1 # Predicting tile drainage discharge using machine learning

2 # algorithms

3 Saghar Khodadad Motarjemi[1], Anders Bjørn Møller[1], Finn Plauborg[1] and Bo Vangsø Iversen[1]

4 [1]Department of Agroecology, Aarhus University, Blichers Alle 20, 8830 Tjele, Denmark

5 *Correspondence to*: Saghar K. Motarjemi (sa.m@agro.au.dk)

6 **Abstract**

7 Drainage systems can significantly improve the water management in agricultural fields. However, they may transport

8 contaminants originating from fertilizers and pesticides and threaten ecosystems. Determining the quantity of drainage

9 water is an important factor for constructed wetlands and other drainage mitigation techniques. This study was carried out

10 in Denmark where tile drainage systems are implemented in more than half of the agricultural fields. The first aim of the

11 study was to predict the annual discharge of tile drainage systems using machine-learning methods, which have been highly

12 popular in recent years. The second objective was to assess the importance of the parameters and their impact on the

13 predictions. Data from 53 drainage stations distributed in different regions of Denmark were collected and used for the

14 analysis. The covariates contained 35 parameters including the calculated percolation and geographic variables such as

15 drainage probability, clay content in different depth intervals, and elevation, all extracted from existing national maps.

16 Random Forest and Cubist were selected as predictive models. Both models were trained on the dataset and used to predict

17 yearly drainage discharge. Results highlighted the importance of the cross-validation methods and indicated that both

18 Random Forest and Cubist can perform as predictive models with a low complexity and good correlation between predicted

19 and observed discharge. Covariate importance analysis showed that among all of the used predictors, the percolation and

20 elevation have the largest effect on the prediction of tile drainage discharge. This work opens up for a better understanding

21 of the dynamics of tile drainage discharge and proves that machine-learning techniques can perform as predictive models

22 in this specific concept. The developed models can be used in regard to a national mapping of expected tile drain discharge.

23 **Keywords:** Tile drainage discharge, Random Forest, Cubist, Cross-validation

24 **1. Introduction**

25 Artificial subsurface drainage has a huge impact on the hydrology, nutrient cycling, and sediment dynamics in

26 agricultural systems (Blann et al. 2009). In temperate climates with fine-textured soils as well as semi-arid regions with

27     irrigated fields (Ayars et al. 2006), tile drainage is a crucial water management system to control runoff, prevent

28     waterlogging, and to increase water use efficiency. On the other hand, tile drainage affects both the quantity and the

29     quality of water resources (Schilling et al. 2012). Nutrient losses and chemical remnants can either be transported

30     through drains to surface water bodies such as lakes and rivers (Stenberg et al. 2012) or be leached to the groundwater,

31     and this fresh-water contamination can threaten both human and ecosystem health (Kuzmanovski et al. 2015).

32     Constructed wetlands are a means to eliminate excessive amounts of nitrogen from drainage water benefiting from

33     natural nitrate reducing processes in a controlled environment (Messer et al. 2017). These systems are mainly installed to

34     reduce the pollution from drainage water from agricultural fields and run-off from industrial areas (Magmedov et al.

35     1996). In order to design constructed wetlands with appropriate sizes, it is necessary to quantify artificial drainage

36     discharge. Physically-based hydrological models have been developed either to estimate the drainage discharge or to

37     include it as a component (De Schepper et al. 2017). These models have a common use in academic research and may as

38     well be used to evaluate various scenarios (Zia et al. 2015). However, they depend on numerous parameters and require

39     calibration to individual areas (Basha et al. 2008), which makes them complicated and time consuming to apply on a

40     national scale. Another disadvantage of these models is the conceptualization as the fundament, which leads to invalid

41     predictions when new empirical data are introduced (Bredehoeft 2005). Beside physically based models, many statistical

42     approaches have been used to model and to predict state variable such as discharge, but there are limited number of

43     literature predicting tile drainage discharge with the means of machine learning approaches. This type of data-driven

44     modelling requires fewer parameters and can perform as an accurate estimation technique and these models have proved

45     to be flexible and robust enough for many regression applications (Park et al., 2016).

46     Machine learning is related to computational statistics and is commonly used for predictions based on learning from

47     historical relationships and trends in the data. Classification and Regression Trees (CART) are a frequently used form of

48     machine learning models. They work by searching through the covariates of a dataset to find the best splitting single

49     value. This creates two different groups of data. The process is repeated for the both created groups until a decision tree

50     forms. Zia et al. (2015) predicted drainage discharge utilizing an M5 decision tree modelling technique on a 17 ha

51     drained farmland in southern Ireland. Predictions were carried out on a daily basis for a 12-month period. They validated

52     the suitability of a simplified discharge prediction model for implementation on a system with limited resources.

53     Kuzmanovski et al. (2015) evaluated machine-learning models in predicting sub-surface tile drainage discharge and

54     surface runoff on an experimental site in La Jaillière, France using daily data from eleven fields including a reference

55     field. The dataset was based on meteorological measurements, agricultural practices, and crop management. By

56   comparing the results from these models with the performance of two physically based models, they found an

57   improvement in the sub-surface discharge predictions.

58   In the present study, two different machine-learning models were used to predict yearly tile drainage discharge, Random

59   Forest (RF) (Breiman, 2001) and Cubist (CB) (Quinlan, 1993). RF is an ensemble approach based on CART (Breiman,

60   2001). It trains a number of regression trees from bootstrap samples drawn from the original dataset and averages the

61   results from each tree for the final prediction. The algorithm furthermore introduces randomness into the splitting process

62   by selecting the optimal split from a random subset of the covariates in each split. CB is a rule-based regression

63   technique, which does not retrieve one final model like RF but a set of rules related to multivariate models (Walton,

64   2008). A specific set of covariates will choose an actual prediction model based on the rule that best fits the predictors.

65   As a commercial and proprietary product, CB has the least algorithmic documentation comparing to random forest.

66   However, Kuhn et al. (2013) ported it into R, which led to its popularity and it is currently being widely used as a

67   regression method.

68   Both RF and CB have been used widely in the recent decades to predict different climatic or environmental parameters.

69   However, there are few studies, which aim to compare RF and CB models. Walton (2008) estimated urban forest canopy

70   cover and impervious surface cover using three different models including CB and RF and compared their performances.

71   They concluded that CB was the best choice for predicting urban impervious surface cover. Noi et al. (2017) compared

72   the results of Multiple Linear Regression, Cubist Regression, and Random Forest Algorithms in estimation of daily air

73   surface temperature. They concluded that using different combinations of data, RF or CB algorithms resulted in high

74   accuracies.

75   In this study, the chosen methodology is based on machine learning, which is considered as a promising modelling

76   method in the fields of agriculture and environmental science (Debeljak and Dzeroski 2011). Here we aim to assess the

77   performance of RF and CB in predicting yearly tile-drainage discharge, to compare the results achieved by both RF and

78   CB, and to analyze and rank the importance of the covariates.

79   **2.   Materials and Methods**

80      **2.1.   Study Area**

81   Denmark is located in northern Europe with a total area of 42,895 km$^2$, of which 66% are used for agricultural purposes

82   (Statistics Denmark, n.d.). The climate is temperate with an approximate mean annual precipitation (P) of 770 mm

83   (Wong, 2013).  The mean temperature is 7.7˚C  ranging from 1.5°C in January to 16.3°C in July. The mean elevation is

84   31 m above sea level and the landscape is generally flat. The geology divides Denmark into two main areas. An eastern

85    part with loamy Weichselian moraines and a western part with sandy glacial outwash plains and Saalian moraines.

86    According to historical maps, wetlands originally covered more than over 20% of the country but due to drainage

87    activities, they have been reduced in extent during the 19th and 20th centuries.

88    **2.2.  Data**

89    Data from 53 drainage stations in different locations and regions of Denmark were used in this study (Fig. 1). It included

90    data from 18 stations established between 2012 and 2016 and historical data from 34 older stations established between

91    1971 and 2009, of which some are still running and some had been shut down (Hansen & Pedersen 1975; Hansen 1981;

92    Simmelsgaard 1994; Grant et al. 2009; Kjær et al. 2011; Kjærgaard et al. 2016). Some data originates from ongoing

93    unpublished drain discharge stations, which have been established in relation to the monitoring of constructed mini-

94    wetlands. Other data belongs to a former project, iDræn (www.idraen.dk, 2011) where data for some of the stations have

95    been published earlier (Hansen et al. 2018a,b; Varvaris et al. 2019a,b). For many stations, drainage discharge (Q) was

96    measured on a daily basis but for some, Q was only measured on a weekly, monthly, or yearly basis. Based on the drain

97    catchment area, yearly values were converted to a water height per year (mm $y^{-1}$) based on the period from 1 July to the

98    end of June to incorporate a full hydrological year. Most of the old stations had available data for a range of 19 to 23

99    years, whereas for some of the new stations there was only data for a few years (1 to 5 years). The lowest discharge (0

100    mm $y^{-1}$) was recorded in southeast Funen during the year 1995 – 1996, whereas the maximum discharge (1183 mm $y^{-1}$)

101    was recorded in eastern Jutland during the year 2015 – 2016. The mean discharge for all the stations was 228 mm $y^{-1}$.

102    The catchment sizes varied from 1 to 164 ha with a mean of 9 ha.

103

104    Thirty-seven different covariates were used as predictors (Table 1). Percolation out of the root zone (Db) was calculated

105    with the simple water balance model EVACROP (Olesen and Heidmann, 1990) driven by input of daily precipitation (P)

106    and reference evapotranspiration ($ET_0$). This was done since it was not expected that P during the growing season would

107    contribute to Q due to the high ET during this period minimizing the percolation out of the root zone. However, the

108    calculated Db is in general closely related to P and Q (Fig. 2). The average Q and the average Db were calculated for

109    each station to determine the ratio between Q and Db (Fig 3). As shown in Figure 3, for seven stations out of 53, the tile

110    drainage discharge is more than the percolated water. These stations are located in large catchments often in stream

111    valleys where external sources (such as regional groundwater) probably flow to the tile drains from outside the

112    catchment. The absolute amounts of discharged water in all the stations is normalized based on catchment area.

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

113    Thirty-three out of 37 covariates were extracted from existing national maps. Topographical variables were calculated by

114    (Møller et al. 2018) based on a digital elevation model (DEM, Fig. 4A) with a 30.4-meter grid size aggregated from a

115    DEM with a 1.6-meter resolution. Adhikari et al. (2013) predicted maps of clay contents for the upper two meter of the

116    soil at a resolution of 30.4 m. These were aggregated by Møller et al. (2018) producing input data in form of maps of clay

117    content in four depth intervals (Clay A%, Clay B%, Clay C%, Clay D%, Table 1, Fig. 4B). Values of clay content were

118    also obtained from a national soil profile database using values from the nearest excavated soil profile. Depth to

119    groundwater (Gwd_model, Table 1, Fig. 4C) was first calculated based on a model at a 500-meter resolution (Henriksen

120    et al., 2012) and then the groundwater table was resampled to a 30.4-meter resolution using bilinear interpolation (Møller

121    et al. 2018). Topographic Wetness Index (TWI, Table 1, Fig. 4D) that quantifies topographic controls of basic

122    hydrological processes (Schillaci et al., 2015) was derived through interactions of fine-scale landform coupled to the up-

123    gradient contributing land surface area by Møller et al. (2018). A map of soil drainage classes (Møller et al., 2017), a

124    rasterized choropleth map of geology (Jacobsen et al., 2015), and a map of wetland areas (Wetlands, Table 1, Greve et al.

125    2014) were also used in the analysis. Horizontal and vertical distances to surface waterbodies were included based on

126    Møller et al. (2018), who calculated horizontal distances to waterbodies as the two-dimensional Euclidean distance to

127    vector layers of waterbodies. Hereafter, they calculated the slope to channel as the angle to the hydrologically nearest

128    waterbody taking into account the surface flow direction. Møller et al. (2018) predicted artificially drained areas

129    (D_DK_New, Table 1) in Denmark by means of a selective model ensemble including number of geographic variables.

130    All 37 covariates were used as input to the statistical models.

131    **2.3.   Models and Measures of Accuracy**

132    As mentioned earlier, two machine-learning algorithms Cubist (CB) (Quinlan, 1993) and Random Forest (RF) (Breiman,

133    2001) were used to predict tile drainage discharge.  Cross-validation was used to adjust the parameters of the models and

134    to assess their predictive accuracy. Cross-validation is a resampling procedure used to evaluate machine-learning models

135    on a given dataset. For CB, the parameters were adjusted to *committees* and *neighbors*. The parameter *committees* sets

136    the number of boosting iterations while the parameter *neighbors* set a number of nearby cases, which can be used for

137    interpolation in order to adjust the predictions. For RF, the parameter *mtry* was adjusted, which sets the number of

138    randomly selected covariates that are available in each split.

139    For both algorithms, three different cross-validation procedures were used. Firstly, in order to assess the ability of each

140    model to predict the tile drain discharge at a new location, *leave-station-out* (LSO) cross-validation was performed. In

141    this procedure, all the measurements were removed from one station in the data sample and a model was trained from the

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

142    remaining measurements and used to predict Q for the excluded cases. This process was repeated for all stations and

143    resulting accuracy was calculated.

144    The stations used in this study are highly clustered in geographic space (Fig. 1). Spatial autocorrelation may therefore

145    affect the accuracy of the LSO procedure as stations may show similar patterns only because they are located close to

146    each other. Therefore, a second cross validation procedure as *leave-cluster-out* was used as well, in which the clusters of

147    stations were left out instead of individual stations. To achieve this, clusters were generated based on the distances

148    between the stations. Stations located less than 10 km from each other were therefore grouped into clusters. This

149    procedure resulted in 23 clusters with 1 – 10 stations each. These clusters were later used for cross-validation.

150    Finally, k-fold cross-validated (KF) RF and CB models were trained on the whole dataset. In this procedure, the dataset

151    were randomly divided into k disjoint folds, which are approximately equal in size. Each of the folds is used to test the

152    generated model from the rest of k-1 folds. The performance of the algorithm was evaluated by the average of the

153    resulting k accuracies from the cross-validation. When a specific value for k was chosen, it could be used in place of k in

154    the reference to the model, which in this case k = 10 and it could therefore be referred as 10-fold cross-validation (Wong

155    2015).

156    In total, six models were trained as the CB and RF models were trained separately with *leave-station-out* (LSO),  *leave-*

157    *cluster-out* (LCO), and *k-fold* (KF) cross validations. The accuracy of all five models were assessed with root mean

158    square error (RMSE):

159    $$RMSE=\sqrt{\frac{\sum_{i=1}^{n}(Q_{m_i} - Q_{o_i})^2}{n}}$$    (1)

160    where $Q_{mi}$ is the predicted value of yearly drainage discharge for the *i*-th instance, $Q_{oi}$ is the observed or measured value

161    of yearly drainage discharge for the *i*-th instance, and n is the total number of instances.

162    The Nash-Sutcliffe efficiency (NSE) was used for validation as well:

163    $$NSE = 1 - \frac{\sum_{i=1}^{n}(Q_{m_i} - Q_{o_i})^2}{\sum_{i=1}^{n}(Q_{o_i} - \bar{Q}_o)^2}$$    (2)

164    where $\bar{Q}_o$ is the mean of observed discharges.

165    Furthermore, to analyze the effect of the covariates in each model, the covariate importance was extracted from all six

166    models. The covariate importance measures were scaled to 100% for the most important covariate in each model. In the

167    beginning, all of the 37 parameters were introduced as covariates to the model. However, the purpose of using machine-

Hydrology and
Earth System
Sciences
Discussions

Open Access

168    learning is to find a simpler way to predict the target and to determine the most effective parameters, which helps to

169    reduce the number of covariates and exclude the ineffective ones.

170

171    **3.    Results and Discussion**

172        **3.1.    Model accuracy**

173    The most accurate predictions were obtained by 10-Fold (KF) cross-validated Cubist (CB) and 10-Fold (KF) cross-

174    validated random forest (RF) with RMSE of 75 and 77 mm/y and NSE of 0.73 and 0.74, respectively (Fig. 5, Table 2).

175    According to Singh et al. (2005), an acceptable value for RMSE in hydrological modelling would normally be half of the

176    standard deviation of training data, which for the current data set was 166 (mm/y).  Therefore, leave-station-out (LSO)

177    cross-validated random forest (RF) with an RMSE of 110 mm/y and LCO cross-validated CB with an RMSE of 113

178    mm/y could be considered as acceptable models regarding the prediction of Q.

179    The purpose of performing three different cross-validations was to test the model accuracy with and without the effect of

180    geological biases. In LSO, a single station containing an entire data set is removed from the training dataset as the target

181    of prediction. However, the model is still trained on the neighbor stations, which are regionally close to the target. That

182    could cause overfitting issues. On the other hand, the LCO ensures that on each run of the model, one of the 23 clusters is

183    excluded as the prediction target, which diminishes the possibility of overfitting caused by geo-regional similarities.

184    Finally, KF randomly divided the whole dataset into 10 fold with equal size, which does not consider the distribution of

185    the stations. Data is sampled based on the rows and the difference in size between the training set used in each fold and

186    the entire dataset is only a single pattern. Each fold contains 41 rows that are selected randomly and each time one of the

187    10 folds is the validation or test data set. The repeated cross-validation guarantees that different combinations of

188    randomly selected stations are in different training folds limiting the possibility of overfitting.

189    With all three cross-validation methods, the accuracies with RF and CB were quite similar. Furthermore, the accuracies

190    calculated with LSO and LCO are relatively similar, compared to KF, which had a substantially higher NSE and lower

191    RMSE than the two other cross-validation methods.

192        **3.2.    Covariate importance**

193    Results of all the six models indicate that the percolation or discharge out of the root zone (Db) has the largest effect on

194    the tile-drainage discharge prediction with 100% importance (Fig. 6). The analyses show that elevation (DEM) follows the

195    Db as the second most important covariate in all the models with more than 80 % importance in LSO-CB and LCO-CB

196    (Fig. 6 a and b) and between approximately 40 to 50% effectiveness for the other four models (Fig. 6c to f). The clay

197    content in the D horizon was the third most important covariate in KF-CB and KF-RF (Fig. 6c to f). For the LCO-CB and

198    LSO-RF models, horizontal distance to the nearest waterbody appears as the third most important covariate with 45% and

199    21% importance, respectively (Fig. 6a and d). Whereas for the LSO-CB model, clay content in the C horizon and the LCO-

200    RF model clay content in the B horizon where the third most important covariates (Fig. 6 b and c). The rest of the list

201    differs between the different models. However, it is observable that for the RF models (Fig 6c to e) only the first covariates

202    have a significant effect where the rest have less than 20% importance. Nevertheless, for all CB models (Fig 6 a, b, and f)

203    the top 10 covariate all have more than 20% importance. As previously stated, percolation and elevation have the largest

204    importance to all of the trained models for the prediction of discharge. Based on the analyses of covariate importance, the

205    results of the predictions for the two most effective covariates were compared to their measurements (Fig. 7). This

206    comparison demonstrates how well the models can simulate the relationship between the most important covariates (Db

207    and elevation) and the prediction target (Q). The open black circles represent the predictors on the x-axis against the

208    measured drainage discharge (Q) on the y-axis. The red open circles represent the predictors on x axis and predicted

209    drainage discharge (Q) on the y-axis by each of the six models mentioned on top of the plots. The best match could be

210    observed on the k-fold cross-validated CB (Fig. 7 e and f).

211    **3.3. Discussion**

212    Similar studies targeting the prediction of discharge with machine learning models developed their models in a catchment

213    scale for time series and chose the daily meteorological data, agricultural practices, and crop management as covariates

214    (Kuzmanovski et al. 2015, Zia et al. 2015). Also in these studies, they used 10-fold cross-validation to evaluate the

215    robustness of their model performance. The present study was carried out on a larger scale with catchments of different

216    sizes distributed in different regions. Along with the percolation, a number of different geological features were used as

217    input parameters to assess if it is possible to predict the tile drainage discharge based on spatially variable geophysical

218    characteristics of the different sites. In the few similar studies (Rasouli et al. 2012, Kuzmanovski et al. 2015, Zia et al.

219    2015), the study area was either one specific catchment or few fields or catchments very close to each other. This means

220    that the geological features were similar. Being able to train machine-learning models on different catchments in very

221    different locations had enabled us to make use of differing geographical characteristics as predictor variables. Predictions

222    were carried out in a yearly basis and were cross-validated with three different methods.

223    The accuracies of RF and CB models in comparison to each other for all the cross-validation methods were quite similar.

224    On the other hand, the obtained accuracies from LSO and LCO are relatively similar but lower compared to KF, which

Hydrology and
Earth System
Sciences
Discussions

Open Access
EGU

225  had a substantially higher NSE and lower RMSE than the two other cross-validation methods. The higher accuracies

226  achieved by KF is most likely results from having the observations of a given station from other years during the

227  prediction procedure. The accuracy obtained with KF could be considered as the internal accuracy of the model, while

228  LSO and LCO better represent the accuracies at new locations without previous measurements of tile drain discharge at

229  the same station. The proposed tile-drainage discharge predictive model is not dependent on the climatic and constantly

230  measured data and makes it possible to use different geographical properties as predictive parameters.

231  Logically, Db is the main driving variable since it takes into account water lost by evaporation from the soil surface,

232  transpiration of water by the crop, and the increase of water stored in the soil. During the growing season, a high value of

233  P will not necessarily lead to a corresponding high value of Q since it is only the part of P that infiltrate out of the root

234  zone that potentially can flow into the tile drains. It is also expected that the clay content in the soil, especially the clay

235  content in the lower horizons below tile drain depths, would have an effect on the drain discharge. A high clay content in

236  the subsoil would lead to a secondary groundwater table building up outside the growing season to the level of the tile

237  drains. That the clay content not play a more important role as a covariate might be explained by the relatively high

238  prediction error of the clay content especially at lower depths for the used soil maps.

239  The position of the tile-drained field in the landscape will have an effect on the tile drain discharge. At low positions in

240  the landscape, the flow of water to the drains is expected to be relatively high due to a high contributing area of expected

241  incoming regional groundwater generated from a larger area outside the tile-drained field. Such areas are also indicated

242  in Figure 3 corresponding to high values of Q/Db. On the other hand, at higher positions in the landscape with no or only

243  a minor contribution of regional groundwater, a proportional part of the water infiltrating into the drains is generated

244  mainly locally from water percolating out of the root zone (Db). It was expected that DEM derived indices such as TWI

245  or SagaWI (Table 1) would describe more precisely the contribution of water in the tile drains and therefore supposed to

246  be important covariates. Both indices attempt to describe the hydrological flow paths in the landscape and should be able

247  to identify areas with a high contribution of water flowing to the drains. However, only for the k-fold cross-validated RF

248  model (Fig. 6E), TWI is found within the list of the top 10 most import covariates. On the other hand, DEM is placed as

249  the second most or the most important covariate for all models. This proves that the position in the landscape does have

250  an effect on the tile drain discharge. That the derived topographical indices only play a minor role in the statistical

251  models might be related to the fact that it can vary considerably within the individual drained catchments. On the other

252  hand, other derived DEM indices such as valley depth (Valldepth), vertical distance to the nearest waterbody (Vdtochn),

253  horizontal distance to the nearest waterbody (Hdtochn), and downhill gradient to the nearest waterbody (Slptochn) are all

254  found in the top 10 list.

255    By applying input from a distributed model predicting Db it is possible to apply the developed model on a national scale

256    developing maps that can be used as a tool to predict the yearly drain discharge. National water resource models in

257    Denmark exists that can be used for such purposes (e.g. Højberg et al. 2013). Outputs from the model can be based on

258    averages for a certain period. Also, the possible variation between years as well as outputs in relation to future climatic

259    scenarios can be studied.

260    **4.    Conclusion**

261    For the current study, two different machine-learning models (RF and CB) were applied on a relatively big dataset

262    containing measured yearly drainage discharge (Q) and 37 parameters as covariates and the results indicated a successful

263    implementation. The predictive models were trained on 53 drainage stations distributed all over Denmark with different

264    characteristics and multiple years of data and cross-validated with three different methods. The best results were

265    achieved by k-fold (KF) cross-validated Cubist (CB) and random forest (RF) and the performance measures certifies the

266    results. RMSE and NSE of both models indicates a good accuracy of the predictive models based on the hydrological

267    modelling standards. Instead of physically-based models that acquire numerous parameters, machine learning models

268    could perform as strong tools for quantifying the tile-drainage discharge with lower complexity. In this study, percolation

269    or discharge out of the root zone (Db) calculated with the simple water balance model EVACROP, and elevation (DEM)

270    where the most important covariate for predicting yearly discharge. Finally, it was concluded that considering the

271    distribution of stations, the method of sampling and the cross-validation has a large effect on estimates of model

272    accuracies. The developed model can be used in relation to a national mapping of yearly tile drain discharge.

273    **Acknowledgments**

277    **References**

278    Adhikari, K., R.B. Kheir, M.B. Greve, P.K. Bocher, B.P. Malone, B. Minasny, et al. 2013. High-Resolution 3-D

279    Mapping of Soil Texture in Denmark. Soil Sci Soc Am J 77: 860-876. doi:10.2136/sssaj2012.0275.

280    Ayars, J.E., E.W. Christen and J.W. Hornbuckle. 2006. Controlled drainage for improved water management in and

281    regions irrigated agriculture. Agr Water Manage 86: 128-139. doi:10.1016/j.agwat.2006.07.004.

282   Basha, E.A., S. Ravela and D. Rus. 2008. Model-Based Monitoring for Early Warning Flood Detection. Sensys'08:

283   Proceedings of the 6th Acm Conference on Embedded Networked Sensor Systems: 295-308.

284   Blann, K.L., J.L. Anderson, G.R. Sands and B. Vondracek. 2009. Effects of Agricultural Drainage on Aquatic

285   Ecosystems: A Review. Crit Rev Env Sci Tec 39: 909-1001. doi:10.1080/10643380801977966.

286   Bredehoeft, J. 2005. The conceptualization model problem—surprise. Hydrogeology Journal 13: 37-46.

287   doi:10.1007/s10040-004-0430-5.

288   Breiman, L. 2001. Random forests. Mach Learn 45: 5-32. doi:Doi 10.1023/A:1010933404324.

289   De Schepper, G., R. Therrien, J.C. Refsgaard, X. He, C. Kjaergaard and B.V. Iversen. 2017. Simulating seasonal

290   variations of tile drainage discharge in an agricultural catchment. Water Resour Res 53: 3896-3920.

291   doi:10.1002/2016wr020209.

292   Debeljak, M. and S. Dzeroski. 2011. Decision Trees in Ecological Modelling. Modelling Complex Ecological Dynamics:

293   An Introduction into Ecological Modelling for Students, Teachers & Scientists: 197-209. doi:10.1007/978-3-642-05029-

294   9_14.

295   Grant, R., G. Blicher-Mathiesen, P.G. Jensen, B. Hansen, L. Thorling 2010. Landovervågningsoplande 2009. NOVANA.

296   Danmarks Miljøundersøgelser, Aarhus Universitet. Faglig rapport fra DMU nr. 802. 124 pp.

297   (https://www2.dmu.dk/Pub/FR802.pdf).

298   Greve, M.H., O.F. Christensen, M.B. Greve and R.B. Kheir. 2014. Change in Peat Coverage in Danish Cultivated Soils

299   During the Past 35 Years. Soil Sci 179: 250-257. doi:10.1097/Ss.0000000000000066.

300   Hansen, A.L., R. Jakobsen, J.C. Refsgaard, A.L. Højberg, B.V. Iversen and C. Kjærgaard 2018a. Groundwater dynamics

301   and effect of tile drainage on water flow across the redox interface in a Danish Weichsel till area. Advances in Water

302   Resources 123:23-39.

303   Hansen, A.L., A. Storgaard, X. He, A.L. Højberg, J.C. Refsgaard, B.V. Iversen, C. Kjærgaard 2018b. Importance of

304   geological information for assessing drain flow in a Danish till landscape, Hydrological Processes 33:450-462.

305   Hansen, B. 1981. Drænvandskvantitet og –kvalitet i Susåens opland. Suså-projekt. Dansk Komite for Hydrologi, Rapport

306   Nr. Suså H 19. København, Denmark. 67 pp.

307   (https://soeg.kb.dk/permalink/45KBDK_KGL/fbp0ps/alma99122878363105763)

308 Hansen, L.; E.F. Pedersen 1975. Drænvandsundersøgelser 1971-74. Tidsskr. Planteavl 79 (670-688).

309 (http://agris.fao.org/agris-search/search.do?recordID=US201303003414)

310 Henriksen, H.J., Højberg, A.L., Olsen, M., Seaby, L.P., van der Keur, P., Stisen, S., Troldborg, L., Sonnenborg, T.O.,

311 Refsgaard, J.C., 2012. Klimaeffekter på hydrologi og grundvand - Klimagrundvandskort. Aarhus University.

312 (https://www.klimatilpasning.dk/media/340310/klimagrundvandskort.pdf)

313 Jakobsen, P.R., Hermansen, B., Tougaard, L., 2015. Danmarks digitale jordartskort 1:25000 version 4.0. GEUS.

314 (http://pubs.geus.net/Danmark/jordartskort/Jordart_25000_beskriv.pdf)

315 Kjær, J., A.E. Rosenbom, W. Brüsch, R.K. Juhler, L. Gudmundsson, F. Plauborg, R. Grant, P. Olsen 2011. The Danish

316 Pesticide Leaching Assessment Programme - Monitoring results May 1999–June 2010. Geological Survey of Denmark

317 and Greenland and Aarhus University. Copenhagen, Denmark. 110 pp. (http://pesticidvarsling.dk/xpdf/vap-results-99-

318 10.pdf)

319 Kjærgaard, C.; Iversen, B.V.; Højberg, A.L.; Mathiesen, G.B. 2016. Drænmålinger som grundlag for emissionsbaseret

320 kvælstofregulering. In Hvid, S.K. Måling af kvælstofudledning emmisionsbaseret kvælstofregulering på bedriftsniveau.

321 SEGES, Aarhus, Denmark. Delrapport C. 67 pp.

322 (https://www.landbrugsinfo.dk/Afrapportering/planter_og_miljoe/2016/Sider/pl_po_999_3682_b3_Delrapport_C_Maali

323 nger_i_draenra.pdf?download=true)

324 Kuhn, M., 2008. Building predictive models in R using the caret package. J. Stat. Softw. 28 (5), 1–26.

325 Kuzmanovski, V., A. Trajanov, F. Leprince, S. Dzeroski and M. Debeljak. 2015. Modeling water outflow from tile-

326 drained agricultural fields. Sci Total Environ 505: 390-401. doi:10.1016/j.scitotenv.2014.10.009.

327 Magmedov, V.G., M.A. Zakharchenko, L.I. Yakovleva and M.E. Ince. 1996. The use of constructed wetlands for the

328 treatment of run-off and drainage waters: The UK and Ukraine experience. Water Sci Technol 33: 315-323. doi:Doi

329 10.1016/0273-1223(96)00247-8.

330 Messer, T.L., M.R. Burchell, F. Birgand, S.W. Broome and G. Chescheir. 2017. Nitrate removal potential of restored

331 wetlands loaded with agricultural drainage water: A mesocosm scale experimental approach. Ecol Eng 106: 541-554.

332 doi:10.1016/j.ecoleng.2017.06.022.

333 Møller, A.B., A. Beucher, B.V. Iversen and M.H. Greve. 2018. Predicting artificially drained areas by means of a

334 selective model ensemble. Geoderma 320: 30-42. doi:10.1016/j.geoderma.2018.01.018.

335     Møller, A.B., Beucher, A., Iversen, B.V., Greve, M.H., 2017. Prediction of soil drainage classes in Denmark by means of

336     decision tree classification. Geoderma. http://dx.doi.org/10.1016/j.geoderma.2017.10.015.

337     Noi, P.T., J. Degener and M. Kappas. 2017. Comparison of Multiple Linear Regression, Cubist Regression, and Random

338     Forest Algorithms to Estimate Daily Air Surface Temperature from Dynamic Combinations of MODIS LST Data.

339     Remote Sens-Basel 9. doi:ARTN 398 10.3390/rs9050398.

340     Olesen, J.E., Heidmann, T., 1990. EVACROP. Et program til beregning af aktuel fordampning og afstrømning fra

341     rodzonen. Version 1.01. AJMET Arbejdsnotat 9, Statens Planteavlsforsøg.

342     (https://docplayer.dk/storage/26/8897995/1577277319/r5LW4nhUDA12pzTOA7v1PA/8897995.pdf)

343     Park, S., J. Im, E. Jang and J. Rhee. 2016. Drought assessment and monitoring through blending of multi-sensor indices

344     using machine learning approaches for different climate regions. Agr Forest Meteorol 216: 157-169.

345     doi:10.1016/j.agrformet.2015.10.011.

346     Rasouli, K., Hsieh, W.W., Cannon, A.J., 2012. Daily streamflow forecasting by machinelearning methods with weather

347     and climate inputs. J. Hydrol. 414, 284–293.

348     Schillaci, C., A. Braun, and J. Kropacek. 2015. Terrain analysis and landform recognition; Chapter 2.4.2, in

349     Geomorphological Techniques; British Society for Geomorphology. 18 pp.

350     Schilling, K.E., P. Jindal, N.B. Basu and M.J. Helmers. 2012. Impact of artificial subsurface drainage on groundwater

351     travel times and baseflow discharge in an agricultural watershed, Iowa (USA). Hydrol Process 26: 3092-3100.

352     doi:10.1002/hyp.8337.

353     Simmelsgaard, S.E. 1994. Nitratkvælstof i drænvand 1971-91. Statens Planteavlsforsøg, SP-rapport nr. 47. Statens

354     Planteavlsforsøg, Lyngby, Denmark. 67 pp.

355     (https://soeg.kb.dk/permalink/45KBDK_KGL/13io73r/faoagrisUS201300291267)

356     Singh, J., H.V. Knapp, J.G. Arnold and M. Demissie. 2005. Hydrological modeling of the iroquois river watershed using

357     HSPF and SWAT. J Am Water Resour As 41: 343-360. doi:DOI 10.1111/j.1752-1688.2005.tb03740.x.

358     Stenberg, M., B. Ulen, M. Soderstrom, B. Roland, K. Delin and C.A. Helander. 2012. Tile drain losses of nitrogen and

359     phosphorus from fields under integrated and organic crop rotations. A four-year study on a clay soil in southwest

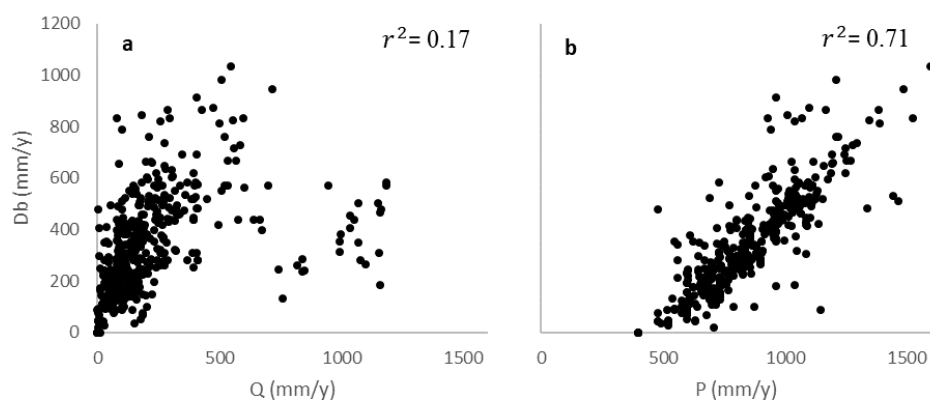360     Sweden. Sci Total Environ 434: 79-89. doi:10.1016/j.scitotenv.2011.12.039.

361    Varvaris, I., C.D. Børgesen, C. Kjærgaard, B.V. Iversen 2019a. Three two-dimensional approaches for simulating the

362    water flow dynamics in a heterogeneous tile-drained agricultural field in Denmark. Soil Science Society of American

363    Journal 82:1367–1383.

364    Varvaris, I, P. Moldrup, Z. Pittaki-Chrysodonta, L. W. de Jonge, and B.V. Iversen 2019b. Coupling vis-NIR and

365    pedotransfer functions for predicting hydraulic properties to simulate water flow dynamics in a tile-drained agricultural

366    field. Vadose Zone Journal (accepted).

367    Walton, J.T. 2008. Subpixel urban land cover estimation: Comparing Cubist, Random Forests, and support vector

368    regression. Photogramm Eng Rem S 74: 1213-1222. doi:Doi 10.14358/Pers.74.10.1213.

369    Wong, T.T. 2015. Performance evaluation of classification algorithms by k-fold and leave-one-out cross validation.

370    Pattern Recogn 48: 2839-2846. doi:10.1016/j.patcog.2015.03.009.

371    Zia, H., N. Harris, G. Merrett and M. Rivers. 2015. Predicting discharge using a low complexity machine learning model.

372    Comput Electron Agr 118: 350-360. doi:10.1016/j.compag.2015.09.012.

373

374

375

**Figure 1. Study area and the location of the 53 drainage stations throughout Denmark**
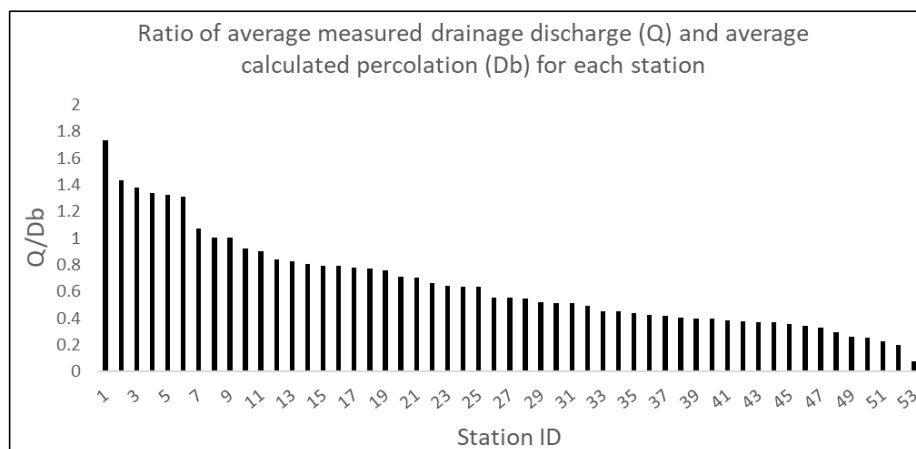


377

378  **Figure 2. a) Measured yearly drainage discharge (Q) against calculated percolation (Db) b) Observed precipitation (P) against**
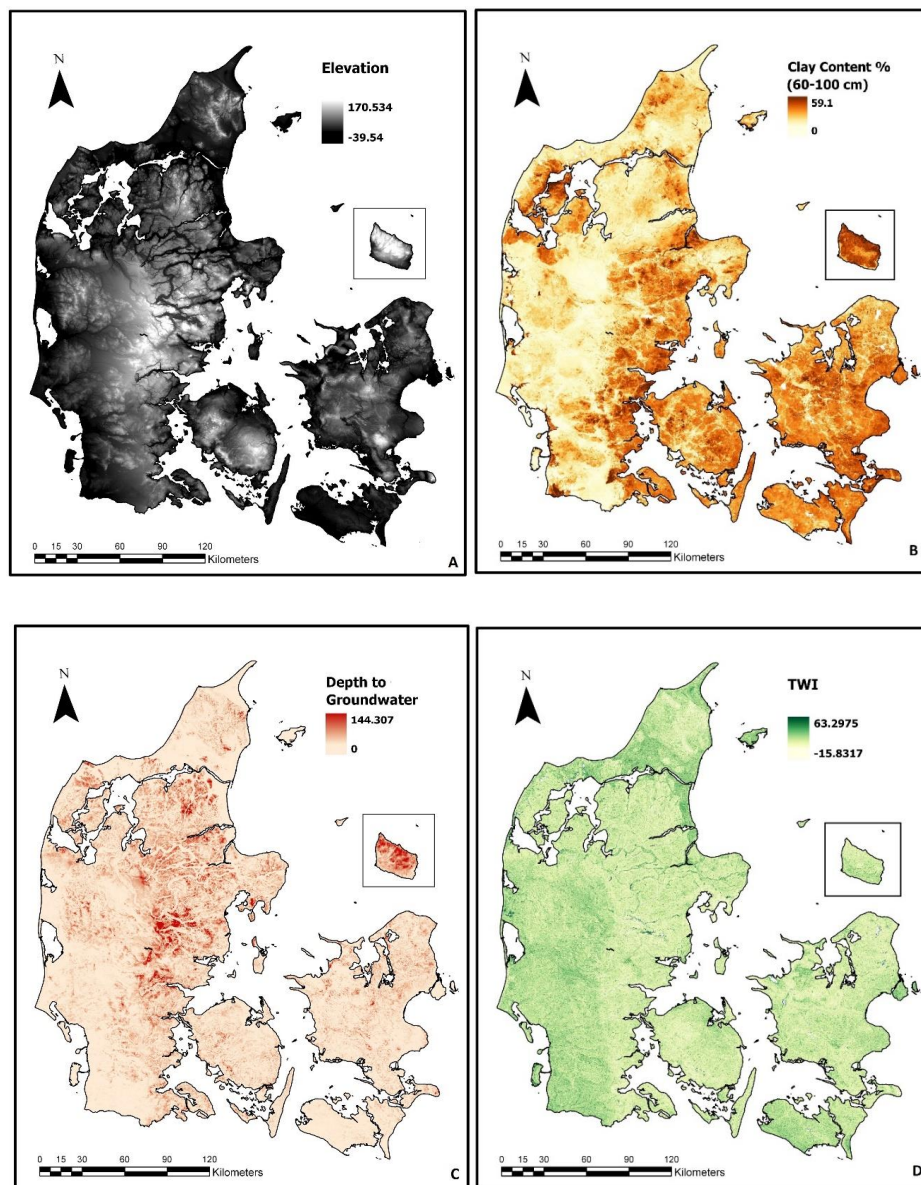379  **calculated percolation (Db)**

380

381

Figure 3. Ratio of average measured drainage discharge (Q) and average calculated percolation (Db) for each station
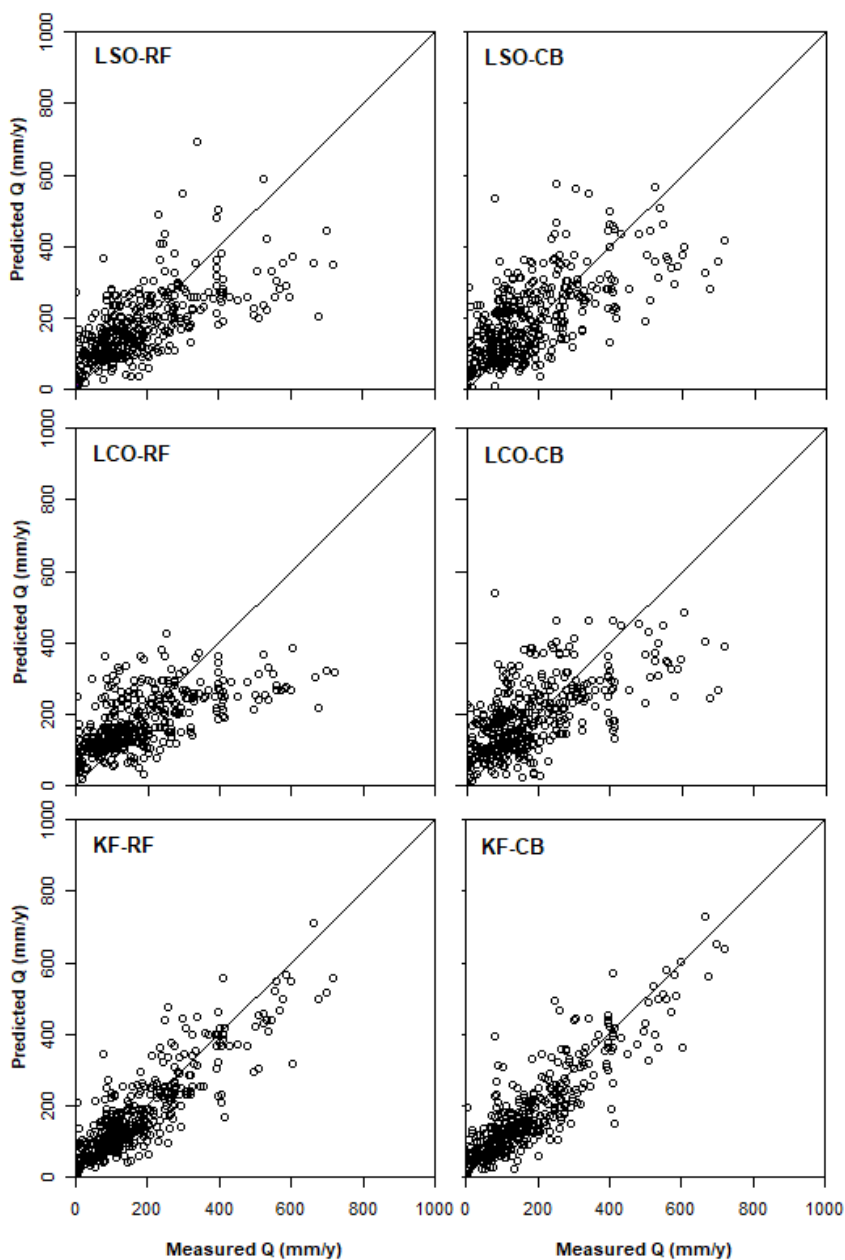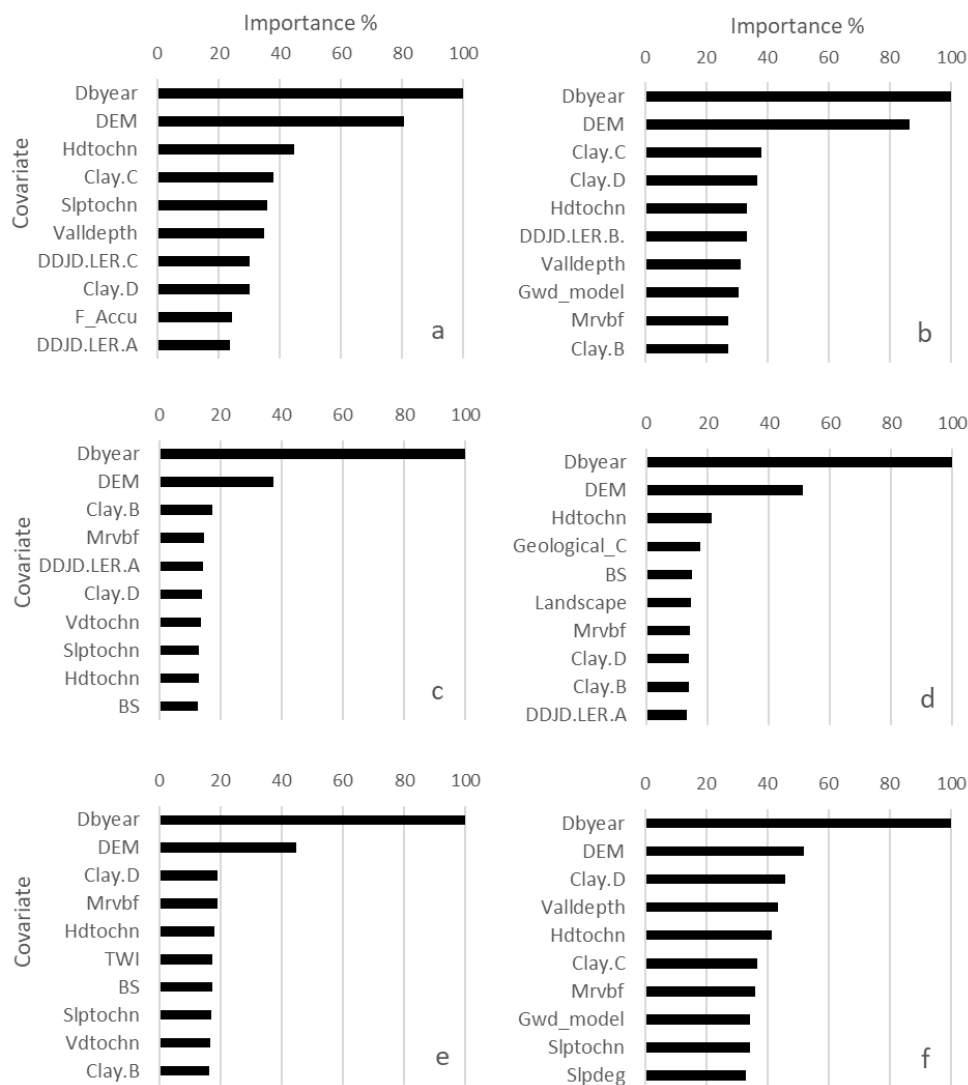
384



385

386     Figure 4. A. Elevation based on a Digital Elevation Map (DEM).  B. Aggregated clay content in the C-horizon (Møller et al.,

387     2018) C. Interpolated depth to groundwater (Møller et al., 2018) D. Topographical wetness index (Møller et al., 2018)
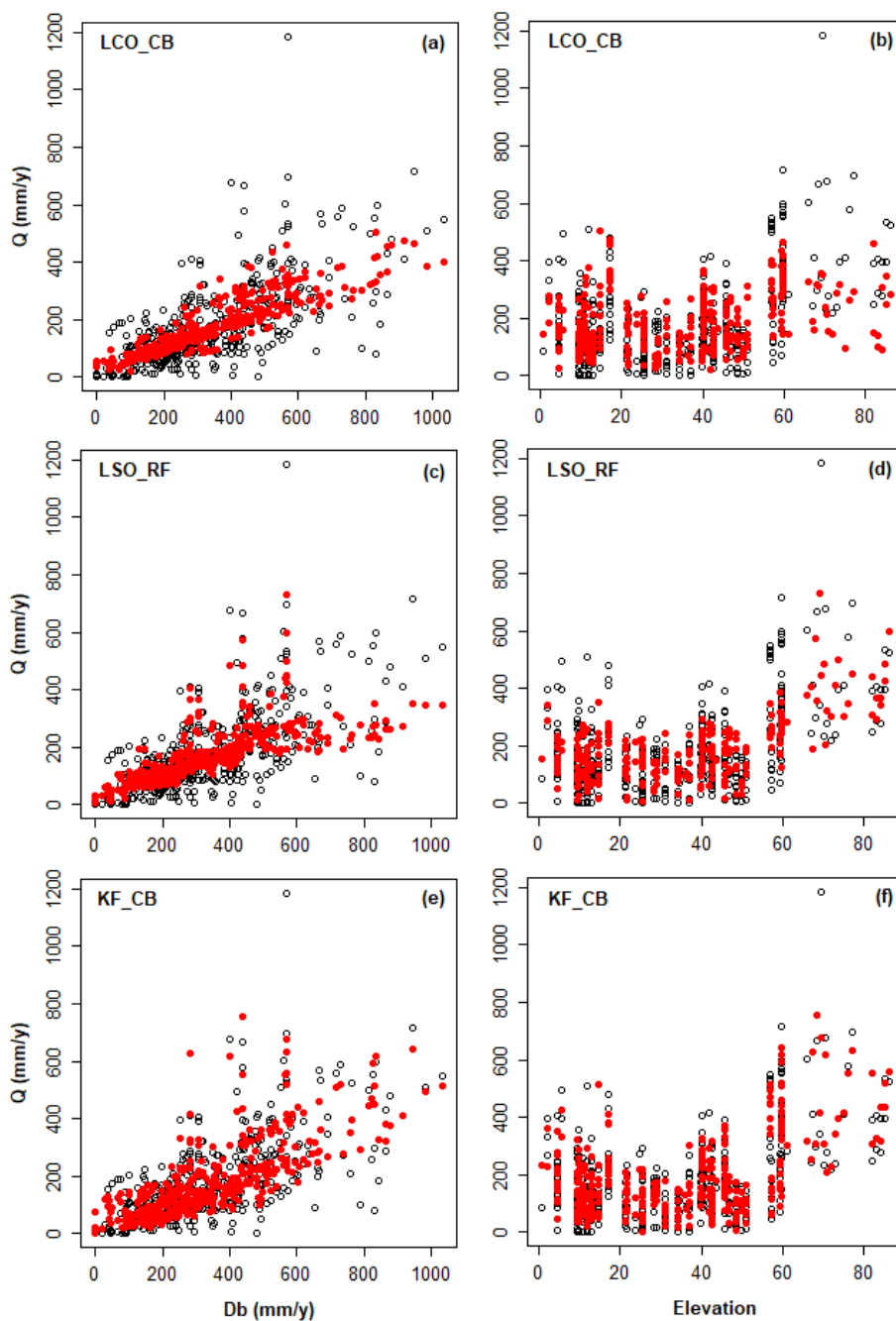
Hydrology and
Earth System
Sciences
Discussions

Open Access

388



**Figure 5. LSO-RF: Leave station out cross-validated random forest model. LSO-CB: Leave station out cross-validated cubist**

**model. LCO-RF: Leave cluster out cross-validated random forest model. LCO-RF: Leave cluster out cross-validated cubist**

**model. KF-RF: K-Fold cross-validated random forest model. KF-CB: k-fold cross-validated Cubist model.**

392

**Figure 6. a) Top 10 most important covariates of the leave-cluster-out cross-validated CB model b) Top 10 most important covariates of the leave-station-out cross-validated CB model c) Top 10 most important covariates of the leave-cluster-out cross-validated RF model d) Top 10 most important covariates of the leave-station-out cross-validated RF model E) Top 10 most important covariates of k-fold cross-validated RF model F) Top 10 most important covariates of the k-fold cross-validated CB model.**

398

**Figure 7. a, c, and e) Measured discharge against calculated percolation in black open circles, predicted discharge against**

**calculated percolation in red open circles for the selected models with the best prediction. b, d, and f) Measured discharge against**

401   elevation in black open circles, predicted discharge against elevation in red open circles for selected models with the best

402   prediction

403   **Table 1. List of covariates used to predict the discharge including a description of the parameter and a range specifying the**

404   **type of covariate.**

| Predictors | Description | Range/ Class |
|---|---|---|
| Db | Percolation/Discharge out of the root zone (mm y$^{-1}$) | 0 – 1033 |
| Geological_R | Geological region | 7 classes |
| DEM | Elevation (m) | 0.74 – 83.16 |
| Geological_C | Geology of the area | 10 classes |
| F_Accu | Flow Accumulation/Number of unslope cells | 1 – 1108 |
| SagaWI | SAGA Wetness Index | 12.16 - 16.58 |
| TWI | Topographic Wetness Index | 3.47 – 12.33 |
| BS | Depth of Sink (m) | 0 – 2.17 |
| D_Class | Drainage class | 5 classes |
| Clay A %† | Clay content 0-30 cm soil depth | 3 – 20.3 |
| Clay B %† | Clay content 30-60 cm soil depth | 2 – 29.1 |
| Clay C %† | Clay content 60-100 cm soil depth | 1.5 – 31 |
| Clay D %† | Clay content 100-200 cm soil depth | 2.2 – 32.6 |
| DDJD LER-A%‡ | Clay content in A horizon | 3 – 24.8 |
| DDJD LER-B%‡ | Clay content in B horizon | 0 – 31.97 |
| DDJD LER-C%‡ | Clay content in C horizon | 0 – 29.1 |
| JB | Danish soil classification for the A horizon | 12 classes |
| Gwd_Int | Depth to groundwater table interpolated from well observations and surface water (m) | 0 – 25.31 |
| Wetlands | 0: Non-wetlands; 1: Wetlands; 2: Central wetlands; 3: Peatlands. | 4 classes |
| D_DK_New | Artifical drainage-new map | 2 classes |
| DP_New | Drainage probability-new map | 0 – 0.86 |

| D_DK | Artifcial drainage-old map | 2 classes |
| DP | Drainage probability-old map | 0 – 0.82 |
| Demdetrend | Elevation minus the mean elevation in a 4 km radius (m) | -11.4 – 26.04 |
| Dirinsola | Direct insolation (kWh/year) | 1150.08 – 1348.61 |
| Gwd_model | Depth to groundwater from the model (m) | 0 – 32.42 |
| Hdtochn | Horizontal distance to the nearest waterbody (m) | 0 – 1114.89 |
| Midslppos | Mid-slope position | 0 – 0.7 |
| Mrvbf | Multi-resolution index of valley bottom flatness | 0.07 – 8.68 |
| Slpdeg | Surface slope gradient (degrees) | 0.09 – 7.53 |
| Slptochn | Downhill gradient to the nearest waterbody (m) | 0 – 3.48 |
| Vdtochn | Vertical distance to the nearest waterbody (m) | 0 – 19.28 |
| Valldepth | Valley depth (m) | 2.43 – 21.35 |
| Landscape | Landform types | 11 classes |

405 † From the map of Adhikari et al. (2013); ‡from the national soil database

406 **Table 2. Error summary of six trained models**

| Model \\ Error | LSO-CB | LCO-CB | LSO-RF | LCO-RF | KF-RF | KF-CB |
|---|---|---|---|---|---|---|
| RMSE | 116.53 | 115.04 | 110.65 | 115.82 | 76.05 | 70.98 |
| NSE | 0.37 | 0.39 | 0.44 | 0.38 | 0.73 | 0.74 |

407

408