

1 **Prediction of volume of shallow landslides due to rainfall using data-driven models**

2

3 Tuganishuri Jérémie<sup>1</sup>, Chan-Young Yune<sup>2</sup>, Gihong Kim<sup>3</sup>, Seung Woo Lee<sup>4</sup>, Manik Das Adhikari<sup>5</sup>,  
4 Sang-Guk Yum<sup>6\*</sup>

5 Department of Civil and Environmental Engineering, Gangneung-Wonju National University,

6 \*Corresponding author: Sang-Guk Yum; [skyeom0401@gwnu.ac.kr](mailto:skyeom0401@gwnu.ac.kr)

7

8 **Abstract**

9 Landslides due to rainfall are among the most destructive natural disasters that cause property  
10 damages, huge financial losses, and human deaths in different parts of the World. To plan for  
11 mitigation and resilience, the prediction of the volume of rainfall-induced landslides is essential to  
12 understand the relationship between the volume of soil materials debris and their associated  
13 predictors. Objectives of this research are to construct a model by utilizing advanced data-driven  
14 algorithms (i.e., ordinary least square or Linear regression (OLS), random forest (RF), support  
15 vector machine (SVM), extreme gradient boosting (EGB), generalized linear model (GLM),  
16 decision tree (DT), deep neural network (DNN), *k*-nearest neighbor (KNN) and Ridge regression  
17 (RR)) for the prediction of the volume of landslides due to rainfall considering geological,  
18 geomorphological, and environmental conditions. Models were trained and tested on the Korean  
19 landslide dataset to obtain the most efficient predictions. The EGB predictions exhibited optimal  
20 predictions with the highest coefficient of determination ( $R^2=0.8841$ ) and lowest mean absolute  
21 error (MAE=146.6120 m<sup>3</sup>), followed by RF ( $R^2=0.8435$ , MAE=330.4876 m<sup>3</sup>) for the holdout set.  
22 The results indicated that the DNN, EGB, and RF models exhibited  $R^2>0.8$  on both the training  
23 and test sets. The difference in coefficient of determination  $R^2$  on the training and holdout set were  
24 1.75, 7.72, and 12.17% for RF, EGB and DNN, respectively, signifying that the model could yield  
25 reliable volume estimates in adjacent areas with similar geomorphological and environmental  
26 settings. The volume of landslides was strongly influenced by slope length, maximum hourly  
27 rainfall, slope angle, aspect, and altitude. The anticipated volume of landslides can be important  
28 for land use allocation and efficient landslide risk management.

29

30 **Keywords:** Data-driven models, volume of landslides, optimal predictive model, rainfall, South  
31 Korea

32

### 33 **1. Introduction**

34 Landslides due to rainfall are phenomena that dislocate a mass of soil from its natural position and  
35 slide downward along a slope due to gravity forces. Intense or long-duration rainfall infiltrates the  
36 soil and increases the pore pressure, resulting in soil saturation that leads to slope failure. The  
37 saturated soil becomes weak and loses cohesion, and the slope fails when rainfall crosses a certain  
38 threshold (Bernardie et al., 2014; Martinović et al., 2018; Lee et al., 2021). The heavy rainfall  
39 saturates a slope and triggers a landslide due to the reduction of the soil's shear strength and the  
40 increase of pore water pressure (Tsai and Chen, 2010; Lacerda et al., 2014; Chatra et al., 2019;  
41 Chen et al., 2021; Luino et al., 2022). For example, steep slopes with loose soils and even moderate  
42 rainfall can lead to the displacement of an enormous quantity of soil mass. On the contrary,  
43 in slopes with more stable, cohesive soils, the surface failure might be smaller (Tsai and Chen,  
44 2010). The rainfall quantity and duration influence the volume of the landslides; the higher the  
45 intensity and the longer the duration of rainfall, the larger the resulting surface failure (Chang and  
46 Chiang, 2009; Bernardie et al., 2014; Chen et al., 2017). The landslide occurrences can also be  
47 influenced by human activities that weaken the slope, such as excavation at the slope toe and  
48 loading caused by construction and land use such as agriculture, mining etc. (Rosi et al., 2016).  
49 The rapid urbanization activities in mountainous regions affect the topography through hill cutting,  
50 deforestation and water drainage (Rahman et al., 2017); these activities disturb the slope structure  
51 and change the water flow, which exacerbates the effect of landslides in regions where human  
52 engineering activities are mostly located (Holcombe et al., 2016; Chen et al., 2019). Therefore, to  
53 mitigate landslide-induced risks in the runout regions, estimation of the volume of landslides due  
54 to rainfall (VLDR) plays a crucial role.

55 The quantification of the VLDR is essential for effective risk management (Tacconi  
56 Stefanelli et al., 2020), emergency response, engineering design (Cheung, 2021), economic  
57 assessment and environmental protection (Alcántara-Ayala and Sassa, 2023). With the estimates  
58 of VLDR, the morphologist can update hazard maps (Van Westen, 2000) to reflect the scale of  
59 potential mass movement in various regions to obtain regions with similar likelihood of landslides  
60 of similar soil mass to highlight risk zone levels, i.e., low, moderate and high. These classifications  
61 help engineers to apply appropriate slope stabilization techniques depending on the level of risk (

62 Dahal and Dahal, 2017). Additionally, enhancing the precision of VLDR estimations and  
63 improving the predictive capabilities is essential for understanding and monitoring landscape  
64 evolution. Montgomery (2009) emphasized that the volume of landslides is a key factor in  
65 determining the extent of downstream damage, particularly for large debris flows or rock  
66 avalanches, which can drastically alter the landscape and affect surrounding ecosystems and  
67 infrastructure. Similarly, Korup (2004) further explored the long-term geomorphological effects  
68 of large-volume landslides, highlighting their importance in reshaping mountainous terrains and  
69 influencing sediment transport, which is critical for understanding both immediate and future  
70 landscape changes. However, the existing landslide susceptibility models mostly used for the  
71 identification of regions susceptible to landslides (i.e., landslide zonation) (Kim et al., 2014;  
72 Gutierrez-Martin, 2020; Chen et al., 2021; Li et al., 2022), which are essential in emergency  
73 management because they provide a general overview of zones with a higher probability of  
74 landslide occurrence; however, they do not emphasize the determination of the approximate value  
75 of the volume of failing mass in relation to excessive rainfall events.

76 Numerous researchers used landslide inventory, remote sensing data and numerical  
77 techniques to establish the relationship between landslide geometry and the influencing factors to  
78 determine the landslide volume quantitatively. For example, Saito et al. (2014) studied the  
79 relationship between rainfall-triggered landslides to test whether the volume of landslides across  
80 Japan that occurred between 2001 and 2011 can be directly predicted from rainfall metrics. The  
81 findings revealed that larger landslides occurred when rainfall exceeded certain thresholds, but  
82 there were significant discrepancies between peaks of rainfall metrics and maximum landslide  
83 volumes, and the total rainfall was the suitable predictor of landslides. Dai and Lee (2001)  
84 established the frequency-volume relation for landslides in Hong Kong and noticed that the  
85 relation for shallow landslides above  $4\text{m}^3$  followed the power law. The 12-hour rolling rainfall  
86 contributed most to the prediction of the volume of landslides. Jaboyedoff et al. (2012) contributed  
87 by demonstrating the value of remote sensing technologies such as Light Detection and Ranging  
88 (LiDAR) in conjunction with field data to improve the accuracy of volume estimates and capture  
89 the geomorphological changes associated with landslides. Ju et al. (2023) constructed an area-  
90 volume power law model for the estimation of the volume of landslides using high-resolution  
91 LiDAR data collected between 2010 and 2020 in Hong Kong. The aim was to estimate accurately  
92 the volume of landslides on small-scale landslides. The reliance on localized datasets limits the

93 model's applicability in regions with different geological settings, and the model does not consider  
94 all variabilities of landslide characteristics. Razakova et al. (2020) calculated landslide volume  
95 using remote sensing data to assess the efficiency of aerial photographs in environmental impact  
96 assessment and ground-based measurement. The study did not consider the effect of vegetation  
97 and topography and only focused on a single landslide case, which may be a source of bias due to  
98 differences in soil composition and environmental factors. Hovius et al. (1997) analyzed multiple  
99 sets of aerial photos and frequency-magnitude relations for landslides in New Zealand. The finding  
100 pinpointed that the landslides frequency-magnitude followed power law and infrequent large  
101 magnitude contributed to the landscape change. The study also noticed the importance of soil  
102 composition in the size of the landslides. This work had a limitation due to the reliance on aerial  
103 photos only, which cannot provide accurate measurement in regions of dense forest, and the  
104 climatic conditions, which are landslide triggering factors, were not considered, and this may affect  
105 the generality of the findings. Guzzetti et al. (2008) applied statistical methods on regional  
106 landslide inventories and antecedent rainfall data ranging between 10 min to 35 days. The findings  
107 revealed that the slope angle and soil type significantly influence landslide volume estimates, and  
108 the rainfall intensity is more important than duration. Chatra et al. (2019) applied numerical  
109 methods to study the effect of rainfall duration and intensity on the generation of pore pressure in  
110 the soil; the finding revealed a higher instability in loose soil compared to medium soil slopes.  
111 Huang et al. (2020) introduced a hybrid machine-learning model combining support vector  
112 regression (SVR) with a genetic algorithm to estimate debris-flow volumes. The model was tested  
113 on real-world case studies, showing improved accuracy in volume predictions compared to  
114 traditional methods. However, a notable weakness of the study is its reliance on a limited dataset,  
115 which may reduce the model's generalizability to environmental contexts. Shirzadi et al. (2017)  
116 compared the effectiveness of statistical and machine-learning models in simulating landslide  
117 volumes-areal relations, demonstrating that machine-learning techniques outperform traditional  
118 statistical methods in terms of accuracy. This method did not consider the climatic and geomorphic  
119 factors such as rainfall, vegetation, soil type, etc., triggering and influencing factors for the  
120 landslide occurrence. It was noted that existing models only treated the interaction of soil and  
121 rainfall without considering the environmental factors, human activity, and non-linear behavior of  
122 the triggering and influencing factors.

123 In the present study, the volume of landslides due to rainfall is predicted using OLS, RF,

124 SVM, EGB, GLM, DT, DNN, KNN and RR algorithms, considering the details of triggering  
125 factors (i.e., rainfall) and predisposing factors (i.e., geomorphological, soil and environmental).  
126 Here, we aim to construct a data-driven algorithm that combines input parameters for physical-  
127 based and empirical models and incorporates more complex non-linear features of input variables  
128 to predict the occurrence of associated events more accurately. The main assumption behind the  
129 data-driven algorithm is that the considered feature input of the model produces a similar volume  
130 of landslides due to rainfall and follows the same pattern at a particular region with the same  
131 features under the same quantity of rainfall. Here, we examine different machine learning (ML)  
132 algorithms and compare their performance using the coefficient of determinations ( $R^2$ ), mean  
133 square errors (MAE), Root mean square error (RMSE), Mean absolute percentage error (MAPE),  
134 and symmetric mean absolute percentage errors (SMAPE) of the predicted volume of landslides.  
135 The focus is to optimize the predictions of the volume of landslides due to rainfall, taking into  
136 account triggering and influencing factors with higher accuracy.

137

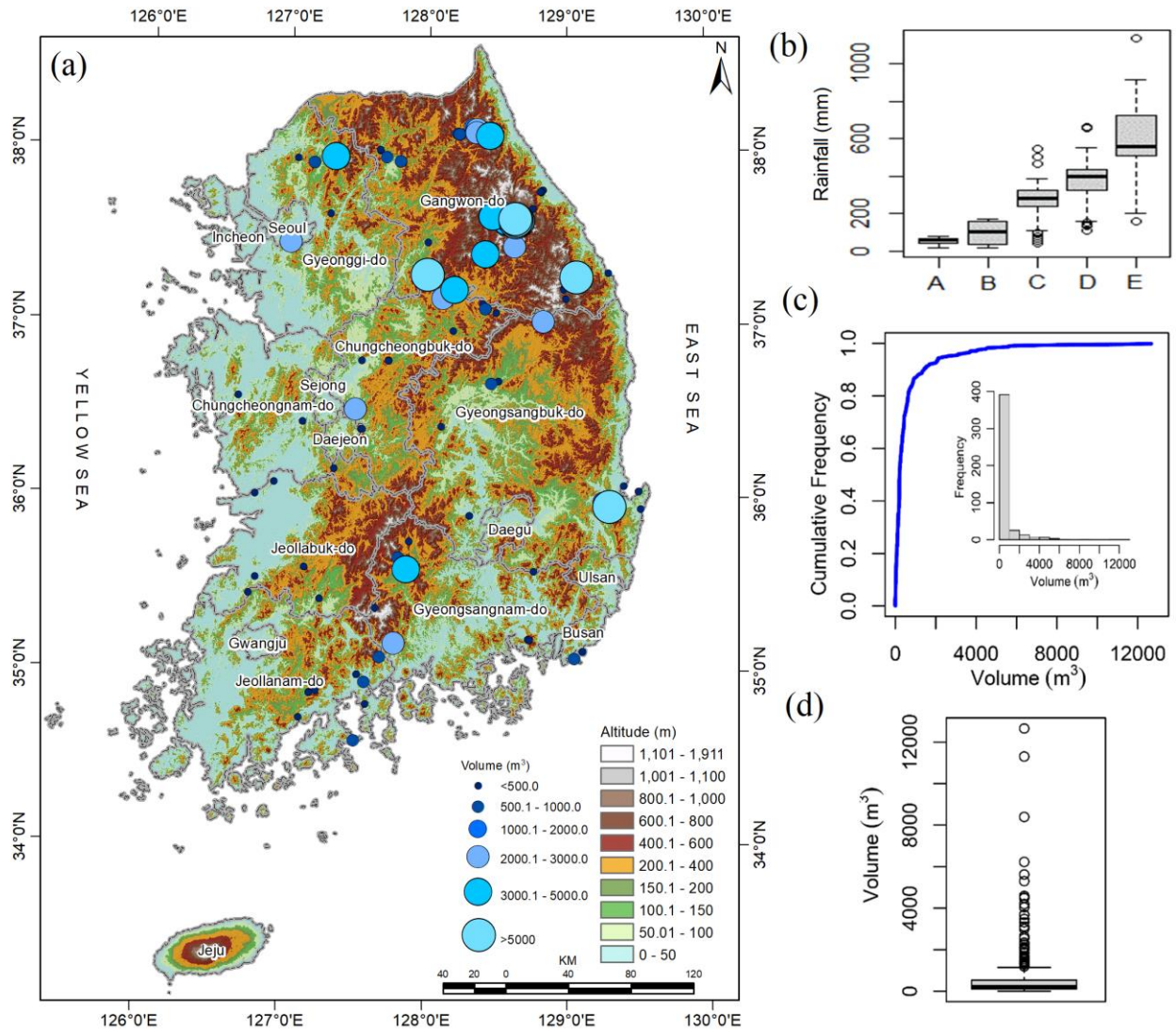
## 138 **2. Data and Study Region**

### 139 ***2.1. Study Region***

140 The region for testing the model is South Korea, characterized by mountainous (63% of total land)  
141 relief, especially in the eastern part of the country (Lee et al., 2022). South Korea is located on the  
142 southern part of the Korean Peninsula, bordered by the Yellow Sea to the west coast and the East  
143 Sea (Sea of Japan) to the East. According to the Korean Meteorological Administration  
144 (<https://www.kma.go.kr/>), the country has a temperate climate characterized by four distinct  
145 seasons: hot and humid summers, cold winters, and springs and falls with moderate temperatures.  
146 The annual rainfall varies between 1000 mm to 1400 mm and 1800 mm for the central region and  
147 southern region, respectively (Jung et al., 2017; Alcantara and Ahn, 2020). During the summer,  
148 heavy rainfall from June to September leads to significant surface runoff, increases landslide risk,  
149 and causes approximately 95% of all landslides each year (Lee et al., 2020; Park and Lee, 2021).  
150 In addition, the landslides may be aggravated by typhoons, which mostly occur in August and  
151 September, and it is anticipated that frequency will increase due to climate change (Kim and Park,  
152 2021). The rainfall trend analysis from 1971 to 2100 predicted an increase in rainfall of 271.23mm,  
153 which indicates the growing risk of landslides associated with climate change (Lee, 2016).  
154 Temperature variations are influenced by its geographical location; the average summer

155 temperatures vary between 25 and 30°C, while winter temperatures can drop to -10°C in some  
156 parts of the country (<https://web.kma.go.kr/>). The South Korean geologically is mainly composed  
157 of granitic and metamorphic rocks, such as gneiss, schist, and granite, which influence the stability  
158 of the landscape (Jung et al., 2024). The geomorphology is characterized by rugged mountains,  
159 river valleys, and coastal plains, with the Taebaek Mountains running along the eastern edge (Kim  
160 et al., 2020). In addition, the influence of rainfall, environmental, geomorphology, and geological  
161 factors increase the vulnerability to landslides across the country, especially in the northeastern  
162 mountainous region, as depicted in Figure 1. The predominant soil types in South Korea include  
163 clay, sandy, and loamy soils, each with different characteristics affecting water infiltration,  
164 retention and erosion (Kang et al., 2022; Lee et al., 2023). Clay soils, being more stable, can  
165 become highly saturated, increasing landslide risk during heavy rains. On the other hand, sandy  
166 soils are more prone to shallow landslides due to fast saturation, leading to instability. Regions  
167 with steep topography and poorly consolidated soil (loose) are mostly at risk, especially after  
168 prolonged rainfalls (Kim et al., 2015).

169         The combination of heavy summer rainfall, geological composition, and geomorphological  
170 factors makes South Korea particularly vulnerable to shallow landslides. Thus, continuous  
171 monitoring and research are vital to understanding the complex interactions between climate,  
172 geology, soil types, and landslide occurrences in this region. Understanding the collective effects  
173 of meteorological, environmental, geological stability, and geomorphological features is crucial  
174 for developing effective disaster management strategies and enhancing public safety in landslide-  
175 prone areas. As climate change continues to impact rainfall patterns, South Korea faces ongoing  
176 challenges in mitigating landslide risks and protecting vulnerable communities.



177  
 178 Figure 1. (a) Spatial distribution of landslides in South Korea, (b) Temporal variation of rainfall,  
 179 i.e., A: Maximum hourly rainfall, B: Four weeks rainfall, C: Three hours rainfall, D:  
 180 Three days rainfall and E: Two weeks rainfall, (c) Cumulative frequency distribution of  
 181 the volume of landslides, and (d) Box plot of the volume of landslides.

182  
 183 **2.2 Data**

184 The landslide inventory dataset contains 455 landslide record information from 2011 to 2012,  
 185 collected from different locations in South Korea by Korean Forest Services. This dataset tabulates  
 186 information on landslide geometry, such as runout length, width, depth, and volume of the affected  
 187 area, along with geomorphological composition, vegetation, and antecedent rainfall prior to  
 188 landslide events. The details regarding landslide predisposing and triggering factors are

189 summarized in Table 1.

190 The majority of landslides in this region were shallow, translational slope failures (Kim et  
 191 al., 2001). The occurred landslides had a volume varying between 1.5m<sup>3</sup> to 12,663m<sup>3</sup> and  
 192 predominantly occurred in the northeastern and southeastern region (Figs.1a,c-d). The occurred  
 193 landslides were hallowed and skewed to the right with 2570.7m<sup>3</sup> as 95<sup>th</sup> quantile, largest volume  
 194 was 12,663m<sup>3</sup>, and the aggregate mass of landslide due to rainfall was 276,986.62m<sup>3</sup>. The  
 195 estimation of the volume of removed material by landslides is important as it helps to assess risks  
 196 the estimated damage can cause down at the toe of the failed slope, such as blocking transportation  
 197 network, burying crops or farmland, the damage-built environment near landslide risks area, and  
 198 post-disaster recovery planning (Evans et al., 2007; Rotaru et al., 2007; Intrieri et al., 2019).

199

200 Table 1. Landslide influencing and triggering factors.

Group	Features	Feature Relevance	References
Vegetation	Fire history	The burning of the vegetation intensifies the mass movement of soil near the uncovered burned stem of trees and free movement on uncovered soil due to post-fire rainfall and storms. The sliding may also be due to loss of vegetation and altered soil property and structure, which lead to soil degradation and infiltration, which increase pore pressure, and change in hydrology by concentrating water flow in places that exacerbate landslides.	Highland and Bobrowsky, 2008; Stoof et al., 2012; Hyde et al., 2016; Culler et al., 2021
	Age of tree	Mature forests have more resistance to shallow landslides due to highly developed roots, which improve soil cohesion and leaves that prevent direct contact of raindrops with the soil surface.	Sato et al., 2023; Lann et al., 2024
	Forest density	The presence of forest reduces the likelihood of landslides about three times compared to grassland. Grassland has been revealed to be three times more vulnerable to shallow landslides than broadleaf, coniferous, and secondary forests.	Greenwood et al., 2004; Turner et al., 2010; Scheidl et al., 2020; Asada and Minagawa, 2023; Lann et al., 2024



Group	Features	Feature Relevance	References
	Timber diameter (m)	Tree spacing and size were used to investigate the effect of root and tree in shallow landslide control. High root density generally enhances slope stability, and specific tree placement and root sizes between 5 to 20 mm effectively prevent landslides.	Wang et al., 2016; Cohen and Schwarz, 2017
Geomorphology	Drainage	The drainage significantly affects slope stability and promotes efficient control of rainfall's influence on groundwater fluctuation. The presence of drainage increases the threshold of landslides due to rainfall.	Korup et al., 2007; Sun et al., 2010; Yan et al., 2019; Wei et al., 2019
	Slope angle (°)	The steeper slopes have a lower presence of landslides due to the low transportable materials. Slopes between 20-40 degrees are most vulnerable to greater landslides as rainfall intensity and duration increase. Generally, the average angle of the terrain at the landslide location provides valuable insight into the region's overall steepness and geomorphic characteristics, which are crucial factors influencing landslide susceptibility and risk modeling.	Donnarumma et al., 2013; Duc, 2013; Qiu et al., 2016
	Slope aspect	The effect of rainfall on slope differs by slope angle and slope aspect, which leads to unevenly distributed landslides.	Panday and Dong, 2021; Cellek, 2021
	Slope length (m)	The volume increases as the slope length increases. A complex interplay exists between rainfall, length of slope and slope angle in the occurrence of landslides.	Turner et al., 2010
	Soil depth (m)	Soil properties, depth, and texture have significant differences in infiltration rates, which have different influences on the occurrence of landslides.	Kitutu et al., 2009; McKenna et al., 2012
	Soil type	Soil types, namely, Sandy loam, silt loam and loam, with their coefficient of permeability 1.7, 1.65 and 1.5, respectively, retain water differently, leading to different saturation	Chen et al., 2015a; Liu et al., 2021a

Group	Features	Feature Relevance	References
		times. The soil with higher permeability tends to drain water more efficiently, making it less prone to saturation. In contrast, the soil with lower permeability, the pore pressure rapidly increases, which leads to shallow landslide initiation during intense rainfall events.	
Location	Altitude	Regional variability of elevation and mountain steepness affect the quantity of rainfall and associated landslides.	Um et al., 2010; Hyun et al, 2010; Yoon and Bae, 2013; Park, 2015
	Maximum hourly rainfall	The rainfall infiltrates the slope and increases pore water pressure, which reduces soil shear strength and leads to soil saturation, that causes surface failure.	Wieczorek, 1987; Dai and Lee, 2001; Smith et al., 2023
Rainfall	Continuous rainfall	Sudden intense rainfall concentrated in short periods is responsible for shallow landslides and debris flow.	Zhang et al., 2019
	Three hours rainfall		
	Three days rainfall	The antecedent rainfalls increase moisture in the soil and weaken soil cohesion.	Bernardie et al., 2014; Chen et al., 2015a; Gariano et al., 2017; Zhang et al., 2019; Ran et al., 2022
	Two weeks rainfall		
Four weeks rainfall			

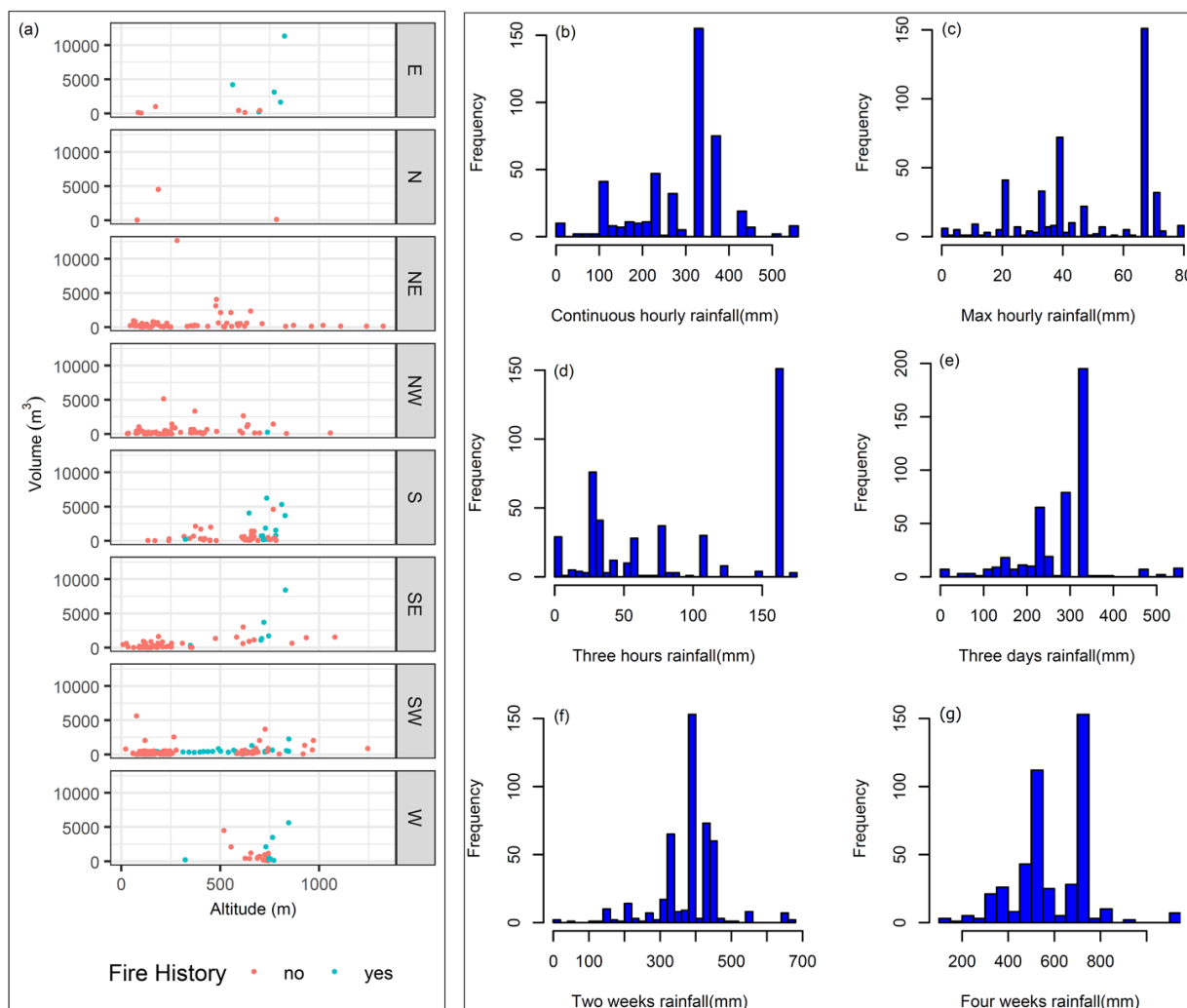
201

202           Location parameters such as altitude, latitude and longitude are essential elements that  
203 determine the microclimate of a given region, influencing rainfall patterns (Hyun et al., 2010; Yoon  
204 and Bae, 2013; Park, 2015). The northeastern region is characterized by high-elevation terrain,  
205 such as the Taebaek and Sobaek ranges, which dry air and lead to orographic precipitation (Yun et  
206 al., 2009). The windward mountain versants receive a substantial amount of rainfall, which can  
207 increase the likelihood of landslides (Jin et al., 2022). This variation of rainfall with respect to the  
208 direction highlights the importance of including slope aspect variables in landslide studies (Kunz  
209 and Kottmeier, 2006). Figure 2(a) depicts the relationship between the slope aspect and the volume  
210 of landslides and slope aspect, altitude and fire history and shows that larger volumes were  
211 localized in regions that faced forest fire and altitudes between 500 and 1000m. Additionally, the

212 topographical features such as slope length and slope angle affect the size of the landslide (Panday  
213 and Dong, 2021), slope failure due to over-saturation from groundwater and rainfall infiltration  
214 that destabilize the slope (Kafle et al., 2022). Furthermore, slope length, slope angle and slope  
215 aspect play an important role in the determination of the volume of geological material uprooted  
216 by landslides (Zaruba and Mencl, 2014; Khan et al., 2021). The slope stability depends on soil  
217 composition properties, including soil permeability indices that affect water infiltration and  
218 saturation level (Chen et al., 2015a). In the study regions, three main soil types, namely, sandy  
219 loam, loam, and silt loam, were observed, and their coefficient of permeability is 1.7, 1.65 and 1.5,  
220 respectively (Lee et al., 2013). Moreover, to reduce the infiltration drainage network that  
221 channeling rainwater terrain drains soil and reduces the saturation, which minimizes the likelihood  
222 of landslide occurrence as a result of groundwater discharge and rainfall water flow (Hovius et al.,  
223 1997; Wei et al., 2019). Furthermore, the vegetation protects the topsoil from the direct impact of  
224 raindrops hitting the ground, which causes erosion due to the force of gravity and reduces  
225 infiltration (Omwega, 1989; Keefer, 2000). The absence of vegetation allows rainwater to seep  
226 away fine topsoil, causing shallow landslides (Gonzalez-Ollauri and Mickovski, 2017). On the  
227 contrary, vegetation improves soil cohesion and prevents potential shallow landslides due to soil-  
228 root interaction (Gong et al., 2021; Phillips et al., 2021). The density of vegetation (forest) and  
229 leafage type (broad, pines or mixture) directly affects the quantity of raindrops intercepted and  
230 prevented from directly hitting the soil, which emphasizes the contributions of vegetation in the  
231 landslides mitigation. Further, the occurrence of forest fires can contribute to the occurrence of  
232 landslides due to the burning of vegetation covering the area, changing soil properties and  
233 increasing soil pH (Lee et al., 2013).

234         The rainfall, a triggering factor of landslides, is the immediate cause of slope instability  
235 and failure due to infiltration that leads to saturation resulting from increased pore water pressure  
236 that reduces soil shear strength (Yune et al., 2010; Khan et al., 2012; Kim et al., 2021; Lee et al.,  
237 2021). The antecedent rainfall increases the moisture in the soil, which accelerates the soil  
238 saturation; the cumulative effect is essential to understand the saturation levels (Ran et al., 2022).  
239 In this study, rainfall variables are grouped based on time, namely, continuous rainfall, which is  
240 the accumulative value of rainfall on the day of a landslide from rainfall start hour to the landslide  
241 event, maximum hourly rainfall, rainfall during the fixed period such as three hours, one day, three  
242 days, two weeks etc. (Fig. 1b). The histograms for rainfall considered in this study are depicted in

243 Figure 2(b-g). The descriptive statistics for all continuous variables are illustrated in Table 2.



244  
 245 Figure 2. (a) The scatter plot showing the variation of landslide volumes with respect to slope  
 246 aspect, fire history and altitude, and (b-g) Histograms of rainfall distribution.

247  
 248 Table 2. Summary statistics for continuous variables.

Variables	Units	N	Min	Mean	Median	Max	Std dev
Max Hourly rain	mm	455	0	48	48	78	20
Continuous rainfall	mm	455	0	285	327	550	106
Three hours rainfall	mm	455	0	88	80	171	60
Twelve Hours rainfall	mm	455	0	150	99	447	95
One day rainfall	mm	455	0	202	162	538	112

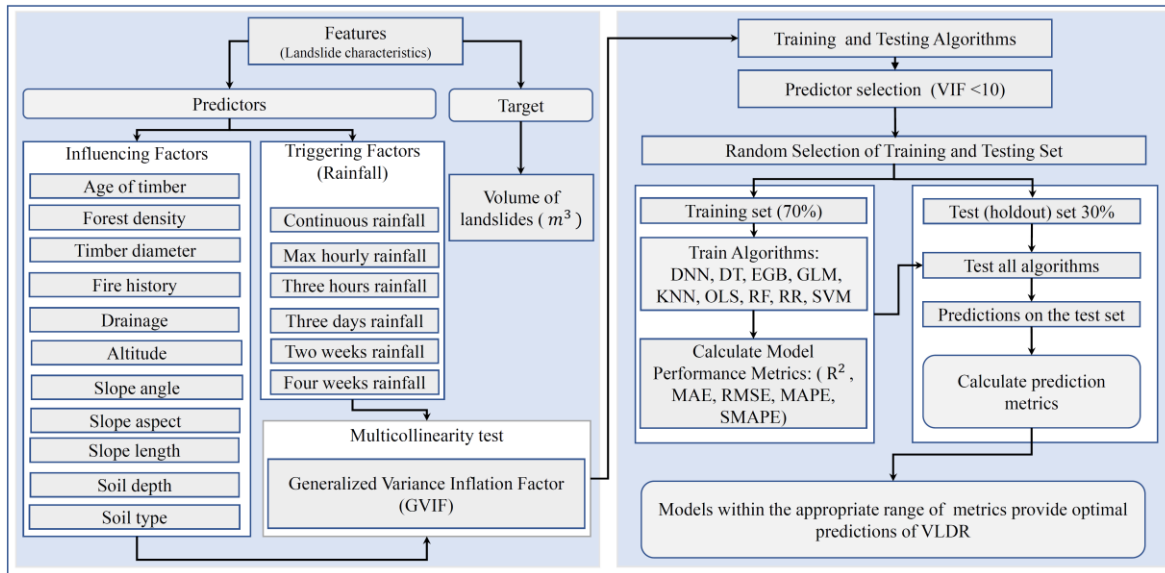
<b>Variables</b>	<b>Units</b>	<b>N</b>	<b>Min</b>	<b>Mean</b>	<b>Median</b>	<b>Max</b>	<b>Std dev</b>
Three days rain	mm	455	0	280	284	550	86
Seven days rain	mm	455	0.5	323	330	634	88
Two weeks rain	mm	455	0.5	385	400	663	90
Three weeks rain	mm	455	86	504	533	914	115
Four weeks rain	mm	455	108	587	561	1135	160
Soil depth	m	455	0.2	0.6	0.75	0.75	0.19
Soil type	-	455	1.5	1.6	1.5	1.7	0.087
Timber diameter	m	455	0.15	0.27	0.23	0.35	0.086
Age of tree	Years	455	10	34	35	60	14
Slope length	m	455	1.8	21	13	180	23
Slope angle	Degree (°)	455	10	34	34	65	7.9
Altitude	m	455	9	391	272	1324	273

249

### 250 **3. Methods**

251 In this paper, we consider nine data-driven models, namely OLS, RF, SVM, EGB, GLM, DT,  
252 DNN, KNN and RR, to predict the volume of landslides due to rainfall. The model is tested on the  
253 South Korean landslides inventories and predisposing factors coupled with triggering factors, i.e.,  
254 rainfall data. The detailed workflow is summarized in Figure 3. The steps for construction of these  
255 models can be briefly summarized as follows: a) the dataset for landslide inventories is cleaned  
256 and combined with rainfall dataset, b) the collinearity analysis is made using variance inflation  
257 factor, c) continuous feature are scaled (Z-score) (Bonamutial and Prasetyo, 2023) to facilitate  
258 algorithms to converge fast, d) the dataset is split into training and test set, e) all models are tested  
259 on the same training set, and the model evaluation on the test set using mean absolute error (MAE),  
260 coefficient of determination ( $R^2$ ), root mean square error (RMSE), symmetric mean absolute  
261 percentage error (SMAPE) and mean absolute percentage error (MAPE) for the comparison of  
262 actual and predicted volume by each model, f) variable importance is calculated for the optimal  
263 model, and g) the distance correlation is calculated for each continuous feature, and Kruskal-Wallis  
264 and Dunn test are conducted to examine the similarity of the effect of each category on the  
265 landslide volume.

266



267

268 Figure 3. Workflow for the prediction of the volume of landslides due to rainfall.

269

270 **3.1 Model Construction**

271 In the present investigation, we aimed to predict landslide volume using models that minimize  
 272 error with interpretability and scalability. Since one model can not have all properties  
 273 simultaneously, we selected some widely used models due to their inherent interpretability and  
 274 scalability properties. The OLS, GLM, and DT were widely used for their high interpretability,  
 275 which helps to understand the influence of individual features on predictions (Gelman, 2007;  
 276 Breiman, 2017). On the other hand, the EGB, RF, SVM, RR, and KNN were used due to their  
 277 robust performance in capturing complex patterns in data, which is essential for accurate  
 278 predictions of landslide volumes (Liaw and Wiener, 2002; Hastie, 2009; Chen and Guestrin, 2016).  
 279 Additionally, considering that the model will be used on a regional scale, which will require big  
 280 data, the EGB, RF, and DNN are designed to efficiently handle large datasets, making them  
 281 suitable for the regional scale analysis. These last models can be scaled to incorporate more data  
 282 from different geographical areas without significant adjustments, enhancing their applicability in  
 283 future research (Krizhevsky et al., 2012). Accordingly, nine data-driven methods were selected and  
 284 tested on a Korean dataset to predict VLDR.

285 The first considered method is OLS, which is applied to estimate parameters of multilinear  
 286 regression that yield the minimum residual sum of squares errors from the data (Kotsakis, 2023)  
 287 under assumptions of no correlation in independent variables and error term, constant variance in

288 error terms, non-linear collinearity of predictors, and normal distribution of error terms. The RF-  
289 regression is a supervised data-driven technique based on ensemble learning, which constructs  
290 many decision trees during the training time of a model by combining multiple decision trees to  
291 produce an improved overall result of the model outcome. The RF-regression is more efficient in  
292 the analysis of multidimensional datasets (Borup et al., 2023). RF is an effective predictive model  
293 due to non-overfitting characteristics based on the law of large numbers (Breiman, 2001). The DT  
294 regression is a predictive modeling technique in the form of a flowchart-like tree structure that  
295 includes all possible results, output, predictor costs, and utility. The DT simplifies the decision-  
296 making due to its algorithm that mimics human brain decision-making patterns (Rathore and  
297 Kumar, 2016). The KNN technique draws an imaginary boundary in which prediction outcomes  
298 are allocated as the average of  $k$ -nearest point predictors and averaging their output variable  
299 (response). The KNN calculates Euclidian distances to identify the likeness between datapoints,  
300 and then it groups points that have smaller distances between them (Kramer and Kramer, 2013).  
301 The RR is an improved form of ordinary least squares, which serves to respond to cases where  
302 collinearity is found in predictor variables. The estimated coefficients of ridge are biased  
303 estimators of true coefficients and are generated after adding a penalty on the OLS model. The RR  
304 has always lower variances compared to OLS (Saleh et al., 2019). The advantage of the GLM over  
305 OLS is that the dependent variable need not follow the normal distribution. The GLM is composed  
306 by random and systematic components and the link function that links the two. In this study, the  
307 GLM with Gaussian link function was applied. GLM is fitted using maximum likelihood  
308 estimation (Dobson and Barnett, 2018). The DNN is among data-driven models that revolutionized  
309 different fields; the DNN learns via multi-processing layers and identifies intricate patterns in the  
310 data to predict the outcome (LeCun et al., 2015). Here, the backpropagation algorithm was used to  
311 predict the estimated outcome. The advantage of DNN is that it can discover the complex structures  
312 in the data using a back propagation algorithm capable of changing the internal parameter (weight  
313 update). The SVM is popular for balanced predictive performance which makes it capable to train  
314 model on small sample size (Pisner and Schnyer, 2020). Subsequently, SVM has been applied in  
315 many different landslide studies (Pham et al., 2018; Miao et al., 2018). SVM methods identify the  
316 optimal hyperplane in multidimensional space that separates different groups in the output values.  
317 The EGB is the most powerful and leading supervised machine learning method in solving  
318 regression problems. It can perform parallel processing on Windows and Linux (Chen et al.,

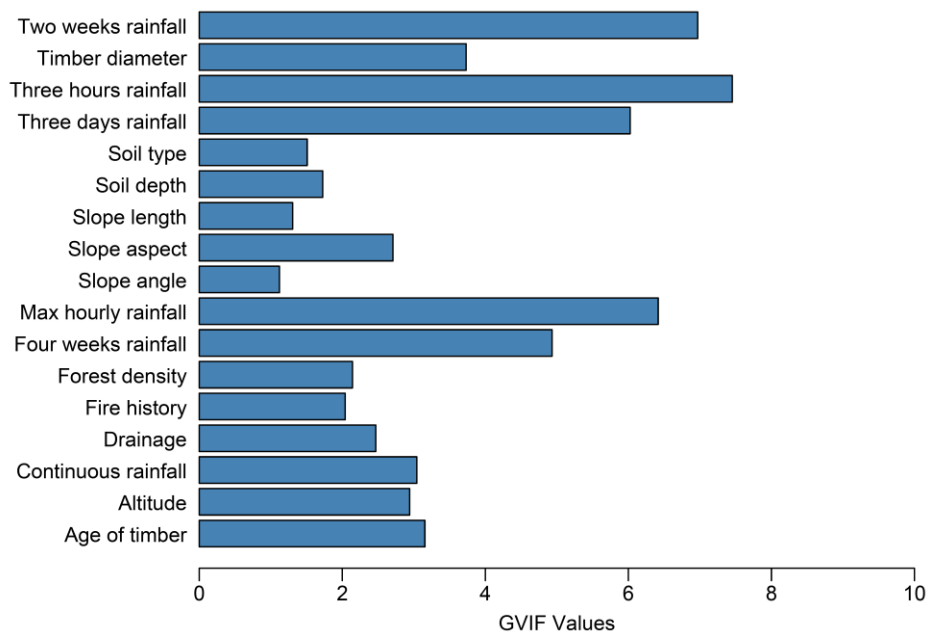
2015b). The gradient boosting trains of differentiable loss function, and the model fits when the gradient is minimized. In this paper, both traditional statistical predictive models and ML models were used. The firsts are known for high clarity and explainability, and the second is famous for handling non-linearity in features. In some cases, the performance of advanced data-driven algorithms is almost similar (Chowdhury et al., 2023).

324

### 3.2 Feature Selection and Data Splitting

The variable selection procedure was based on previous literature and applied in the model using generalized variance inflation factor (GVIF) (O'Brien, 2007) to eliminate collinear variables. The variable with  $GVIF < 10$  was considered non-collinear and used in the model. Figure 4 depicts retained features and corresponding GVIF values. The retained features have GVIF less than 10 (O'brien, 2007). Accordingly, all depicted variables were considered for the model training. Further, to train the model, the datasets were split randomly, with 70% of the data for the training set and 30% for testing (Nguyen et al., 2021). The 10-fold cross-validation was performed to obtain an optimal model. The training and test set was scaled (Z-score or variance stability scaling) to solve convergence issues that are associated with running the model without feature scaling (Singh and Singh, 2022). To run the model on the data using driven methods that accept numerical features only, the test and training set was one-hot-encoded to create a feature matrix (Seger, 2018).

337



338

339

340 Figure 4. Generalized Variance Inflation Factor (GVIF) bar plot for features.



341 **3.3 Model Evaluation Metrics**

342 The model performance evaluation is a process of quantifying the difference between the  
 343 observed value not used in the modeling process and the predicted value by the model. Different  
 344 metrics are applied depending on the type of task, whether it is a classification or a regression  
 345 problem. Subsequently, the widely used evaluation metrics for regression models, namely,  $R^2$ ,  
 346 MAE, RMSE, MAPE and SMAPE, were utilized to evaluate the model performances. The metric  
 347 formulae and evaluation criteria are summarized in Table 3.

348

349 Table 3. Model evaluation metrics.

Metrics	Evaluation	References
$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$	<ul style="list-style-type: none"> <li>• Measures the square root of the average squared differences between predicted and actual values.</li> <li>• Lower values indicate better model performance.</li> </ul>	Hyndman and Koehler, 2006
$MAE = \frac{1}{n} \sum_{i=1}^n  y_i - \hat{y}_i $	<ul style="list-style-type: none"> <li>• The average of the absolute differences between predicted and actual values.</li> <li>• Lower values indicate better model performance.</li> </ul>	Willmott and Matsuura, 2005
$MAPE = \frac{100}{n} \sum_{i=1}^n \left  \frac{y_i - \hat{y}_i}{y_i} \right $	<ul style="list-style-type: none"> <li>• Measures the accuracy of a model as a percentage, which can be more interpretable.</li> <li>• Lower values indicate better model performance.</li> </ul>	Armstrong, 2001
$SMAPE = \frac{100}{n} \sum_{i=1}^n \frac{ y_i - \hat{y}_i }{ y_i  +  \hat{y}_i }$	<ul style="list-style-type: none"> <li>• Unlike MAPE, which can be skewed by very small actual values, SMAPE accounts for both the actual and predicted values, making it symmetric.</li> <li>• SMAPE is expressed as a percentage</li> <li>• Mitigates the impact of small actual values on the error metric, providing a more balanced assessment.</li> <li>• Lower values indicate better model performance.</li> </ul>	Hyndman and Koehler, 2006

Metrics	Evaluation	References
$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$	<ul style="list-style-type: none"> <li>• Represents the proportion of variance in the dependent variable that can be explained by the independent variables.</li> <li>• Values closer to 1 indicate a better fit</li> </ul>	Darlington, 1990; Chicco et al., 2021

350 \* $y_i$  and  $\hat{y}_i$  representing the actual and predicted value and,  $\bar{y}$  and  $n$  standing for the mean of actual value and number  
351 of observations in the dataset, respectively.

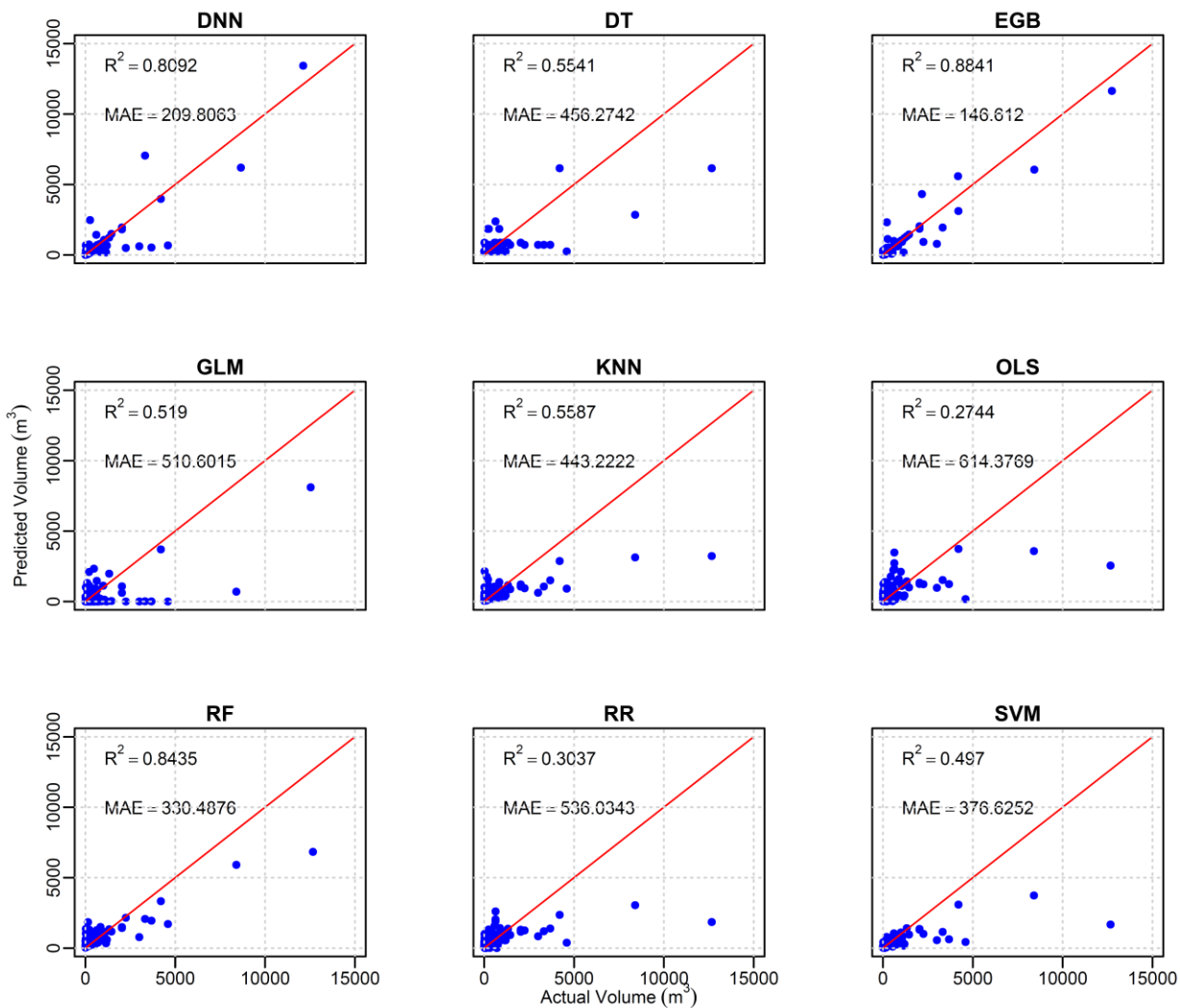
352

#### 353 4. Results

354 The model was developed in R with different libraries, as discussed below. The DNN  
355 regression model was constructed using `dnn()` function from the `cito` library (Amesoeuder et al.,  
356 2023), with two hidden layers of (50, 50) nodes. The model was trained on 1500L epochs, learning  
357 rate ( $lr = 0.01$ ), and  $loss = "mae"$ . The DT regression model was constructed with `tree()` function  
358 from the `tree` library, with the recursive-partition method. The RR model was constructed using  
359 `glmnet()` from the `glmnet` package (Friedman et al., 2010), with ridge penalty ( $alpha=0$ ). The  
360 optimal lambda was obtained by performing 10-fold cross-validation. The EGB model was built  
361 using `xgboost()` function in `xgboost` package (Chen et al., 2022). The optimal model was obtained  
362 at 524<sup>th</sup> boosting iteration with max depth =5 and other parameters set to default. The GLM  
363 regression model was constructed using `glm()` function (R core Team, 2022) with family Gaussian  
364 and log link to constrain the model of predicting positive outcomes. The KNN regression was  
365 constructed using `knnreg()` function from the `caret` package (Kuhn, 2022), with number of  
366 neighbors,  $k=17$ . The OLS model was constructed `lm()` from the `stats` package (R core Team,  
367 2022). The RF model was run using `randomForest()` from the `randomforest` package (Liaw and  
368 Wiener, 2002) with default parameters and the optimal model was reached at 256<sup>th</sup> iteration. The  
369 SVM regression model with linear kernel was built using `e1071` package (Meyer et al., 2021) and  
370 other parameters set to default.

371 The predictive performance of all tested models on the holdout dataset is depicted by the  
372 scatterplot (Fig. 5) of actual volume as recorded in the test set and predicted outcome values of  
373 each model. The red line represents the perfect prediction. The scatter plot of actual and predicted  
374 values of tested models shows that OLS performed least compared to other models with  
375  $R^2=0.2744$ , that is, 27% of variances in the model were explained by predictors. The second least  
376 performing was the RR with  $R^2= 0.3034$ , which is 3.6% improvement compared to OLS. Among  
377 all models, three out of nine, namely, OLS, SVM, and RR, performed below 50%; however, these

378 models predicted well small values of volume (below 2000m<sup>3</sup>). The MAE of these three models  
 379 was higher than the remaining six models, namely DNN, DT, GLM, KNN, RF, and EGB. Among  
 380 these lasts, the most performing was EGB with R<sup>2</sup>= 0.88 of variance explained by predictors and  
 381 MAE=146.6 m<sup>3</sup>. The evaluation metrics for the training and tested models are summarized in Table  
 382 4. Considering the R<sup>2</sup>, the three models, namely EGB, RF, and DNN, had a value of R<sup>2</sup> above 80%  
 383 on the holdout set.



384  
 385 Figure 5. Scatterplot of actual and predicted values for the nine tested models.

386  
 387 Regarding the prediction on the training set, the GLM had an R<sup>2</sup> of 83%. Nevertheless, the  
 388 prediction on the holdout set was 51.9%; this large variation in variance explained by predictors  
 389 indicates that the GLM model did not catch all non-linear patterns in the holdout set. Notably, the  
 390 prediction difference in R<sup>2</sup> on both training and test for the random forest exhibited a very small

391 difference compared to EGB and DNN, that is, 1.75% compared to 12.17% and 7.72% for DNN  
 392 and EGB, respectively. Despite the stable prediction of RF, the performance in terms of SMAPE,  
 393 the DNN was the second lowest symmetric mean absolute percentage error, 43.83m<sup>3</sup> and 39.79 m<sup>3</sup>  
 394 on training and test sets, respectively. According to Chicco et al. (2021), the R<sup>2</sup> is more informative  
 395 in regression modeling; thus, RF had better predictions than the DNN.

396

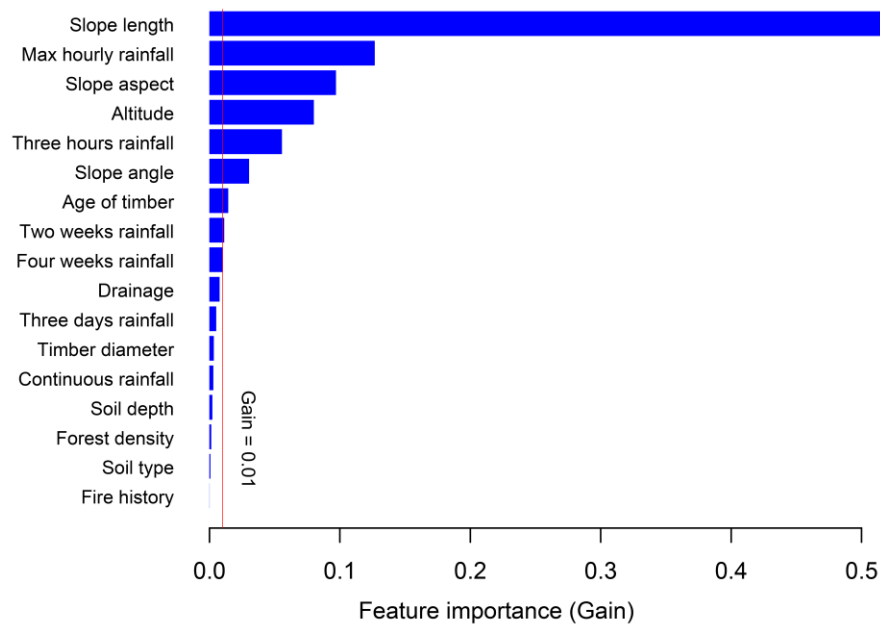
397 Table 4. Summary of prediction metrics for tested models on the training and test set.

Metrics		Models								
		DNN	DT	EGB	GLM	KNN	OLS	RF	RR	SVM
R <sup>2</sup>	Train	0.9309	0.4514	0.9613	0.8380	0.3470	0.3775	0.8610	0.3382	0.5510
	Test	0.8092	0.5822	0.8841	0.5190	0.5587	0.2744	0.8435	0.3037	0.4970
MAE	Train	132.7429	407.0814	75.1250	308.9700	410.2945	502.0053	236.9516	470.1633	276.2000
	Test	209.8063	435.5836	146.6120	510.6015	443.2222	614.3769	330.4876	536.0343	376.6252
RMSE	Train	348.6190	940.4850	113.4940	570.0070	1027.3730	1001.7620	574.9720	1042.9110	916.5471
	Test	646.5438	1047.4880	501.8960	1055.9190	1115.5270	1234.1220	737.0857	1237.9420	1176.9410
MAPE	Train	0.5240	0.7930	0.1540	76.3530	0.6280	5.2310	0.3810	1.5330	1.1588
	Test	0.5623	0.8892	0.3132	1819.2220	0.6623	4.1277	0.4939	5.8428	1.0421
SMAPE	Train	43.8375	79.8680	13.1780	150.4262	67.4715	103.0555	52.3359	93.4002	67.3221
	Test	39.7998	81.4539	22.7237	152.4991	73.6498	106.9756	63.7582	93.9244	76.9794

398

399 To dive deep into the prediction performance of the EGB model, we analyzed variables  
 400 importance in the prediction of the volume. It was observed that slope length was the most  
 401 contributing predictor in the performance of the EGB model, followed by maximum hourly rainfall  
 402 and slope aspect. The altitude, three hours rainfall, slope angle and age of timber contributed  
 403 moderately to the prediction of the outcome volumes with gain above 0.01 and less than 0.2. The  
 404 antecedent rainfall from three days and above and continuous rainfall had a minor contribution,  
 405 with a gain of less than 0.01 for each. The presence of rainwater drainage channels had a moderate  
 406 contribution, with a gain close to 0.01. On the other hand, the contribution of soil depth and forest  
 407 density in the models was insignificant and far below 0.01. Though Figure 2(a) depicted the  
 408 association between larger volumes and fire history, the variable importance indicates that the  
 409 relation was not significant. Even though some variables had minor contributions, depending on  
 410 the case, the contribution of those variables may also increase depending on other regional settings.  
 411 Therefore, all variables with GVIF below 10 were kept in the model. Figure 6 illustrates the  
 412 variables importance for the EGB model. The vertical red line split the variables into two groups,

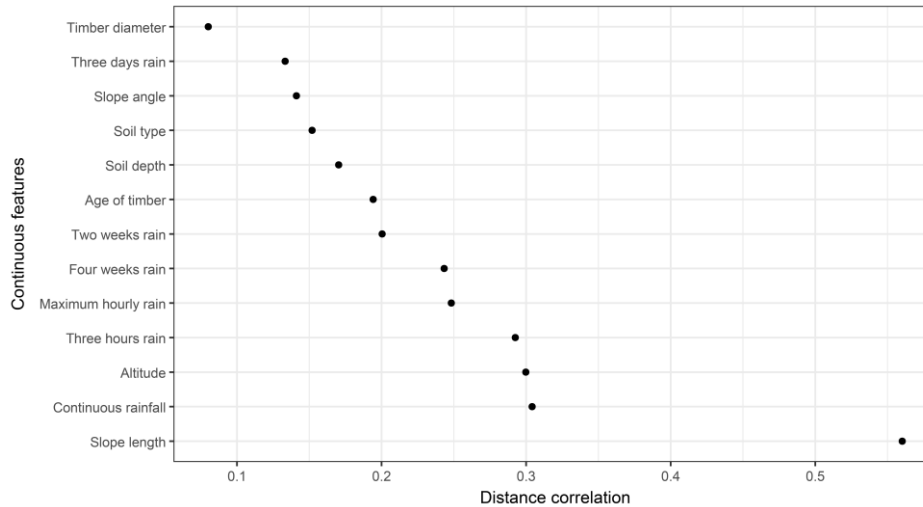
413 the first containing variables that contributed a gain above 0.01 and others with minor  
414 contributions.



415  
416 Figure 6. Variable importance for the EGB model.

417  
418 The variable importance plot depicts the overall contribution of a given variable; however,  
419 it does not provide detailed information. To get more insight into the relationship between the  
420 volume of landslides and predictors, statistical tests for normality, namely, Shapiro-Wilk's test,  
421 and Dunn's test were conducted. The Shapiro-Wilk's test (Dudley, 2023) results revealed that the  
422 distribution of volume was non-normal ( $W = 0.40642$ ,  $p\text{-value} < 0.001$ ). Noting that the volume  
423 distribution was non-normal, we opted for the non-parametric tests, which do not rely on normality  
424 to conduct the distance correlation (Székely et al., 2007) test (dcor) for continuous independent  
425 features. Figure 7 illustrates that the slope length exhibited a higher value (dcor=0.56) followed  
426 by continuous rainfall altitude and three hours rainfall and kept decreasing up to timber diameter  
427 with a distance correlation of 0.08. Overall, the distance correlation between the volume of  
428 landslides shows a moderate strength of association between continuous predictors.

429

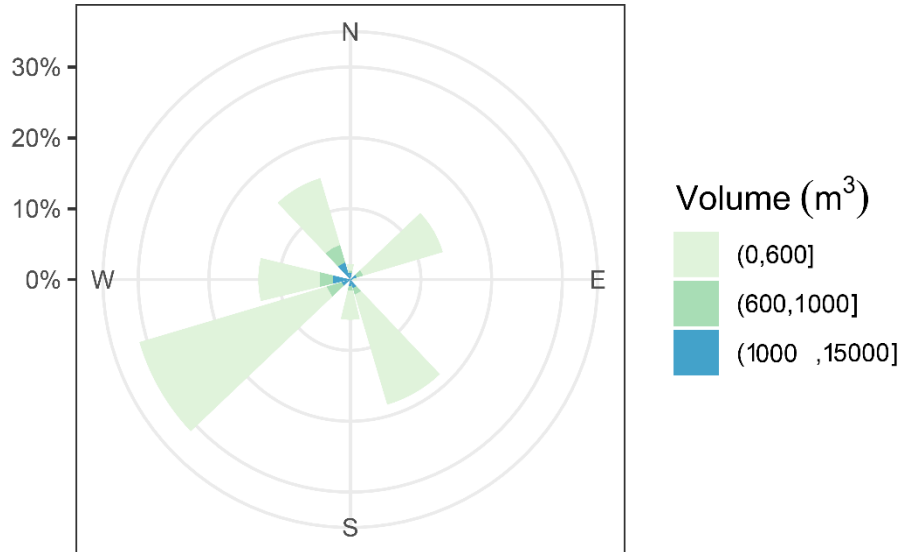


430

431 Figure 7. Distance correlation plot for the volume and continuous features.

432

433 Furthermore, to test for categorical features, Kruskal-Wallis test (McKight and Najab,  
 434 2010) was used to check whether the volume of the landslide was different in each category and  
 435 Dunn's tests (Dinno, 2015) were applied to examine which categories had similar means of the  
 436 volume of landslides due to rainfall in different categories. The  $H_0$  (null hypothesis) was that the  
 437 mean volume of landslides in different categories is the same, and the  $H_1$  (alternative hypothesis)  
 438 was that the means of landslides are different in some categories. For the slope aspect, the second  
 439 most significant predictor for the EGB model, the results of Kruskal-Wallis test (chi-squared =  
 440 20.889,  $df = 7$ ,  $p\text{-value} = 0.003938$ ) showed that there is a significant difference in median of  
 441 volume in some categories of slope aspects. To know which classes of slope aspects had  
 442 significantly different mean volumes, the Dunn's test results at 95% confidence interval, pairs  
 443 (East-South west, East-South East, East-South, East-North West and North West-South East) had  
 444 significantly different means of landslides' volume (with  $p\text{-value} < 0.05$ ). Figure 8 depicts that the  
 445 southwest and southeast aspects had a higher frequency of landslides.



446  
 447 Figure 8. The distribution of the volume of landslides due to rainfall with respect to the slope  
 448 aspect.  
 449

450 The Kruskal-Wallis test for the difference in mean of drainage classes showed the result  
 451 was: chi-squared = 15.792, df = 2, p-value = 0.000372, which shows that the means of volume per  
 452 class were different. This was clarified by Dunn's test results, p-values were less than 0.05 in all  
 453 pairwise mean difference comparisons. The results of these tests highlighted that drainage has a  
 454 remarkable influence on the occurrence of rainfall-induced landslides in the Korean Peninsula.

455  
 456 **5. Discussion**

457 Numerical models have traditionally been employed due to their foundation in physical principles  
 458 such as slope stability and hydrological dynamics (Glade et al., 2005). These models are valuable  
 459 for understanding the underlying mechanisms of landslide processes but often face limitations  
 460 when applied to regions with complex or heterogeneous terrain, as they require detailed, high-  
 461 quality input data that may not always be available (Caine, 1980). In the same way, statistical  
 462 models, which use historical rainfall and landslide data to establish correlations, can offer useful  
 463 predictions of VLDR in regions with extensive historical records (Chung and Fabbri, 2003).  
 464 However, these models may struggle to account for local variations in topography or rapidly  
 465 changing weather patterns, limiting their general applicability. Additionally, ML techniques have  
 466 shown significant promise in improving predictive accuracy at the regional level due to the  
 467 capability of processing large, diverse datasets and capturing complex, non-linear relationships

468 that traditional models might fail to capture (Pourghasemi and Rahmati, 2018). Further, ML  
469 models can adapt to regional variations and continuously improve as new data is introduced,  
470 offering a more flexible and dynamic approach to predict VLDR on a regional scale (Liu et al.,  
471 2021b). Subsequently, the aim of this study was to construct a data-driven algorithm that accurately  
472 predicts the VLDR. The result of nine different tested algorithms revealed a tremendous difference  
473 between classical regression models (OLS, RR, and GLM) and other data-driven machine learning  
474 models. In this study, apart from SVM regression, DT and KNN, other machine learning models  
475 (DNN, DT, RF, and EGB) exhibited high prediction capability with  $R^2$  above 50% (Fig. 5). The  
476 DNN, EGB, and RF models achieved  $R^2 > 0.8$  on both training and test set with accuracy reduced  
477  $R^2$  by 1.75, 7.72, and 12.17% for RF, EGB and DNN respectively, on the holdout set, indicating  
478 that the model could yield reliable volume estimates in adjacent areas with similar geological and  
479 environmental conditions. The random forest model performed well in predicting smaller volume;  
480 however, as the volume increased, the model underpredicted volume values. The DNN model  
481 performed quite well with low MAE compared to random forest; however, the model did not  
482 perform well on moderate volume values, resulting in reduced  $R^2$ . The EGB model tested on South  
483 Korean landslide inventory coupled with rainfall data at the time of landslide events and antecedent  
484 rainfall within one month of the event exhibited more accurate predictions compared to other  
485 constructed algorithms. The difference in performance may be due to the internal structure of each  
486 algorithm; the RF builds multiple decision trees and averages predictions to improve accuracy  
487 (Breiman, 2001), while the EGB builds sequential trees in a recursive order where the new built  
488 tree improves error occurred while building the previous decision tree and optimizes the loss  
489 function through a gradient descent (Chen and Guestrin, 2016).

490 The slope aspect played an important role in the prediction of the volume, and the landslide  
491 mostly occurred in locations oriented toward south-southwest and southeast. That may be due to  
492 the direction taken by typhoons, which hit the southwest versants of mountains upon landfall on  
493 the Korean peninsula toward the North East Pacific (Lee et al., 2013; Ha, 2022). The findings of  
494 this research are congruent with those of Lee et al. (2013), who also highlighted that the mountain  
495 versant oriented to strong wind direction may face more landslides. The study also highlighted that  
496 a moderate rainwater drainage channel plays an important role in the prevention of landslides due  
497 to its stabilizing effect. The landslide location and pattern follow the rainfall climate scenario,  
498 which highlighted a higher intensity of rainfall in the northeastern region of South Korea (Lee,



499 2016). In addition, the findings of this study are congruent with Zhang et al. (2019) observations  
500 that highlighted the low influence of soil type in landslide modeling and the maximum rainfall and  
501 cumulative three hours of rainfall were the most contributing rainfall, which indicated that these  
502 shallow landslides may have been triggered by sudden rainfall concentrated in few hours before  
503 the occurrence of the event. The occurrence of landslides triggered by rainfall is a complex  
504 phenomenon that involves many interrelated environmental settings, human activity, geological  
505 conditions and climatic conditions. Moreover, the occurrence of typhoons is known to aggravate  
506 the landslides impacts on communities (Chang et al., 2008); incorporating typhoon variables in  
507 future studies to customize for regional settings may improve the accuracy of the model. The  
508 advantage of his research is that the constructed model has high predictive accuracy and can handle  
509 the non-linearity of predisposing factors. The model came to fill the gap in a few literatures related  
510 to the prediction of the volume of landslides using data-driven techniques. This model can be a  
511 good tool to help policy-makers integrate the landslides volume risks in policy to protect  
512 infrastructure and inhabitants dwelling near the foot of mountains with high risks of being buried  
513 by geological materials resulting from landslides.

514 To understand the applicability of the developed models, the trained model was tested  
515 using unknown data (test data), with volume predictions generated solely based on the predictor  
516 variables; actual volume values were utilized only for evaluating model prediction accuracy. The  
517 outcome exhibited that the difference in  $R^2$  on the training and holdout set of 7.72% for the optimal  
518 model (i.e., EGB) highlights that the model can be applied to another region of a similar setting. It  
519 was noted that without proper model calibration with the independent data set, it's difficult to  
520 determine whether these discrepancies in performance are due to model limitations or data  
521 differences in different regions (Huang et al., 2020). Therefore, in future work, we plan to develop  
522 an independent database based on collecting the extensive recent landslide geometry at different  
523 parts of the Korean Peninsula to improve the models further by calibrating region-specific  
524 parameters to ensure the transferability of the model to other regions.

525 The major limitation of this study is that the analysis is solely focused on shallow-seated  
526 landslides, specifically translational slope failures with volumes below 13,000m<sup>3</sup>. Thus, the  
527 analysis may not fully capture the variability in landslide characteristics across different  
528 geomorphological and geological contexts. Deep-seated landslides, for instance, often exhibit  
529 distinct failure mechanisms, material compositions, and depositional patterns that influence their

530 volumetric characteristics, which were not considered in this investigation. Similarly, debris flows,  
531 known for their unique channelization and entrainment behaviors, were not included, potentially  
532 limiting the applicability of the optimized models to other landslide types. Further, this study was  
533 also performed using point-based landslide inventory data, which may not capture all variability  
534 of influencing factors and their exact state. The incorporation of high-resolution data from remote  
535 sensing and other sources may also improve the efficiency of the predictions. These limitations  
536 may impact the broader applicability of the proposed model; however, future studies will aim to  
537 address this by conducting separate analyses for deep-seated landslides and debris flows, allowing  
538 for a more comprehensive understanding of landslide volume predictions across diverse landslide  
539 types and geomorphological settings.

540

## 541 **6. Conclusions**

542 In this paper, the aim was to construct a data-driven model that predicts the volume of landslides  
543 due to rainfall. To this, nine different classical regression models and machine learning algorithms  
544 were tested on South Korean landslide data set containing features of landslides that occurred  
545 between 2011 and 2012. Among the tested models, the EGB model produced the most accurate  
546 prediction. This is proven by the evaluation of the difference between actual and predicted values,  
547 such as  $R^2= 88.41\%$  and  $MAE=146.6120m^3$  on the holdout set. The analysis of feature variables  
548 in the contribution to the prediction of the model revealed that the slope length was the most  
549 influencing predictor. The EGB model can be a promising tool for the prediction of the volume of  
550 landslides due to its high predictive performance. The model can be customized in different  
551 environmental settings. The model can be applied to estimate the expected volume of landslides  
552 based on forecasted rainfall once the model is well-adjusted to fit the geomorphological and  
553 environmental settings of the region of interest after re-training on the regional historical data to  
554 include regional variability. Therefore, this model can be a good tool for planning for resilience  
555 and infrastructure pre-construction risk assessment to ensure the new infrastructure is placed in  
556 stable regions free from severe landslides.

557

## 558 **Acknowledgments**

559 This research was supported by the Korean government (MSIT) (2021R1C1C2003316) and Basic  
560 Science Research Program through the National Research Foundation of Korea (NRF) funded by  
561 Ministry of Education (2021R1A6A1A03044326).

562 The authors highly appreciate both anonymous reviewers and editor for their constructive  
563 suggestions that helped us improve the preprint version.

564

## 565 **Reference**

566 Alcántara, A. L., and Ahn, K. H. (2020). Probability distribution and characterization of daily  
567 precipitation related to tropical cyclones over the Korean Peninsula. *Water*, 12(4), 1214.  
568 <https://doi.org/10.3390/w12041214>

569 Alcántara-Ayala, I., and Sassa, K. (2023). Landslide risk management: from hazard to disaster risk  
570 reduction. *Landslides*, 20(10), 2031-2037. <https://doi.org/10.1007/s10346-023-02140-5>

571 Amesoeder, C., Hartig, F., and Pichler, M. (2023). cito: An R package for training neural networks  
572 using torch. arXiv e-prints, arXiv-2303. <https://doi.org/10.1111/ecog.07143>

573 Armstrong, J. S. (2001). Combining forecasts (pp. 417-439). Springer US.  
574 [https://doi.org/10.1007/978-0-306-47630-3\\_19](https://doi.org/10.1007/978-0-306-47630-3_19)

575 Asada, H., and Minagawa, T. (2023). Impact of vegetation differences on shallow landslides: a  
576 case study in Aso, Japan. *Water*, 15(18), 3193. <https://doi.org/10.3390/w15183193>

577 Bernardie, S., Desramaut, N., Malet, J.-P., Gourlay, M., and Grandjean, G. (2014). Prediction of  
578 changes in landslide rates induced by rainfall. *Landslides*, 12(3), 481–494.  
579 <https://doi.org/10.1007/s10346-014-0495-8>

580 Bonamutial, M., and Prasetyo, S. Y. (2023). Exploring the Impact of Feature Data Normalization  
581 and Standardization on Regression Models for Smartphone Price Prediction. In 2023  
582 International Conference on Information Management and Technology (ICIMTech) (pp.  
583 294-298). IEEE. <https://doi.org/10.1109/ICIMTech59029.2023.10277860>

584 Borup, D., Christensen, B. J., Mühlbach, N. S., and Nielsen, M. S. (2023). Targeting predictors in  
585 random forest regression. *International Journal of Forecasting*, 39(2), 841-868.  
586 <https://doi.org/10.1016/j.ijforecast.2022.02.010>

587 Breiman, L. (2001). Random forests. *Machine learning*, 45, 5-32. <https://doi.org/10.1023/A:1010933404324>

588  
589 Breiman, L. (2017). Classification and regression trees. Routledge. <https://doi.org/10.1201/9781315139470>

590  
591 Caine, N. (1980). The rainfall intensity-duration control of shallow landslides and debris flows.  
592 *Geografiska annaler: series A, physical geography*, 62(1-2), 23-27.  
593 <https://doi.org/10.1080/04353676.1980.11879996>

594 Cellek, S. (2021). The effect of aspect on landslide and its relationship with other parameters. In  
595 *Landslides*. IntechOpen.

- 596 Chang, K. T., and Chiang, S. H. (2009). An integrated model for predicting rainfall-induced  
597 landslides. *Geomorphology*, 105(3-4), 366-373. [https://doi.org/10.1016/](https://doi.org/10.1016/j.geomorph.2008.10.012)  
598 [j.geomorph.2008.10.012](https://doi.org/10.1016/j.geomorph.2008.10.012)
- 599 Chang, K. T., Chiang, S. H., and Lei, F. (2008). Analysing the relationship between typhoon-  
600 triggered landslides and critical rainfall conditions. *Earth Surface Processes and*  
601 *Landforms: The Journal of the British Geomorphological Research Group*, 33(8), 1261-  
602 1271. <https://doi.org/10.1002/esp.1611>
- 603 Chatra, A. S., Dodagoudar, G. R., and Maji, V. B. (2019). Numerical modelling of rainfall effects  
604 on the stability of soil slopes. *International Journal of Geotechnical Engineering*.  
605 <https://doi.org/10.1080/19386362.2017.1359912>
- 606 Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., Chen, K., Mitchell, R., Cano, I.,  
607 Zhou, T., Li, M., Xie, J., Lin, M., Geng, Y., Li, Y., and Yuan, J. (2022). *\_xgboost: Extreme*  
608 *Gradient Boosting\_*. R package version 1.6.0.1, <[https://CRAN.R-](https://CRAN.R-project.org/package=xgboost)  
609 [project.org/package=xgboost](https://CRAN.R-project.org/package=xgboost)>.[Accessed 2025-01-25]
- 610 Chen, C. W., Oguchi, T., Hayakawa, Y. S., Saito, H., and Chen, H. (2017). Relationship between  
611 landslide size and rainfall conditions in Taiwan. *Landslides*, 14, 1235-1240.  
612 <https://doi.org/10.1007/s10346-016-0790-7>
- 613 Chen, L., Guo, Z., Yin, K., Shrestha, D. P., and Jin, S. (2019). The influence of land use and land  
614 cover change on landslide susceptibility: a case study in Zhushan Town, Xuan'en County  
615 (Hubei, China). *Natural Hazards and Earth System Sciences*, 19(10), 2207-2228.  
616 <https://doi.org/10.5194/nhess-19-2207-2019>
- 617 Chen, T., and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the*  
618 *22nd acm sigkdd international conference on knowledge discovery and data mining* (pp.  
619 785-794). <https://doi.org/10.1145/2939672.2939785>
- 620 Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., ... and Zhou, T. (2015b). Xgboost:  
621 extreme gradient boosting. R package version 0.4-2, 1(4), 1(4).
- 622 Chen, X., Zhang, L., Zhang, L., Zhou, Y., Ye, G., and Guo, N. (2021). Modelling rainfall-induced  
623 landslides from initiation of instability to post-failure. *Computers and Geotechnics*, 129,  
624 103877. <https://doi.org/10.1016/j.compgeo.2020.103877>
- 625 Chen, Z., Luo, R., Huang, Z., Tu, W., Chen, J., Li, W., ... and Ai, Y. (2015a). Effects of different  
626 backfill soils on artificial soil quality for cut slope revegetation: Soil structure, soil  
627 erosion, moisture retention and soil C stock. *Ecological Engineering*, 83, 5-12.  
628 <https://doi.org/10.1016/j.ecoleng.2015.05.048>
- 629 Cheung, R. W. (2021). Landslide risk management in Hong Kong. *Landslides*, 18(10), 3457-3473.  
630 <https://doi.org/10.1007/s10346-020-01587-0>
- 631 Chicco, D., Warrens, M. J., and Jurman, G. (2021). The coefficient of determination R-squared is  
632 more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis  
633 evaluation. *Peerj Computer Science*, 7, e623. <https://doi.org/10.7717/peerj-cs.623>
- 634 Chowdhury, M. Z. I., Leung, A. A., Walker, R. L., Sikdar, K. C., O'Beirne, M., Quan, H., and  
635 Turin, T. C. (2023). A comparison of machine learning algorithms and traditional

636 regression-based statistical modeling for predicting hypertension incidence in a Canadian  
637 population. *Scientific Reports*, 13(1), 13. <https://doi.org/10.1038/s41598-022-27264-x>

638 Chung, C. J. F., and Fabbri, A. G. (2003). Validation of spatial prediction models for landslide  
639 hazard mapping. *Natural Hazards*, 30, 451-472. [https://doi.org/10.1023/  
640 B:NHAZ.0000007172.62651.2b](https://doi.org/10.1023/B:NHAZ.0000007172.62651.2b)

641 Cohen, D., and Schwarz, M. (2017). Tree-root control of shallow landslides. *Earth Surface  
642 Dynamics*, 5(3), 451-477. <https://doi.org/10.5194/esurf-5-451-2017>

643 Culler, E. S., Livneh, B., Rajagopalan, B., and Tiampo, K. F. (2021). A data-driven evaluation of  
644 post-fire landslide susceptibility. *Natural Hazards and Earth System Sciences  
645 Discussions*, 2021, 1-24. <https://doi.org/10.5194/nhess-23-1631-2023>

646 Dahal, B. K., and Dahal, R. K. (2017). Landslide hazard map: tool for optimization of low-cost  
647 mitigation. *Geoenvironmental Disasters*, 4, 1-9. [https://doi.org/10.1186/s40677-017-  
648 0071-3](https://doi.org/10.1186/s40677-017-<br/>
648 0071-3)

649 Dai, F. C., and Lee, C. F. (2001). Frequency–volume relation and prediction of rainfall-induced  
650 landslides. *Engineering Geology*, 59(3-4), 253-266. [https://doi.org/10.1016/S0013-  
651 7952\(00\)00077-6](https://doi.org/10.1016/S0013-<br/>
651 7952(00)00077-6)

652 Darlington, R. B. (1990). *Regression and linear models*. McGraw-Hill College.

653 Dinno, A. (2015). Nonparametric pairwise multiple comparisons in independent groups using  
654 Dunn's test. *The Stata Journal*, 15(1), 292-300. [https://doi.org/10.1177/  
655 /1536867X1501500117](https://doi.org/10.1177/<br/>
655 /1536867X1501500117)

656 Dobson, A. J., and Barnett, A. G. (2018). *An introduction to generalized linear models*. CRC press.  
657 <https://doi.org/10.1201/9781315182780>

658 Donnarumma, A., Revellino, P., Grelle, G., and Guadagno, F. M. (2013). Slope angle as indicator  
659 parameter of landslide susceptibility in a geologically complex area. *Landslide Science  
660 and Practice: Volume 1: Landslide Inventory and Susceptibility and Hazard Zoning*, 425-  
661 433. [https://doi.org/10.1007/978-3-642-31325-7\\_56](https://doi.org/10.1007/978-3-642-31325-7_56)

662 Duc, D. M. (2013). Rainfall-triggered large landslides on 15 December 2005 in Van Canh district,  
663 Binh Dinh province, Vietnam. *Landslides*, 10(2), 219-230. [https://doi.org/10.1007/  
664 /s10346-012-0362-4](https://doi.org/10.1007/<br/>
664 /s10346-012-0362-4)

665 Dudley, R. (2023). The Shapiro–Wilk test for normality. Available at  
666 <https://math.mit.edu/~rmd/46512/shapiro.pdf> [Accessed 2025-01-25]

667 Evans, S. G., Mugnozza, G. S., Strom, A., and Hermanns, R. L. (2007). *Landslides from massive  
668 rock slope failure (Vol. 49)*. Springer Science and Business Media.

669 Friedman, J. H., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear  
670 models via coordinate descent. *Journal of statistical software*, 33, 1-22. Available at  
671 <https://pmc.ncbi.nlm.nih.gov/articles/PMC2929880/>

672 Gariano, S. L., Rianna, G., Petrucci, O., and Guzzetti, F. (2017). Assessing future changes in the  
673 occurrence of rainfall-induced landslides at a regional scale. *Science of the Total  
674 Environment*, 596, 417-426. <https://doi.org/10.1016/j.scitotenv.2017.03.103>

675 Gelman, A. (2007). Data analysis using regression and multilevel/hierarchical models. Cambridge  
676 University Press.

677 Glade, T., Anderson, M. G., and Crozier, M. J. (2005). Landslide hazard and risk (Vol. 807). John  
678 Wiley & Sons. <https://doi.org/10.1002/9780470012659>

679 Gong, Q., Wang, J., Zhou, P., and Guo, M. (2021). A regional landslide stability analysis method  
680 under the combined impact of rainfall and vegetation roots in south China. *Advances in*  
681 *Civil Engineering*, 2021, 1-12. <https://doi.org/10.1155/2021/5512281>

682 Gonzalez-Ollauri, A., and Mickovski, S. B. (2017). Hydrological effect of vegetation against  
683 rainfall-induced landslides. *Journal of Hydrology*, 549, 374-387.  
684 <https://doi.org/10.1016/j.jhydrol.2017.04.014>

685 Greenwood, J. R., Norris, J. E., and Wint, J. (2004). Assessing the contribution of vegetation to  
686 slope stability. *Proceedings of the Institution of Civil Engineers-Geotechnical*  
687 *Engineering*, 157(4), 199-207. <https://doi.org/10.1680/geng.2004.157.4.199>

688 Gutierrez-Martin, A. (2020). A GIS-physically-based emergency methodology for predicting  
689 rainfall-induced shallow landslide zonation. *Geomorphology*, 359, 107121.  
690 <https://doi.org/10.1016/j.geomorph.2020.107121>

691 Guzzetti, F., Peruccacci, S., Rossi, M., and Stark, C. P. (2008). The rainfall intensity–duration  
692 control of shallow landslides and debris flows: an update. *Landslides*, 5, 3-17.  
693 <https://doi.org/10.1007/s10346-007-0112-1>

694 Ha, K. M. (2022). predicting typhoon tracks around Korea. *Natural Hazards*, 113(2), 1385-1390.  
695 <https://doi.org/10.1007/s11069-022-05335-6>

696 Hastie, T. (2009). *The elements of statistical learning: data mining, inference, and prediction*. 2nd  
697 edition. <https://doi.org/10.1111/j.1541-0420.2010.01516.x>

698 Highland, L. and Bobrowsky, P. (2008). *The Landslide Handbook: A Guide to Understanding*  
699 *Landslides*, United States Geological Survey, Reston, VA, Circular 1325, Available at  
700 <https://pubs.usgs.gov/circ/1325/> [ Accessed: 2025-01-25]

701 Holcombe, E. A., Beesley, M. E., Vardanega, P. J., and Sorbie, R. (2016). Urbanisation and  
702 landslides: hazard drivers and better practices. In *Proceedings of the Institution of Civil*  
703 *Engineers-Civil Engineering* (Vol. 169(3), pp. 137-144). Thomas Telford Ltd.  
704 <https://doi.org/10.1680/jcien.15.00044>

705 Hovius, N., Stark, C. P., and Allen, P. A. (1997). Sediment flux from a mountain belt derived by  
706 landslide mapping. *Geology*, 25(3), 231-234. [https://doi.org/10.1130/0091-7613\(1997\)025<0231:SFFAMB>2.3.CO;2](https://doi.org/10.1130/0091-7613(1997)025<0231:SFFAMB>2.3.CO;2)

707

708 Huang, J., Hales, T. C., Huang, R., Ju, N., Li, Q., and Huang, Y. (2020). A hybrid machine-learning  
709 model to estimate potential debris-flow volumes. *Geomorphology*, 367, 107333.  
710 <https://doi.org/10.1016/j.geomorph.2020.107333>

711 Hyde, K. D., Riley, K., and Stoof, C. (2016). Uncertainties in predicting debris flow hazards  
712 following wildfire. *Natural Hazards*. <https://doi.org/10.1002/9781119028116.ch19>

713 Hyndman, R. J., and Koehler, A. B. (2006). Another look at measures of forecast accuracy.  
714 International Journal of Forecasting, 22(4), 679-688. [https://doi.org/10.1016/](https://doi.org/10.1016/j.ijforecast.2006.03.001)  
715 [j.ijforecast.2006.03.001](https://doi.org/10.1016/j.ijforecast.2006.03.001)  
716

717 Hyun, Y. K., Kar, S. K., Ha, K. J., and Lee, J. H. (2010). Diurnal and spatial variabilities of  
718 monsoonal CG lightning and precipitation and their association with the synoptic weather  
719 conditions over South Korea. Theoretical and Applied Climatology, 102, 43-60.  
720 <https://doi.org/10.1007/s00704-009-0235-5>

721 Intrieri, E., Carlà, T., and Gigli, G. (2019). Forecasting the time of failure of landslides at slope-  
722 scale: A literature review. Earth-science reviews, 193, 333-349.  
723 <https://doi.org/10.1016/j.earscirev.2019.03.019>

724 Jaboyedoff, M., Choffet, M., Derron, M. H., Horton, P., Loye, A., Longchamp, C., Mazotti, B.,  
725 Michoud, C., and Pedrazzini, A. (2012). Preliminary slope mass movement susceptibility  
726 mapping using DEM and LiDAR DEM. In Terrigenous mass movements: Detection,  
727 modelling, early warning and mitigation using geoinformation technology, 109-170.  
728 Springer, Berlin Heidelberg, [https://doi.org/10.1007/978-3-642-25495-6\\_5](https://doi.org/10.1007/978-3-642-25495-6_5)

729 Jin, H. G., Lee, H., and Baik, J. J. (2022). Characteristics and possible mechanisms of diurnal  
730 variation of summertime precipitation in South Korea. Theoretical and Applied  
731 Climatology, 148(1), 551-568. <https://doi.org/10.1007/s00704-022-03965-1>

732 Ju, L. Y., Zhang, L. M., and Xiao, T. (2023). Power laws for accurate determination of landslide  
733 volume based on high-resolution LiDAR data. Engineering Geology, 312, 106935.  
734 <https://doi.org/10.1016/j.enggeo.2022.106935>

735 Jung, M. J., Jeong, Y. J., Shin, W. J., and Cheong, A. C. S. (2024). Isotopic distribution of  
736 bioavailable Sr, Nd, and Pb in Chungcheongbuk-do Province, Korea. Journal of  
737 Analytical Science and Technology, 15(1), 46. [https://doi.org/10.1186/s40543-024-](https://doi.org/10.1186/s40543-024-00460-2)  
738 [00460-2](https://doi.org/10.1186/s40543-024-00460-2)

739 Jung, Y., Shin, J. Y., Ahn, H., and Heo, J. H. (2017). The spatial and temporal structure of extreme  
740 rainfall trends in South Korea. Water, 9(10), 809. <https://doi.org/10.3390/w9100809>

741 Kafle, L., Xu, W. J., Zeng, S. Y., and Nagel, T. (2022). A numerical investigation of slope stability  
742 influenced by the combined effects of reservoir water level fluctuations and precipitation:  
743 A case study of the Bianjiazhai landslide in China. Engineering Geology, 297, 106508.  
744 <https://doi.org/10.1016/j.enggeo.2021.106508>

745 Kang, M. W., Yibeltal, M., Kim, Y. H., Oh, S. J., Lee, J. C., Kwon, E. E., and Lee, S. S. (2022).  
746 Enhancement of soil physical properties and soil water retention with biochar-based soil  
747 amendments. Science of the Total Environment, 836, 155746. [https://doi.org/10.1016/](https://doi.org/10.1016/j.scitotenv.2022.155746)  
748 [j.scitotenv.2022.155746](https://doi.org/10.1016/j.scitotenv.2022.155746)

749 Keefer, R. F. (2000). Handbook of soils for landscape architects. Oxford University Press.

750 Khan, M. A., Basharat, M., Riaz, M. T., Sarfraz, Y., Farooq, M., Khan, A. Y., Pham, Q. B., Ahmed,  
751 K. S., and Shahzad, A. (2021). An integrated geotechnical and geophysical investigation

752 of a catastrophic landslide in the Northeast Himalayas of Pakistan. *Geological Journal*,  
753 56(9), 4760-4778. <https://doi.org/10.1002/gj.4209>

754 Khan, Y. A., Lateh, H., Baten, M. A., and Kamil, A. A. (2012). Critical antecedent rainfall  
755 conditions for shallow landslides in Chittagong City of Bangladesh. *Environmental Earth*  
756 *Sciences*, 67, 97-106. <https://doi.org/10.1007/s12665-011-1483-0>

757 Kim, D. E., Seong, Y. B., Weber, J., and Yu, B. Y. (2020). Unsteady migration of Taebaek Mountain  
758 drainage divide, Cenozoic extensional basin margin, Korean Peninsula. *Geomorphology*,  
759 352, 107012. <https://doi.org/10.1016/j.geomorph.2019.107012>

760 Kim, H. G., and Park, C. Y. (2021). Landslide susceptibility analysis of photovoltaic power stations  
761 in Gangwon-do, Republic of Korea. *Geomatics, Natural Hazards and Risk*, 12(1), 2328-  
762 2351. <https://doi.org/10.1080/19475705.2021.1950219>

763 Kim, J., Lee, K., Jeong, S., and Kim, G. (2014). GIS-based prediction method of landslide  
764 susceptibility using a rainfall infiltration-groundwater flow model. *Engineering Geology*,  
765 182, 63-78. <https://doi.org/10.1016/j.enggeo.2014.09.001>

766 Kim, M. S., Onda, Y., Kim, J. K., and Kim, S. W. (2015). Effect of topography and soil  
767 parameterisation representing soil thicknesses on shallow landslide modelling.  
768 *Quaternary International*, 384, 91-106. <https://doi.org/10.1016/j.quaint.2015.03.057>

769 Kim, S. W., Chun, K. W., Kim, M., Catani, F., Choi, B., and Seo, J. I. (2021). Effect of antecedent  
770 rainfall conditions and their variations on shallow landslide-triggering rainfall thresholds  
771 in South Korea. *Landslides*, 18, 569-582. <https://doi.org/10.1007/s10346-020-01505-4>

772 Kitutu, M. G., Muwanga, A., Poesen, J., and Deckers, J. A. (2009). Influence of soil properties on  
773 landslide occurrences in Bududa district, Eastern Uganda. *African Journal of Agricultural*  
774 *Research*, 4(7), 611-620. Available at <https://lirias.kuleuven.be/retrieve/78489> [Accessed  
775 2025-01-25]

776 Korup, O. (2004). Geomorphometric characteristics of New Zealand landslide dams. *Engineering*  
777 *Geology*, 73(1-2), 13-35. <https://doi.org/10.1016/j.enggeo.2003.11.003>

778 Korup, O., Clague, J. J., Hermanns, R. L., Hewitt, K., Strom, A. L., and Weidinger, J. T. (2007).  
779 Giant landslides, topography, and erosion. *Earth and Planetary Science Letters*, 261(3-  
780 4), 578-589. <https://doi.org/10.1016/j.epsl.2007.07.025>

781 Kotsakis, C. (2023). Ordinary Least Squares. In *Encyclopedia of Mathematical Geosciences* (pp.  
782 1032-1038). Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-  
783 030-85040-1\\_237](https://doi.org/10.1007/978-3-030-85040-1_237)

784 Kramer, O., and Kramer, O. (2013). K-nearest neighbors. Dimensionality reduction with  
785 unsupervised nearest neighbors, 13-23. [https://doi.org/10.1007/978-3-642-38652-7\\_2](https://doi.org/10.1007/978-3-642-38652-7_2)

786 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep  
787 convolutional neural networks. *Advances in Neural Information Processing Systems*, 25.  
788 Available at  
789 [https://proceedings.neurips.cc/paper\\_files/paper/2012/file/c399862d3b9d6b76c8436e92  
790 4a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf) [Accessed 2025-01-25].



791 Kuhn, M. (2022). caret: Classification and Regression Training\_. R package version 6.0-92,  
792 Available at <https://CRAN.R-project.org/package=caret> [Accessed 2025-01-25]

793 Kunz, M., and Kottmeier, C. (2006). Orographic enhancement of precipitation over low mountain  
794 ranges. Part II: Simulations of heavy precipitation events over southwest Germany.  
795 Journal of applied meteorology and climatology, 45(8), 1041-1055.  
796 <https://doi.org/10.1175/JAM2390.1>

797 Lacerda, W. A., Palmeira, E. M., Netto, A. L. C., and Ehrlich, M. (Eds.). (2014). Extreme rainfall  
798 induced landslides: an international perspective. Oficina de Textos. ISBN 978-85-7975-  
799 150-9.

800 Lann, T., Bao, H., Lan, H., Zheng, H., and Yan, C. (2024). Hydro-mechanical effects of vegetation  
801 on slope stability: A review. Science of the Total Environment, 171691.  
802 <https://doi.org/10.1016/j.scitotenv.2024.171691>

803 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.  
804 <https://doi.org/10.1038/nature14539>

805 Lee, D. B., Kim, Y. N., Sonn, Y. K., and Kim, K. H. (2023). Comparison of Soil Taxonomy (2022)  
806 and WRB (2022) Systems for classifying Paddy Soils with different drainage grades in  
807 South Korea. Land, 12(6), 1204. <https://doi.org/10.3390/land12061204>

808 Lee, D. H., Cheon, E., Lim, H. H., Choi, S. K., Kim, Y. T., and Lee, S. R. (2021). An artificial  
809 neural network model to predict debris-flow volumes caused by extreme rainfall in the  
810 central region of South Korea. Engineering Geology, 281, 105979.  
811 <https://doi.org/10.1016/j.enggeo.2020.105979>

812 Lee, D. H., Kim, Y. T., and Lee, S. R. (2020). Shallow landslide susceptibility models based on  
813 artificial neural networks considering the factor selection method and various non-linear  
814 activation functions. Remote Sensing, 12(7), 1194. <https://doi.org/10.3390/rs12071194>

815 Lee, J. U., Cho, Y. C., Kim, M., Jang, S. J., Lee, J., and Kim, S. (2022). The effects of different  
816 geological conditions on landslide-triggering rainfall conditions in South Korea. Water,  
817 14(13), 2051. <https://doi.org/10.3390/w14132051>

818 Lee, M. J. (2016). Rainfall and landslide correlation analysis and prediction of future rainfall base  
819 on climate change. In Geohazards Caused by Human Activity. IntechOpen.

820 Lee, S. W., Kim, G., Yune, C. Y., and Ryu, H. J. (2013). Development of landslide-risk assessment  
821 model for mountainous regions in eastern Korea. Disaster Advances, 6(6), 70-79.

822 Li, C. J., Guo, C. X., Yang, X. G., Li, H. B., and Zhou, J. W. (2022). A GIS-based probabilistic  
823 analysis model for rainfall-induced shallow landslides in mountainous areas.  
824 Environmental Earth Sciences, 81(17), 432. [https://doi.org/10.1007/s12665-022-10562-  
825 y](https://doi.org/10.1007/s12665-022-10562-<br/>
825 y)

826 Liaw, A., and Wiener, M., (2002). Classification and regression by randomForest. R News 2(3),  
827 18--22. Available at <https://journal.r-project.org/articles/RN-2002-022/RN-2002-022.pdf>  
828 [Accessed 2025-01-24].

829 Liu, Y., Deng, Z., and Wang, X. (2021a). The effects of rainfall, soil type and slope on the processes  
830 and mechanisms of rainfall-induced shallow landslides. *Applied Sciences*, 11(24), 11652.  
831 <https://doi.org/10.3390/app112411652>

832 Liu, Z., Gilbert, G., Cepeda, J. M., Lysdahl, A. O. K., Piciullo, L., Hefre, H., and Lacasse, S.  
833 (2021b). Modelling of shallow landslides with machine learning algorithms. *Geoscience*  
834 *Frontiers*, 12(1), 385-393. <https://doi.org/10.1016/j.gsf.2020.04.014>

835 Luino, F., De Graff, J., Biddoccu, M., Faccini, F., Freppaz, M., Roccati, A., Ungaro, F., D'Amico,  
836 M., and Turconi, L. (2022). The Role of soil type in triggering shallow landslides in the  
837 alps (Lombardy, Northern Italy). *Land*, 11(8), [https://doi.org/1125.](https://doi.org/1125.10.3390/land11081125)  
838 [10.3390/land11081125](https://doi.org/1125.10.3390/land11081125)

839 Martinović, K., Gavin, K., Reale, C., and Mangan, C. (2018). Rainfall thresholds as a landslide  
840 indicator for engineered slopes on the Irish Rail network. *Geomorphology*, 306, 40-50.  
841 <https://doi.org/10.1016/j.geomorph.2018.01.006>

842 McKenna, J. P., Santi, P. M., Amblard, X., and Negri, J. (2012). Effects of soil-engineering  
843 properties on the failure mode of shallow landslides. *Landslides*, 9, 215-228.  
844 <https://doi.org/10.1007/s10346-011-0295-3>

845 McKight, P. E., and Najab, J. (2010). Kruskal-wallis test. *The corsini encyclopedia of psychology*,  
846 1-1. <https://doi.org/10.1002/9780470479216.corpsy0491>

847 Meyer D, Dimitriadou E, Hornik K, Weingessel A, Leisch F (2021). `e1071: Misc Functions of`  
848 `the Department of Statistics, Probability Theory Group (Formerly: E1071), TU Wien`.  
849 `R package version 1.7-9`. <https://doi.org/10.32614/CRAN.package.e1071>

850 Miao, F., Wu, Y., Xie, Y., and Li, Y. (2018). Prediction of landslide displacement with step-like  
851 behavior based on multialgorithm optimization and a support vector regression model.  
852 *Landslides*, 15, 475-488. <https://doi.org/10.1007/s10346-017-0883-y>

853 Montgomery, D. R., Schmidt, K. M., Dietrich, W. E., and McKean, J. (2009). Instrumental record  
854 of debris flow initiation during natural rainfall: Implications for modeling slope stability.  
855 *Journal of Geophysical Research: Earth Surface*, 114(F1).  
856 <https://doi.org/10.1029/2008JF001078>

857 Nguyen, Q. H., Ly, H. B., Ho, L. S., Al-Ansari, N., Le, H. V., Tran, V. Q., Prakash, I., and Pham,  
858 B. T. (2021). Influence of data splitting on performance of machine learning models in  
859 prediction of shear strength of soil. *Mathematical Problems in Engineering*, 2021(1),  
860 4832864. <https://doi.org/10.1155/2021/4832864>

861 O'brien, R. M. (2007). A caution regarding rules of thumb for variance inflation factors. *Quality*  
862 *and quantity*, 41, 673-690. <https://doi.org/10.1007/s11135-006-9018-6>

863 Omwega, A. K. (1989). Crop cover, rainfall energy and soil erosion in Githunguri (Kiambu  
864 District), Kenya. The University of Manchester (United Kingdom). Available at  
865 [https://www.proquest.com/openview/dd7c169f804775d18041ec262d03e4c1/1?cbl=202](https://www.proquest.com/openview/dd7c169f804775d18041ec262d03e4c1/1?cbl=2026366&diss=y&pq-origsite=gscholar)  
866 [6366&diss=y&pq-origsite=gscholar](https://www.proquest.com/openview/dd7c169f804775d18041ec262d03e4c1/1?cbl=2026366&diss=y&pq-origsite=gscholar) [Accessed 2025-01-24].

867 Panday, S., and Dong, J. J. (2021). Topographical features of rainfall-triggered landslides in Mon  
868 State, Myanmar, August 2019: Spatial distribution heterogeneity and uncommon large

869 relative heights. *Landslides*, 18(12), 3875-3889. [https://doi.org/10.1007/s10346-021-](https://doi.org/10.1007/s10346-021-01758-7)  
870 [01758-7](https://doi.org/10.1007/s10346-021-01758-7)

871 Park, C. Y. (2015). The classification of extreme climate events in the Republic of Korea. *Journal*  
872 *of the Korean Association of Regional Geographers*, 21(2), 394-410.  
873 Available at <https://koreascience.kr/article/JAKO201507740043627.page>. [Accessed:  
874 2025-01-24]

875 Park, S. J., and Lee, D. K. (2021). Predicting susceptibility to landslides under climate change  
876 impacts in metropolitan areas of South Korea using machine learning. *Geomatics,*  
877 *Natural Hazards and Risk*, 12(1), 2462-2476.  
878 <https://doi.org/10.1080/19475705.2021.1963328>

879 Pham, B. T., Tien Bui, D., and Prakash, I. (2018). Bagging based support vector machines for  
880 spatial prediction of landslides. *Environmental Earth Sciences*, 77, 1-17.  
881 <https://doi.org/10.1007/s12665-018-7268-y>

882 Phillips, C., Hales, T., Smith, H., and Basher, L. (2021). Shallow landslides and vegetation at the  
883 catchment scale: A perspective. *Ecological Engineering*, 173, 106436.  
884 <https://doi.org/10.1016/j.ecoleng.2021.106436>

885 Pisner, D. A., and Schnyer, D. M. (2020). Support vector machine. In *Machine learning* (pp. 101-  
886 121). Academic Press. <https://doi.org/10.1016/B978-0-12-815739-8.00006-7>

887 Pourghasemi, H. R., and Rahmati, O. (2018). Prediction of the landslide susceptibility: Which  
888 algorithm, which precision?. *Catena*, 162, 177-192. [https://doi.org/10.1016/j.catena.](https://doi.org/10.1016/j.catena.2017.11.022)  
889 [2017.11.022](https://doi.org/10.1016/j.catena.2017.11.022)

890 Qiu, H., Regmi, A. D., Cui, P., Cao, M., Lee, J., and Zhu, X. (2016). Size distribution of loess  
891 slides in relation to local slope height within different slope morphologies. *Catena*, 145,  
892 155-163. <https://doi.org/10.1016/j.catena.2016.06.005>

893 R Core Team (2022). R: A language and environment for statistical computing. R Foundation for  
894 Statistical Computing, Vienna, Austria. Available at <https://www.R-project.org/>  
895 [Accessed 2025-01-24]

896 Rahman, M. S., Ahmed, B., and Di, L. (2017). Landslide initiation and runout susceptibility  
897 modeling in the context of hill cutting and rapid urbanization: a combined approach of  
898 weights of evidence and spatial multi-criteria. *Journal of Mountain Science*, 14(10),  
899 1919-1937. <https://doi.org/10.1007/s11629-016-4220-z>

900 Ran, Q., Wang, J., Chen, X., Liu, L., Li, J., and Ye, S. (2022). The relative importance of antecedent  
901 soil moisture and precipitation in flood generation in the middle and lower Yangtze River  
902 basin. *Hydrology and Earth System Sciences*, 26(19), 4919-4931.  
903 <https://doi.org/10.5194/hess-26-4919-2022>

904 Rathore, S. S., and Kumar, S. (2016). A decision tree regression-based approach for the number of  
905 software faults prediction. *ACM SIGSOFT Software Engineering Notes*, 41(1), 1-6.  
906 <https://doi.org/10.1145/2853073.2853083>

907 Razakova, M., Kuzmin, A., Fedorov, I., Yergaliev, R., and Ainakulov, Z. (2020). Methods of  
908 calculating landslide volume using remote sensing data. In E3S Web of Conferences (Vol.  
909 149, p. 02009). EDP Sciences. <https://doi.org/10.1051/e3sconf/202014902009>

910 Rosi, A., Peternel, T., Jemec-Auflič, M., Komac, M., Segoni, S., and Casagli, N. (2016). Rainfall  
911 thresholds for rainfall-induced landslides in Slovenia. *Landslides*, 13, 1571-1577.  
912 <https://doi.org/10.1007/s10346-016-0733-3>

913 Rotaru, A., Oajdea, D., and Răileanu, P. (2007). Analysis of the landslide movements. *International*  
914 *journal of geology*, 1(3), 70-79. Available at [https://naun.org/multimedia/NAUN/  
915 geology/ijgeo-10.pdf](https://naun.org/multimedia/NAUN/geology/ijgeo-10.pdf). [Accessed: 2025-01-24]

916 Saito, H., Korup, O., Uchida, T., Hayashi, S., and Oguchi, T. (2014). Rainfall conditions, typhoon  
917 frequency, and contemporary landslide erosion in Japan. *Geology*, 42(11), 999-1002.  
918 <https://doi.org/10.1130/G35680.1>

919 Saleh, A. M. E., Arashi, M., and Kibria, B. G. (2019). *Theory of ridge regression estimation with*  
920 *applications*. John Wiley and Sons.

921 Sato, T., Katsuki, Y., and Shuin, Y. (2023). Evaluation of influences of forest cover change on  
922 landslides by comparing rainfall-induced landslides in Japanese artificial forests with  
923 different ages. *Scientific reports*, 13(1), 14258. [https://doi.org/10.1038/s41598-023-  
924 41539-x](https://doi.org/10.1038/s41598-023-41539-x)

925 Scheidl, C., Heiser, M., Kamper, S., Thaler, T., Klebinder, K., Nagl, F., Lechner, L., Markart, G.,  
926 Rammer, W., and Seidl, R. (2020). The influence of climate change and canopy  
927 disturbances on landslide susceptibility in headwater catchments. *Science of the Total*  
928 *Environment*, 742, 140588. <https://doi.org/10.1016/j.scitotenv.2020.140588>

929 Seger, C. (2018). An investigation of categorical variable encoding techniques in machine  
930 learning: binary versus one-hot and feature hashing. Available at [https://www.diva-  
931 portal.org/smash/get/diva2:1259073/FULLTEXT01.pdf](https://www.diva-portal.org/smash/get/diva2:1259073/FULLTEXT01.pdf). [last accessed: 2025-01-24]

932 Shirzadi, A., Shahabi, H., Chapi, K., Bui, D. T., Pham, B. T., Shahedi, K., and Ahmad, B. B.  
933 (2017). A comparative study between popular statistical and machine learning methods  
934 for simulating volume of landslides. *Catena*, 157, 213-226. [https://doi.org/10.1016/  
935 j.catena.2017.05.016](https://doi.org/10.1016/j.catena.2017.05.016)

936 Singh, D., and Singh, B. (2022). Feature wise normalization: An effective way of normalizing data.  
937 *Pattern Recognition*, 122, 108307. <https://doi.org/10.1016/j.patcog.2021.108307>

938 Smith, H. G., Neverman, A. J., Betts, H., and Spiekermann, R. (2023). The influence of spatial  
939 patterns in rainfall on shallow landslides. *Geomorphology*, 437, 108795.  
940 <https://doi.org/10.1016/j.geomorph.2023.108795>

941 Stoof, C. R., Vervoort, R. W., Iwema, J., Van Den Elsen, E., Ferreira, A. J. D., and Ritsema, C. J.  
942 (2012). Hydrological response of a small catchment burned by experimental fire.  
943 *Hydrology and Earth System Sciences*, 16(2), 267-285. [https://doi.org/10.5194/hess-16-  
944 267-2012](https://doi.org/10.5194/hess-16-267-2012)

945 Sun, H. Y., Wong, L. N. Y., Shang, Y. Q., Shen, Y. J., and Lü, Q. (2010). Evaluation of drainage  
946 tunnel effectiveness in landslide control. *Landslides*, 7, 445-454.  
947 <https://doi.org/10.1007/s10346-010-0210-3>

948 Székely, G. J., Rizzo, M. L., and Bakirov, N. K. (2007). Measuring and testing dependence by  
949 correlation of distances. <https://doi.org/10.1214/009053607000000505>

950 Tacconi Stefanelli, C., Casagli, N., and Catani, F. (2020). Landslide damming hazard susceptibility  
951 maps: a new GIS-based procedure for risk management. *Landslides*, 17, 1635-1648.  
952 <https://doi.org/10.1007/s10346-020-01395-6>

953 Tsai, T. L., and Chen, H. F. (2010). Effects of degree of saturation on shallow landslides triggered  
954 by rainfall. *Environmental Earth Sciences*, 59, 1285-1295. <https://doi.org/10.1007/s12665-009-0116-3>

956 Turner, T. R., Duke, S. D., Fransen, B. R., Reiter, M. L., Kroll, A. J., Ward, J. W., Bach, J. L.,  
957 Justice, T. E., and Bilby, R. E. (2010). Landslide densities associated with rainfall, stand  
958 age, and topography on forested landscapes, southwestern Washington, USA. *Forest  
959 Ecology and Management*, 259(12), 2233-2247. [https://doi.org/10.1016/  
960 j.foreco.2010.01.051](https://doi.org/10.1016/j.foreco.2010.01.051)

961 Um, M. J., Yun, H., Cho, W., and Heo, J. H. (2010). Analysis of orographic precipitation on Jeju-  
962 Island using regional frequency analysis and regression. *Water Resources Management*,  
963 24, 1461-1487. <https://doi.org/10.1007/s11269-009-9509-z>

964 Van Westen, C. J. (2000). The modelling of landslide hazards using GIS. *Surveys in geophysics*,  
965 21(2), 241-255. <https://doi.org/10.1023/A:1006794127521>

966 Wang, D., Hollaus, M., Schmaltz, E., Wieser, M., Reifeltshammer, D., and Pfeifer, N. (2016). Tree  
967 stem shapes derived from TLS data as an indicator for shallow landslides. *Procedia Earth  
968 and Planetary Science*, 16, 185-194. <https://doi.org/10.1016/j.proeps.2016.10.020>

969 Wei, Z. L., Shang, Y. Q., Sun, H. Y., Xu, H. D., and Wang, D. F. (2019). The effectiveness of a  
970 drainage tunnel in increasing the rainfall threshold of a deep-seated landslide. *Landslides*,  
971 16, 1731-1744. <https://doi.org/10.1007/s10346-019-01241-4>

972 Wieczorek, G. (1987). *Debris flows/avalanches: process, recognition, and mitigation*, Volume VII,  
973 The Geological Society of America, Boulder, Colorado.

974 Willmott, C. J., and Matsuura, K. (2005). Advantages of the mean absolute error (MAE) over the  
975 root mean square error (RMSE) in assessing average model performance. *Climate  
976 Research*, 30(1), 79-82. <https://doi.org/10.3354/cr030079>

977 Yan, L., Xu, W., Wang, H., Wang, R., Meng, Q., Yu, J., and Xie, W. C. (2019). Drainage controls  
978 on the Donglingxing landslide (China) induced by rainfall and fluctuation in reservoir  
979 water levels. *Landslides*, 16, 1583-1593. <https://doi.org/10.1007/s10346-019-01202-x>

980 Yoon, S. S., and Bae, D. H. (2013). Optimal rainfall estimation by considering elevation in the  
981 Han River Basin, South Korea. *Journal of Applied Meteorology and Climatology*, 52(4),  
982 802-818. <https://doi.org/10.1175/JAMC-D-11-0147.1>

- 983 Yun, H. S., Um, M. J., Cho, W. C., and Heo, J. H. (2009). Orographic precipitation analysis with  
984 regional frequency analysis and multiple linear regression. *Journal of Korea Water*  
985 *Resources Association*, 42(6), 465-480. <https://doi.org/10.3741/JKWRA.2009.42.6.465>
- 986 Yune, C. Y., Jun, K. J., Kim, K. S., Kim, G. H., and Lee, S. W. (2010). Analysis of slope hazard-  
987 triggering rainfall characteristics in Gangwon Province by database construction. *Journal*  
988 *of the Korean Geotechnical Society*, 26(10), 27-38. [https://doi.org/10.7843/kgs.](https://doi.org/10.7843/kgs.2010.26.10.27)  
989 [2010.26.10.27](https://doi.org/10.7843/kgs.2010.26.10.27)
- 990 Zaruba, Q., and Mencl, V. (2014). *Landslides and their control*. Elsevier. ISBN 0444600760,  
991 9780444600769
- 992 Zhang, K., Wang, S., Bao, H., and Zhao, X. (2019). Characteristics and influencing factors of  
993 rainfall-induced landslide and debris flow hazards in Shaanxi Province, China. *Natural*  
994 *Hazards and Earth System Sciences*, 19(1), 93-105. [https://doi.org/10.5194/nhess-19-93-](https://doi.org/10.5194/nhess-19-93-2019)  
995 [2019](https://doi.org/10.5194/nhess-19-93-2019)