	R^2	<i>MAE</i>	<i>RMSE</i>	<i>RRMSE</i> (%)
1 (Proposed model)	<i>P</i> , <i>LSPPT</i> , <i>CS</i>	0.87	0.81	1.34	35.9	0.64	1.30	1.95	55.9
2	<i>P</i> , <i>LSPPT</i>	0.84	0.97	1.51	39.2	0.58	1.52	2.17	61.4
3	<i>P</i>	0.51	1.86	2.66	70.1	0.09	2.30	3.05	91.4

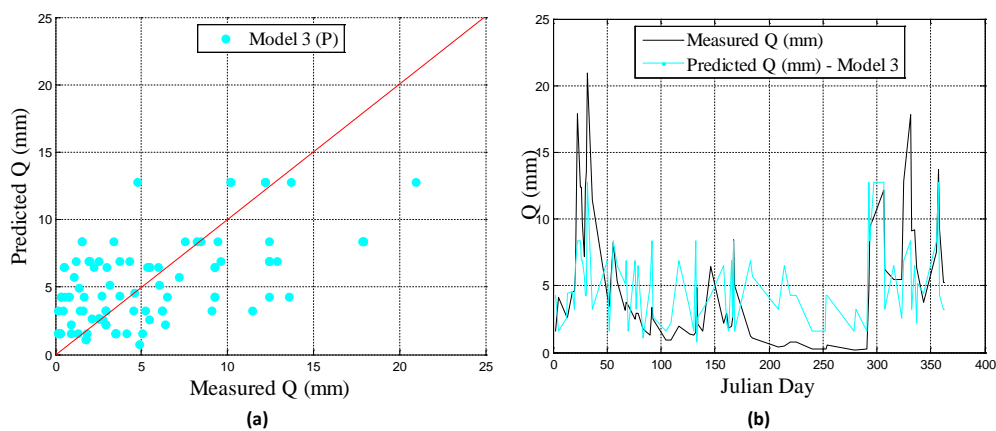
317
 318 To illustrate the comparison among models further, test data is used to predict *Q* values for model 2 and
 319 model 3. The predicted values are plotted against the known measured values and against days of the year in
 320 Figure 6 and Figure 7. Model 2 show a good fit (as shown in Figure 6 (a) with R^2 of 0.844. The predicted values
 321 for test data and measured values are also plotted against days of the year in Figure 6 (b) to represent the
 322 difference between the two curves. Figure 6 (b) indicates that model 2 under-predicts during the initial 50 days

323 as well as last 50 days of the year during which highest discharges were observed (as shown in Figure 1 (b)).
324 Furthermore, the predicted Q for test data using Model 3 is plotted in Figure 7 (a), which shows a poor fit, with
325 an R^2 value of 0.505. This validates that reliance on only precipitation values, as an input for developing the Q -
326 predictive model, will result in weak learning of the model hence poor prediction. Figure 7 (b) further illustrates
327 a plot of the deviation of predicted Q compared to measured Q . The model seems to over-predict during summer
328 and under-predict during winters leading to an unreliable system. To compare the prediction accuracy for these
329 three models, we plot the results on a single graph as shown in Figure 8. This clearly illustrates the deviation of
330 predicted results of these models from the measured results; model 1 has the closest similarity with the
331 measured values as compared to other models. Hence, it is concluded that the proposed model parameters
332 cannot be further simplified without degrading the performance.

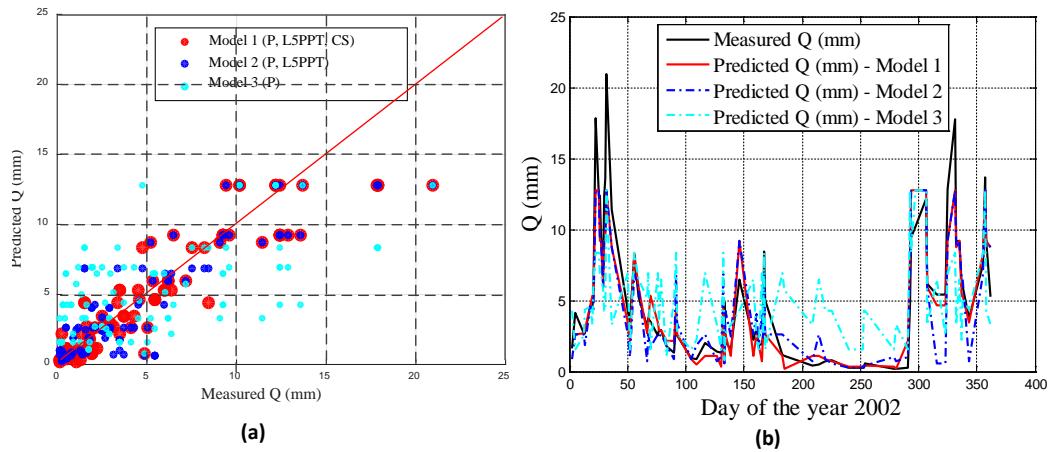


333
334 Figure 6: (a) Scatter plot of predicted Q using Model-2 against measured Q , (b) Curve plot for predicted and measured
335 Q plotted against days of the year

336



337
338 Figure 7: (a) Scatter plot of predicted Q using Model-3 against measured Q , (b) Curve plot for predicted and measured
339 Q plotted against day of the year

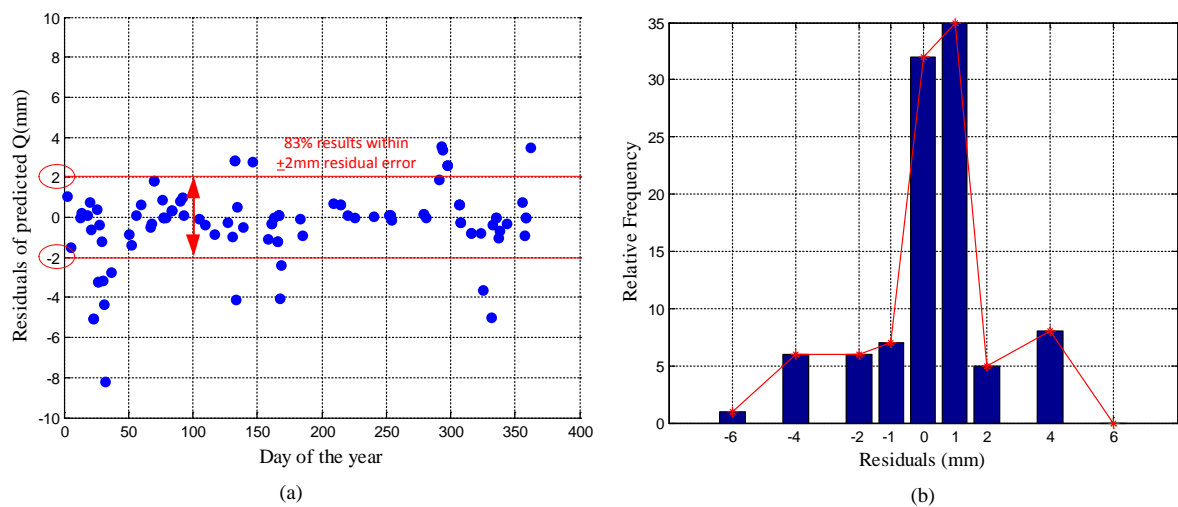


340

341 Figure 8: Comparison of predicted and measured Q for model 1, 2 & 3 in a (a) scatter plot and (b) curve plot

342 4.2 Uncertainty Analysis

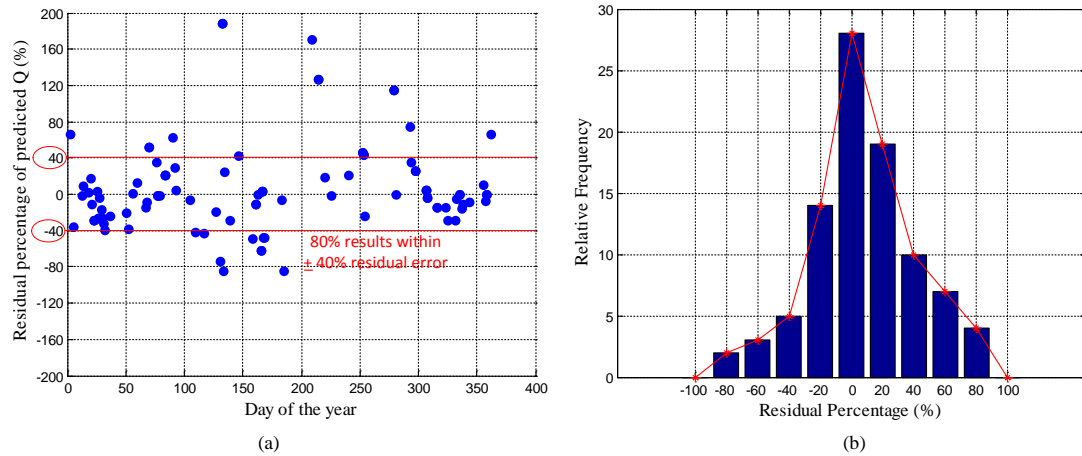
343 To determine if any time dependency exists for the prediction error over summer and winter, residuals are
 344 plotted against day of the year. This provides us with a time-dependent confidence interval for the predicted
 345 values. Figure 9 (a) shows that 83% of the residuals for the predicted test values fall within a range of ± 2 mm. A
 346 prediction error of this scale is not significant because this does not yield substantial outflows for this small
 347 catchment (17 ha). Therefore, incorrect estimates at this scale will not adversely impact decision making based
 348 on modelled results. Furthermore, as already pointed out in the earlier discussion of the results, there is a
 349 seasonal variation in the predicted data. This is linked to the performance of the model for predicting high Q
 350 events, and these high Q events tend to occur in winter, thus concentrating the uncertainty in this period. This is
 351 a feature of the current model structure that will be investigated in more detail in future work, but is currently
 352 limited by the availability of real data.



353

354 Figure 9: (a) Trend of residual error of predicted Q , (b) relative frequency of the residual error

355 A frequency plot for the residual error illustrates an approximately normal distribution of residuals
356 produced by the model with highest frequency corresponding to 0 and 1 mm errors (Figure 9 (b)). To explain
357 this further, Figure 10 (a) shows the residuals percentile of predicted Q against days of the year. This illustrates
358 that 80% of the predicted values are within $\pm 40\%$ residual errors. A histogram for this shows a normal
359 distribution having maximum number of values with 0% residual error (Figure 10 (b)).



360 (a) (b)
361 Figure 10: (a) Trend of percentile value of residual error versus, (b) relative frequency of per cent residual error

362 5. Comparative Assessment – Q -Predictive Model versus similar models

363 In order to evaluate if the proposed model has acceptable (or comparable) performance, we compare its
364 results with the results of the most relevant existing models. For this, work by *Kuzmanovski* (Kuzmanovski, 2012),
365 *Solomatine et al.* (Solomatine and Xue, 2004b), *Corzo et al.* (Corzo et al., 2007) and *Galelli et al.* (Galelli and
366 Castelletti, 2013) has been selected. All of these models use M5 decision trees. Although, all of these works were
367 aimed at predicting discharges either in the form of drainage from a small field plot, or flow volumes and rates
368 in a river or stream drained by large catchments, it is important to note that none of these works are entirely
369 similar, in objective and methodology with the proposed Q -predictive model. Table 3 lists the experimental
370 details and performance metrics for all the models including the Q -predictive model. The model listed first in
371 the table is the Q -predictive model proposed in this paper.

372 Overall, it is clear that the Q -predictive model has a comparable performance compared to other models.
373 More specifically, work done by *Kuzmanovski* (Kuzmanovski, 2012) is closer in objective to the Q -predictive
374 model, hence provides better comparison. In this work, the aim was to measure drainage discharges from fields,
375 ranging from 0.83 ha to 1ha, in order to control pollutant outflows. The model uses 10 parameters related to
376 crop stage, day of the season, slope of the field, rainfall, temperature, runoff, drainage etc., and used 22 years of

377 daily data to train the model. Without cross validation, the performance measures for the models developed for
 378 various field plots was calculated; R^2 ranges between 0.75 and 0.89, and RMSE is between 45% and 65.9%. The
 379 performance (without cross validation) of the Q -predictive model is acceptable when compared with the best
 380 performing model by *Kuzmanovski*, although the former model is developed with only 3 input parameters. The
 381 reason for good performance of the Q -predictive model can be possibly attributed to the use of the most relevant
 382 parameters, especially the 5-day rainfall value. It is believed that results would further improve if actual soil
 383 moisture measurements are used. However, the models developed by *Kuzmanovski* use only 2 – day rainfall
 384 value, which may not accurately represent soil moisture conditions. The performance of the Q -predictive model
 385 would need to be validated with more data samples.

386 The models developed by *Solomatine et al.*, *Corzo et al.* and *Galelli et al.*, were all predicting discharges in
 387 a river and stream draining very large catchments, and hence used only 2 to 3 parameters related to climate and
 388 flow. This is because field conditions (soil moisture, vegetation cover) for such large heterogeneous catchments
 389 can vary tremendously; hence a single average value may not represent the field conditions for the whole
 390 catchment. Furthermore, these models are developed for hourly predictions; field conditions do not change in an
 391 hourly manner, however previous flow and rainfall intensity do. Hence, these models rely only on climatic and
 392 flow parameters for predictions and use thousands of samples so that the model could learn all possible samples
 393 over many years (decades in some cases). These models have slightly better performance when compared with
 394 the Q -predictive model for obvious difference in the training set size. For instance, the R^2 values of the model
 395 proposed by *Solomatine et al.* and *Corzo et al.* are 0.97 and 0.89 respectively as compared to the 0.86 for the Q -
 396 predictive model. Similarly, RRMSE of the model proposed by *Galelli et al.* is 48% as compared to 55%
 397 RRMSE of the Q -predictive model.

398 The proposed Q -predictive model is between the models developed by *Kuzmanovski* (*Kuzmanovski*, 2012)
 399 and the others (*Solomatine and Xue*, 2004b, *Corzo et al.*, 2007, *Galelli and Castelletti*, 2013) in terms of the
 400 number of parameters used. The Q -predictive model has climatic conditions as used in the latter models.
 401 However it also uses field conditions but keeps the number of parameters simpler and fewer compared to
 402 *Kuzmanovski*.

403

404 Table 3: Performance comparison of Q -predictive models with the existing models developed using M5 trees

S. No.	Predictive	Output	No of Input	No of	Drainage	Cross	RRMSE	R2
--------	------------	--------	-------------	-------	----------	-------	-------	----

	Models	Value	variables	Training samples	area	Validation for training	(%)	
1	Proposed Q -predictive Model	Daily Discharge	3 (rainfall, last 5-day rainfall, crop stage)	200 samples (2002)	17 ha	10- fold	55	0.63
		s from a farm				Not done	35	0.86
2	<i>Kuzmanovski</i> (Kuzmanovski, 2012)	Daily Drainage volume from farms	10 parameters	7000 samples (22 years data, 1987-2011)	0.83-1 ha field plots	Not done	45 - 65.9	0.75-0.89
3	<i>Solomatine et al.</i> (Solomatine and Xue, 2004b)	Flood discharge	3 parameters (11 sub values)	5000 samples (> 13 years data, 1976-1989)	106 ha	Not done	-	0.97
3	<i>Corzo et al.</i> (Corzo et al., 2007)	Hourly discharge in a stream	Previous runoff, effective rainfall	2000 (8 years data, 1988-1996)	37,000 ha	Not done	-	0.89
4	<i>Galelli et al.</i> (Galelli and Castelletti, 2013)	Hourly Stream flow	2 (rainfall & inflow value), 3 time lag sets	24120 (2.5 years data, 2009-2011)	10,000 ha	10-fold	48	-

405

406 **Conclusions**

407 In this paper, we have successfully validated a discharge (Q) predictive model for the proposed simplified
408 parameters by employing M5 decision tree learning algorithms. The input parameters included daily
409 precipitation, vegetation cover, and last-5-day rainfall value. The model is assessed for its prediction accuracy in
410 comparison to major data-driven hydrological models. The Q -predictive model was evaluated by training it on a

411 year-long discharge dataset measured for a sub-catchment in Ireland. The significance of the proposed Q -
412 predictive model is in the fact that it uses just a year-long training sample set which includes climatic as well as
413 field parameters to develop accurate predictive model, and that it can adequately predict flows for smaller sub-
414 catchments with simpler parameters and acceptable (comparable) prediction accuracy.

415 Results for the Q -predictive model show that:

- 416 i) Performance measures for the proposed Q -predictive model provide good performance; R^2 is 0.823,
417 RMSE is 1.33, and RRMSE value is 35.9%. 10-fold cross-validated results yield R^2 of 0.63, RMSE of
418 1.95, and RRMSE of 55.7%.
- 419 ii) 83% of the residuals for the predicted test values fall within +2 mm range. It has been argued here that
420 prediction errors of this scale are not significant, and a plot of the residuals percentile values illustrates
421 that 80% of the predicted values are within $\pm 40\%$ residual errors.
- 422 iii) Investigation for further simplification of model parameters results in poor performance. 10-fold cross
423 validated results for model developed using rainfall and last-5-day rainfall value as input parameters
424 (and omitting the crop cover) result in RRMSE of 61.4%. For the model developed using only rainfall
425 as the input parameter, RRMSE is 91.4%.
- 426 iv) In comparison to the existing models developed for discharge predictions using M5 decision trees, the
427 proposed model presents great promise and comparable performance by only using a year-long dataset
428 and simplifying the model parameters tremendously for small-scale land areas.

429 **References**

- 430 ADELMAN, D. D. 2000. Simulation of irrigation reuse system nitrate losses and potential corn yield
431 reductions. *Environmental Science & Policy*, 3, 213-217.
- 432 BASHA, E. A., RAVELA, S. & RUS, D. Model-based monitoring for early warning flood detection.
433 6th ACM Conference on Embedded network sensor systems 2008. ACM, 295-308.
- 434 BHATTACHARYA, B., PRICE, R. & SOLOMATINE, D. 2005. Data-driven modelling in the
435 context of sediment transport. *Physics and Chemistry of the Earth, Parts A/B/C*, 30, 297-302.
- 436 BREIMAN, L. 1996. Bagging predictors. *Machine learning*, 24, 123-140.
- 437 CARR, G., POTTER, R. B. & NORTCLIFF, S. 2011. Water reuse for irrigation in Jordan:
438 Perceptions of water quality among farmers. *Agricultural Water Management*, 98, 847-854.
- 439 CASTELLETTI, A., GALELLI, S., RESTELLI, M. & SONCINI-SESSA, R. 2010. Tree-based
440 reinforcement learning for optimal water reservoir operation. *Water Resources Research*, 46.
- 441 COMPUTING, W. N. 2004. Model-Independent Parameter Estimation. *User Manual*.
- 442 CORZO, G., SIEK, M., SOLOMATINE, D. & PRICE, R. Modular data-driven hydrologic models
443 with incorporated knowledge: neural networks and model trees. International congress
444 IAHR, ASCE, Italy, 2007.
- 445 DAWSON, C. W. & WILBY, R. 1998. An artificial neural network approach to rainfall-runoff
446 modelling. *Hydrological Sciences Journal*, 43, 47-66.
- 447 EISENHAEUER, D. E. 2011. Irrigation Efficiency and Uniformity, and Crop Water Use Efficiency.
448 The Board of Regents of the University of Nebraska - Lincoln Extension.

- 449 FORTIN, J., MORAIS, A., ANCTIL, F. & PARENT, L. 2014. Comparison of Machine Learning
450 Regression Methods to Simulate NO₃ Flux in Soil Solution under Potato Crops. *Journal of*
451 *Applied Mathematics*, 5, 832-841.
- 452 GALELLI, S. & CASTELLETTI, A. 2013. Assessing the predictive capability of randomized tree-
453 based ensembles in streamflow modelling. *Hydrology & Earth System Sciences Discussions*,
454 10.
- 455 HARPER, H. H. 2012. Impacts of Reuse Irrigation on Nutrient Loadings and Transport in Urbanized
456 Drainage Basins. *Florida Stormwater Association Annual Meeting*. Environmental Research
457 & Design, Inc.
- 458 JEKABSONS, G. 2010. *M5PrimeLab: M5' regression tree and model tree toolbox for Matlab*
459 [Online]. Institute of Applied Computer Systems Riga Technical University, Latvia.
460 Available: <http://www.cs.rtu.lv/jekabsons/Files/M5PrimeLab.pdf>.
- 461 KEILY, G. 2003. *Phosphorus, Nitrogen and Suspended Sediment loss from Soil to Water from*
462 *Agricultural Grassland* [Online]. Environmental Protection Agency (EPA). Available:
463 <http://erc.epa.ie/safer/resource?id=ad1f3acf-5035-102a-90c6-0593d266866d> [2014].
- 464 KHANDOKAR, F. 2003. *Phosphorus in soils of pasture farms subject to chemical and slurry*
465 *fertilization*. NUI, 2003 at Department of Civil and Environmental Engineering, UCC.
- 466 KOHAVI, R. A study of cross-validation and bootstrap for accuracy estimation and model selection.
467 IJCAI, 1995. 1137-1145.
- 468 KUZMANOVSKI, V. 2012. Integration of expert knowledge and predictive learning: Modelling
469 water flows in agriculture - MS Thesis. Jožef Stefan International Postgraduate School,
470 Ljubljana, Slovenia
- 471 LEWIS, C. 2003. *Phosphorus, nitrogen and suspended sediment loss from soil to water from*
472 *agricultural grassland*. Department of Civil and Environmental Engineering, University
473 College Cork.
- 474 LIU, X., JU, X., ZHANG, F., PAN, J. & CHRISTIE, P. 2003. Nitrogen dynamics and budgets in a
475 winter wheat-maize cropping system in the North China Plain. *Field Crops Research*, 83,
476 111-124.
- 477 MINNS, A. & HALL, M. 1996. Artificial neural networks as rainfall-runoff models. *Hydrological*
478 *Sciences Journal*, 41, 399-417.
- 479 MITCHELL, T. M. 1999. Machine learning and data mining. *Communications of the ACM*, 42, 30-36.
- 480 NATRELLA, M. 2010. *NIST/SEMATECH: e-Handbook of Statistical Methods*, Available online:
481 <http://www.itl.nist.gov/div898/handbook>.
- 482 OSTER, J. & GRATAN, S. 2002. Drainage water reuse. *Irrigation and Drainage Systems*, 16, 297-
483 310.
- 484 PIÑEROS GARCET, J. D., ORDOÑEZ, A., ROOSEN, J. & VANCLOOSTER, M. 2006.
485 Metamodelling: Theory, concepts and application to nitrate leaching modelling. *Ecological*
486 *modelling*, 193, 629-644.
- 487 QUINLAN, J. R. Learning with continuous classes. 5th Australian joint Conference on Artificial
488 Intelligence, 1992. Singapore, 343-348.
- 489 RASOULI, K., HSIEH, W. W. & CANNON, A. J. 2012. Daily streamflow forecasting by machine
490 learning methods with weather and climate inputs. *Journal of Hydrology*, 414-415, 284-293.
- 491 SINGH, J., KNAPP, H. V., ARNOLD, J. & DEMISSIE, M. 2005. Hydrological modeling of the
492 iroquois river watershed using HSPF and SWAT1. *JAWRA Journal of the American Water*
493 *Resources Association*, 41, 343-360.
- 494 SOLOMATINE, D. & OSTFELD, A. 2008. Data-driven modelling: some past experiences and new
495 approaches. *Journal of hydroinformatics*, 10, 3-22.
- 496 SOLOMATINE, D. & XUE, Y. 2004a. M5 Model Trees and Neural Networks: Application to Flood
497 Forecasting in the Upper Reach of the Huai River in China. *Journal of Hydrologic*
498 *Engineering*, 9, 491-501.
- 499 SOLOMATINE, D. P. & DULAL, K. N. 2003. Model trees as an alternative to neural networks in
500 rainfall—runoff modelling. *Hydrological Sciences Journal*, 48, 399-411.
- 501 SOLOMATINE, D. P. & XUE, Y. 2004b. M5 model trees and neural networks: application to flood
502 forecasting in the upper reach of the Huai River in China. *Journal of Hydrologic Engineering*,
503 9, 491-501.

- 504 TINDULA, G., ORANG, M. & SNYDER, R. 2013. Survey of Irrigation Methods in California in
505 2010. *Journal of Irrigation and Drainage Engineering*, 139, 233-238.
- 506 VELLIDIS, G., SAVELLE, H., RITCHIE, R., HARRIS, G., HILL, R. & HENRY, H. 2011. NDVI
507 response of cotton to nitrogen application rates in Georgia. *Precision Agriculture*, 359.
- 508 VELLIDIS, G., TUCKER, M., PERRY, C., KVIEN, C. & BEDNARZ, C. 2008. A real-time wireless
509 smart sensor array for scheduling irrigation. *Computers and Electronics in Agriculture*, 61,
510 44-50.
- 511 VILLA-VIALANEIX, N., FOLLADOR, M., RATTO, M. & LEIP, A. 2012. A comparison of eight
512 metamodeling techniques for the simulation of N₂O fluxes and N leaching from corn crops.
513 *Environmental Modelling & Software*, 34, 51-66.
- 514 WILBY, R., ABRAHART, R. & DAWSON, C. 2003. Detection of conceptual model rainfall—runoff
515 processes inside an artificial neural network. *Hydrological Sciences Journal*, 48, 163-181.
- 516 WILLARDSON, L., BOELS, D. & SMEDEMA, L. 1997. Reuse of drainage water from irrigated
517 areas. *Irrigation and Drainage Systems*, 11, 215-239.
- 518 ZHAO, Y. & ZHANG, Y. 2008. Comparison of decision tree methods for finding active objects.
519 *Advances in Space Research*, 41, 1955-1959.
- 520 ZIA, H., HARRIS, N. R. & MERRETT, G. V. Empirical Modelling and Simulation for Discharge
521 Dynamics Enabling Catchment-Scale Water Quality Management. The 26th European
522 Modeling & Simulation Symposium, 2014a Bordeaux, France.
- 523 ZIA, H., HARRIS, N. R. & MERRETT, G. V. Water Quality Monitoring, Control and Management
524 (WQMCM) Framework using Collaborative Wireless Sensor Networks. 11th International
525 Conference on Hydroinformatics 2014b New York City, USA.
- 526 ZIA, H., HARRIS, N. R., MERRETT, G. V., RIVERS, M. & COLES, N. 2013. The impact of
527 agricultural activities on water quality: A case for collaborative catchment-scale management
528 using integrated wireless sensor networks. *Computers and Electronics in Agriculture*, 96,
529 126-138.
- 530