

Journal of Hydrology

Machine Learning Techniques as an Alternative Approach to Rainfall-Runoff Inundation Models for Flood Susceptibility Prediction

--Manuscript Draft--

Manuscript Number:	HYDROL46050
Article Type:	Research paper
Keywords:	machine learning; random forest; LightGBM; CatBoost; flooding susceptibility mapping; rainfall-runoff inundation model
Corresponding Author:	Mohamed Saber Uji, JAPAN
First Author:	Mohamed Saber
Order of Authors:	Mohamed Saber Tayeb Boulmaiz Mawloud Guermoui Karim I. Abdrado Sameh A. Kantoush Tetsuya Sumi Hamouda Boutaghane Tomoharu Hori Doan Van Binh Binh Quang Nguyen Thao T. P. Bui Ngoc Duong Vo Emad Mabrouk
Abstract:	Vietnam has experienced many natural disasters, particularly typhoons. This study aims to examine three machine learning (ML) approaches—random forest (RF), LightGBM, and CatBoost—for flooding susceptibility maps (FSMs) in the Vu Gia-Thu Bon (VGTB) River Basin of Vietnam. The results of ML are compared with those of the rainfall-runoff model, and different training dataset sizes are utilized in the performance assessment. Ten independent factors that influence the FSMs in the study area, namely, aspect, rainfall, curvature, DEM, horizontal distance from the river, geology, hillshade, land use, slope, and stream power index, are assessed. An inventory map that includes approximately 850 flooding sites is considered based on several post-flood surveys after the typhoons in 1999, 2006, 2007, 2009, 2013, and 2020. The inventory dataset is randomly divided into two sets: training (70%), and testing (30%). The AUC-ROC results are 97.9%, 99.5%, 99.5% for CatBoost, LightGBM, and RF, respectively. The FSMs developed by the ML methods show good agreement with flood inundation mapping developed using the rainfall-runoff model. The FSMs show that downstream areas (both urbanized and agricultural) are under “high” and “very high” levels of susceptibility. Additionally, different sizes of the input datasets (i.e., 30, 60, 90, 200, 400, 600, 800, 1000, and 1250 data points) are tested to determine the least number of data points having an acceptable reliability. The results show that the ML methods can reasonably predict FSMs, regardless of the number of training samples, although the final FSMs show some spatial differences when changes in susceptibility level are seen. The developed FSMs for such typhoon-prone regions can be used by decision-makers and planners in Vietnam to propose effective mitigation measures for community resilience and development.
Suggested Reviewers:	Hisham Eldardiry dardiry@uw.edu

<p>Expert in hydrology</p>
<p>Hadir Abd-El Moneim hadir_eng@yahoo.com Expert in Hydrological Modeling</p>
<p>Ashraf Elmoustafa elmoustafa010@yahoo.co.uk</p>
<p>Ali Al-Maktoom ali4530@squ.edu.om</p>
<p>Takahiro Sayama sayama.takahiro.3u@kyoto-u.ac.jp</p>

1 **Machine Learning Techniques as an Alternative Approach to Rainfall-Runoff
2 Inundation Models for Flood Susceptibility Prediction**

3 *Mohamed Saber^{1,*}, Tayeb Boulmaiz², Mawloud Guermout³, Karim I. Abdrado⁴, Sameh A.
4 Kantoush¹, Tetsuya Sumi¹, Hamouda Boutaghane⁵, Tomoharu Hori¹, Doan Van Binh⁶,
5 Binh Quang Nguyen^{1,7}, Thao T. P. Bui¹, Ngoc Duong Vo⁷, Emad Mabrouk^{8,9}*

6

7 *¹Disaster Prevention Research Institute (DPRI), Kyoto University, Kyoto 611-0011, Japan;
8 *mohamedmd.saber.3u@kyoto-u.ac.jp, kantoush.samehahmed.2n@kyoto-u.ac.jp,
9 sumi.tetsuya.2s@kyoto-u.ac.jp*

10 *²Materials, Energy Systems Technology and Environment Laboratory, Ghardaia
11 University, Ghardaia, Algeria; boulmaiz.tayeb@univ-ghardaia.dz*

12 *³Unité de Recherche Appliquée en Energies Renouvelables, URAER, Centre de
13 Développement des Energies Renouvelables, CDER, 47133 Ghardaïa, Algeria;
14 gue.mouloud@gmail.com*

15 *⁴Faculty of Urban and Regional Planning, Cairo University, Giza 12613, Egypt;
16 m.karim.ibrahim@cu.edu.eg*

17 *⁵Hydraulic Department, Badji Mokhtar-Annaba University, P.O. Box 12, Annaba, Algeria;
18 hamouda.boutaghane@univ-annaba.dz*

19 *⁶Master Program in Water Technology, Reuse, and Management, Faculty of Engineering,
20 Vietnamese German University, 2-Le Lai Street, Hoa Phu Ward, Thu Dau Mot City, Binh
21 Duong Province 820000, Vietnam; binh.dv@vgu.edu.vn*

22 *⁷The University of Danang -University of Science and Technology, 54 Nguyen Luong Bang,
23 Danang, Vietnam; nqbinh@dut.udn.vn*

24 *⁸College of Engineering and Technology, American University of the Middle East, Egaila
25 54200, Kuwait; emad.mabrouk@aum.edu.kw*

26 *⁹Department of Mathematics, Faculty of Science, Assiut University, Assiut 71516, Egypt;
27 mabrouk@aun.edu.eg*

28 ***Corresponding:** *mohamedmd.saber.3u@kyoto-u.ac.jp*

29

31 **Abstract**

32 Vietnam has experienced many natural disasters, particularly typhoons. This study aims to
33 examine three machine learning (ML) approaches—random forest (RF), LightGBM, and
34 CatBoost—for flooding susceptibility maps (FSMs) in the Vu Gia-Thu Bon (VGTB) River Basin
35 of Vietnam. The results of ML are compared with those of the rainfall–runoff model, and different
36 training dataset sizes are utilized in the performance assessment. Ten independent factors that
37 influence the FSMs in the study area, namely, aspect, rainfall, curvature, DEM, horizontal
38 distance from the river, geology, hillshade, land use, slope, and stream power index, are assessed.
39 An inventory map that includes approximately 850 flooding sites is considered based on several
40 post-flood surveys after the typhoons in 1999, 2006, 2007, 2009, 2013, and 2020. The inventory
41 dataset is randomly divided into two sets: training (70%), and testing (30%). The AUC-ROC
42 results are 97.9%, 99.5%, 99.5% for CatBoost, LightGBM, and RF, respectively. The FSMs
43 developed by the ML methods show good agreement with flood inundation mapping developed
44 using the rainfall-runoff model. The FSMs show that downstream areas (both urbanized and
45 agricultural) are under “high” and “very high” levels of susceptibility. Additionally, different
46 sizes of the input datasets (i.e., 30, 60, 90, 200, 400, 600, 800, 1000, and 1250 data points) are
47 tested to determine the least number of data points having an acceptable reliability. The results
48 show that the ML methods can reasonably predict FSMs, regardless of the number of training
49 samples, although the final FSMs show some spatial differences when changes in susceptibility
50 level are seen. The developed FSMs for such typhoon-prone regions can be used by decision-
51 makers and planners in Vietnam to propose effective mitigation measures for community
52 resilience and development.

53 **Keywords**

54 Machine learning, random forest, LightGBM, CatBoost, flooding susceptibility mapping,
55 rainfall-runoff inundation model.

56 **1. Introduction**

58 Globally, floods are the most damaging natural disaster. Flash floods are more devastating than
59 any other flooding type because of their short lag times (Bui et al., 2019b; Vinet, 2008). The disastrous
60 impacts of flash floods have been documented in both developed and developing countries (Bisht et
61 al., 2018); however, flood events are more destructive in developing countries, such as Vietnam. The
62 observed increase in flash flood frequency is driven mainly by extreme changes in storm patterns and
63 global climate change (Hirabayashi et al., 2013; Pachauri et al., 2014). Causes of severe flooding in
64 Vietnam include tropical cyclones, typhoons, dense river networks, and extended coastal areas. The
65 area is also highly vulnerable to floods caused by extreme storms. Vietnam is ranked eighth among
66 the 10 countries with the highest number of weather events (Thao et al., 2020), where densely
67 populated areas are more vulnerable to floods. Consequently, continuous risk in terms of human life
68 and assets will always exist (Luu et al., 2021). Flash flood mitigation for risk reduction requires
69 efficient monitoring measures for sustainable support of flood hazard management (Arora et al.,
70 2020). Mapping flash flood susceptibility (FFS) areas is critical for scientists and governments
71 worldwide (Ali et al., 2020).

72

73 Several studies have been performed to predict the likelihood of flood events. These studies can be
74 clustered into three major groups: conventional analysis, rainfall-runoff, and pattern categorization
75 (Tien Bui and Hoang, 2017). Conventional analysis uses time-series data for an extended period
76 obtained from rainfall stations to produce regression models. The rainfall-runoff models (e.g., MIKE,
77 PCSWMM 2D, HEC-RAS, etc.) focus on determining the relationship between rainfall and runoff to
78 calculate temporal and spatial floods (Nguyen et al., 2015). This task is generally difficult because of
79 difficulties in accessing affected areas, especially in developing countries; consequently, the
80 performance of hydrological models can be affected, and detailed observational datasets are necessary
81 for calibration and validation (Abdrabo et al., 2020; Abushandi and Merkel, 2011). Both groups have
82 a major deficiency: the lack of data required often limits their applications and incurs a considerable
83 cost for collecting the data (Fenicia et al., 2014). On the other hand, the last group (pattern
84 classification) employs machine learning (ML) models that utilize historical geological,
85 environmental, and flood data. Accordingly, flood-prone areas are defined as “flood” and “non-flood”
86 classes (Bui et al., 2019b). However, comparative studies and integration between these groups are
87 lacking (Demirel et al., 2009; Hsu et al., 1995; Humphrey et al., 2016; Kratzert et al., 2019; Yang et
88 al., 2020).

89

90 Globally, the application of ML approaches for flood susceptibility prediction have been extensively
91 assessed over the past 20 years. Therefore, the recent development of ML methods has resulted in
92 substantial enhancements in flood modeling. Such practices have become widespread because of their
93 ability to capture information without applying predefined assumptions and process complex datasets
94 with high levels of accuracy in short periods (Arabameri et al., 2020; Costache et al., 2020b). Several
95 studies have been conducted to develop reliable flooding susceptibility maps (FSMs) using GIS and
96 remote sensing techniques. Currently, ML models are being adopted in combination with GIS to
97 address many hydrological and environmental challenges (Akay and Taş, 2020). ML models used to
98 predict flood susceptibility include artificial neural networks (ANNs), logistic regression (LR) (Arora
99 et al., 2020; Shahabi et al., 2020), adaptive neuro-fuzzy inference system (ANFIS) (Costache, Hong,
100 et al. 2020; Arora et al. 2020), genetic algorithms (GAs) (Darabi et al., 2019; Shirzadi et al., 2020),
101 support vector machines (SVMs) (Choubin et al., 2019; Dodangeh et al., 2020), and random forest
102 (RF) models. The RF model has been extensively used for flood risk assessments (Chen et al., 2020;
103 Esfandiari et al., 2020). Several ML models and ensemble methods have been used to predict FFS
104 (Shahabi et al., 2020). ML approaches involve multiple steps (Arora et al., 2020), including the
105 preparation of inventory and influencing factors, as well as the assessment of the accuracy of the ML
106 model. However, few studies have discussed the effect of inventory dataset size on the accuracy of
107 results (Catal and Diri, 2009; Meadows and Wilson, 2021; Tiwari and Chatterjee, 2010).

108

109 Recently, ensemble and hybrid ML models have appeared, and they outperform single models in their
110 prediction accuracy (Zenggang et al., 2021). Several ensemble ML techniques, such as the alternating
111 decision tree, bagging, dagging, reduced-error pruning tree, logistic model tree, J48 decision tree,
112 naïve Bayes tree, AdaBoost, and random subspace ensembles, have been applied to enhance the
113 predictive accuracy of the FSM (Luu et al., 2021; Pham et al., 2021a; Tuyen et al., 2021). In Vietnam,
114 several studies have developed FFS maps using ML, which can be categorized into three groups. The
115 first group evaluates the utilization of new ML models and their ability to detect areas prone to floods.
116 For instance, the AdaBoost, bagging, dagging, and random subspace ensemble learning techniques
117 were combined with the Partial Decision Tree (PART) classifier to develop new GIS-based ensemble
118 models for FSM in the Quang Binh Province (Luu et al., 2021). The second group attempts to
119 overcome the limitations in the number of studies that utilize remote sensing data to generate input
120 variables for FSM despite the merits of using such available data (Pham et al., 2019, pp. 2010–2018).
121 As such, (Dhara et al., 2020; Ngo et al., 2021; V.-N. Nguyen et al., 2020; Nhu et al., 2020, p. 202)

122 suggested a hybrid approach using remotely sensed data with ML models for the FFS. The third group
123 introduced a novel deep learning neural network (DLNN) algorithm for FSM (Tien Bui et al., 2020),
124 integrating extreme learning machines (ELMs) and particle swarm optimization (PSO) (Bui et al.,
125 2019b, 2020) along with a comparison between ML and deep learning techniques (Pham et al., 2021b)
126 for the same study area.

127

128 In this study, we examined two methods, categorical boosting (CatBoost) and the light gradient
129 boosting machine (LightGBM), for FSM for the first time in humid regions after successful
130 application in arid regions (Saber et al., 2021). Previously, both methods have been applied to
131 LightGBM and CatBoost. For instance, LightGBM has been used in previous studies owing to its
132 accurate predictions, short computational time, and outstanding ability to avoid overfitting issues.
133 Accordingly, our primary objectives are (1) to evaluate how practical the two ML approaches
134 (LightGBM and CatBoost) are for predicting FFS in humid environments (Vu Gia-Thu Bon basin in
135 Vietnam); (2) to compare the performance of these methods to the typical RF method; (3) to test the
136 effect of the inventory datasets (number of points) on the accuracy of the results in the study area;
137 and (4) to compare the rainfall-runoff inundation (RRI) 2D hydrological model with the proposed
138 ML integrated models.

139

140 The subsequent sections of this study are as follows. In Section 2, a description of the study area is
141 presented. In Section 3, an introduction to the datasets and methodology is developed. In Section 4,
142 the results are presented and discussed. Finally, our conclusions are given in Section 5.

143

144 **2. Study Area**

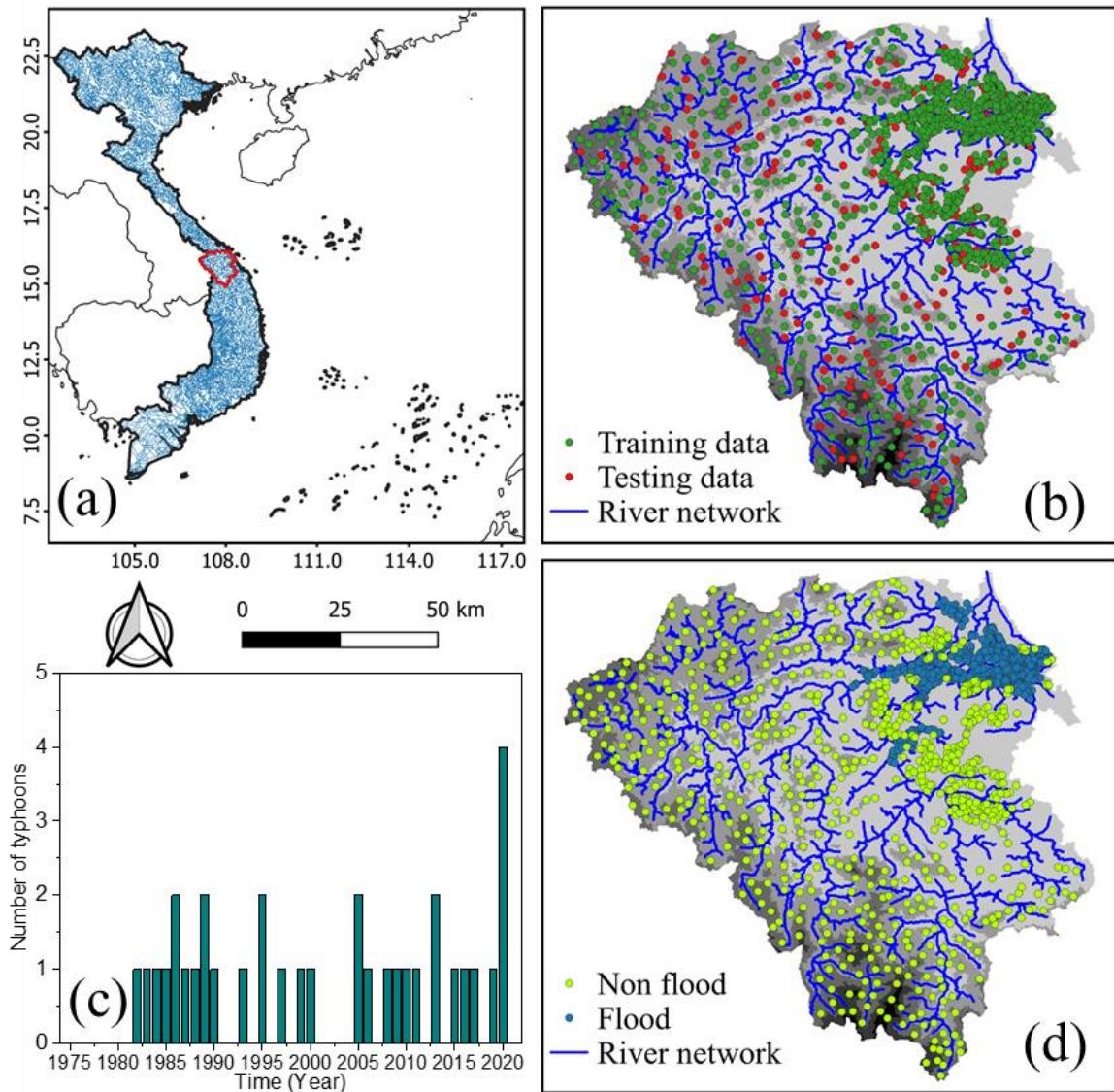
145 The Vu Gia-Thu Bon (VGTB) River Basin (Fig. 1) is one of the largest river basins in Vietnam,
146 with a surface area of 10350 km² (RETA, 2011). The land use types in the basin are forest (47%),
147 cropland (26%), and grassland (20%) (Avitabile et al., 2016). The basin has a tropical monsoon
148 climate with two seasons: dry summer (January–August) and rainy (September–December) seasons.
149 The topographic features of the basin are hilly mountainous areas with approximately 60% of the
150 basin having an elevation over 552 m. The average annual rainfall varies significantly, from 2000
151 mm in the downstream regions, to more than 4000 mm in mountainous areas. There are seasonal
152 differences, with 65%–80% of the annual rainfall concentrated between September and December

153 (RETA, 2011). The rainfall in eight months of the dry season is only 20%–35% of the annual rainfall
154 (Nauditt et al., 2017).

155 Owing to the difference in the spatial distribution of rainfall, the runoff in the VGTB basin varies
156 significantly across seasons. The flooding season lasts from October to December and is accompanied
157 by rains. River flow in this period accounts for approximately 62.5% to 69.2% of the total annual
158 flow. The impacts of both intense rain and steep terrain usually lead to flooding with a high intensity
159 and short occurrence times. Approximately 4–8 floods occur annually. The maximum flood peak
160 occurs in October and November owing to different weather patterns, such as tropical depressions,
161 cold air, and typhoons (Vu et al., 2011). The number of deaths and property damage caused by storms
162 and floods is increasing over the years, especially in 2020, based on the report of the Commanding
163 Committee for Disaster Prevention and Search and Rescue in Quang Nam Province (Fig. 2).

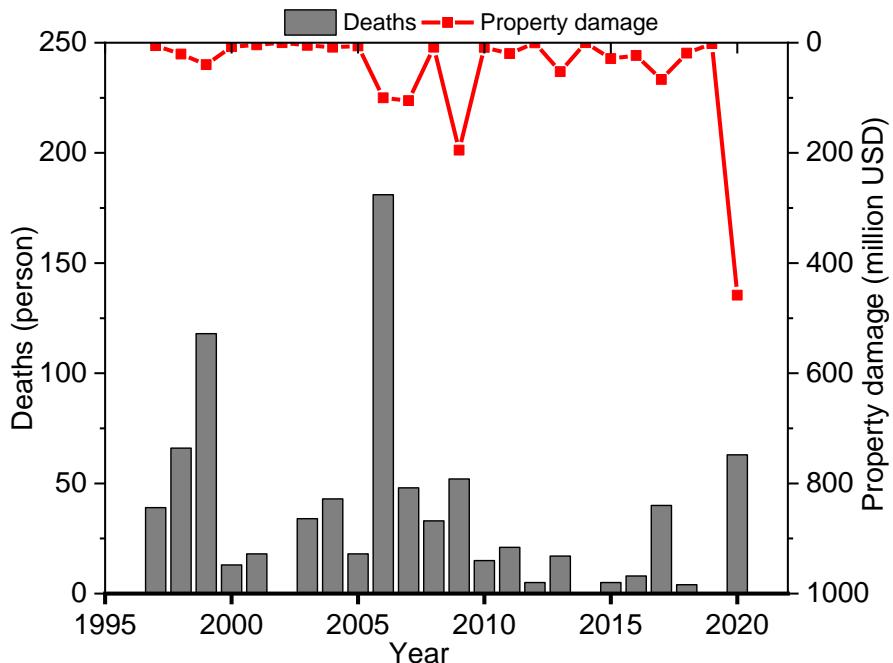
164 There are two main sub-basins in the VGTB river system: the Vu Gia Basin and Thu Bon Basin.
165 The Quang Hue River is connecting the Vu Gia and Thu Bon Rivers. The Vu Gia River originates
166 from the western slope in Kon Tum and flows towards Quang Nam province and Danang city. It
167 connects with the sea at the Cua Han estuary. The length of the main river from the source to the Cua
168 Han estuary was 204 km. The Thu Bon River originates from a mountain at an elevation of 1500 m
169 in Kon Tum province. The length of the river from the source to the Cua Dai estuary was 198 km.

170



171

172 Fig. 1. Location map of the VGTB river basin: a) The Vietnam Map, b) flood inventory map
 173 for training and validation datasets, (c) total annual precipitation of the entire basin, and (d)
 174 flooded and non-flooded locations,



175

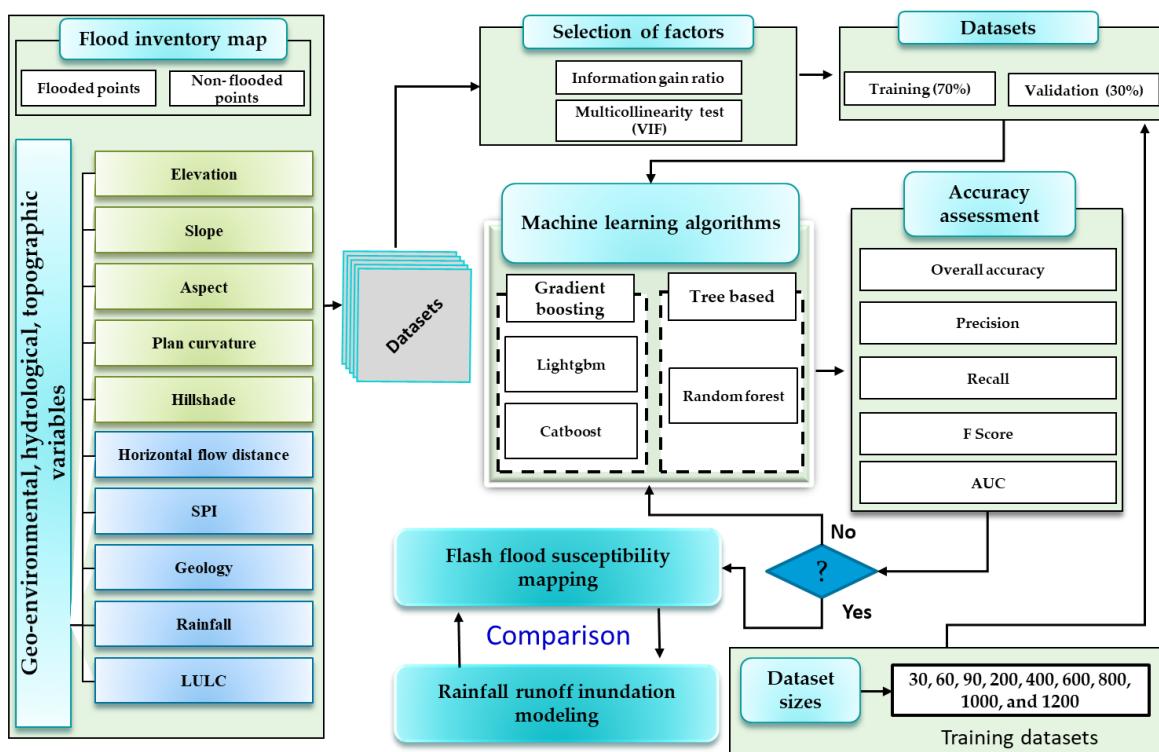
176 Fig. 2. Number of deaths and property damage caused by storms and floods from 1997 to
 177 2020 in VGTB River basin (source: Commanding Committee for Natural Disaster Prevention
 178 and Control, Search and Rescue in Quang Nam Province).

179

180 3. Methodology

181 The methodology of this study consists of several steps, as illustrated in the flowchart in Fig.
 182 3. There are two main parts to this methodology. First, a flood inventory map based on 850
 183 flooded points is developed. These points were identified based mainly on post-flood surveys
 184 after typhoons in 1999, 2006, 2007, 2009, 2013, and 2020. Non-flooded points (850) over
 185 the catchment were randomly selected using GIS tools. Additionally, a total of 10
 186 independent flood susceptibility factors (FSFs) were considered for modeling based on
 187 topographical, hydrological, geological, and landform characteristics, which have been used
 188 in several studies. The FFSFs, namely elevation, aspect, slope, hillshade, plan curvature,
 189 horizontal flow distance, stream power index (SPI), geology, rainfall, and land use/land cover
 190 were used to determine the linear relationship with other factors. In the later stages, the
 191 dataset was divided into two sets: (70%) for training and (30 %) for testing using the random

192 selection scheme. Spatial maps for each FFSF were produced using ArcGIS, considering the
 193 consistency of spatial resolution. Following that, two approaches, the information gain ratio
 194 (IGR) and variance inflation factor (VIF), were applied to examine the importance of the
 195 influencing factors in the FFS. Subsequently, the ML algorithms employed, namely, RF,
 196 LightGBM, and CatBoost were employed. The final results of the ML models were assessed
 197 for accuracy by using different statistical measures, including the most dominant, area under
 198 the curve (AUC). Moreover, as we have very high-quality observational flood locations, we
 199 tested the models to check the different sizes of the training datasets (Fig. 3). The final flood
 200 susceptibility maps developed by the ML models were then compared with the flood
 201 inundation maps from the physical hydrological model.
 202
 203



204
 205 Fig. 3. Methodology flowchart for flash flood susceptibility mapping.

206 **3.1. Datasets**

207 **3.1.1. Flash Flood inventory data**

208

209 The first step in FFS mapping involves identifying flood points or locations based on
210 historical records of previous floods. This is the most important data input for FFS mapping
211 (Tehrany et al., 2014). From the records of past occurrences, the locations of future hazard
212 events could be estimated (Devkota et al., 2013; Tehrany and Kumar, 2018). Therefore,
213 analysis of previous historical flood events and their influencing factors is the primary step
214 in flood susceptibility analysis (Masood and Takeuchi, 2012). The flood inventory map
215 shows the sites of the flooded areas in any flood-prone basin (Bellu et al., 2016). An inventory
216 map can be developed from several sources, such as field surveys, flood forecasting records,
217 and remote sensing data (Band et al., 2020; Esfandiari et al., 2020; Wang et al., 2019). The
218 accuracy in selecting flood points will enhance the model accuracy for FSM (Arora et al.,
219 2019; Tehrany et al., 2013). In this study, 1700 ground control points (Fig. 1) were identified
220 for flooded (850) and non-flooded points (850). Approximately 1250 were used for training,
221 and 450 for testing the models. The flooded locations were compiled from historical flood
222 records and post-flood field surveys in 1999, 2006, 2007, 2009, 2013, and 2020 (Fig. 1). The
223 non-flood points were randomly selected. The flood points and non-floods were assigned
224 with values 1 and 0, respectively, and the points were divided into 70% and 30% for training
225 to set the flash flood prediction model and validation to verify the model performance and
226 generalization ability by using the random selection method.

227 **3.1.2. Geospatial database (flood controlling parameters)**

228

229 The selection of flood controlling factors for the FFS mapping is important and affects the
230 accuracy of the models (Kia et al., 2012). During floods in a drainage system, runoff depends
231 on the characteristics of the watershed, topography, catchment area, and land use/land cover
232 types (Hölting and Coldewey, 2019). There are no standard and universal criteria for
233 selecting the controlling factors for FSM; therefore, according to the previous review and

study area characteristics, as well as the data availability, 10 flood triggering factors were prepared, including topographic, hydrological, geologic, and landform factors, namely, elevation, plan curvature, slope, aspect, horizontal flow distance, hillshade, SPI, rainfall, geology, and land use/land cover. Using ArcGIS, all data were resampled and prepared in spatial raster formats with a spatial resolution of 90 m (Fig. 4). All topographic factors were constructed based on the MERIT digital elevation model (Yamazaki et al., 2017). The spatial resolution of the terrain elevation was 3 s (~90 m at the equator). It was developed by eliminating the error components from existing digital elevation models (DEMs), such as Aw3D-30m v1 and SRTM3 v2.1. These data are freely available and accessible at http://hydro.iis.u-tokyo.ac.jp/~yamadai/MERIT_DEM/.

244

245 Below are the details of all the considered flood influencing parameters in this study.

246 **Elevation:** There is a direct relationship between altitude and flooding (Tehrany et al.,
247 which means that lowland surfaces are more vulnerable to floods than elevated areas
248 (Khosravi et al., 2016). This means that the higher the topographic elevation, the lower the
249 flood possibility (Tehrany et al., 2014; Youssef et al., 2016). The study area has complex
250 topographic features, with very high elevations up to 2600 m, and the downstream region of
251 the basin and coastal area has low elevations ranging from -3 m to 200 m; these regions are
252 mainly residential and agricultural areas (Fig. 4a).

253 **Slope:** This is a significant factor influencing flooding (Khosravi et al., 2016; Meraj et
254 al., 2018; Tien Bui et al., 2016), because of its effect on water velocity and surface flow
255 (Torabi Haghghi et al., 2018). The steeper slopes contribute to a high velocity and increased
256 volume of water in downstream areas (Chen et al., 2020). Slope influences hydrological
257 features that directly affect water runoff (Tehrany et al., 2019). In the study area, the slope
258 varied from 0° to 70° (Fig. 4b).

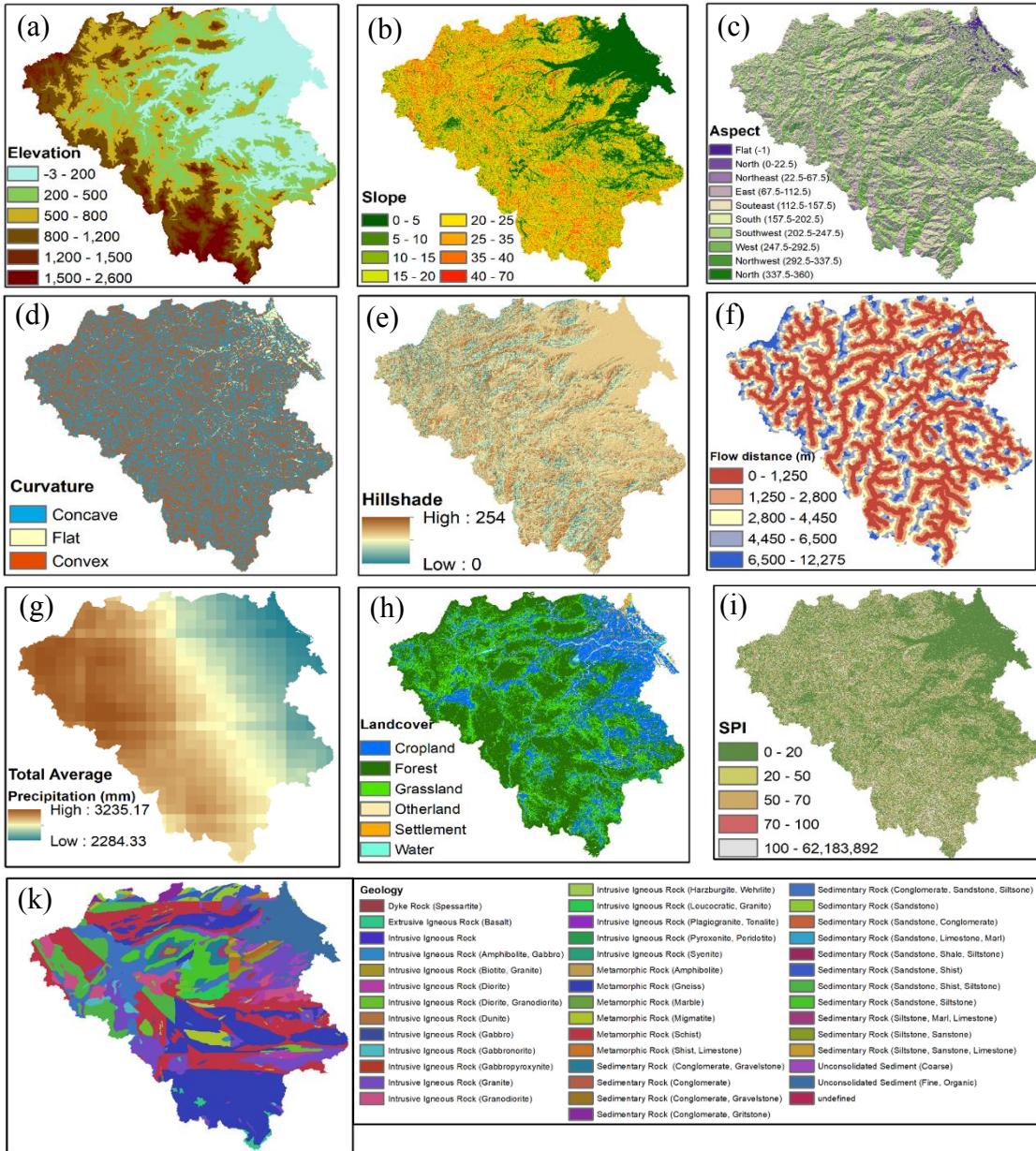
259 **Aspect:** Aspect is important for flooding (Choubin et al., 2019) as stated by Yates et
260 al. (2002); the hydrological parameters are influenced by this aspect. When aspects receive
261 low intensity of sunlight, which means more soil moisture, the moist slope will be more likely
262 to increase runoff, leading to increased flooding risk (Yariyan et al., 2020). There is an

263 indirect relationship between the aspects of floods owing to their control over several geo-
264 environmental factors, such as soils, rainfall, and vegetation (Rahmati et al., 2016). In this
265 study, the aspect map was categorized into 10 classes ranging from flat to northwest (Fig.
266 4c).

267 **Plan Curvature:** This is considered to be an important and essential flood influencing
268 factor by many researchers (Hong et al., 2018) and also affects hyporheics and heterogeneity
269 (Cardenas et al., 2004). The different values of curvatures differentiate between the areas of
270 the faster runoff from those with a slower runoff. While negative values cause an increase in
271 runoff, and the positive values decrease it. The runoff is affected by the slope shape, as the
272 flat form (zero curvature) and concave form (negative) have more potential for flooding than
273 the convex form (positive) (Shahabi et al., 2020; Tehrany et al., 2015, 2014). For instance,
274 concave slopes decelerate surface flow and increase filtration losses (Young and Mutchler,
275 1969), but convex slopes accelerate flow discharge and perhaps reverse filtration losses (Cao
276 et al., 2016). The curvature map was developed from the DEM with three main classes
277 (concave, flat, and convex), and the flat class was more dominant in the downstream area, as
278 shown in Fig. 4d.

279 **Hillshade:** Hillshade or toposhade is directly related to the length and shade of
280 hillslopes, which may affect the convergence of overland flow (Aryal et al., 2003). However,
281 the toposhade is limited in previous studies (Bui et al., 2019b), and it was found to be the
282 most important factor in predicting FFS mapping after slope and elevation (Bui et al., 2019a).
283 Therefore, toposhade was selected as a flood influencing factor, as shown in Fig. 4e.

284



285

286 Fig. 4. Flood influencing factors: a) elevation, b) slope, c) aspect, d) plan curvature, e)
287 hillshade, f) horizontal flow distance, g) rainfall, h) land use/land cover, i) SPI, j) geology.

288

289 **Flow distance:** The distance from the main rivers or streams has a significant impact on
290 flooding occurrence in any area (Glenn et al., 2012). The areas adjacent to streams are

291 typically more prone to flooding(Chapi et al., 2017). The risk of flooding is proportional to
292 the distance from rivers. It has been emphasized as a significant influencing factor for
293 flooding (Predick and Turner, 2008), and flooding is frequent in areas adjacent to rivers (Bui
294 et al., 2018; Darabi et al., 2019). The distance from streams signifies the distance from stream
295 channels, which are the main means of overland flow (Gigović et al., 2017; González-
296 Arqueros et al., 2018). In the present study, the horizontal flow distance was estimated (Fig.
297 4f.) using ArcGIS from the flow accumulation, flow direction, and DEM in ArcGIS.

298

299 **Rainfall:** Precipitation is one of the triggering factors for flooding as no rainfall indicates
300 a lack of flooding. The total average rainfall was estimated between 2001 and 2019 using the
301 PERSIANN Dynamic Infrared–Rain rate model (PDIR). Estimation of precipitation was
302 done using remotely sensed information that utilizes ANNs (P. Nguyen et al., 2020, p.). It is
303 a real-time global dataset with a high resolution of approximately $(0.04^\circ \times 0.04^\circ$, or 4 km \times
304 4 km, and it is freely available and accessible for download (<https://chrsdata.eng.uci.edu/>).
305 The spatial maps show that the average precipitation ranges from 2235 mm in the
306 downstream area to 3284 mm in the upstream region and mountainous area, as shown in Fig.
307 4g.

308

309 **Land use and land cover:** The influence of this factor was confirmed using the global
310 cover map developed by the geospatial information authority of Japan
311 (https://www.gsi.go.jp/kankyochoiri/gm_global_e.html), and mainly from this website
312 (<https://globalmaps.github.io/glcnmo.html>). Land use/land cover type was also considered as
313 a controlling factor because of its influence on filtration and runoff velocity. The study area
314 has approximately six classes (Fig. 4h), including cropland, forest, grassland, other land,
315 settlement, and water. The forest is the more dominant type of land cover in the mountainous
316 area, especially upstream of the basins, and agricultural land and urban areas are located in
317 the downstream region.

318 **Stream Power Index:** This parameter indicates the power of erosion and discharge
319 within a specific area of the river system (Poudyal et al., 2010). Several researchers have

320 considered the SPI as a flooding contributor because of its indication of surface flow. The
321 highest values of SPI imply a fast flow of downstream water, which reveals lower flooding
322 susceptibility, and low values imply slow flow leading to more inundation (Tehrany and
323 Kumar, 2018). The SPI was calculated based on a method derived from Jebur et al. (2014).
324 In the study area, SPI was classified into five classes, as shown in Fig. 4i.

325 **Geology:** This parameter is a crucial factor because of its effect on the infiltration and
326 flow velocity. The lithology was obtained from the Land Use and Climate Change Interaction
327 in Central Vietnam (LUCCI) (Nauditt and Ribbe, 2017). It has been classified into several
328 geological types (Fig. 4k) with high variation in sedimentary, igneous, and metamorphic rock
329 types.

330 **3.1.3. Selection of flood influencing factors**

331 Feature selection is a basic step in ML modeling for FSM. Removing redundant features may
332 decrease in the training speed caused by a larger number of features and the loss of ill-
333 conditioned landscape issues (Öztürk and Akdeniz, 2000). The latter occurs when highly
334 correlated features exist in the training dataset, causing a problem in determining the
335 learning-rate hyperparameter. Therefore, an incorrect choice of this hyperparameter value
336 may affect the estimation ability of the model. Therefore, the feature selection process was
337 based on three analyses to detect irrelevant factors, Spearman's rank correlation, the
338 multicollinearity test, and the IGR.

339

340 **3.1.4. Spearman's rank correlation coefficient**

341 Spearman's rank correlation coefficient is a nonparametric parameter used to evaluate
342 the strength of the monotonic link between two parameters, X and Y. The value of the
343 coefficient varies between -1 and 1, representing negative and positive degrees of association,
344 respectively. The closer the coefficient value is to zero, the weaker the relationship between
345 X and Y. A Pearson correlation coefficient value greater than 0.7 indicates a high level of
346 collinearity (Tien Bui et al., 2016). The correlation coefficient is estimated as follows:

347

$$r(x, y) = 1 - \frac{6\sum(x - y)^2}{n(n^2 - 1)} \quad (1)$$

348 where r is the correlation coefficient, x and y are the two variables, and n is the length of
 349 each variable.

350

351 **3.1.5. Multicollinearity Test**

352 In addition to the correlations between two features carried out using Spearman's
 353 coefficient, multicollinearity was checked in this study between all influencing factors.
 354 Multicollinearity analysis aims to detect the existing interrelatedness between variables and
 355 was performed (in this study) using the VIF. This factor is commonly used in flood
 356 susceptibility assessment studies (Bui et al. 2019; Khosravi et al. 2019; Rahman et al. 2019)
 357 which suggests a threshold > 5 to consider multicollinearity. However, in other studies, if the
 358 VIF value is greater than 10, the corresponding predictors are considered collinear; thus, it is
 359 recommended to exclude them from the models (Dou et al., 2019; Wang et al., 2019).
 360 Therefore, in this study, we considered a value of five as the threshold for selection. The
 361 independent predictors are defined as $X = \{X_1, X_2, \dots, X_n\}$ and R_j^2 , and refer to the coefficient
 362 of determination when the j th independent predictor X_j is regressed on the other predictors.
 363 VIF is calculated using the following equation:

364

$$VIF = \frac{1}{1-R_j^2} \quad (2)$$

365 **3.1.6. Information gain ratio**

366 Conditioning factors were evaluated to identify their relative importance in flood
 367 occurrence using the IGR test (Quinlan, 1986; Xu et al., 2013). The latter is among the feature
 368 selection methods that have been considered in many classification studies (Khosravi et al.,
 369 2019; Shahabi et al., 2020). The use of an input with an IGR equal to zero means that no
 370 relationship exists between this factor and the output. This situation indicates that the use of
 371 such an input in the model will not add any information to the applied model; in contrast, it

372 will generate noise that will decrease its predictive capability. Therefore, removing these
373 factors from the inputs is highly recommended. The IGR is calculated using Eq. (3):

374

$$375 \quad IGR(x, Z) = \frac{Entropy(Z) - \sum_{i=1}^n \frac{|Z_i|}{|Z|} Entropy(Z_i)}{\sum_{i=1}^n \frac{|Z_i|}{|Z|} \log \frac{|Z_i|}{|Z|}} \quad (3)$$

376 **3.2. Machine learning methods**

377 ML approaches are the basic concept of employing algorithms to analyze and learn
378 from the data to produce forecasting or classification systems. These techniques can be
379 learned from previous experience or a given historical database. These methods can
380 generalize the learning examples provided in the training phase to identify the main tasks that
381 must be performed.

382 Several ML algorithms have been developed. These techniques can be classified
383 according to their learning mechanisms (i.e., supervised, unsupervised, and semi-supervised
384 learning). The choice of a suitable ML model and training method depends on the problems
385 to be solved or the available data and its types. In the current study, we focused on using
386 supervised ML techniques for FFS assessment. According to previous research, various ML
387 techniques have been proposed recently to deal with FFS assessment (i.e., SVMs, ELMs,
388 ANNs, Gaussian process regression (GPR), and classification and regression trees (CART)).
389 In addition, few studies have addressed FFS using ensemble-learning approaches based on
390 decision trees. These algorithms are based on boosting techniques that concentrate on
391 misclassified data during the training phase. In this respect, the objective of this study is to
392 evaluate the performance of two new modeling techniques, CatBoost and LightGBM,
393 benchmarked with the conventional RF approach.

394 **3.2.1. Random forest**

395 RF models have proven efficiency when dealing with prediction and classification
396 problems (Esfandiari et al., 2020; Schoppa et al., 2020). RF is an ensemble learning approach
397 based on a decision tree model. It was developed by Breiman (2001), who combined bagging
398 (Breiman, 2001) and random subspace (Ho 1998) techniques. This ML algorithm has proven
399 to be reliable in many fields (Izquierdo-Verdiguier and Zurita-Milla, 2020; Pourghasemi et
400 al., 2020; Zahedi et al., 2018). In this study, we aimed to predict flood or non-flood regions
401 according to several conditioning factors; therefore, the RF model was used as a classifier
402 method.

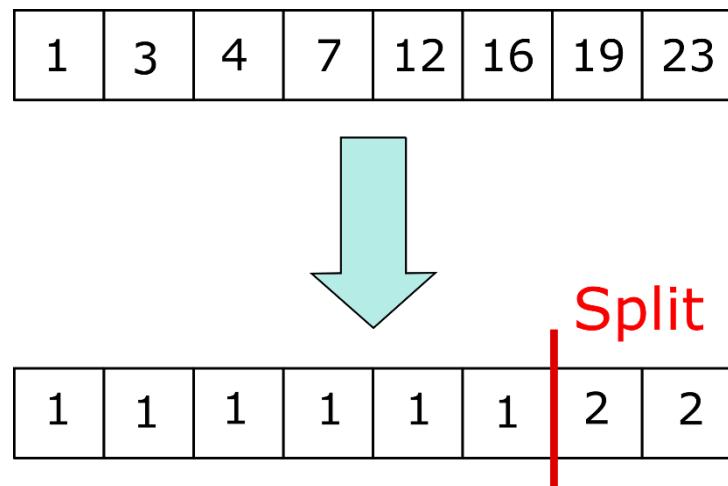
403 The weakness of decision trees is their sensitivity to training data, which may result in
404 very different tree structures. In the RF method, the original training set is used to randomly
405 generate several training sets, thereby allowing the creation of different trees (bagging
406 method). The inputs of the decision trees have the same data size as the initial training, and
407 because the data are randomly generated, the samples may be repeated two or more times. In
408 addition, each tree in the RF is trained with a subset of features that allows the development
409 of diversified trees that are not correlated. The final result (classification) was obtained by
410 performing a majority voting method on each decision tree results (Pal, 2005).

411 Decision tree models are simple to use and easy to interpret; however, their
412 performance is not always better than that of other classification methods (Malekipirbazari
413 and Aksakalli, 2015). On the other hand, RF outperformed other ML algorithms, such as
414 ANNs (Bachmair et al., 2017).

415 **3.2.2. Light Gradient Boosting Machine**

416 LightGBM is a variant of the gradient boosting decision tree (GBDT) that was
417 developed by Microsoft (Ke et al., 2017). It uses a combination of weak learners to generate
418 a robust model. The new variant includes algorithms such as histograms, leafwise tree growth,
419 gradient-based one-side sampling (GOSS), and exclusive feature bundling (EFB).

420 In GBDT models, the commonly used algorithm for split operations is the presorted
421 algorithm. To determine the optimal split, all possible split points are tested based on the
422 information gain, which is a time-consuming operation. A new histogram algorithm was
423 adopted in the LightGBM method. To reduce the time and complexity of the operation, the
424 data are grouped into a histogram, and the split point is chosen based on it (Fig. 5).

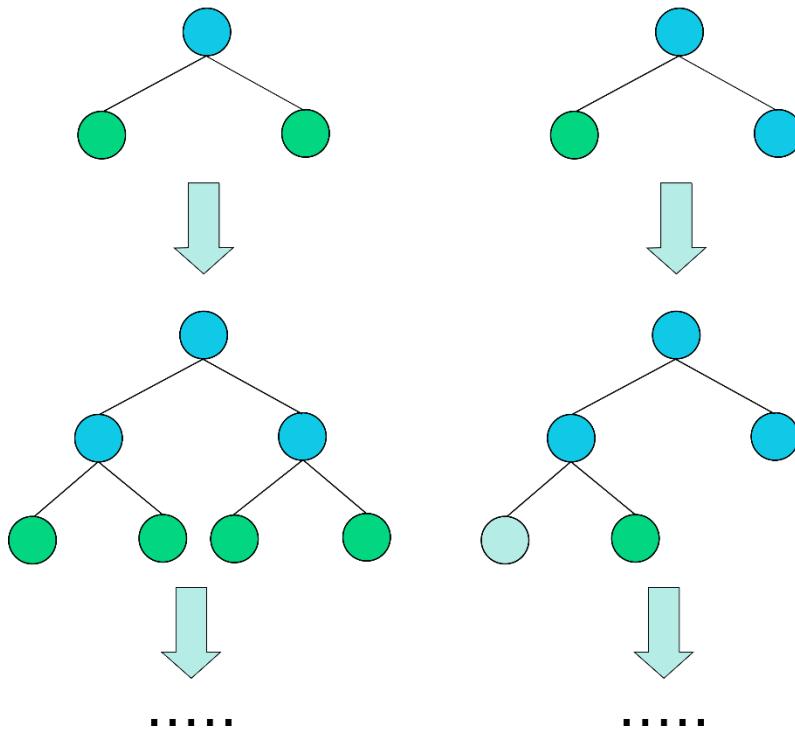


425

426 Fig. 5. Split operation example based on histogram algorithm.

427 In LightGBM, the decision tree growth strategy was changed by replacing the level-wise
428 approach with the leafwise tree growth approach. When finding the best node to split, the
429 former approach of the GBDT splits one level down, forming symmetric trees (Fig. 6). In
430 LightGBM, only the leaves that reduced the maximum error were split (Fig. 6). Ge et al.
431 (2019) recommended defining a maximum leaf-wise depth to avoid deep growth of trees and
432 provoke overfitting of the model.

Level-wise tree growth Leafwise tree growth



433

.....

.....

434

Fig. 6. Level-wise and leafwise tree growth strategies.

435 The LightGBM model also uses two algorithms (GOSS and EFB), making it faster than
 436 GBDT models while maintaining a high performance (Saber et al., 2021).

437

438 3.2.3. Categorical boosting

439 The CatBoost model is another enhanced boosting decision-tree learning technique, which
 440 was proposed by (Dorogush et al., 2018). It employs a gradient boosting scheme to construct
 441 a regression model through adjusted estimation. Furthermore, various refinements were
 442 performed to minimize overfitting of the model. The gradient boosting model is a useful ML
 443 tool that has yielded accurate results in many disciplines, including environmental parameter
 444 estimation, geospatial ecosystem factor dispersion, and meteorological forecasting. The
 445 CatBoost model operates well in terms of categorical attributes. Typically, the absence of

446 categorical characteristics increases the accuracy of the model. It is primarily dependent on
447 the use of gradient boosting, which employs a binary-tree classification scheme. The
448 following points outline the differences between CatBoost and the other boosting techniques.

- A sophisticated method was incorporated to convert category characteristics into numerical information. As mentioned by (Prokhorenkova et al. 2017), target statistics are very effective for dealing with categorical attributes with minimal information errors.
 - CatBoost combines categorical variables to take advantage of the existing relationship between different parameters.
 - To reduce the overfitting problem and improve the classification performance, a symmetrical tree strategy is used.

457

458 Let us suppose we have a dataset:

$$459 \quad D = \{(X_J, Y_J)\} \quad J = 1, \dots, m \quad (4)$$

460 where $X_j = (x_j^1, x_j^2, \dots, x_j^n)$ is a combination of attributes, and $Y_j \in R$, denotes the desired
 461 target. Input-output data are dispersed independently and identically, depending on an
 462 unknown function $\rho(\cdot, \cdot)$. The main objective of the learning scheme is to train a function
 463 $H: R^n \rightarrow R$ that can decrease information loss, that is, $L(H) := EL(y, H(X))$, where L is the
 464 smoothness error function and (X, y) denotes the testing samples from D. The gradient
 465 boosting approach builds a greedy series of approximations $H_t: R \rightarrow R$,
 466 $t = 0, 1, 2, \dots$, $H_t = H((t-1)) + g_t$ is the final function produced from prior approximation using
 467 an additive process $H_t = H((t-1)) + g_t$.

468

$$469 \quad g^t = \arg \min_{g \in G} L(H^{t-1} + g) = \arg \min_{g \in G} E L(y, H^{t-1}(X)) \quad (5)$$

470 In general, greedy techniques, such as Newton's method, employing a second-order approach
471 of $L(H(t-1) + g)$ at $H(t-1)$ or adopting (negative) gradient stages, are used to address the
472 optimization issue.

473 **3.3. Rainfall-runoff inundation model**

474 The RRI model was developed by the International Center for Water Hazard and Risk
475 Management in Japan. It is a 2D model distributed hydrological capable of simulating the
476 rainfall-runoff and flood inundation simultaneously (Sayama et al., 2012). The RRI model
477 has been successfully applied in many studies worldwide (Abdel-Fattah et al., 2018; Perera
478 et al., 2017; Saber et al., 2020; Tam et al., 2019; Try et al., 2020). In this study, the model
479 was calibrated and validated based on the typhoon of 2020, showing acceptable results with
480 the actual flood discharge and good agreement with flood inundation maps. The final flood
481 inundation map produced by the model was used for comparison with the ML FSMs.

482

483 **3.4. Evaluation of the model's performance validation**

484 The receiver operating characteristic (ROC) curve measure is a commonly used and
485 validated strategy for assessing the reliability of a model in geospatial research (Chen et al.,
486 2020; Tehrany et al., 2013). The ROC curve is the most commonly used mechanism for
487 analyzing flood susceptibility and landslide approaches. The classification performance of a
488 given technique was evaluated using the AUC in several previous studies (Bui et al., 2012;
489 Youssef et al., 2016; Youssef and Hegab, 2019). A high classification efficiency for a given
490 classification model should have an AUC-ROC value of 0.5 to 1, and the model's
491 performance is enhanced by boosting the AUC-ROC scores. When the AUC-ROC value was
492 close to 1.0, the models offered the best rate of precision and consistency. This demonstrates
493 the ability of the model to forecast the occurrence of disasters without bias (Bui et al. 2012).
494 In this study, the ROC score was determined using the following formula (Chang et al.,
495 2018):

496 Other quantitative metrics (accuracy, recall, precision, and F1-score) were employed to
497 check the model performance and compare its classification ability with that of its counterpart
498 models in the literature. Accuracy is defined as the ratio of correctly classified data to total
499 observations [Eq. (6)]; precision can be defined by the ratio of properly positive classified
500 data to total positive data [Eq. (7)]. Recall, which is also known as sensitivity, is defined by
501 the ratio of positive observations to the total observations [Eq. (8)]. F1-score uses weighted
502 averaging for both the precision and recall [Eq. (9)].

503

504
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

505
$$Precision = \frac{TP}{TP + FP} \quad (7)$$

506

507
$$Recall = \frac{TP}{TP + FN} \quad (8)$$

508

509
$$F1\ score = \frac{2(Recall * Precision)}{Recall + Precision} \quad (9)$$

510 where true positive (TP) represents a properly categorized flooded pixel, true negative (TN)
511 represents a correctly categorized non-flood pixel, false positive (FP) indicates the number
512 of pixels categorized incorrectly as flood pixels, and false negative (FN) indicates the number
513 of pixels categorized incorrectly as non-flood pixels.

514

515 **4. Results and Discussion**

516

517 **4.1. Multicollinearity assessment and feature selection**

518 According to (Chen et al., 2020), a value greater than 0.7 indicates a strong correlation
519 between variables. This value was adopted in this study to detect the existence of a correlation

520 between the flood-influencing factors. Ten of the conditioning factors (DEM, NDVI, flow
521 accumulation, vertical distance from the river, and slope) were identified as correlated with
522 each other (Table 1). The VIF of the vertical distance from the river (= 12), DEM (= 10.5),
523 SPI (= 7.7), and flow accumulation (= 7.4) factors were greater than the threshold value (>
524 5), which indicates a problem of multicollinearity (Fig. 7a).

525 To formulate an opinion on the importance of influencing factors in relation to flood
526 generation, the IGR scores were computed and are illustrated in Fig. 7b. According to the
527 results, most of the factors had an IGR greater than 0.05, and only four of them had an inferior
528 IGR, that is, flow accumulation, flow direction, rainfall, and aspect.

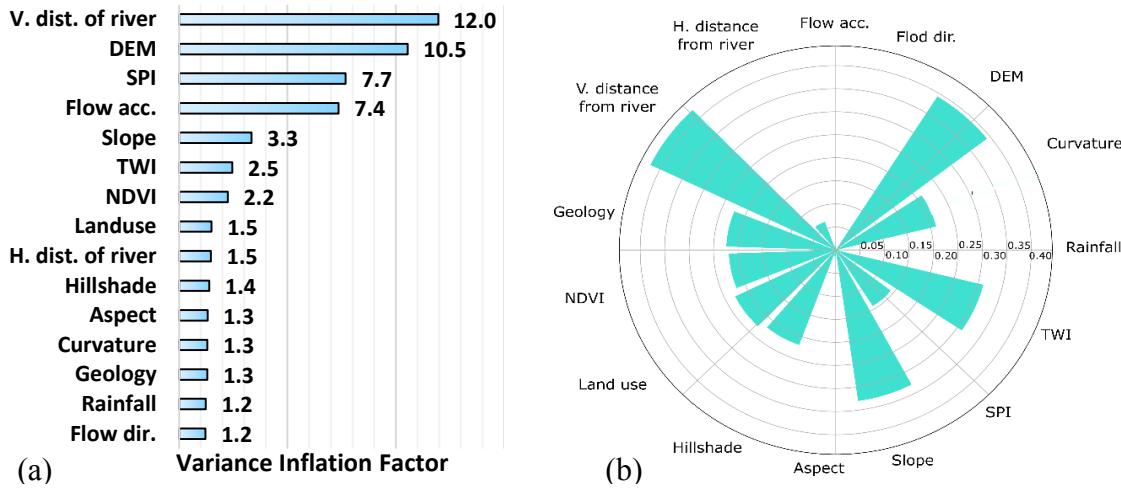
529 The selection of conditioning factors was performed as follows:

- 530 1) Based on multicollinearity analysis, the vertical distance from the river, DEM, SPI, and
531 flow accumulation factors were removed from the selection list.
- 532 2) Using the IGR as a selection criterion, flow direction, rainfall, and aspect were also
533 removed because their IGR was almost equal to zero.
- 534 3) After removing the aforementioned factors, only the slope and the topographic wetness
535 index (TWI) remained as correlated variables. By comparing the IGR (Fig. 7), we find
536 that the slope factor is more important than the TWI in relation to flood generation.
537 Therefore, the slope factor was selected for flood prediction based on the normalized
538 difference vegetation index (NDVI), land use, curvature, geology, hillshade, and
539 horizontal distance from the river.

540 Table 1. Spearman's correlation coefficients for FFS mapping.

	Aspect	Rainfall	Curvature	DEM	NDVI	Flow acc.	Flow dir.	H. Dist.	V. Dist.	Geology	Hillshade	Land use	Slope	SPI	TWI
Aspect	1.00														
Rainfall	-0.03	1.00													
Curvature	0.08	0.01	1.00												
DEM	0.03	0.20	0.27	1.00											
NDVI	0.04	-0.03	0.20	0.72	1.00										
Flow acc.	-0.07	0.01	-0.48	0.05	0.06	1.00									
Flow dir.	0.06	-0.07	0.11	0.06	0.06	-0.11	1.00								
H. Dist.	-0.04	0.03	0.02	0.49	0.39	0.08	0.02	1.00							
V. Dist.	0.04	0.17	0.28	0.97	0.72	0.00	0.06	0.52	1.00						
Geology	0.00	0.29	-0.02	0.08	-0.02	-0.01	-0.03	-0.09	0.06	1.00					
Hillshade	0.38	-0.05	0.08	0.03	0.02	-0.05	0.51	-0.06	-0.03	-0.06	1.00				
Land use	-0.03	-0.10	-0.11	0.59	-0.53	-0.13	-0.05	-0.25	-0.58	-0.07	0.00	1.00			
Slope	0.10	0.07	0.28	0.85	0.69	0.03	0.08	0.38	0.86	0.06	0.01	-0.53	1.00		
SPI	0.04	0.03	-0.41	0.32	0.27	0.86	-0.06	0.21	0.30	-0.01	-0.05	-0.27	0.38	1.00	
TWI	-0.10	-0.07	-0.49	0.76	-0.61	0.35	-0.12	-0.31	-0.79	-0.07	-0.04	0.45	-0.90	0.00	1.00

541



542

543 Fig. 7. Analysis of influencing factors: (a) VIF and (b) IGR for flood susceptibility.

544

545

546 4.2. Evaluation of the models

547 This section provides a detailed evaluation and comparison of all models developed in this
 548 study with respect to various classification criteria. The learning phase was performed using
 549 the K-fold cross-validation strategy. The examined data were partitioned into a learning set
 550 (60%), and the rest were deployed to evaluate the precision. The learning set was separated
 551 into two sets of training data (80%), which were used to adjust model weights and minimize
 552 classification errors, and validation data, which were used for hyperparameter tuning.
 553 Through the grid search approach, the appropriate hyperparameters for each classification
 554 technique were chosen. During the procedure, a wide variety of hyperparameter values was
 555 examined. Table 2 lists the optimal designs for each classifier.

556

557 The accuracy rates of all the studied models are listed in Table 3. As can be seen, all
 558 developed classification techniques achieved approximately identical results in terms of
 559 statistical metrics. The LightGBM model slightly outperformed the others in terms of speed
 560 convergence and classification metrics. The ROC curve of the generated models on the test

561 ensembles is displayed in Fig. 8, which reveals that the three suggested boosting strategies
562 have similar qualities and provide significant accuracy. The maximum AUC was reached by
563 LightGBM and RF models with the same score (99.5%), and CatBoost was ranked as the
564 worst model with an AUC of 97.9%.

565 Furthermore, CatBoost scored the first rank in terms accuracy performance accuracy equal
566 to 97.8 %, precision equal to 96%, accompanied by LightGBM classifier with an accuracy
567 of 97.3% and precision of 95%. Finally, the RF model was ranked as the last classifier model
568 with an accuracy equal to 95.5%, and precision of 96.2%. In comparison with previous
569 studies, RF in this study outperformed many of the previous applications, including (e.g.,
570 AUC = 0.925, Chen et al. (2020); AUC = 0.886, Tang et al. (2020); AUC = 0.7878, Lee
571 et al. (2017); AUC = 0.972, Achour and Pourghasemi (2020)).

572 The confusion matrix in Fig. 9. shows the performance of the used models in the study area,
573 where CatBoost shows better prediction followed by LightGBM and, finally, the RF
574 methods; however, all of them display acceptable prediction.

575 Two novel boosting classification models were examined in this study for FFS assessment
576 in the VGTB River Basin. From the evaluation statistics, we can conclude that the LightGBM
577 and CatBoost models proved their performance for FFS and can be used as essential tools for
578 real-time application compared to their counterpart models because of their high performance
579 and speed convergence.

580

581 Table 2. Parameter values of random forest, CatBoost, and LightGBM models.

METHOD	HYPERPARAMETER	GRID SEARCH VALUES	SELECTED VALUE
RF	Max depth	[2:2:40]	22
	Number of trees	[100:100:2000]	800

LIGHTGBM	Learning rate	[0.01:0.01:0.1]	0.09
	Max depth	[2:2:40]	32
	Number of leaves	[20:20:300]	220
	Min data	[5:5:100]	25
CATBoost	Learning rate	[0.01:0.01:0.1]	0.07
	Max depth	[2:2:40]	6
	Leaf estimation iterations	[1:1:10]	8
	12 leaf reg	Log [10(-21),10(-8)]	Log (10(-18))

582

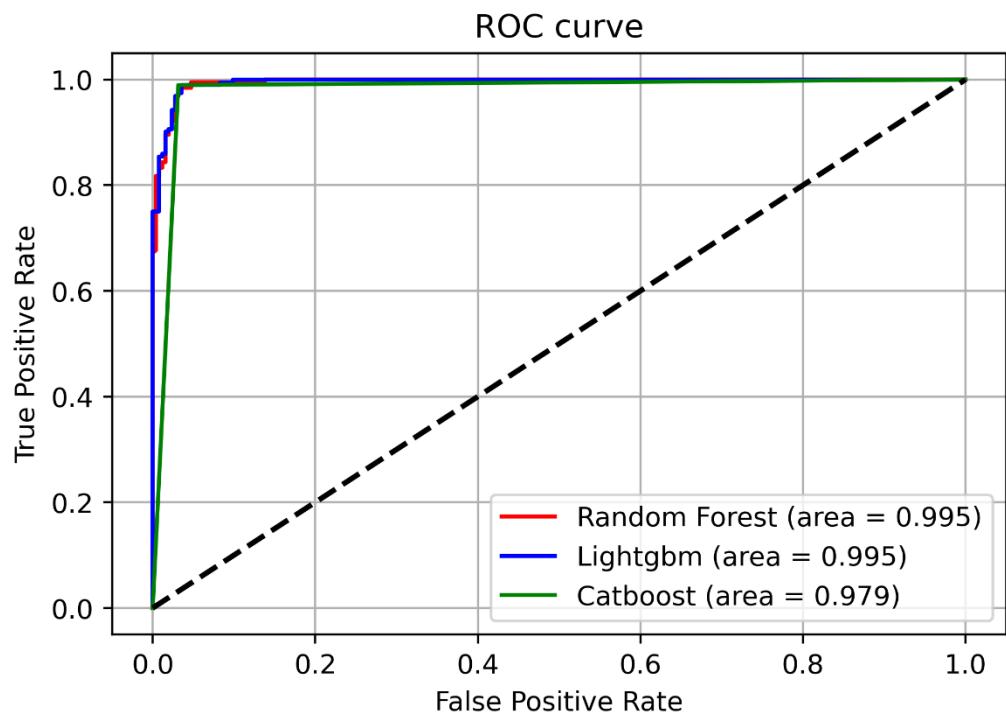
583

584 Table 3. Statistical measures used for the model performance evaluation.

CRITERIA	RANDOM FOREST		LIGHTGBM		CATBOOST	
	Train	Test	Train	Test	Train	Test
ACCURACY	0.999	0.955	0.999	0.973	0.999	0.978
PRECISION	0.998	0.962	0.998	0.950	0.998	0.960
RECALL	1.000	0.932	1.000	0.990	1.000	0.990
F1_SCORE	0.999	0.947	0.999	0.969	0.999	0.974
AUC	1.000	0.995	1.000	0.995	0.999	0.979

585

586 In this study, two new boosting classification techniques were investigated for flood
587 susceptibility projection in the VGTB. This is the first work that investigates the use of
588 CatBoost and LightGBM for flood classification in humid environments against the
589 commonly used RF models. The obtained results revealed that LightGBM outperformed its
590 counterpart ML models, especially in terms of classification metrics and processing time.
591 This agrees with the findings of Saber et al. (2021) that LightGBM has proven its efficiency
592 in flash flood prediction and outperforms the other two methods in terms of classification
593 and processing time. In addition, it was stated that LightGBM outperformed other methods
594 such as the RF, M5Tree, and other empirical models for estimating daily evapotranspiration
595 in China as a humid subtropical region (Fan et al. 2019). Similarly, it was also found that
596 LightGBM performed better than the others in terms of AUC (99.5%). The accuracy of
597 CatBoost (97.9%) was also high compared to the previous studies in other fields. Among
598 other methods, CatBoost, SVM, and RF have been applied to evapotranspiration modeling
599 in China (Huang et al. 2019). They stated that CatBoost presented higher accuracy and lower
600 computational cost than the other approaches (RF and SVM).



601

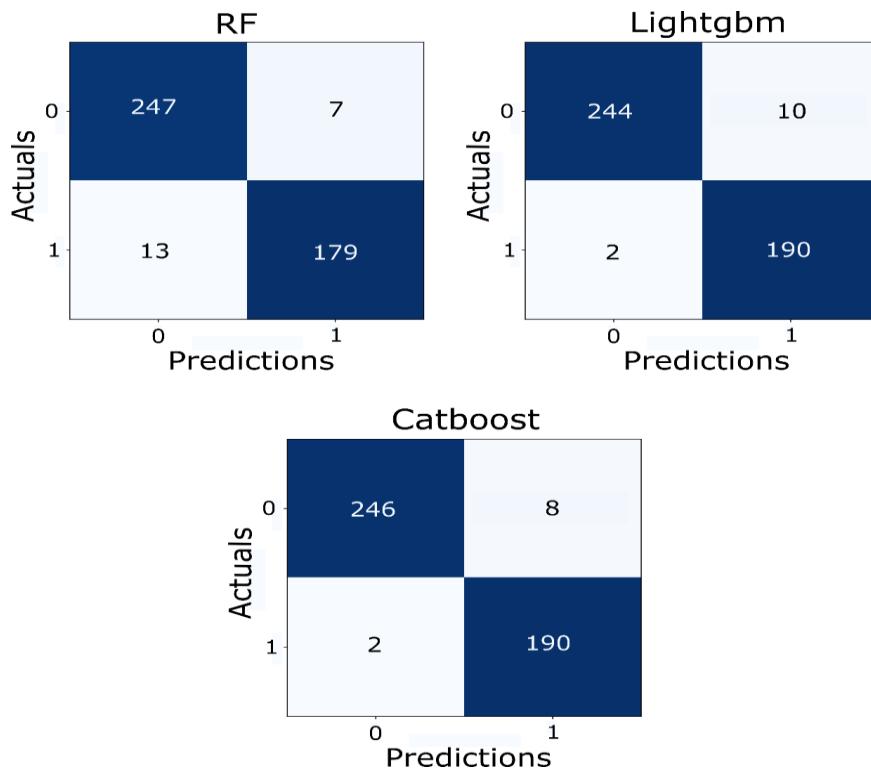
602

603 Fig. 8. Performance of random forest, CatBoost, and LightGBM models based on the area
604 under the receiver operating characteristic (ROC) curves.

605

606

607



608

609 Fig. 9. Confusion matrix showing the performance of the used models in VGTB River Basin.

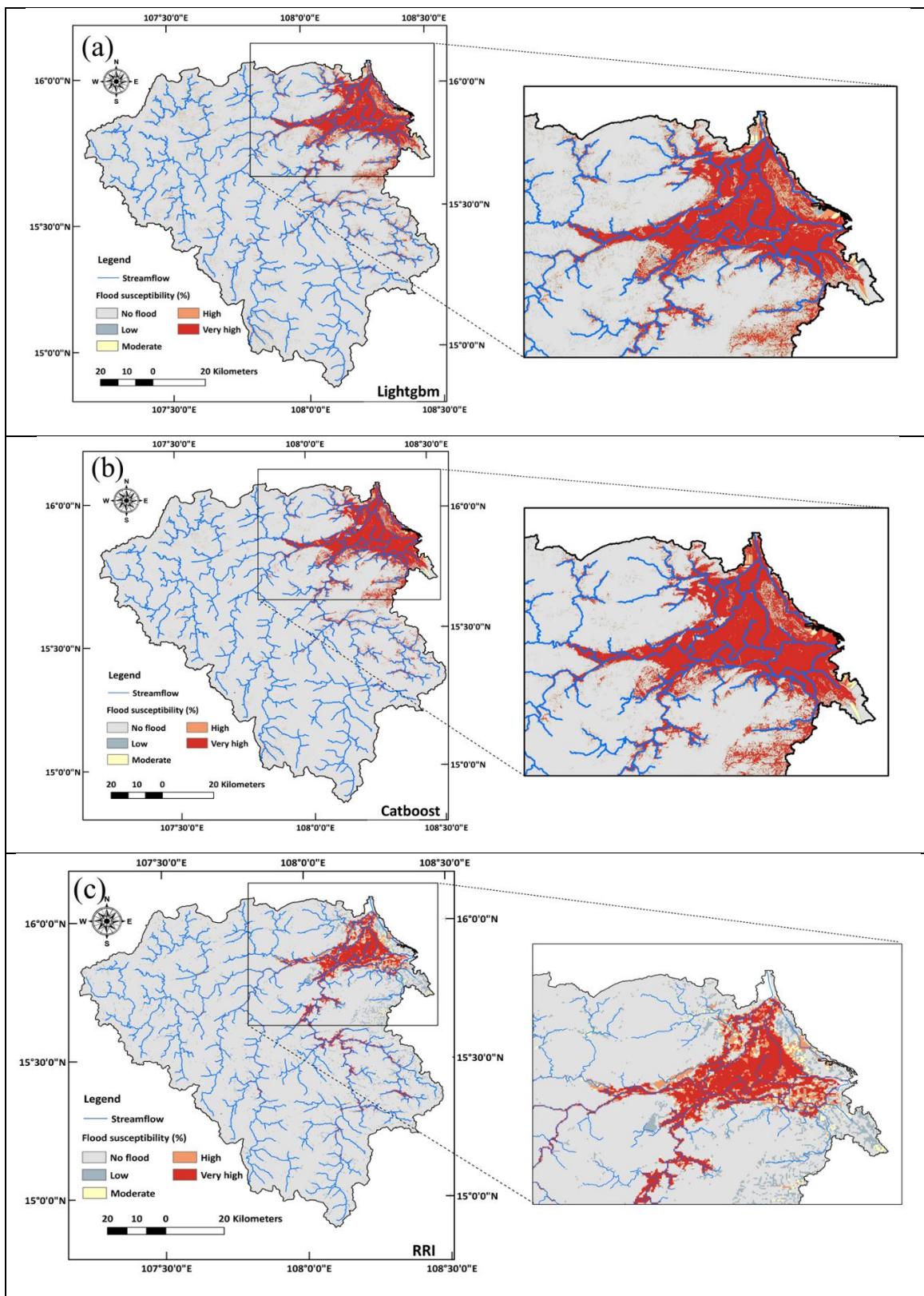
610

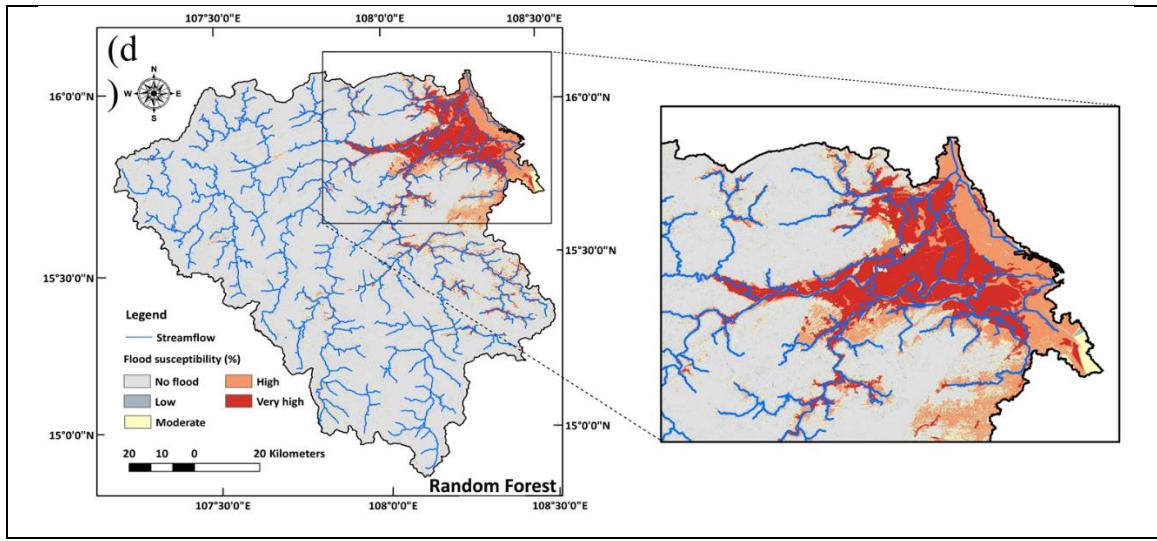
611 **4.3. Flash Flood Susceptibility Modeling**

612 The evaluation metrics of the newly applied boosting methods (CatBoost and LightGBM)
613 along with RF demonstrated their high performance in predicting flooding in a humid climate
614 environment. Accordingly, these methods were employed to estimate the flood susceptibility
615 maps for the entire VGTB river basin. The three FSMs developed using the three methods

616 (RF, LightGBM, and CatBoost) were compared with the FIM of the RRI model, as shown in
617 Fig. 10. The flooding susceptibility values were then mapped under five levels of
618 susceptibility classes: no flood, low, moderate, high, and very high. The regions affected by
619 these different susceptibility levels varied slightly depending on the model.

620 The FSMs by the three models showed that the areas of high and very high susceptibility to
621 flooding to be 13% (RF), 11% (LightGBM), and 10% (CatBoost) of the study area, which
622 agrees with the FIM developed by RRIat approximately 11%. This level of susceptibility is
623 predominant in the coastal and plane areas along the Vu Gia and Thu Bon Rivers (Fig. 10
624 and 11). The spatial distributions of the high and very high levels were similar in all the maps
625 produced by the ML and RRI models. The areas affected by a moderate level of susceptibility
626 to flooding (Fig. 11) were estimated at 10% (RF), 0% (LightGBM), and 1% (CatBoost),
627 indicating that both LightGBM and CatBoost are more similar to the RRI model which shows
628 a value of approximately 1%. The areas affected by the low level of susceptibility to flooding
629 (Fig. 11) were estimated at 36% (RF), 1% (LightGBM), and 2% (CatBoost), also revealed
630 that both LightGBM and CatBoost are performed better, with good agreement with the RRI
631 model showing the value of approximately 3%. It was also found that the areas that were not
632 subjected to the flash floods were approximately 42% (RF), 83% (LightGBM), and 87%
633 (CatBoost) of the total study area (Fig. 11), showing good agreement with RRI model that
634 shows approximately 90%. However, the performances of the employed models are almost
635 the same and the two new methods of LightGBM and CatBoost outperform RF in terms of
636 the spatial coverage of the flood susceptibility levels in comparison with the RRI model. The
637 RF overestimated the low flood susceptibility in the study area. All methods show an
638 agreement in terms of the same spatial pattern of FSM, highlighting that the coastal areas are
639 the areas prone to flooding where most of the residential and agricultural regions are located.



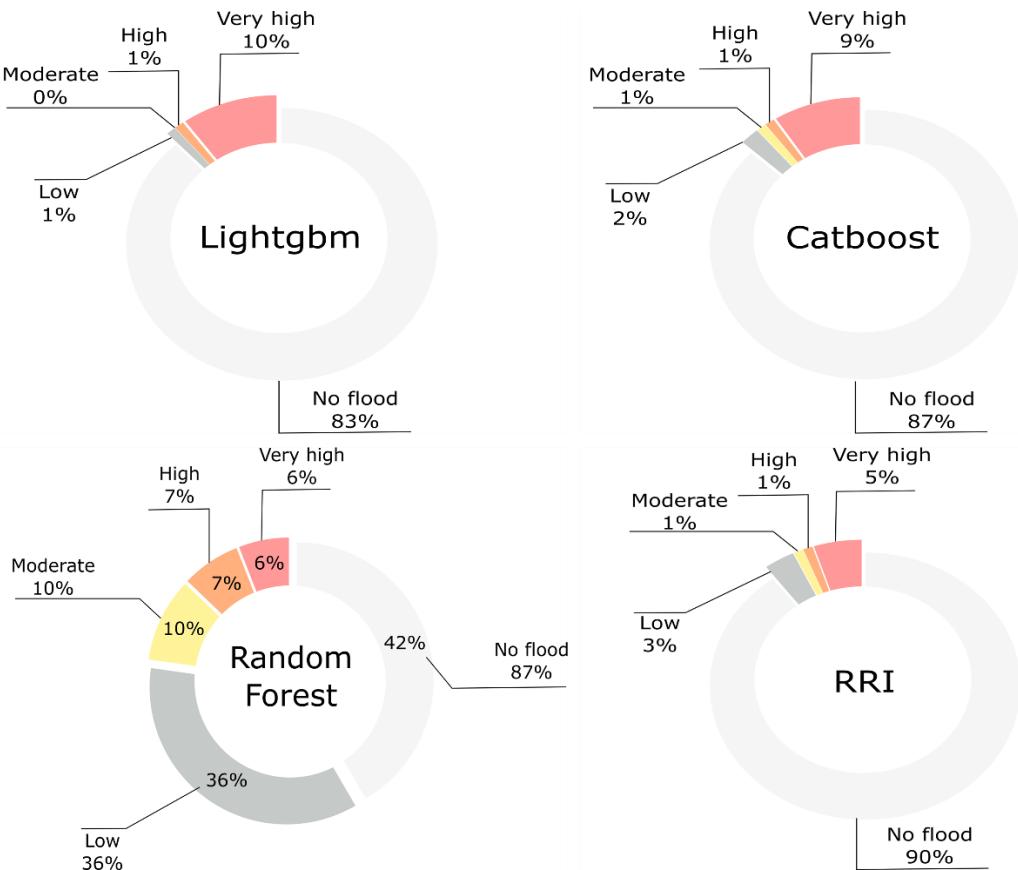


640

641 Fig. 10. Flood susceptibility maps by LightGBM (a), CatBoost (b), RF (c), and RRI (d),
 642 respectively, from the top to the bottom.

643

644



645
646 Fig. 11. Affected area of the flood susceptibility levels for the three applied ML methods
647 and flood inundation map of RRI model.

648

649 **4.4. Testing different sizes of the datasets**

650

651 In this section, we tested different sizes of the datasets, including flooded and non-flooded
652 points (1250, 1000, 800, 600, 400, 90, 60, and 30) of the training model (Fig. 12). The
653 training datasets classified as 50% and 50% for flooded and non-flooded points, respectively;
654 however, the testing datasets were the same during the simulation (Fig. 12). We found that
655 accuracy scores for all the models and all the tested cases were greater than 90% (Fig. 13),
656 except for dataset sizes of 60 and 30 points in LightGBM. The accuracy score slightly
657 decreases with the decrease in the datasets in both the LightGBM and CatBoost models, but

is not consistent in the RF model. This implies that the ML approaches employed in this study can effectively work with very limited training datasets with a slight decrease in accuracy, which will be applicable for ungauged or regions with very low monitoring and observations of flooding occurrences and impacts. The FSMs developed based on different training datasets show that most of the spatial maps are acceptable as overall spatial coverage; however, there are some small spatial differences in the susceptibility flooding levels (Figs. 14 and 15). For instance, the affected areas (Figs. 14 and 15) in the very high flood susceptibility flood category are almost the same about 6% for all datasets (1250, 1000, 800, 600, 400, 90, 60, and 30), except for dataset of 200, which was approximately 5%. On the other hand, the affected areas by the high flood susceptibility level are also varying but the highest percentage was 9% for the 200 and 60 datasets and the lowest was 6% for 30, 90, 600, and 800 datasets. The variation in moderate flood occurrence ranged from 17% to 9%. The dataset size of 30 was the highest among the others, about 17%. The range of the low flood susceptibility category was highly variable from 20% to 41%, the lowest was for the dataset of 30, and the highest was for the dataset of 800. The reasons for such variation are probably the random selection of the flooded samples, which in some cases are not representative of all the influencing factors. In general, we noticed that the spatial coverage was not extremely different, but some differences based on the categories were observed. The areas with no flood levels are also changeable by about 42%, 42%, 38%, 45%, 44%, 45%, 42%, 44%, and 51% for the datasets of 1250, 1000, 800, 600, 400, 200, 90, 60, and 30, respectively. Interestingly, the highest percentage was recorded by the dataset of 30 points and the lowest was recorded by dataset 800. The results of the analysis of different data sizes for ML training show that ML can effectively predict the flood susceptibility maps in the study area regardless of the number of samples, with the condition of the used data being observational flooded sites.

683

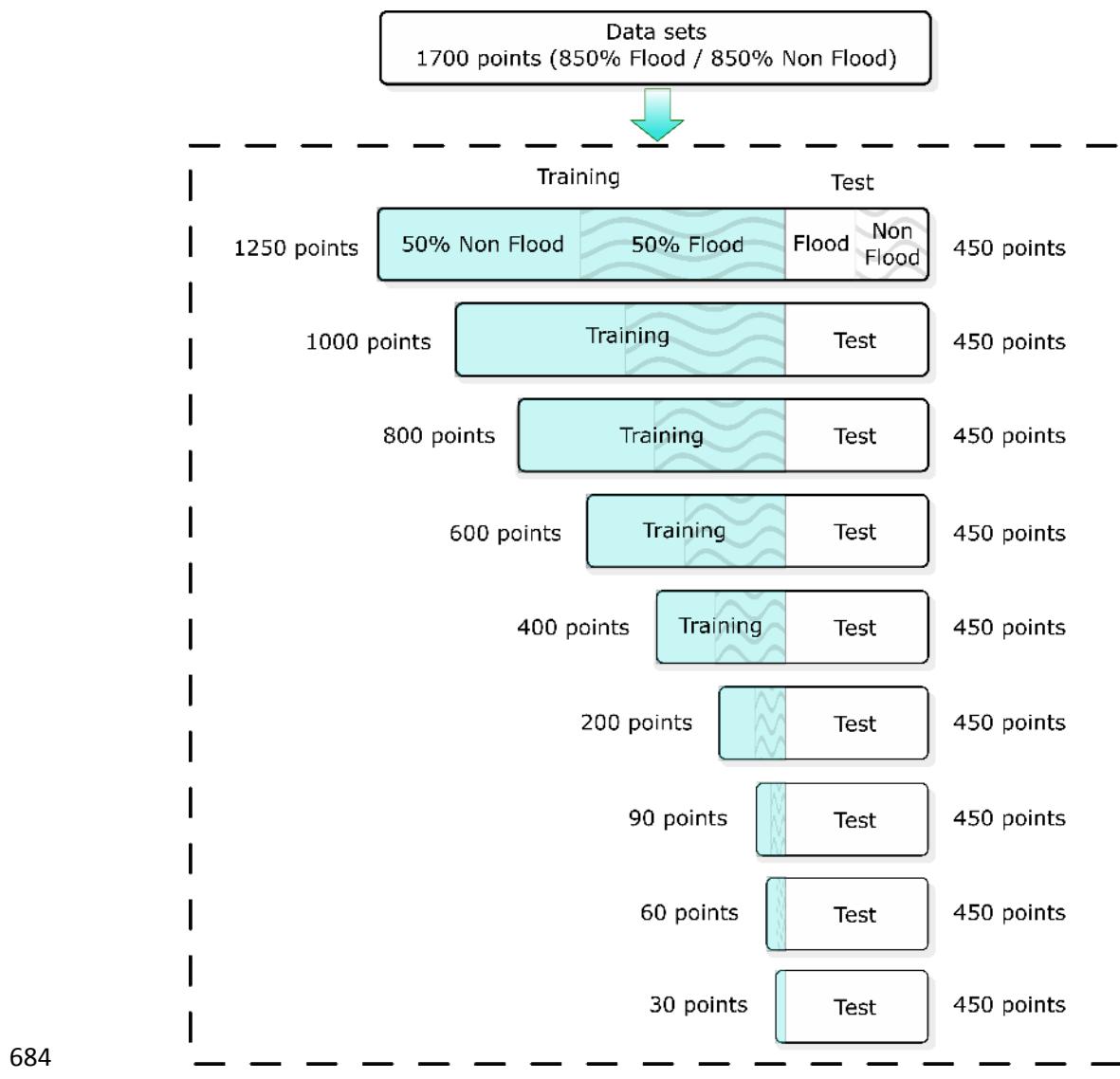
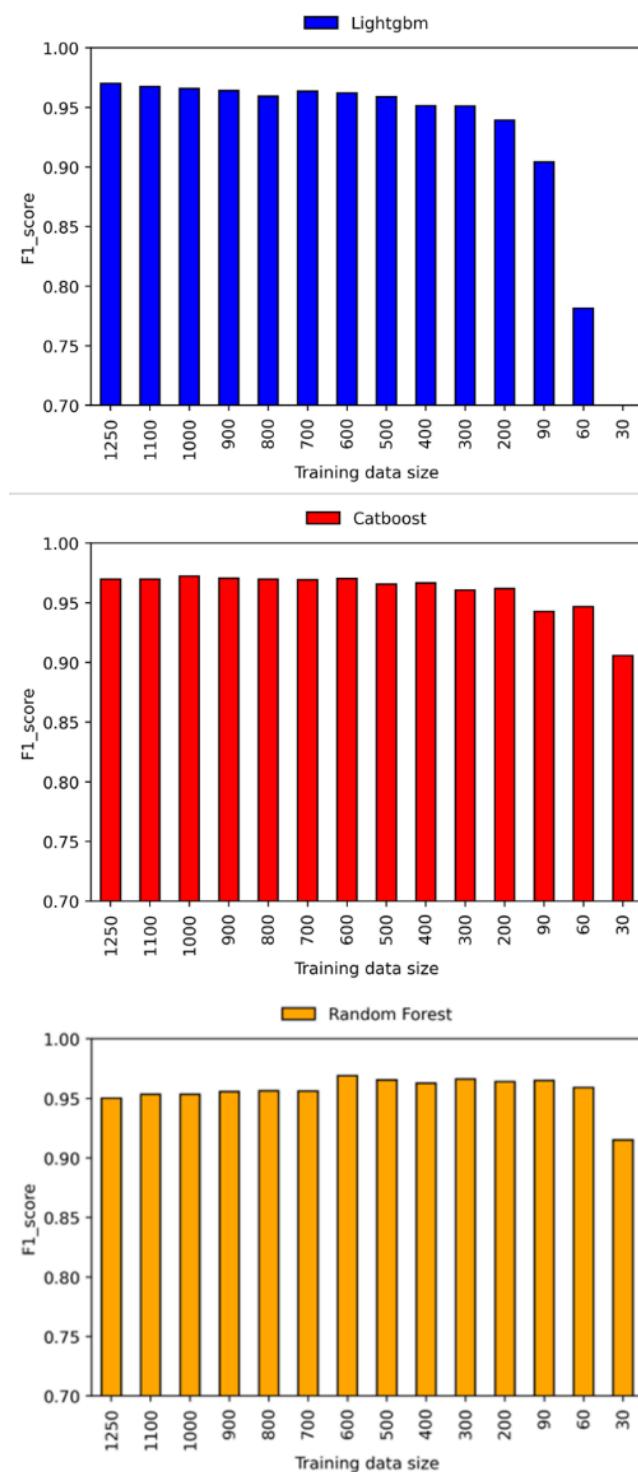


Fig. 12. Datasets used in the training and testing of the ML models.

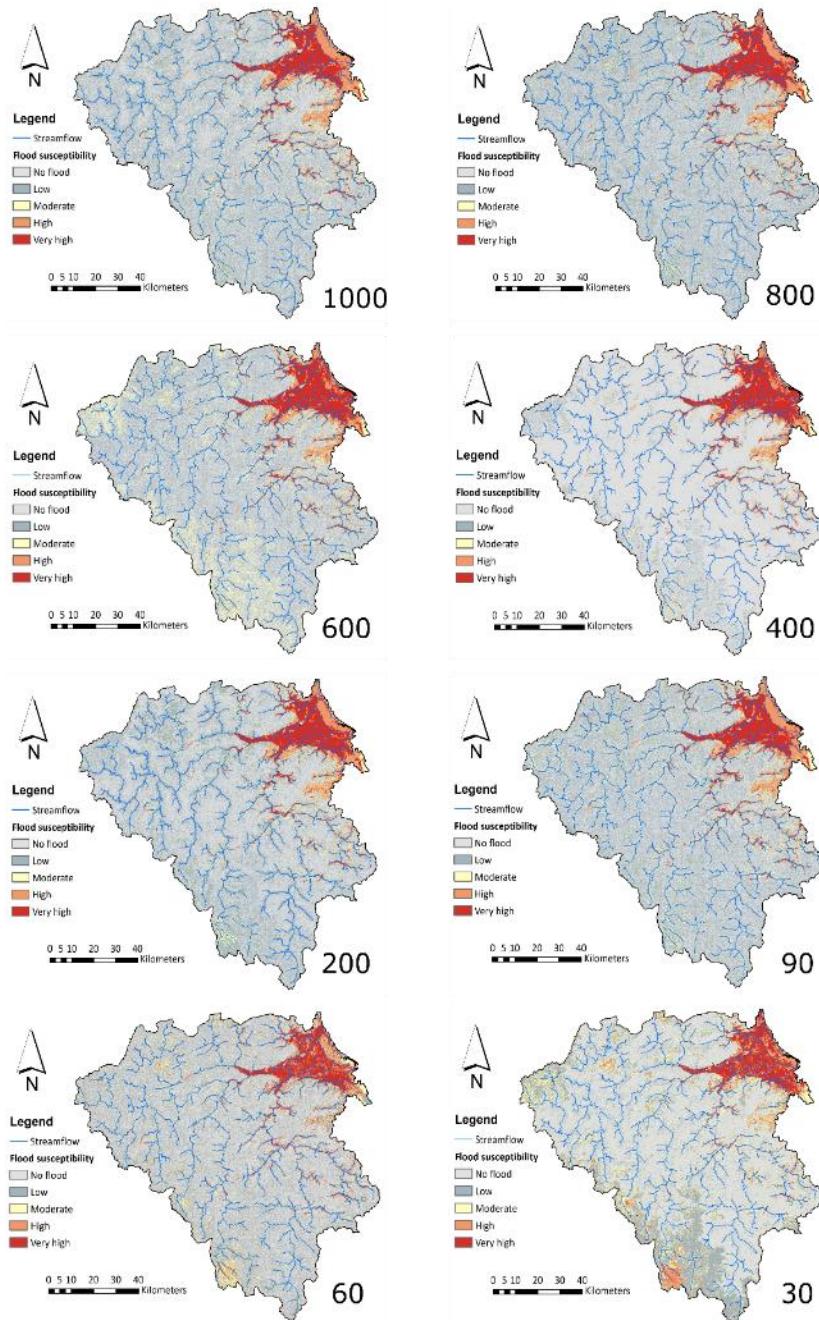


687

688

Fig. 13. Accuracy of the models based on different training datasets.

689



690

691

Fig. 14. Impact of data size on flood susceptibility map.

692

693

694

695

696

697

698

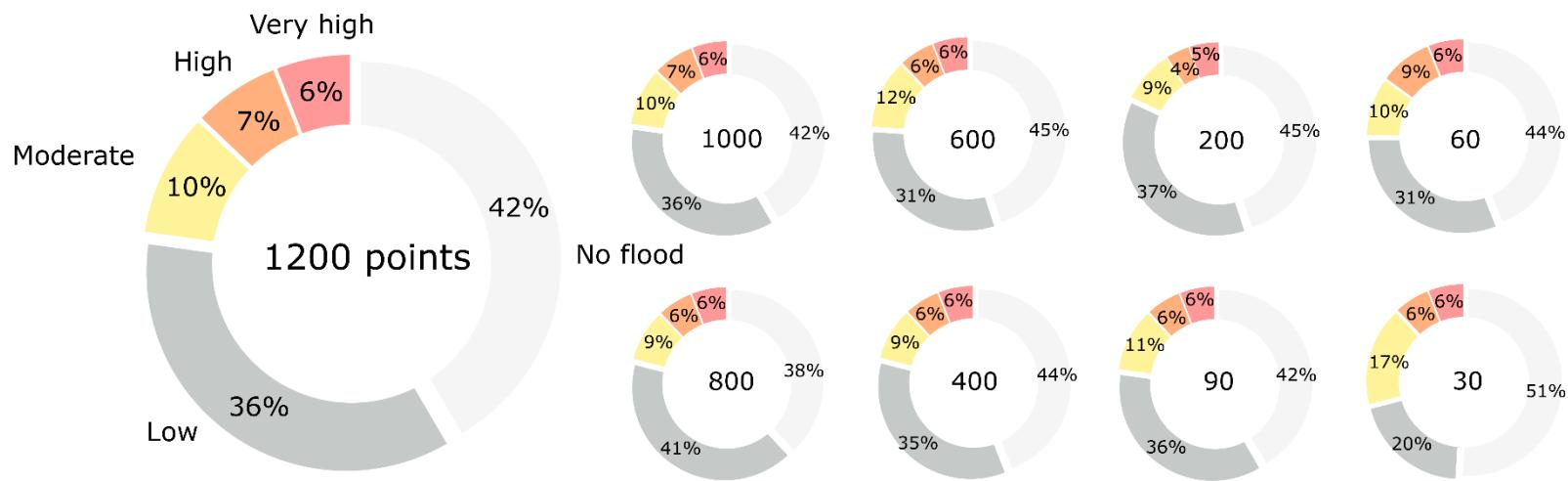


Fig. 15. Percentage of the affected areas under different flood susceptibility classes using different dataset sizes (RF method).

699 **4.5. Discussion and comparison with results of RRI model**

700 Currently, there is increasing interest in the use of data-driven methods as appropriate
701 alternatives to hydrological and hydraulic models. Therefore, the scientific community of
702 flood risk assessment is endeavoring to develop much more logical and mathematical
703 methods to predict flood-susceptible areas at different catchment scales (Arora et al., 2020).
704 In the study area, there have been some previous studies on flood susceptibility mapping that
705 use ML approaches and deep learning. Therefore, testing different methods is highly
706 recommended, particularly in regions where data are limited and hydrological models are
707 challenging. This study provides the application of three ML methods: RF, LightGBM, and
708 CatBoost. The latter two methods were tested for FFS mapping for the first time in this humid
709 region with a high frequency of typhoon events. The results of the FFS maps prove that both
710 methods can predict flood-prone areas with acceptable accuracy in comparison with the RF
711 method, which is widely applicable. RFs are applied in many other related studies with
712 different accuracies, e.g., AUC = 78% (Band et al., 2020), 99.3% (Li et al., 2019), 94.5%
713 (Talukdar et al., 2020), 93.8% (Park and Lee, 2020), and 89.4% (Nguyen et al., 2018). In this
714 study AUC = 99% for RF was higher than most of the previous studies. Additionally, the
715 newly applied methods of LightGBM and CatBoost showed almost the same accuracy of
716 99% and 98%, respectively, revealing better performance than most of the previous studies.
717 These three methods have been tested in Hurghada, Egypt (Saber et al., 2021), stating that
718 LightGBM has the advantage of better classification metrics and fast processing time, and
719 outperforms other methodologies such as CatBoost and RF. In addition, their results showed
720 that LightGBM and CatBoost have proven their efficiency in flash flood prediction in arid
721 regions.

722 Additionally, the three methods exhibited better performance than the average of the
723 previously applied methods for flood susceptibility mapping, which was 90% as an average
724 of approximately 140 previous applications from more than 30 publications that have been
725 reviewed. The performance of the previous methods applied for FSM based on AUC varies
726 from 64% (Shafizadeh-Moghadam et al., 2018) to 99.3% (Li et al., 2019). LightGBM has

727 also been applied in other applications, such as the real-time prediction of regional urban
728 runoff (<https://doi.org/10.1016/j.jhydrol.2021.127124>), showing high accuracy and fast
729 processing predictions. CatBoost was also applied in Germany, with better performance in
730 comparison with other methods, showing good accuracy with an AUC of 0.816 (Kaiser, M.
731 H. E. (2021)).

732 The maps of flood susceptibility developed using ML techniques (Fig. 16) showed an
733 acceptable fit with the developed flood inundation map by the RRI model, showing that the
734 ML approaches are promising for flood prediction and can be used without detailed
735 observations and challenges of model calibrations as alternative tools for hydrological
736 models. The results of LightGBM and CatBoost are more comparable to the flood inundation
737 map developed by the physical RRI model, indicating that they are more acceptable than RF,
738 which overestimates the low flood susceptibility level in the study area.

739 Furthermore, we tested different datasets for the training of the three ML models, with the
740 conclusion that datasets greater than 90 points can be sufficiently accurate for reasonable
741 prediction of the FSM. LightGBM and CatBoost showed a slightly declining trend in the
742 accuracy of the results based on the dataset sizes; however, RF did not show such a trend.
743 These results are highly valuable for the application of ML to ungauged basins with very
744 limited datasets.

745

746 **5. Conclusions**

747 Flooding resulting from typhoons is one of the most threatening disasters in Asian
748 countries and worldwide. Therefore, the present study introduced three ML methods to
749 accurately predict flooding susceptibility in humid areas in Vietnam. The first method is RF,
750 which is well known and widely applicable in many applications including FSM, and the
751 other two methods (LightGBM and CatBoost) were examined for the first time for FSM in
752 this humid region. The methods were trained and validated based on a flood inventory map
753 and 10 influencing flood factors. Owing to the availability of high-quality observations, we

754 also tested different datasets for the training (i.e., 30, 60, 90, 200, 400, 600, 800, 1000, and
755 1250 data points) to determine the least data points that provided acceptable reliability, as
756 well as to understand the differences in the spatial FSMs in the study area. Interestingly, we
757 found that the accuracy of results based on all the tested datasets was higher than 90%,
758 indicating that a limited number of observations can be used efficiently in terms of model
759 accuracy, although the final FSMs showed some differences spatially from one susceptibility
760 level to the others. This finding is highly important to demonstrate that ML methods can
761 work efficiently with an acceptable level of accuracy within a small number of actual training
762 datasets. The conclusions of this study can be summarized as follows:

- 763 ② We applied three ML models—RF, LightGBM, and CatBoost—to predict flood
764 susceptibility in humid areas that experienced successive extreme typhoons.
 - 765 ② The LightGBM and CatBoost models were tested for the first time in this specific
766 climatic region, and showed high performance in comparison with the RF method.
 - 767 ② The results of the ML methods showed good agreement with the rainfall-runoff model
768 for flood inundation mapping, especially the LightGBM and CatBoost models in
769 terms of coverage areas of the flood susceptibility levels.
 - 770 ② Different training datasets were examined to determine the lowest acceptable number
771 of observations for flood susceptibility in ML.
 - 772 ② The FSMs demonstrated that downstream areas with high residential and agricultural
773 activity are highly susceptible to flooding.
 - 774 ② The outcomes of this study could be used as guidance and reference for flood risk
775 mitigation and management in this region and, consequently, to assist managers,
776 decision makers, and planners in effectively managing and mitigating floods in high-
777 flood-susceptible areas.
- 778 As an extension of this research, we are examining ML for predicting flood depth, as a
779 comparison between ML approaches and RRI models. The study concludes that the ML
780 approach can be an alternative for hydrological models in terms of flood susceptibility
781 mapping; however, we are still working to prove this with flood depth prediction.

782 **Acknowledgment:**

783 Data sources, such as the Geospatial Information Authority of Japan, Chiba University, and
784 collaborating organizations.

785

786 **Funding:**

787 This work was funded by the Asia-Pacific Network for Global Change Research (APN)
788 under project reference number CRRP2020-09MYKantoush (Funder ID:
789 <https://doi.org/10.13039/100005536>).

790

791 **Declaration of Competing Interest**

792 The authors declare that they have no known competing financial interests or personal
793 relationships that could have influenced the work reported in this study.

794

795 **References**

- 796 [1] Abdel-Fattah, M., Kantoush, S.A., Saber, M., Sumi, T., 2018. Rainfall-runoff modeling for
797 extreme flash floods in wadi samail, oman. J. Jpn. Soc. Civ. Eng. Ser B1 Hydraul. Eng. 74,
798 I_691-I_696.
- 799 [2] Abdrabo, K.I., Kantoush, S.A., Saber, M., Sumi, T., Habiba, O.M., Elleithy, D., Elboshy, B., 2020.
800 Integrated Methodology for Urban Flood Risk Mapping at the Microscale in Ungauged
801 Regions: A Case Study of Hurghada, Egypt. Remote Sens. 12, 3548.
802 <https://doi.org/10.3390/rs12213548>
- 803 [3] Abushandi, E.H., Merkel, B.J., 2011. Application of IHACRES rainfall-runoff model to the
804 Wadi Dhuliel arid catchment, Jordan. J. Water Clim. Change 2, 56–71.
- 805 [4] Akay, A.E., Taş, İ., 2020. Mapping the risk of winter storm damage using GIS-based fuzzy
806 logic. J. For. Res. 31, 729–742.
- 807 [5] Ali, S.A., Parvin, F., Pham, Q.B., Vojtek, M., Vojteková, J., Costache, R., Linh, N.T.T., Nguyen,
808 H.Q., Ahmad, A., Ghorbani, M.A., 2020. GIS-based comparative assessment of flood
809 susceptibility mapping using hybrid multi-criteria decision-making approach, naïve Bayes
810 tree, bivariate statistics and logistic regression: A case of Topľa basin, Slovakia. Ecol. Indic.
811 117, 106620.
- 812 [6] Arabameri, A., Saha, S., Mukherjee, K., Blaschke, T., Chen, W., Ngo, P.T.T., Band, S.S., 2020.
813 Modeling Spatial Flood using Novel Ensemble Artificial Intelligence Approaches in Northern
814 Iran. Remote Sens. 12, 3423. <https://doi.org/10.3390/rs12203423>
- 815 [7] Arora, A., Arabameri, A., Pandey, M., Siddiqui, M.A., Shukla, U.K., Bui, D.T., Mishra, V.N.,
816 Bhardwaj, A., 2020. Optimization of state-of-the-art fuzzy-metaheuristic ANFIS-based
817 machine learning models for flood susceptibility prediction mapping in the Middle Ganga

- Plain, India. Sci. Total Environ. 750, 141565.
<https://doi.org/10.1016/j.scitotenv.2020.141565>
- [8] Arora, A., Pandey, M., Siddiqui, M.A., Hong, H., Mishra, V.N., 2019. Spatial flood susceptibility prediction in Middle Ganga Plain: comparison of frequency ratio and Shannon's entropy models. *Geocarto Int.* 1–32.
- [9] Aryal, S.K., Mein, R.G., O'Loughlin, E.M., 2003. The concept of effective length in hillslopes: assessing the influence of climate and topography on the contributing areas of catchments. *Hydrol. Process.* 17, 131–151.
- [10] Avitabile, V., Schultz, M., Herold, N., De Bruin, S., Pratiast, A.K., Manh, C.P., Quang, H.V., Herold, M., 2016. Carbon emissions from land cover change in Central Vietnam. *Carbon Manag.* 7, 333–346.
- [11] Bachmair, S., Svensson, C., Prosdocimi, I., Hannaford, J., Stahl, K., 2017. Developing drought impact functions for drought risk management. *Nat. Hazards Earth Syst. Sci.* 17, 1947–1960.
- [12] Band, S.S., Janizadeh, S., Chandra Pal, S., Saha, A., Chakrabortty, R., Melesse, A.M., Mosavi, A., 2020. Flash Flood Susceptibility Modeling Using New Approaches of Hybrid and Ensemble Tree-Based Machine Learning Algorithms. *Remote Sens.* 12, 3568. <https://doi.org/10.3390/rs12213568>
- [13] Bellu, A., Fernandes, L.F.S., Cortes, R.M., Pacheco, F.A., 2016. A framework model for the dimensioning and allocation of a detention basin system: The case of a flood-prone mountainous watershed. *J. Hydrol.* 533, 567–580.
- [14] Bisht, S., Chaudhry, S., Sharma, S., Soni, S., 2018. Assessment of flash flood vulnerability zonation through Geospatial technique in high altitude Himalayan watershed, Himachal Pradesh India. *Remote Sens. Appl. Soc. Environ.* 12, 35–47.
- [15] Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5–32.
- [16] Bui, D.T., Hoang, N.-D., Pham, T.-D., Ngo, P.-T.T., Hoa, P.V., Minh, N.Q., Tran, X.-T., Samui, P., 2019a. A new intelligence approach based on GIS-based multivariate adaptive regression splines and metaheuristic optimization for predicting flash flood susceptible areas at high-frequency tropical typhoon area. *J. Hydrol.* 575, 314–326.
- [17] Bui, D.T., Ngo, P.-T.T., Pham, T.D., Jaafari, A., Minh, N.Q., Hoa, P.V., Samui, P., 2019b. A novel hybrid approach based on a swarm intelligence optimized extreme learning machine for flash flood susceptibility mapping. *CATENA* 179, 184–196. <https://doi.org/10.1016/j.catena.2019.04.009>
- [18] Bui, D.T., Panahi, M., Shahabi, H., Singh, V.P., Shirzadi, A., Chapi, K., Khosravi, K., Chen, W., Panahi, S., Li, S., 2018. Novel hybrid evolutionary algorithms for spatial prediction of floods. *Sci. Rep.* 8, 1–14.
- [19] Bui, D.T., Pradhan, B., Lofman, O., Revhaug, I., Dick, O.B., 2012. Spatial prediction of landslide hazards in Hoa Binh province (Vietnam): a comparative assessment of the efficacy of evidential belief functions and fuzzy logic models. *Catena* 96, 28–40.
- [20] Bui, Q.-T., Nguyen, Q.-H., Nguyen, X.L., Pham, V.D., Nguyen, H.D., Pham, V.-M., 2020. Verification of novel integrations of swarm intelligence algorithms into deep learning neural network for flood susceptibility mapping. *J. Hydrol.* 581, 124379. <https://doi.org/10.1016/j.jhydrol.2019.124379>
- [21] Cardenas, M.B., Wilson, J., Zlotnik, V.A., 2004. Impact of heterogeneity, bed forms, and stream curvature on subchannel hyporheic exchange. *Water Resour. Res.* 40.

- 862 [22]Catal, C., Diri, B., 2009. Investigating the effect of dataset size, metrics sets, and feature
863 selection techniques on software fault prediction problem. *Inf. Sci.* 179, 1040–1058.
- 864 [23]Chang, M.-J., Chang, H.-K., Chen, Y.-C., Lin, G.-F., Chen, P.-A., Lai, J.-S., Tan, Y.-C., 2018. A
865 support vector machine forecasting model for typhoon flood inundation mapping and early
866 flood warning systems. *Water* 10, 1734.
- 867 [24]Chapi, K., Singh, V.P., Shirzadi, A., Shahabi, H., Bui, D.T., Pham, B.T., Khosravi, K., 2017. A
868 novel hybrid artificial intelligence approach for flood susceptibility assessment. *Environ.
869 Model. Softw.* 95, 229–245.
- 870 [25]Chen, T.-H.K., Qiu, C., Schmitt, M., Zhu, X.X., Sabel, C.E., Prischepov, A.V., 2020. Mapping
871 horizontal and vertical urban densification in Denmark with Landsat time-series from 1985
872 to 2018: A semantic segmentation solution. *Remote Sens. Environ.* 251, 112096.
873 <https://doi.org/10.1016/j.rse.2020.112096>
- 874 [26]Choubin, B., Moradi, E., Golshan, M., Adamowski, J., Sajedi-Hosseini, F., Mosavi, A., 2019.
875 An ensemble prediction of flood susceptibility using multivariate discriminant analysis,
876 classification and regression trees, and support vector machines. *Sci. Total Environ.* 651,
877 2087–2096.
- 878 [27]Costache, R., Hong, H., Pham, Q.B., 2020a. Comparative assessment of the flash-flood
879 potential within small mountain catchments using bivariate statistics and their novel hybrid
880 integration with machine learning models. *Sci. Total Environ.* 711, 134514.
881 <https://doi.org/10.1016/j.scitotenv.2019.134514>
- 882 [28]Costache, R., Popa, M.C., Tien Bui, D., Diaconu, D.C., Ciubotaru, N., Minea, G., Pham, Q.B.,
883 2020b. Spatial predicting of flood potential areas using novel hybridizations of fuzzy
884 decision-making, bivariate statistics, and machine learning. *J. Hydrol.* 585, 124808.
885 <https://doi.org/10.1016/j.jhydrol.2020.124808>
- 886 [29]Darabi, H., Choubin, B., Rahmati, O., Haghghi, A.T., Pradhan, B., Kløve, B., 2019. Urban flood
887 risk mapping using the GARP and QUEST models: A comparative study of machine learning
888 techniques. *J. Hydrol.* 569, 142–154.
- 889 [30]Demirel, M.C., Venancio, A., Kahya, E., 2009. Flow forecast by SWAT model and ANN in
890 Pracana basin, Portugal. *Adv. Eng. Softw.* 40, 467–473.
- 891 [31]Devkota, K.C., Regmi, A.D., Pourghasemi, H.R., Yoshida, K., Pradhan, B., Ryu, I.C., Dhital, M.R.,
892 Althuwaynee, O.F., 2013. Landslide susceptibility mapping using certainty factor, index of
893 entropy and logistic regression models in GIS and their comparison at Mugling–Narayanghat
894 road section in Nepal Himalaya. *Nat. Hazards* 65, 135–165.
- 895 [32]Dhara, S., Dang, T., Parial, K., Lu, X.X., 2020. Accounting for Uncertainty and Reconstruction
896 of Flooding Patterns Based on Multi-Satellite Imagery and Support Vector Machine
897 Technique: A Case Study of Can Tho City, Vietnam. *Water* 12, 1543.
898 <https://doi.org/10.3390/w12061543>
- 899 [33]Dodangeh, E., Choubin, B., Eigdir, A.N., Nabipour, N., Panahi, M., Shamshirband, S., Mosavi,
900 A., 2020. Integrated machine learning methods with resampling algorithms for flood
901 susceptibility prediction. *Sci. Total Environ.* 705, 135983.
902 <https://doi.org/10.1016/j.scitotenv.2019.135983>
- 903 [34]Dorogush, A.V., Ershov, V., Gulin, A., 2018. CatBoost: gradient boosting with categorical
904 features support. *ArXiv Prepr. ArXiv181011363*.
- 905 [35]Dou, J., Yunus, A.P., Bui, D.T., Merghadi, A., Sahana, M., Zhu, Z., Chen, C.-W., Khosravi, K.,
906 Yang, Y., Pham, B.T., 2019. Assessment of advanced random forest and decision tree

- 907 algorithms for modeling rainfall-induced landslide susceptibility in the Izu-Oshima Volcanic
908 Island, Japan. *Sci. Total Environ.* 662, 332–346.
- 909 [36]Esfandiari, M., Jabari, S., McGrath, H., Coleman, D., 2020. FLOOD MAPPING USING RANDOM
910 FOREST AND IDENTIFYING THE ESSENTIAL CONDITIONING FACTORS; A CASE STUDY IN
911 FREDERICTON, NEW BRUNSWICK, CANADA. *ISPRS Ann. Photogramm. Remote Sens. Spat.*
912 *Inf. Sci.* 5.
- 913 [37]Fenia, F., Kavetski, D., Savenije, H.H., Clark, M.P., Schoups, G., Pfister, L., Freer, J., 2014.
914 Catchment properties, function, and conceptual model representation: is there a
915 correspondence? *Hydrol. Process.* 28, 2451–2467.
- 916 [38]Gigović, L., Pamučar, D., Bajić, Z., Drobnjak, S., 2017. Application of GIS-interval rough AHP
917 methodology for flood hazard mapping in urban areas. *Water* 9, 360.
- 918 [39]Glenn, E.P., Morino, K., Nagler, P.L., Murray, R.S., Pearlstein, S., Hultine, K.R., 2012. Roles of
919 saltcedar (*Tamarix spp.*) and capillary rise in salinizing a non-flooding terrace on a flow-
920 regulated desert river. *J. Arid Environ.* 79, 56–65.
- 921 [40]González-Arqueros, M.L., Mendoza, M.E., Bocco, G., Castillo, B.S., 2018. Flood susceptibility
922 in rural settlements in remote zones: The case of a mountainous basin in the Sierra-Costa
923 region of Michoacán, Mexico. *J. Environ. Manage.* 223, 685–693.
- 924 [41]Hirabayashi, Y., Mahendran, R., Koitala, S., Konoshima, L., Yamazaki, D., Watanabe, S., Kim,
925 H., Kanae, S., 2013. Global flood risk under climate change. *Nat. Clim. Change* 3, 816–821.
- 926 [42]Hölting, B., Coldewey, W.G., 2019. Surface water infiltration, in: *Hydrogeology*. Springer, pp.
927 33–37.
- 928 [43]Hong, H., Tsangaratos, P., Ilia, I., Liu, J., Zhu, A.-X., Chen, W., 2018. Application of fuzzy
929 weight of evidence and data mining techniques in construction of flood susceptibility map
930 of Poyang County, China. *Sci. Total Environ.* 625, 575–588.
- 931 [44]Hsu, K., Gupta, H.V., Sorooshian, S., 1995. Artificial neural network modeling of the rainfall-
932 runoff process. *Water Resour. Res.* 31, 2517–2530.
- 933 [45]Humphrey, G.B., Gibbs, M.S., Dandy, G.C., Maier, H.R., 2016. A hybrid approach to monthly
934 streamflow forecasting: integrating hydrological model outputs into a Bayesian artificial
935 neural network. *J. Hydrol.* 540, 623–640.
- 936 [46]Izquierdo-Verdiguier, E., Zurita-Milla, R., 2020. An evaluation of Guided Regularized Random
937 Forest for classification and regression tasks in remote sensing. *Int. J. Appl. Earth Obs.*
938 *Geoinformation* 88, 102051.
- 939 [47]Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., Liu, T.-Y., 2017. Lightgbm: A
940 highly efficient gradient boosting decision tree. *Adv. Neural Inf. Process. Syst.* 30, 3146–
941 3154.
- 942 [48]Khosravi, K., Nohani, E., Maroufinia, E., Pourghasemi, H.R., 2016. A GIS-based flood
943 susceptibility assessment and its mapping in Iran: a comparison between frequency ratio
944 and weights-of-evidence bivariate statistical models with multi-criteria decision-making
945 technique. *Nat. Hazards* 83, 947–987.
- 946 [49]Khosravi, K., Shahabi, H., Pham, B.T., Adamowski, J., Shirzadi, A., Pradhan, B., Dou, J., Ly, H.-
947 B., Gróf, G., Ho, H.L., Hong, H., Chapi, K., Prakash, I., 2019. A comparative assessment of
948 flood susceptibility modeling using Multi-Criteria Decision-Making Analysis and Machine
949 Learning Methods. *J. Hydrol.* 573, 311–323. <https://doi.org/10.1016/j.jhydrol.2019.03.073>

- 950 [50]Kia, M.B., Pirasteh, S., Pradhan, B., Mahmud, A.R., Sulaiman, W.N.A., Moradi, A., 2012. An
951 artificial neural network model for flood simulation using GIS: Johor River Basin, Malaysia.
952 Environ. Earth Sci. 67, 251–264.
- 953 [51]Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., Nearing, G., 2019.
954 Benchmarking a catchment-aware long short-term memory network (LSTM) for large-scale
955 hydrological modeling. Hydrol. Earth Syst. Sci. Discuss. 1–32.
- 956 [52]Li, X., Yan, D., Wang, K., Weng, B., Qin, T., Liu, S., 2019. Flood Risk Assessment of Global
957 Watersheds Based on Multiple Machine Learning Models. Water 11, 1654.
958 <https://doi.org/10.3390/w11081654>
- 959 [53]Luu, C., Pham, B.T., Phong, T.V., Costache, R., Nguyen, H.D., Amiri, M., Bui, Q.D., Nguyen,
960 L.T., Le, H.V., Prakash, I., Trinh, P.T., 2021. GIS-based ensemble computational models for
961 flood susceptibility prediction in the Quang Binh Province, Vietnam. J. Hydrol. 599, 126500.
962 <https://doi.org/10.1016/j.jhydrol.2021.126500>
- 963 [54]Malekipirbazari, M., Aksakalli, V., 2015. Risk assessment in social lending via random forests.
964 Expert Syst. Appl. 42, 4621–4631.
- 965 [55]Masood, M., Takeuchi, K., 2012. Assessment of flood hazard, vulnerability and risk of mid-
966 eastern Dhaka using DEM and 1D hydrodynamic model. Nat. Hazards 61, 757–770.
- 967 [56]Meadows, M., Wilson, M., 2021. A Comparison of Machine Learning Approaches to Improve
968 Free Topography Data for Flood Modelling. Remote Sens. 13, 275.
- 969 [57]Meraj, G., Khan, T., Romshoo, S.A., Farooq, M., Rohitashw, K., Sheikh, B.A., 2018. An
970 integrated geoinformatics and hydrological modelling-based approach for effective flood
971 management in the Jhelum Basin, NW Himalaya. Multidiscip. Digit. Publ. Inst. Proc. 7, 8.
- 972 [58]Nauditt, A., Firoz, A.B.M., Trinh, V.Q., Fink, M., Stolpe, H., Ribbe, L., 2017. Hydrological
973 drought risk assessment in an anthropogenically impacted tropical catchment, Central
974 Vietnam, in: Land Use and Climate Change Interactions in Central Vietnam. Springer, pp.
975 223–239.
- 976 [59]Nauditt, A., Ribbe, L., 2017. Land use and climate change interactions in central Vietnam.
977 Springer.
- 978 [60]Ngo, P.-T.T., Pham, T.D., Hoang, N.-D., Tran, D.A., Amiri, M., Le, T.T., Hoa, P.V., Bui, P.V., Nhu,
979 V.-H., Bui, D.T., 2021. A new hybrid equilibrium optimized SysFor based geospatial data
980 mining for tropical storm-induced flash flood susceptible mapping. J. Environ. Manage. 280,
981 111858. <https://doi.org/10.1016/j.jenvman.2020.111858>
- 982 [61]Nguyen, H.Q., Degener, J., Kappas, M., 2015. Flash flood prediction by coupling KINEROS2
983 and HEC-RAS models for tropical regions of Northern Vietnam. Hydrology 2, 242–265.
- 984 [62]Nguyen, P., Ombadi, M., Gorooh, V.A., Shearer, E.J., Sadeghi, M., Sorooshian, S., Hsu, K.,
985 Bolvin, D., Ralph, M.F., 2020. PERSIANN Dynamic Infrared–Rain Rate (PDIR-Now): A Near-
986 Real-Time, Quasi-Global Satellite Precipitation Dataset. J. Hydrometeorol. 21, 2893–2906.
- 987 [63]Nguyen, V.-N., Tien Bui, D., Ngo, P.-T.T., Nguyen, Q.-P., Nguyen, V.C., Long, N.Q., Revhaug,
988 I., 2018. An Integration of Least Squares Support Vector Machines and Firefly Optimization
989 Algorithm for Flood Susceptible Modeling Using GIS, in: Tien Bui, D., Ngoc Do, A., Bui, H.-B.,
990 Hoang, N.-D. (Eds.), Advances and Applications in Geospatial Technology and Earth
991 Resources. Springer International Publishing, Cham, pp. 52–64.
992 https://doi.org/10.1007/978-3-319-68240-2_4
- 993 [64]Nguyen, V.-N., Yariyan, P., Amiri, M., Dang Tran, A., Pham, T.D., Do, M.P., Thi Ngo, P.T., Nhu,
994 V.-H., Quoc Long, N., Tien Bui, D., 2020. A New Modeling Approach for Spatial Prediction of

- 995 Flash Flood with Biogeography Optimized CHAID Tree Ensemble and Remote Sensing Data.
996 Remote Sens. 12, 1373. <https://doi.org/10.3390/rs12091373>
- 997 [65]Nhu, V.-H., Thi Ngo, P.-T., Pham, T.D., Dou, J., Song, X., Hoang, N.-D., Tran, D.A., Cao, D.P.,
998 Aydilek, I.B., Amiri, M., Costache, R., Hoa, P.V., Tien Bui, D., 2020. A New Hybrid Firefly–PSO
999 Optimized Random Subspace Tree Intelligence for Torrential Rainfall-Induced Flash Flood
1000 Susceptible Mapping. Remote Sens. 12, 2688. <https://doi.org/10.3390/rs12172688>
- 1001 [66]Öztürk, F., Akdeniz, F., 2000. Ill-conditioning and multicollinearity. Linear Algebra Its Appl.
1002 321, 295–305.
- 1003 [67]Pachauri, R.K., Allen, M.R., Barros, V.R., Broome, J., Cramer, W., Christ, R., Church, J.A.,
1004 Clarke, L., Dahe, Q., Dasgupta, P., 2014. Climate change 2014: synthesis report. Contribution
1005 of Working Groups I, II and III to the fifth assessment report of the Intergovernmental Panel
1006 on Climate Change. Ipcc.
- 1007 [68]Pal, M., 2005. Random forest classifier for remote sensing classification. Int. J. Remote Sens.
1008 26, 217–222.
- 1009 [69]Park, S.-J., Lee, D.-K., 2020. Prediction of coastal flooding risk under climate change impacts
1010 in South Korea using machine learning algorithms. Environ. Res. Lett. 15, 094052.
1011 <https://doi.org/10.1088/1748-9326/aba5b3>
- 1012 [70]Perera, E.D.P., Sayama, T., Magome, J., Hasegawa, A., Iwami, Y., 2017. RCP8.5-based future
1013 flood hazard analysis for the lower Mekong river basin. Hydrology 4, 55.
- 1014 [71]Pham, B.T., Jaafari, A., Nguyen-Thoi, T., Van Phong, T., Nguyen, H.D., Satyam, N., Masroor,
1015 M., Rehman, S., Sajjad, H., Sahana, M., 2021a. Ensemble machine learning models based on
1016 Reduced Error Pruning Tree for prediction of rainfall-induced landslides. Int. J. Digit. Earth
1017 14, 575–596.
- 1018 [72]Pham, B.T., Luu, C., Phong, T.V., Trinh, P.T., Shirzadi, A., Renoud, S., Asadi, S., Le, H.V., von
1019 Meding, J., Clague, J.J., 2021b. Can deep learning algorithms outperform benchmark
1020 machine learning algorithms in flood susceptibility modeling? J. Hydrol. 592, 125615.
1021 <https://doi.org/10.1016/j.jhydrol.2020.125615>
- 1022 [73]Pham, T.D., Xia, J., Ha, N.T., Bui, D.T., Le, N.N., Tekeuchi, W., 2019. A review of remote
1023 sensing approaches for monitoring blue carbon ecosystems: Mangroves, seagrassesand salt
1024 marshes during 2010–2018. Sensors 19, 1933.
- 1025 [74]Pourghasemi, H.R., Kariminejad, N., Amiri, M., Edalat, M., Zarafshar, M., Blaschke, T., Cerdá,
1026 A., 2020. Assessing and mapping multi-hazard risk susceptibility using a machine learning
1027 technique. Sci. Rep. 10, 1–11.
- 1028 [75]Predick, K.I., Turner, M.G., 2008. Landscape configuration and flood frequency influence
1029 invasive shrubs in floodplain forests of the Wisconsin River (USA). J. Ecol. 96, 91–102.
- 1030 [76]Quinlan, J.R., 1986. Induction of decision trees. Mach. Learn. 1, 81–106.
- 1031 [77]Rahman, M., Ningsheng, C., Islam, M.M., Dewan, A., Iqbal, J., Washakh, R.M.A., Shufeng, T.,
1032 2019. Flood susceptibility assessment in Bangladesh using machine learning and multi-
1033 criteria decision analysis. Earth Syst. Environ. 3, 585–601.
- 1034 [78]Rahmati, O., Pourghasemi, H.R., Zeinivand, H., 2016. Flood susceptibility mapping using
1035 frequency ratio and weights-of-evidence models in the Golestan Province, Iran. Geocarto
1036 Int. 31, 42–70.
- 1037 [79]RETA, 2011. Investment, Managing water in Asia's river basins: Charting progress and
1038 facilitating - The Vu Gia-Thu Bon Basin.

- 1039 [80]Saber, M., Abdrabo, K.I., Habiba, O.M., Kantosh, S.A., Sumi, T., 2020. Impacts of Triple
1040 Factors on Flash Flood Vulnerability in Egypt: Urban Growth, Extreme Climate, and
1041 Mismanagement. *Geosciences* 10, 24.
- 1042 [81]Saber, M., Boulmaiz, T., Guermoui, M., Abdrado, K.I., Kantoush, S.A., Sumi, T., Boutaghane,
1043 H., Nohara, D., Mabrouk, E., 2021. Examining LightGBM and CatBoost models for wadi flash
1044 flood susceptibility prediction. *Geocarto Int.* 1–27.
- 1045 [82]Schoppa, L., Disse, M., Bachmair, S., 2020. Evaluating the performance of random forest for
1046 large-scale flood discharge simulation. *J. Hydrol.* 590, 125531.
- 1047 [83]Shafizadeh-Moghadam, H., Valavi, R., Shahabi, H., Chapi, K., Shirzadi, A., 2018. Novel
1048 forecasting approaches using combination of machine learning and statistical models for
1049 flood susceptibility mapping. *J. Environ. Manage.* 217, 1–11.
1050 <https://doi.org/10.1016/j.jenvman.2018.03.089>
- 1051 [84]Shahabi, H., Shirzadi, A., Ronoud, S., Asadi, S., Pham, B.T., Mansouripour, F., Geertsema, M.,
1052 Clague, J.J., Bui, D.T., 2020. Flash flood susceptibility mapping using a novel deep learning
1053 model based on deep belief network, back propagation and genetic algorithm. *Geosci. Front.*
1054 S1674987120302401. <https://doi.org/10.1016/j.gsf.2020.10.007>
- 1055 [85]Shirzadi, A., Asadi, S., Shahabi, H., Ronoud, S., Clague, J.J., Khosravi, K., Pham, B.T., Ahmad,
1056 B.B., Bui, D.T., 2020. A novel ensemble learning based on Bayesian Belief Network coupled
1057 with an extreme learning machine for flash flood susceptibility mapping. *Eng. Appl. Artif.*
1058 *Intell.* 96, 103971.
- 1059 [86]Talukdar, S., Ghose, B., Salam, R., Mahato, S., Pham, Q.B., Linh, N.T.T., Costache, R., Avand,
1060 M., 2020. Flood susceptibility modeling in Teesta River basin, Bangladesh using novel
1061 ensembles of bagging algorithms. *Stoch. Environ. Res. Risk Assess.* 34, 2277–2300.
- 1062 [87]Tam, T.H., Abd Rahman, M.Z., Harun, S., Hanapi, M.N., Kaoje, I.U., 2019. Application of
1063 Satellite rainfall products for flood inundation modelling in Kelantan River Basin, Malaysia.
1064 *Hydrology* 6, 95.
- 1065 [88]Tehrany, M.S., Jones, S., Shabani, F., 2019. Identifying the essential flood conditioning
1066 factors for flood prone area mapping using machine learning techniques. *CATENA* 175, 174–
1067 192. <https://doi.org/10.1016/j.catena.2018.12.011>
- 1068 [89]Tehrany, M.S., Kumar, L., 2018. The application of a Dempster–Shafer-based evidential
1069 belief function in flood susceptibility mapping and comparison with frequency ratio and
1070 logistic regression methods. *Environ. Earth Sci.* 77, 1–24.
- 1071 [90]Tehrany, M.S., Pradhan, B., Jebur, M.N., 2015. Flood susceptibility analysis and its
1072 verification using a novel ensemble support vector machine and frequency ratio method.
1073 *Stoch. Environ. Res. Risk Assess.* 29, 1149–1165.
- 1074 [91]Tehrany, M.S., Pradhan, B., Jebur, M.N., 2014. Flood susceptibility mapping using a novel
1075 ensemble weights-of-evidence and support vector machine models in GIS. *J. Hydrol.* 512,
1076 332–343.
- 1077 [92]Tehrany, M.S., Pradhan, B., Jebur, M.N., 2013. Spatial prediction of flood susceptible areas
1078 using rule based decision tree (DT) and a novel ensemble bivariate and multivariate
1079 statistical models in GIS. *J. Hydrol.* 504, 69–79.
- 1080 [93]Thao, N.T.P., Linh, T.T., Ha, N.T.T., Vinh, P.Q., Linh, N.T., 2020. Mapping flood inundation
1081 areas over the lower part of the Con River basin using Sentinel 1A imagery. *Vietnam J. Earth*
1082 *Sci.* 42, 288–297.

- 1083 [94] Tien Bui, D., Hoang, N.-D., 2017. A Bayesian framework based on a Gaussian mixture model
1084 and radial-basis-function Fisher discriminant analysis (BayGmmKda V1. 1) for spatial
1085 prediction of floods. *Geosci. Model Dev.* 10, 3391–3409.
- 1086 [95] Tien Bui, D., Hoang, N.-D., Martínez-Álvarez, F., Ngo, P.-T.T., Hoa, P.V., Pham, T.D., Samui,
1087 P., Costache, R., 2020. A novel deep learning neural network approach for predicting flash
1088 flood susceptibility: A case study at a high frequency tropical storm area. *Sci. Total Environ.*
1089 701, 134413. <https://doi.org/10.1016/j.scitotenv.2019.134413>
- 1090 [96] Tien Bui, D., Pradhan, B., Nampak, H., Bui, Q.-T., Tran, Q.-A., Nguyen, Q.-P., 2016. Hybrid
1091 artificial intelligence approach based on neural fuzzy inference model and metaheuristic
1092 optimization for flood susceptibility modeling in a high-frequency tropical cyclone area
1093 using GIS. *J. Hydrol.* 540, 317–330. <https://doi.org/10.1016/j.jhydrol.2016.06.027>
- 1094 [97] Tiwari, M.K., Chatterjee, C., 2010. Uncertainty assessment and ensemble flood forecasting
1095 using bootstrap based artificial neural networks (BANNs). *J. Hydrol.* 382, 20–33.
- 1096 [98] Torabi Haghghi, A., Menberu, M.W., Darabi, H., Akanegbu, J., Kløve, B., 2018. Use of remote
1097 sensing to analyse peatland changes after drainage for peat extraction. *Land Degrad. Dev.*
1098 29, 3479–3488.
- 1099 [99] Try, S., Tanaka, S., Tanaka, K., Sayama, T., Oeurng, C., Uk, S., Takara, K., Hu, M., Han, D., 2020.
1100 Comparison of gridded precipitation datasets for rainfall-runoff and inundation modeling in
1101 the Mekong River Basin. *PLoS One* 15, e0226814.
- 1102 [100] Tuyen, T.T., Jaafari, A., Yen, H.P.H., Nguyen-Thoi, T., Phong, T.V., Nguyen, H.D., Van
1103 Le, H., Phuong, T.T.M., Nguyen, S.H., Prakash, I., Pham, B.T., 2021. Mapping forest fire
1104 susceptibility using spatially explicit ensemble models based on the locally weighted
1105 learning algorithm. *Ecol. Inform.* 63, 101292. <https://doi.org/10.1016/j.ecoinf.2021.101292>
- 1106 [101] Vinet, F., 2008. Geographical analysis of damage due to flash floods in southern
1107 France: The cases of 12–13 November 1999 and 8–9 September 2002. *Appl. Geogr.* 28, 323–
1108 336. <https://doi.org/10.1016/j.apgeog.2008.02.007>
- 1109 [102] Vu, T.T.L., Nguyen, L.D., Hoang, T.S., Bui, T.A., Nguyen, M.T., Nguyen, T.H., 2011.
1110 Solutions for flood and drought prevention and mitigation in Quang Nam.
- 1111 [103] Wang, Y., Hong, H., Chen, W., Li, S., Panahi, M., Khosravi, K., Shirzadi, A., Shahabi,
1112 H., Panahi, S., Costache, R., 2019. Flood susceptibility mapping in Dingnan County (China)
1113 using adaptive neuro-fuzzy inference system with biogeography based optimization and
1114 imperialistic competitive algorithm. *J. Environ. Manage.* 247, 712–729.
- 1115 [104] Xu, Y., Dai, Y., Dong, Z.Y., Zhang, R., Meng, K., 2013. Extreme learning machine-
1116 based predictor for real-time frequency stability assessment of electric power systems.
1117 *Neural Comput. Appl.* 22, 501–508.
- 1118 [105] Yamazaki, D., Ikeshima, D., Tawatari, R., Yamaguchi, T., O'Loughlin, F., Neal, J.C.,
1119 Sampson, C.C., Kanae, S., Bates, P.D., 2017. A high-accuracy map of global terrain elevations.
1120 *Geophys. Res. Lett.* 44, 5844–5853.
- 1121 [106] Yang, S., Yang, D., Chen, J., Santisirisomboon, J., Lu, W., Zhao, B., 2020. A physical
1122 process and machine learning combined hydrological model for daily streamflow
1123 simulations of large watersheds with limited observation data. *J. Hydrol.* 590, 125206.
- 1124 [107] Yariyan, P., Janizadeh, S., Van Phong, T., Nguyen, H.D., Costache, R., Van Le, H.,
1125 Pham, B.T., Pradhan, B., Tiefenbacher, J.P., 2020. Improvement of best first decision trees
1126 using bagging and dagging ensembles for flood probability mapping. *Water Resour. Manag.*
1127 34, 3037–3053.

- 1128 [108] Young, R.A., Mutchler, C.K., 1969. Soil movement on irregular slopes. *Water Resour.*
1129 *Res.* 5, 1084–1089.
- 1130 [109] Youssef, A.M., Hegab, M.A., 2019. Flood-hazard assessment modeling using
1131 multicriteria analysis and GIS: a case study—Ras Gharib area, Egypt, in: *Spatial Modeling in*
1132 *GIS and R for Earth and Environmental Sciences*. Elsevier, pp. 229–257.
- 1133 [110] Youssef, A.M., Pradhan, B., Sefry, S.A., 2016. Flash flood susceptibility assessment
1134 in Jeddah city (Kingdom of Saudi Arabia) using bivariate and multivariate statistical models.
1135 *Environ. Earth Sci.* 75, 12.
- 1136 [111] Zahedi, P., Parvandeh, S., Asgharpour, A., McLaury, B.S., Shirazi, S.A., McKinney, B.A.,
1137 2018. Random forest regression prediction of solid particle Erosion in elbows. *Powder*
1138 *Technol.* 338, 983–992.
- 1139 [112] Zenggang, X., Zhiwen, T., Xiaowen, C., Xue-min, Z., Kaibin, Z., Conghuan, Y., 2021.
1140 Research on image retrieval algorithm based on combination of color and shape features. *J.*
1141 *Signal Process. Syst.* 93, 139–146.
- 1142

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Highlights (3 to 5 bullet points (maximum 85 characters including spaces per bullet point)

Highlights

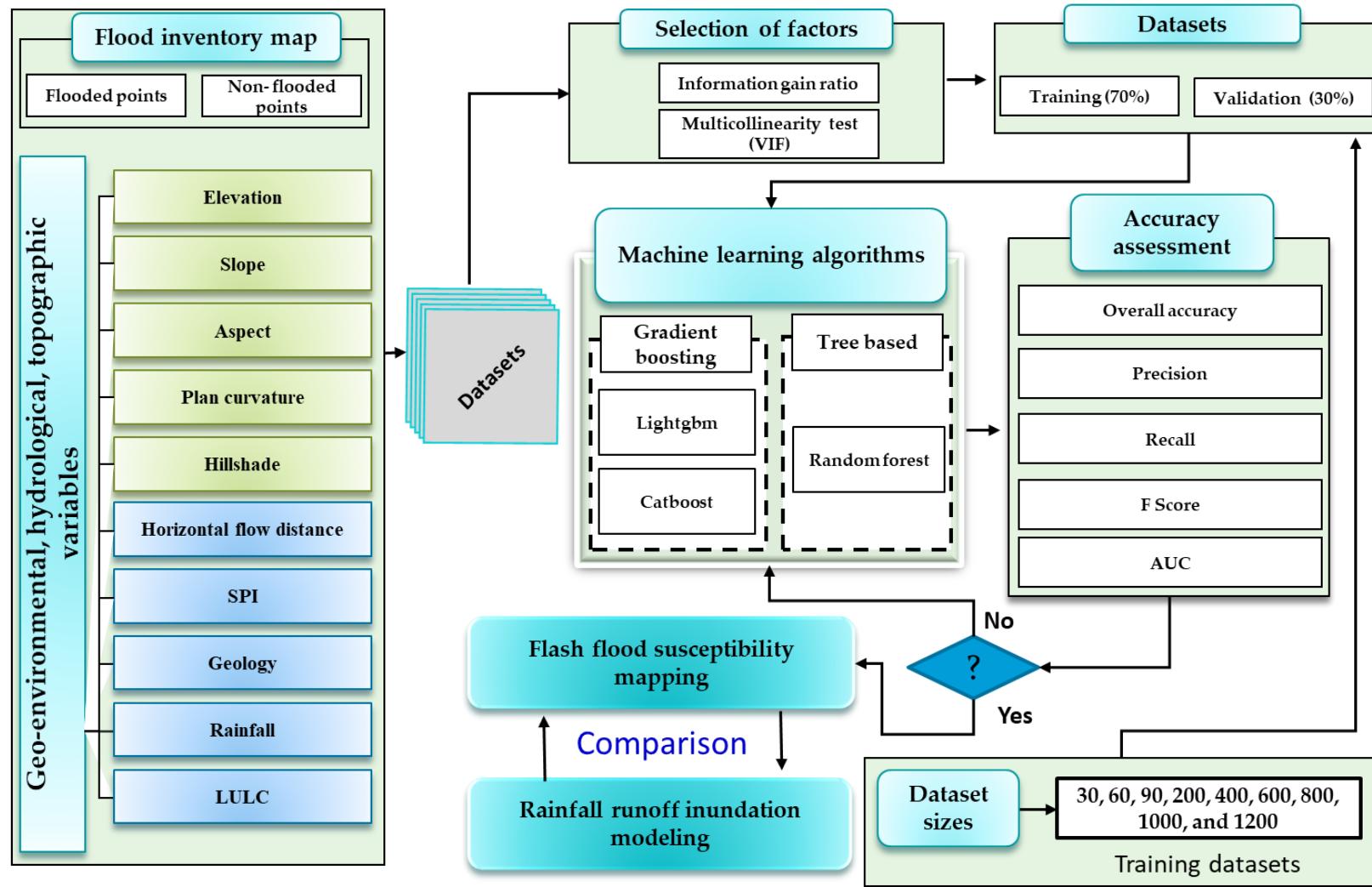
- Random forest, LightGBM, and CatBoost machine learning models were applied to predict flood susceptibility in a humid area that experienced successive extreme typhoons
- LightGBM and CatBoost were tested for the first time in this specific climatic region, with acceptable agreement with the RRI model
- Different training datasets were examined to find the lowest acceptable number of observations for flood susceptibility in machine learning
- The FSMs showed that downstream areas with high residential and agricultural activities were highly susceptible to flooding

Abstract

Vietnam has experienced many natural disasters, particularly typhoons. This study aims to examine three machine learning (ML) approaches—random forest (RF), LightGBM, and CatBoost—for flooding susceptibility maps (FSMs) in the Vu Gia-Thu Bon (VGTB) River Basin of Vietnam. The results of ML are compared with those of the rainfall–runoff model, and different training dataset sizes are utilized in the performance assessment. Ten independent factors that influence the FSMs in the study area, namely, aspect, rainfall, curvature, DEM, horizontal distance from the river, geology, hillshade, land use, slope, and stream power index, are assessed. An inventory map that includes approximately 850 flooding sites is considered based on several post-flood surveys after the typhoons in 1999, 2006, 2007, 2009, 2013, and 2020. The inventory dataset is randomly divided into two sets: training (70%), and testing (30%). The AUC-ROC results are 97.9%, 99.5%, 99.5% for CatBoost, LightGBM, and RF, respectively. The FSMs developed by the ML methods show good agreement with flood inundation mapping developed using the rainfall-runoff model. The FSMs show that downstream areas (both urbanized and agricultural) are under “high” and “very high” levels of susceptibility. Additionally, different sizes of the input datasets (i.e., 30, 60, 90, 200, 400, 600, 800, 1000, and 1250 data points) are tested to determine the least number of data points having an acceptable reliability. The results show that the ML methods can reasonably predict FSMs, regardless of the number of training samples, although the final FSMs show some spatial differences when changes in susceptibility level are seen. The developed FSMs for such typhoon-prone regions can be used by decision-makers and planners in Vietnam to propose effective mitigation measures for community resilience and development.

Keywords

Machine learning, random forest, LightGBM, CatBoost, flooding susceptibility mapping, rainfall-runoff inundation model.



Graphical Abstract

4/26/2022

Prof. Andras Bardossy
Editor
Journal of Hydrology

Dear Prof. Andras Bardossy

I wish to submit an original article for publication in Journal of hydrology, titled “Machine Learning Techniques as an Alternative Approach to Rainfall-Runoff Inundation Models for Flood Susceptibility Prediction”.

This study elaborates on the machine learning (ML) techniques that can be used to predict flood susceptibility as an alternative for the rainfall-runoff inundation model (RRI). Two machine learning methods (CatBoost and LightGBM) along with conventional Random forest (RF) were examined for the first time in a typhoon-prone area in humid environment in Vietnam. We believe that our study makes a significant contribution to the literature because it is the first study that addresses this application with varying sizes of data sets and analyses the affecting on the model accuracy of the used ML techniques, as well as comparing with RRI model.

Further, we believe that this paper will be of interest to the readership of your journal because the employed two ML methods (CatBoost and LightGBM) superior the conventional Method of RF in terms of accuracy and processing speed. Additionally, the flood susceptibility maps developed by the ML techniques agree with the flood inundation map developed by RRI Model.

This manuscript has not been published or presented elsewhere in part or in entirety and is not under consideration by another journal. We have read and understood your journal’s policies, and we believe that neither the manuscript nor the study violates any of these. There are no conflicts of interest to declare.

Thank you for your consideration. I look forward to hearing from you.

Sincerely,
Mohamed Saber, on behalf of the coauthors
Disaster Prevention Research Institute
Kyoto University
Goka-sho, Uji City, Kyoto 611-0011, Japan
Tel.: +81-70-3600-6556
Email: mohamedmd.saber.3u@kyoto-u.ac.jp