

Chapman University

## Chapman University Digital Commons

---

Mathematics, Physics, and Computer Science  
Faculty Articles and Research

Science and Technology Faculty Articles and  
Research

---

3-27-2024

### Enhancing Landslide Susceptibility Modelling Through a Novel Non-landslide Sampling Method and Ensemble Learning Technique

Chao Zhou

Yue Wang

Ying Cao

Ramesh P. Singh

Bayes Ahmed

*See next page for additional authors*

Follow this and additional works at: [https://digitalcommons.chapman.edu/scs\\_articles](https://digitalcommons.chapman.edu/scs_articles)



Part of the [Environmental Indicators and Impact Assessment Commons](#), [Environmental Monitoring Commons](#), [Geology Commons](#), [Geophysics and Seismology Commons](#), and the [Other Earth Sciences Commons](#)

---

---

# Enhancing Landslide Susceptibility Modelling Through a Novel Non-landslide Sampling Method and Ensemble Learning Technique

## Comments

This article was originally published in *Geocarto International*, volume 39, issue 1, in 2024. <https://doi.org/10.1080/10106049.2024.2327463>

## Creative Commons License



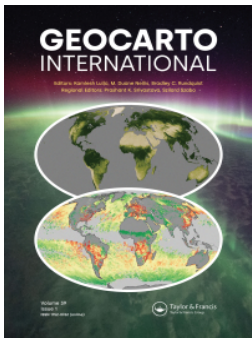
This work is licensed under a [Creative Commons Attribution-Noncommercial 4.0 License](https://creativecommons.org/licenses/by-nc/4.0/)

## Copyright

The authors

## Authors

Chao Zhou, Yue Wang, Ying Cao, Ramesh P. Singh, Bayes Ahmed, Mahdi Motagh, Yang Wang, Ling Chen, Guangchao Tan, and Shanshan Li



## Enhancing landslide susceptibility modelling through a novel non-landslide sampling method and ensemble learning technique

Chao Zhou, Yue Wang, Ying Cao, Ramesh P. Singh, Bayes Ahmed, Mahdi Motagh, Yang Wang, Ling Chen, Guangchao Tan & Shanshan Li

To cite this article: Chao Zhou, Yue Wang, Ying Cao, Ramesh P. Singh, Bayes Ahmed, Mahdi Motagh, Yang Wang, Ling Chen, Guangchao Tan & Shanshan Li (2024) Enhancing landslide susceptibility modelling through a novel non-landslide sampling method and ensemble learning technique, Geocarto International, 39:1, 2327463, DOI: [10.1080/10106049.2024.2327463](https://doi.org/10.1080/10106049.2024.2327463)

To link to this article: <https://doi.org/10.1080/10106049.2024.2327463>



© 2024 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group



Published online: 27 Mar 2024.



Submit your article to this journal [↗](#)



Article views: 27



View related articles [↗](#)



View Crossmark data [↗](#)



# Enhancing landslide susceptibility modelling through a novel non-landslide sampling method and ensemble learning technique

Chao Zhou<sup>a,b</sup>, Yue Wang<sup>c</sup>, Ying Cao<sup>c</sup>, Ramesh P. Singh<sup>d</sup>, Bayes Ahmed<sup>e</sup>,  
Mahdi Motagh<sup>b,f</sup>, Yang Wang<sup>c</sup>, Ling Chen<sup>c</sup>, Guangchao Tan<sup>g</sup> and Shanshan Li<sup>g</sup>

<sup>a</sup>School of Geography and Information Engineering, China University of Geosciences, Wuhan, China;

<sup>b</sup>Department of Geodesy, Section of Remote Sensing and Geoinformatics, Helmholtz Centre Potsdam, GFZ German Research Centre for Geosciences, Potsdam, Germany; <sup>c</sup>Engineering Faculty, China University of Geosciences, Wuhan, China; <sup>d</sup>School of Life and Environmental Sciences, Schmid College of Science and Technology, Chapman University, Orange, CA, USA; <sup>e</sup>Institute for Risk and Disaster Reduction, University College London (UCL), London, UK; <sup>f</sup>Institute for Photogrammetry and Geoinformation, Leibniz University Hannover, Hannover, Germany; <sup>g</sup>Hydrogeology and Engineering Geology Institute of Hubei Geological Bureau, Jingzhou, China

## ABSTRACT



In recent years, several catastrophic landslide events have been observed throughout the globe, threatening to lives and infrastructures. To minimize the impact of landslides, the need of landslide susceptibility map is important. The study aims to extract high-quality non-landslide samples and improve the accuracy of landslide susceptibility modelling (LSM) outcomes by applying a coupled method of ensemble learning and Machine Learning (ML). The Zigui-Badong section of the Three Gorges Reservoir area (TGRA) in China was considered in the present study. Twelve influencing factors were selected as inputs for LSM, and the relationship between each causal factor and landslide spatial development was quantitatively analyzed. A total of 179 landslides have been used in the present study. About 70% of the landslide pixels were randomly considered for training, and the remaining 30% were used for validation. Logistic Regression (LR) model was applied to produce an initial susceptibility map, and the non-landslide samples were selected within the classified low-susceptibility zone. Subsequently, two ML classifiers – the Classification and Regression Tree (CART), and the Multi-Layer Perceptron (MLP), and four coupling models – the CART-Bagging, CART-Boosting, MLP-Bagging, and MLP-Boosting, were utilized for LSM. Finally, the receiver operating characteristics (ROC) curve and statistical analysis were applied for accuracy assessment. The results show that altitude and distance to rivers were the main causal factors of landslides in the study area. The LR-MLP-Boosting performed the best with an accuracy of 0.986 followed by the LR-CART-Bagging, LR-CART-Boosting, and LR-MLP-Bagging. Accuracy comparisons demonstrate that ensemble learning algorithm can

## ARTICLE HISTORY

Received 28 November 2023  
Accepted 1 March 2024

## KEYWORDS

Reservoir landslides;  
susceptibility mapping;  
non-landslide sampling;  
ensemble learning; machine  
learning

**CONTACT** Ying Cao  [caoying@cug.edu.cn](mailto:caoying@cug.edu.cn)  Engineering Faculty, China University of Geosciences, Wuhan, China

© 2024 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited. The terms on which this article has been published allow the posting of the Accepted Manuscript in a repository by the author(s) or with their consent.

notably enhance the LSM performance of ML classifiers, and the Boosting algorithm marginally outperforms the Bagging algorithm. Moreover, the LR model can effectively constrain the selection range of non-landslide samples. The non-landslide sampling method constrained by LR yields higher quality samples compared to traditional random sampling method with no constraints, which develops a more excellent LSM.

## 1. Introduction

Landslides are severe, sudden, and frequent geological disasters that occur throughout the globe, taking huge loss of life, and infrastructures that affect the day-to-day life. The Ministry of Natural Resources of China reported that 4,810 landslide disasters occurred in 2020, resulting in 139 deaths and 730 million US dollars in direct economic losses. The Three Gorges Reservoir area (TGRA) is a highly landslide-prone area, with more than 5,000 landslide occurrences recorded (Zhou et al. 2022a). During the first impoundment of the TGRA in July 2003, the Qianjiangping landslide caused a death toll of 24 and an economic loss of about 11.6 million US dollars. Landslide risk assessment has been widely used as a vital means of disaster prevention and mitigation. It is the foundation for quantitative risk assessment and the final land use and planning maps. However, due to the nonlinear relationship between landslide occurrence and their influencing factors, accurate landslide susceptibility modelling (LSM) is challenging for geoscientists and engineers.

Over the past two decades, numerous qualitative and quantitative methods have been developed for LSM (Sabokbar et al. 2014; Zhou et al. 2018; Yavuz Ozalp et al. 2023; Liu et al. 2024). In the qualitative methods, the weight of various controlling factors is determined by engineering geologists and geotechnical engineers based on their past experiences. These methods require landslide-vulnerable areas based on past landslide events, geology, and slope. The qualitative methods include expert scoring and analytic hierarchy methods based on numerous controlling landslide parameters (Meena et al. 2022; Roy et al. 2023). The quantitative method is divided based on data-driven and physically driven. The advancement of earth observation techniques has markedly enhanced data quality, particularly in landslide catalogues and topographic landforms. This improvement has propelled the widespread adoption of data-driven methods in LSM. The Machine Learning (ML) technique has a strong nonlinear fitting ability that are being used in various fields including LSM (Zhou et al. 2020a; 2016; Yavuz Ozalp et al. 2023). The popular ML methods include support vector machines, artificial neural networks, random forest, extreme gradient Boosting, CatBoost (Yu et al. 2022; Zhou et al. 2022b; Yavuz Ozalp et al. 2023; Liu et al. 2024).

Machine learning methods are reported to surpass traditional approaches in LSM studies. Contemporary studies additionally indicate that a sole machine learning classifier frequently encounters challenges in achieving optimal performance within the context of LSM (Akinci and Zeybek, 2021; Long et al. 2021; Tanyu et al. 2021). Users of models may lack familiarity with the historical performance of an individual machine learning classifier, posing a challenge in assessing whether a modelling method adequately supports specific susceptibility assessments (Hu et al. 2021). This increases the possibility of inappropriate model selection. Ensemble learning techniques address the limitations of individual classifiers by combining the predictions of multiple classifiers, resulting in more robust and accurate predictions. This approach has gained widespread attention

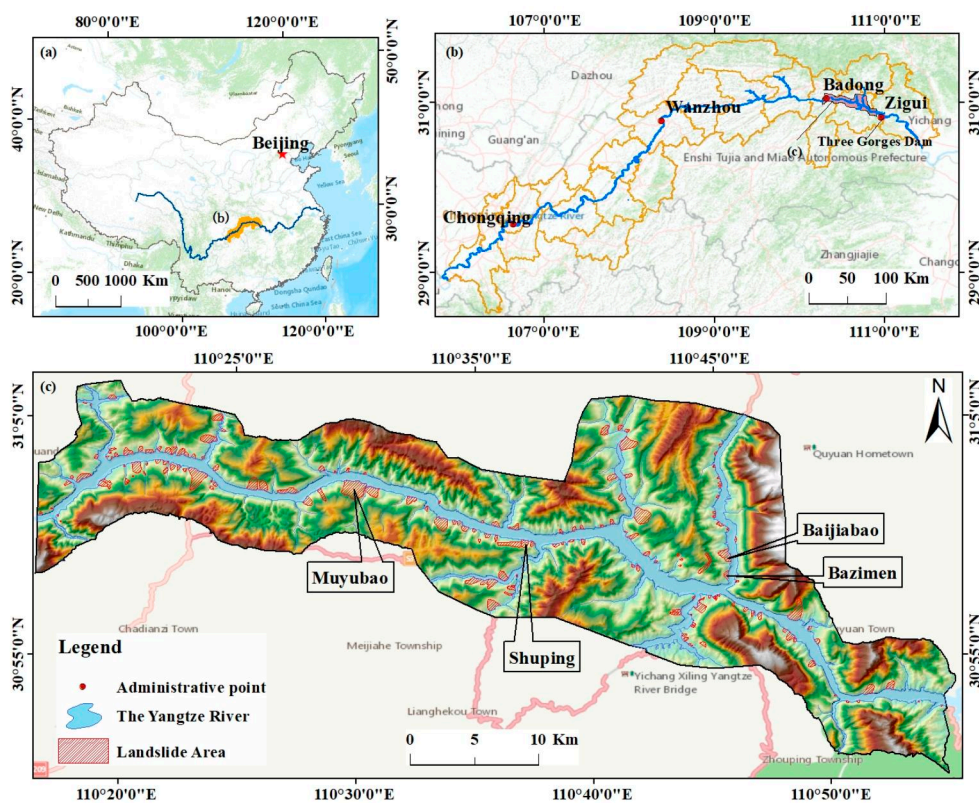
across various fields (Mohammed and Kora, 2023). Recently, ensemble learning has also gained popularity in LSM with impressive research results (Sharma et al. 2024). Song et al. (2020) introduced a stacking ensemble learning framework consisting of multiple tree-based classifiers as base learners and logistic regression as a meta-learner in a two-layer structure. Hong (2023a) applied best-first decision tree as the base classifier and explored the effect of coupling with five ensemble algorithms in LSM. However, there is no consensus on the applicability of ensemble machine learning in LSM. To date, most of the base classifiers for ensemble learning are tree-based algorithms, and less attention is paid to the performance of ensemble models derived from MLP as the base classifier (Hong, 2023b). Many studies concentrate on comparing the prediction accuracy of ensemble algorithms with single models and among different ensemble learning methods. Few studies focus on the impact of key parameters on ensemble learning modelling performance, such as the relationship between the number of base classifiers and LSM accuracy.

The selection of high-quality training samples is also crucial for enhancing LSM. At the regional scale, the non-landslide area is much larger than the landslide area, even if landslides are extremely developed. How to conduct high-quality non-landslide sampling in landslide-free areas is a focus of LSM research (Zhu et al. 2019; Hong et al. 2024). The most used method currently is to randomly sample non-landslide pixels from the whole landslide-free areas with no constraint. The selected non-landslide samples may exhibit similarities to the geo-environmental context of the landslide occurrence area, adversely impacting the accuracy of susceptibility modelling. The buffer-controlled (Gameiro et al. 2021; Lucchese et al. 2021) and slope-controlled (Kavzoglu et al. 2014) methods enable the acquisition of high-quality non-landslide samples. However, various conditions, such as low slope-gradient and hard rock properties, can serve as controlling factors against landslides. This kind of method employs simple rules for non-landslide sampling and is not able to fully capture the non-landslide development pattern across the entire study area. Acquiring high-quality non-landslide samples that encompass various control patterns remains a pressing issue to be addressed for a better LSM.

The Zigui-Badong section which is located at the head area of the TGRA is considered for the detailed LSM study (Figure 1). Over the past two decades, the combination of precipitation and periodic water level fluctuations in the reservoir has led to significant number of landslides. In this study, twelve controlling landslide factors are statistically analyzed and selected as inputs for modelling. An initial landslide susceptibility map is produced using Logical Regression (LR), and the non-landslide training samples are selected in the low susceptibility area. Two single models, namely Classification and Regression Tree (CART), Multi-Layer Perceptron (MLP), and four coupling models (CART-Bagging, CART-Boosting, MLP-Bagging, and MLP-Boosting) were utilized for LSM. Finally, the modelling performance is compared by Receiver Operating Characteristic (ROC) Curve and statistical analysis method. Our results attempt to develop a high-accurate susceptibility model for reservoir landslides in the TGRA.

## 2. Study area

The study area is located within Zigui and Badong counties, (latitude  $30^{\circ}51' - 31^{\circ}4' N$ , and longitude  $100^{\circ}17' - 100^{\circ}52' E$ ) with the total area about  $656 \text{ km}^2$  (Figure 1). It is a high-prone area for landslide disasters with an altitude range of 80–2,020 m. The geological structure in the study area is complex, with developed faults and fragmented rock mass. Triassic and Jurassic dominate the stratum, and the lithology is mostly carbonate, sand



**Figure 1.** (a) Map of China, (b) map shows Three gorges reservoir area, and (c) topography and landslide distribution in the study area. The background maps are made using ArcGIS.

shale, marlstone, and mudstone, which is sensitive to landslide occurrence. Quaternary is widely exposed in the study area and accumulates on the terraces and slope surfaces. In addition, the study area experiences excessive rainfall, with an average annual rainfall of 1,250 mm, mainly during May to September. In 2003, the TGRA was first impounded up to 135 m. After September 2008, the reservoir water level periodically varies between 145 m – 175 m per year, significantly changing the bank slope's hydrogeological conditions (Cao et al. 2016; Zhou et al. 2020b). The water level fluctuations induce the deformation and trigger a large number of reservoir landslides, such as the Qianjiangping landslide, the Muyubao landslide, and the Shuping landslide.

### 3. Methodology

#### 3.1. Developing LSM

The development of LSM includes four parts (Figure 2): influencing factor selection and landslide pixel sampling, non-landslide pixel sampling, model construction, and accuracy evaluation. First, we have considered the influencing factors for LSM; 70% of the landslide pixels were selected as the training data, while the remaining 30% was applied for validation. Second, we produce a preliminary susceptibility map using LR and non-landslide pixels with an equal number of landslide pixels are randomly selected from the low



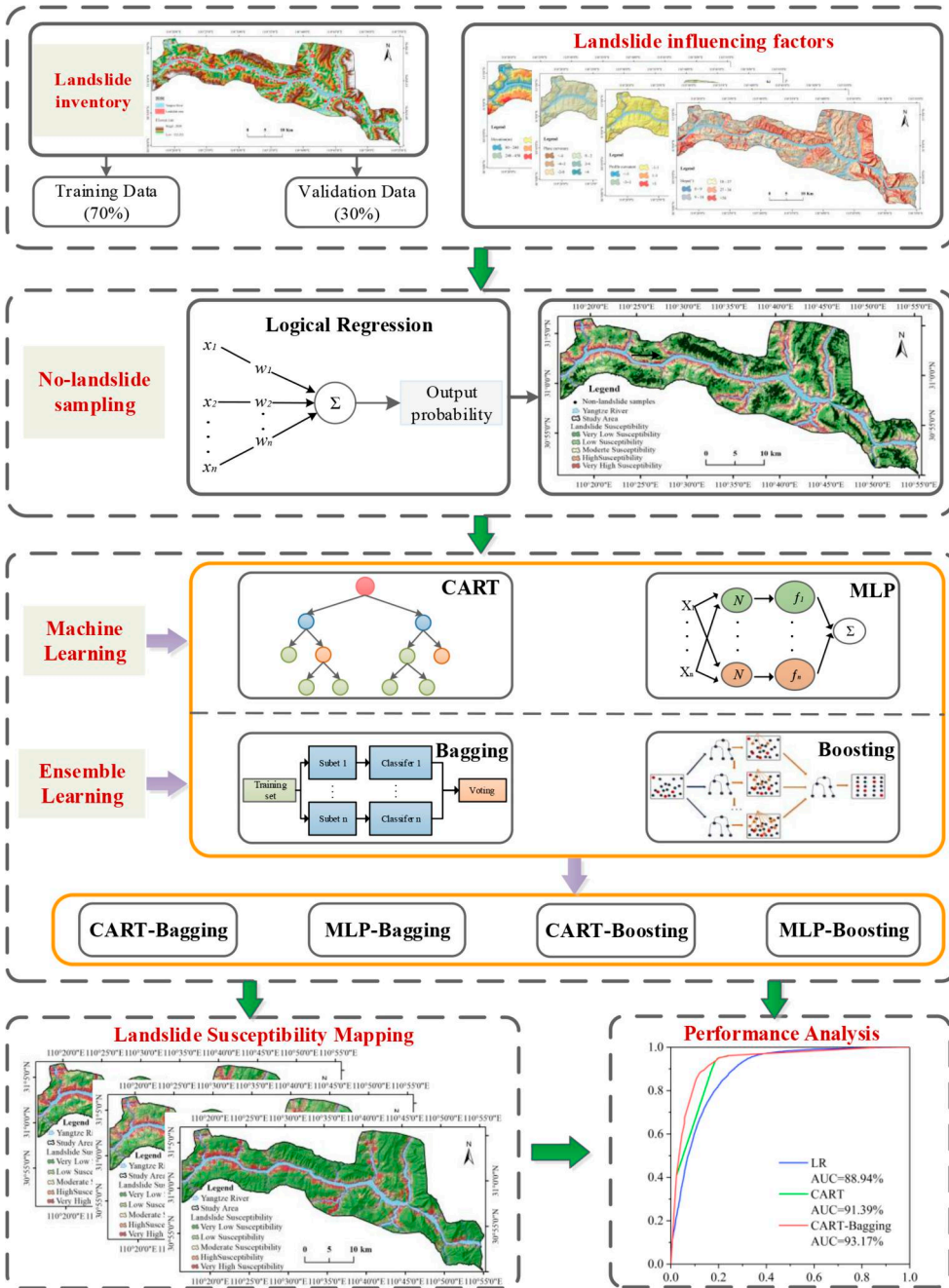


Figure 2. The flowchart of the landslide susceptibility mapping in this study.

susceptibility area. Third, two single classifiers (CART and MLP) and four coupling models (CART-Bagging, CART-Boosting, MLP-Bagging, and MLP-Boosting,) are applied for LSM. Finally, we use statistical analysis method and ROC curves to evaluate the partitioning results and model performance.



### 3.2. Information value method

Information value is a statistical method based on information theory. In this study, we apply this method to quantitatively assess the influence of various factors on landslide occurrence. In LSM, the information value is given as follows:

$$I_i = \sum_{i=1}^n \text{Ln} \frac{S_i/S}{A_i/A} \quad (1)$$

where  $I_i$  is the information value of the  $i$ -th influencing factor;  $S_i$  is the number of landslide pixels within the  $i$ -th influencing factor;  $S$  is the total number of landslide pixels;  $A_i$  is the number of pixels for the  $i$ -th influencing factor;  $A$  is the total number of pixels in the study area;  $n$  is the number of influencing factors (Zhou et al. 2018). When the information value is greater than 0, the factor promotes the occurrence of landslides. Conversely, when the value is less than 0, that indicates the factor inhibits the occurrence of landslides. Moreover, the larger the absolute information value, the stronger the effect.

### 3.3. Classifier

#### 3.3.1. Multi-Layer Perceptron

Multi-Layer Perceptron (MLP) is a feed-forward artificial neural network widely used in many fields. It consists of three layers: input layer, hidden layer, and output layer (Figure 3). MLP with a sufficient number of hidden layer neurons can realize any nonlinear mapping from  $n$ -dimensional to  $m$ -dimensional (Gardner and Dorling 1998). During the calculation process, the input layer neurons receive sample data, and the hidden layer and the output layer neurons deal with the inputs according to the weight value. To build a better model, MPL modifies the weight value through backpropagation. The learning process of MPL models is constantly adjusting the parameters, and the training of the MPL model is the process of constantly adjusting network parameters.

#### 3.3.2. Classification and Regression Tree

Classification and Regression Tree (CART) is a simple but powerful approach to forecast an event. This method is easy to understand and implement for LSM. During modeling, CART does not need to presuppose a relationship between the predictor and target variables. The child nodes are obtained to form a binary tree by recursively dividing the data set. The child nodes are continuously expanded to generate a complete decision tree

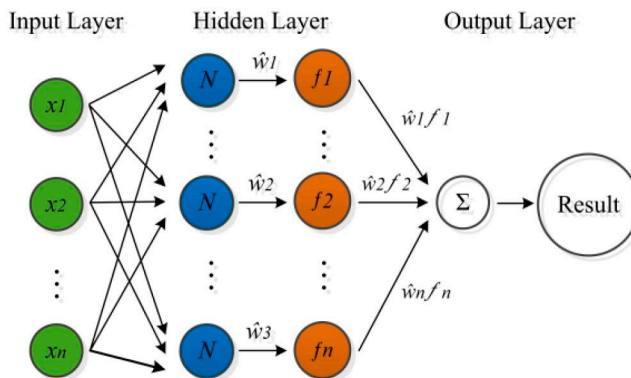


Figure 3. MLP neural network structure.

and perform the necessary pruning to prevent overfitting. The Gini index minimization criterion determines the optimal segmentation point, and the smaller the Gini index is, the better the effect of tree division. The Gini index represents the classification error rate for the binary classification problem. For example, if the sample set  $D$  contains  $k$  categories, the Gini coefficient of the sample set can be expressed as:

$$Gini(D) = 1 - \sum_{i=1}^k \left( \frac{C_i}{D} \right)^2 \quad (2)$$

where  $C_i$  is a subset of class  $i$  samples in  $D$ .

### 3.3.3. Logical regression

The logical regression (LR) model is a statistical analysis model suitable for binomial categorical dependent variables and is widely used in landslide and subsidence prediction (Youssef et al. 2016; Hu et al. 2021). Training testing on known landslide events establishes the nonlinear relationship between the dependent variable and multiple independent variables. Therefore, the occurrence probability of future landslides can be predicted or evaluated using the established formula. The method considers the landslide influencing factor as the independent variable and the occurrence probability of landslides as the dependent variable (landslide is 1, the non-landslide is 0). The independent variable can be continuous or discrete. Assuming that the probability of landslide occurrence is  $P$ , the regression equation can be written as follows:

$$P = \frac{1}{1 + e^{-(b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n)}} \quad (3)$$

where  $b_0$  is a constant,  $n$  is the number of independent variables,  $x_1, x_2, \dots, x_n$  is the landslide influencing factors, and  $b_1, b_2, \dots, b_n$  are the coefficients of LR.

## 3.4. Ensemble learning

### 3.4.1. Boosting

Boosting algorithm is a process of enhancing a simple weak classification algorithm to reduce variance and bias through iterative training to improve the ability to classify model data (Youssef et al. 2016). This algorithm generates a classifier combination through multiple iterations. Each iteration constructs a new training set from a sample returned to the total dataset, and each iteration adjusts the weight of the sample to get higher weight values at the next iteration. After  $T$  iterations, the updated weak classifiers are weighted and superimposed to obtain the strong classifier (Figure 4).

### 3.4.2. Bagging

The Bagging algorithms an ensemble learning technique aimed at improving the stability and accuracy of machine learning models. Its main idea is to repeat the input training sets by Bootstrap sampling to obtain  $n$  subsets and build a weak classifier for each subset. The voting method integrates the weak ( $n$ ) classifiers to form a strong classifier. The Bagging algorithm can observe small changes in the training data, effectively improving the accuracy and stability of the model prediction results, especially for models susceptible to sample disturbances (Figure 5).

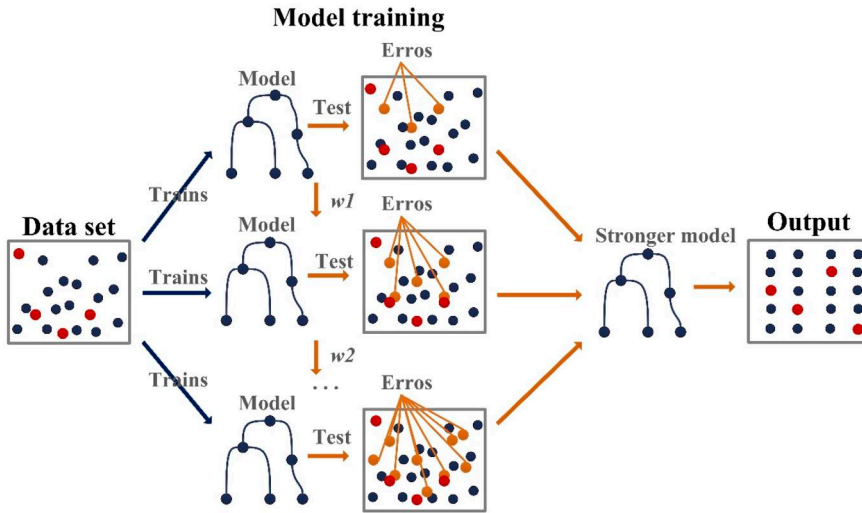


Figure 4. Schematic diagram of Boosting ensemble algorithm.

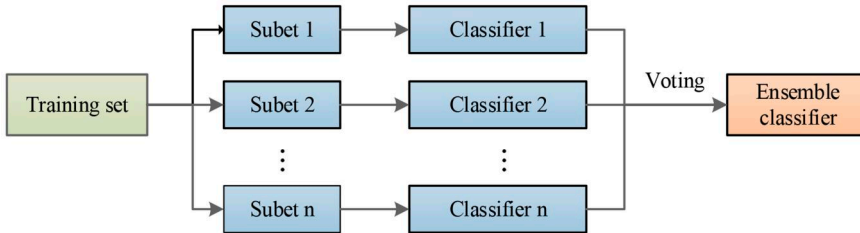


Figure 5. Flowchart of the Bagging ensemble algorithm.

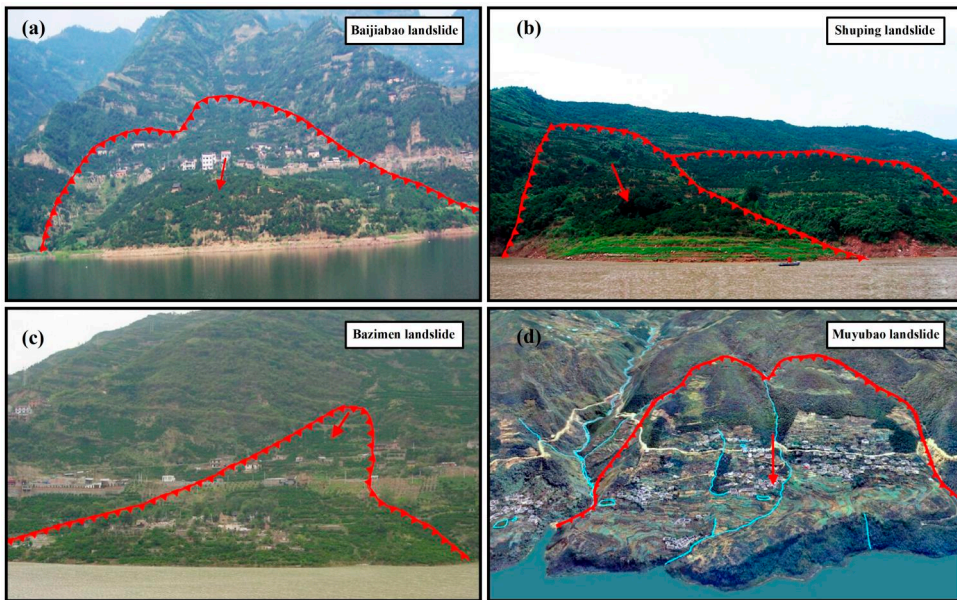
## 4. Modelling and results

### 4.1. Data used

Accurate and reliable landslide inventory data are essential for LSM, here we prepared the landslide inventory using field data, historical landslide inventory, and high-resolution remote sensing images. The data source of this study includes: (a) a topographic map (1:10,000) for extraction of topography, landscape, and rivers; (b) a geological map (1:50,000) for extraction of lithology, geologic structure, faults, and so on; (c) field investigation data; and (d) the historical landslide inventory. A total of 179 landslides with a total area of 22.14 km<sup>2</sup> are obtained in the study area, shown as dots or bands along the Yangtze River (Figure 1). In contrast, the area of individual landslides ranges from 0.13 km<sup>2</sup> to 1.80 km<sup>2</sup>, with four typical reservoir landslides of Baijiabao, Shuping, Bazimen, and Muyubao landslides (Figure 6).

### 4.2. Analysis of influencing factors

Different influencing factors cause landslide occurrence in various regions due to the diverse geological environments. With the consideration of the regional geological conditions, landslide inventory, and earlier studies (Yu et al. 2019), 12 influencing factors were prepared initially for LSM, namely altitude, slope, aspect, terrain roughness index (TRI), topographic relief, slope geometry, slope structure, lithology, topographic wetness index



**Figure 6.** Reservoir landslides in the Three Gorges reservoir area (Figure 1): (a) Baijiabao landslide, (b) Shuping landslide, (c) Bazimen landslide, and (d) Muyubao landslide.

(TWI), landuse, distance to rivers, and distance to faults. According to the Technical requirement for the geo-hazard survey (1:50,000) of China Geological Survey, the raster of  $30\text{ m} \times 30\text{ m}$  is adopted as the basic unit for LSM. All layers of twelve influencing factors are extracted in *ArcGIS 10.2*.

#### 4.2.1. Altitude

In the study area, there are many new infrastructures are developed, such as roads, bridges, electric lines, telephone lines etc. at low altitudes severely affecting the stability of natural slopes. The altitude of the study area varies from 145 to 2,020 m, which is divided into five classes: (145 ~ 240), (240 ~ 450), (450 ~ 650), (650 ~ 1200), (1200 ~ 2020) (Figure 7a). As given in Table 1, landslides mainly occur in the altitude range of 145 ~ 240 m, the information value of which is the highest of 1.49. No landslides occur in regions above an altitude of 1200 m, since the slopes are observed to remain undisturbed.

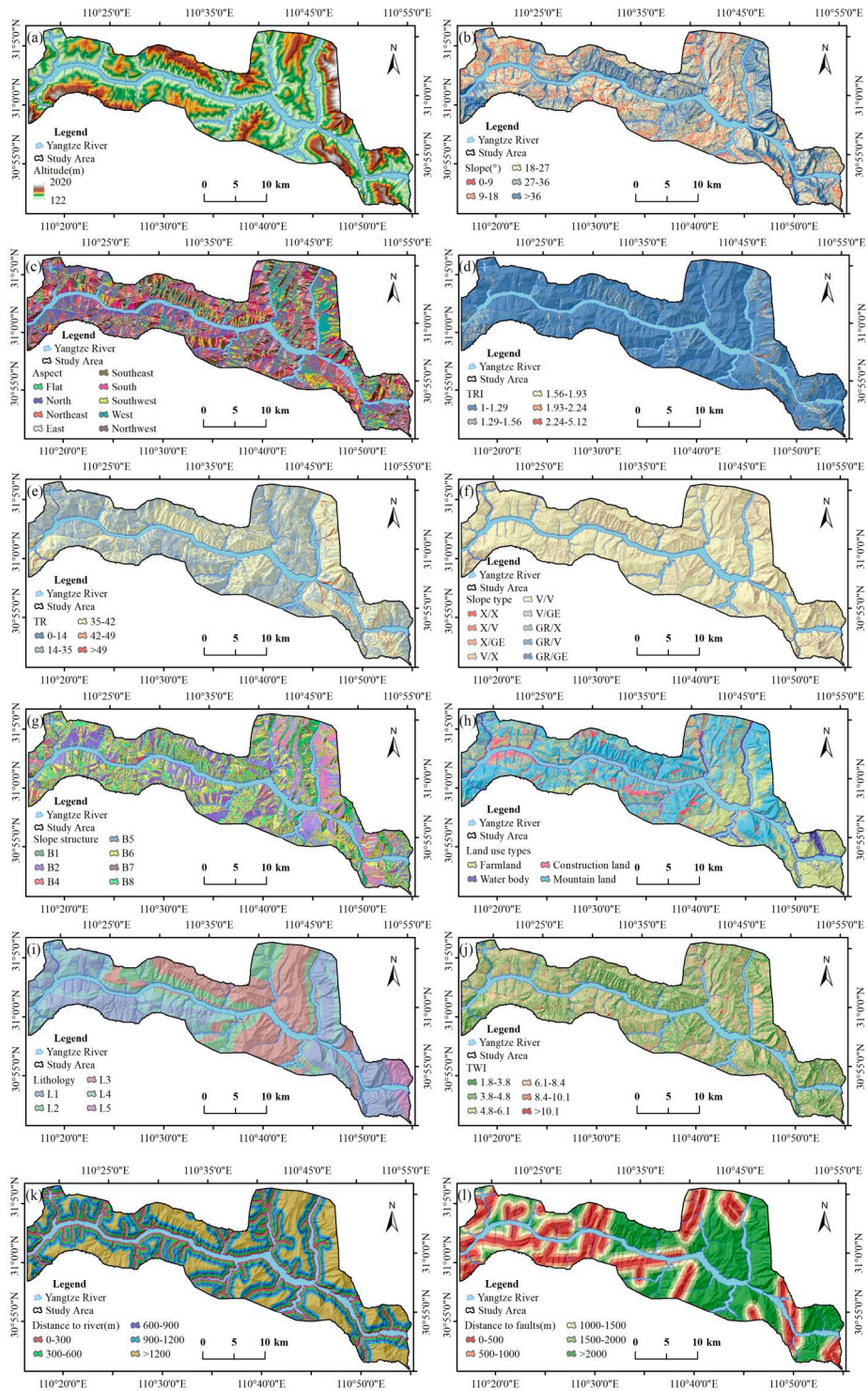
#### 4.2.2. Slope

The slope affects the stress distribution, materials accumulation, and surface runoff. It can be divided into five classes: very gentle ( $0^\circ$ ,  $9^\circ$ ), gentle ( $9^\circ$ ,  $18^\circ$ ), moderate ( $18^\circ$ ,  $27^\circ$ ), steep ( $27^\circ$ ,  $36^\circ$ ), and very steep ( $36^\circ$ ,  $90^\circ$ ) (Figure 7(b)). Landslides generally occur on the gentle slope, whose information value is 0.42. When the slope is more than  $36^\circ$ , the occurrence of a landslide is significantly inhibited, and the information value is the lowest of  $-1.63$ .

#### 4.2.3. Aspect

Rainfall and sunlight exposure vary according to the aspect, leading to contrasting differences in the physical and mechanical properties of sliding masses. These differences may arise from the weathering impacts, ultimately affecting the stability of landslides (Pham et al. 2021). The aspect in this study was divided into nine classes. The slopes with north





**Figure 7.** Shows various landslide factors: (a) altitude, (b) slope, (c) aspect, (d) TRI, (e) topographic relief, (f) slope geometry, (g) slope structure, (h) landuse, (i) lithology, (j) TWI, (k) distance to rivers, and (l) distance to faults.

**Table 1.** Statistics between causal factors and landslides occurrence.

Causal factor	Category	IV	Causal factor	Category	IV
Slope (°)	0 ~ 9	-0.90	Aspect	flat	-1.93
	9 ~ 18	0.37		north	0.47
	18 ~ 27	0.42		northeast	0.23
	27 ~ 36	-0.28		east	-0.27
	>36	-1.63		southeast	-0.31
TRI	1 ~ 1.1	0.37		south	0.43
	1.1 ~ 1.2	-0.04		southwest	-0.29
	1.2 ~ 1.3	-0.88		west	-0.71
	1.3 ~ 1.4	-1.9	northwest	-0.03	
	1.4 ~ 1.5	-2.48	Distance to faults (m)	0-500	-0.04
>1.5	-2.49	500-1,000		0.08	
Distance to rivers (m)	0 ~ 300	0.92		1,000-1,500	0.24
	300 ~ 600	0.29	1,500-2,000	0.23	
	600 ~ 900	-0.57	>2,000	-0.22	
	900 ~ 1,200	-1.7	Altitude (m)	<240	1.49
>1,200	-2.95	240-450		0.54	
Lithology	L1	-3.94		450-650	-1.36
	L2	-0.3		650-1,200	-3.84
	L3	0.17	>1,200	-∞	
	L4	-0.34	Slope geometry	X/X	0.16
	L5	0.42		X/V	-0.96
TWI	1.37 ~ 3	-2.49		X/GE	-1.45
	3 ~ 4.5	-0.37		V/X	0.04
	4.5 ~ 6	-0.20		V/V	-1.74
	6 ~ 7.5	0.47	V/GE	-1.19	
	7.5 ~ 9	0.72	GR/X	0.01	
>9	0.02	GR/V	-1.02		
Topographic relief (m)	0-14	-0.63	GR/GE	-1.46	
	14-35	0.47	Slope structure	B1	-∞
	35-42	0.08		B2	0.14
	42-49	-0.47		B4	0.13
	>49	-1.70		B5	0.08
Landuse	Mountain land	-0.59		B6	-0.17
	Farmland	0.12	B7	-0.34	
	Waterbody	0.43	B8	-0.67	
	Construction land	0.85			

Note: IV refers to Information Value; The meaning of the lithology, slope structure, and slope type abbreviation of the formation is given in Tables 2, 3 and 4.

and northeast aspects are more prone to landslides, their information values are 0.47 and 0.23, respectively (Figure 7(c)).

#### 4.2.4. Topographic relief

The relative height difference in the study area is computed using following equation:

$$D = Hmax - Hmin \tag{4}$$

where,  $D$  is the relief factor;  $Hmax$  is the highest altitude value;  $Hmin$  is the lowest altitude value. The topographic relief in this study is divided into five classes: (0 ~ 14 m), (14 ~ 35 m), (35 ~ 42 m),  $\geq 49$  m (Figure 7(d)). The information values are -0.63, 0.47, 0.08, -0.47 and -1.70, respectively.

#### 4.2.5. Slope geometry

Slope curvature is the microscopic performance of the earth’s surface landforms, divided into plane curvature and profile curvature. It reflects the concavity of the slope along the aspect, which controls the flow speed of surface material and rainfall confluence. We classify the plane and profile curvatures into three classes respectively. The nine slope



**Table 2.** Definition for slope geometry classification.

Profile curvature	Plan curvature		
	Outward slope (X)	Inward slope (V)	Straight slope (GR)
Convex slope (X)	X/X	V/X	GE/X
Concave slope (V)	X/V	V/V	GE/V
Straight slope (GE)	X/GR	V/GR	GE/GR

**Table 3.** Classification of slope structure (Zhou et al. 2018).

Category	Definition (slope: $\theta$ , aspect: $\sigma$ , bed dip angle: $\alpha$ , bed dip direction: $\beta$ )
B1	$a < 10^\circ$
B2	$(( \alpha - \beta  \in (0, 30^\circ)) \vee ( \alpha - \beta  \in (330^\circ, 360^\circ))) \wedge (\alpha > 10^\circ) \wedge (\theta > \alpha)$
B3	$(( \alpha - \beta  \in (0, 30^\circ)) \vee ( \alpha - \beta  \in (330^\circ, 360^\circ))) \wedge (\alpha > 10^\circ) \wedge (\theta = \alpha)$
B4	$(( \alpha - \beta  \in (0, 30^\circ)) \vee ( \alpha - \beta  \in (330^\circ, 360^\circ))) \wedge (\alpha > 10^\circ) \wedge (\theta < \alpha)$
B5	$( \alpha - \beta  \in (30^\circ, 60^\circ)) \vee ( \alpha - \beta  \in (330^\circ, 360^\circ))$
B6	$( \alpha - \beta  \in (60^\circ, 120^\circ)) \vee ( \alpha - \beta  \in (240^\circ, 300^\circ))$
B7	$( \alpha - \beta  \in (90^\circ, 150^\circ)) \vee ( \alpha - \beta  \in (210^\circ, 240^\circ))$
B8	$( \alpha - \beta  \in (120^\circ, 180^\circ)) \vee ( \alpha - \beta  \in (180^\circ, 210^\circ))$

geometries are defined with the combination of plane and profile curvatures (Table 2, Figure 7€). In this study area, landslides mainly occur in the slope type of X/X, with the highest information value of 0.16.

#### 4.2.6. Terrain roughness index

The terrain roughness index (TRI) reflects the degree of surface fluctuation and erosion, which is computed using the following equation:

$$TRI = \sqrt{Abs(\max^2 - \min^2)} \quad (5)$$

The TRI is divided into six classes: (1 ~ 1.1), (1.1 ~ 1.2), (1.2 ~ 1.3), (1.3 ~ 1.4), (1.4 ~ 1.5), and  $\geq 1.5$  (Figure 7(f)). The information values are 0.37, -0.04, -0.88, -1.90, -2.48, and -2.49, respectively (Table 1).

#### 4.2.6. Landuse

Landuse and landslide development are closely related to the triggering of landslides due to the changes in the slope. The landuse is divided into four categories, namely water bodies, construction land, farmland, and mountain land (Figure 7(g)). In the foothills and mountainous areas, landuse and landcover is changing, affecting the slope that triggers the landslides. Further, the construction land is mainly concentrated in the gentle river terraces on both sides of the Yangtze River. A large number of excavations, slope cutting, and other activities in the construction of houses and roads directly impact the slope's stability, and the information value is the highest of 0.85.

#### 4.2.8. Slope structure

Slope structure indicates the intersection relationship between strata and slope, which determines the direction of the sedimentary stack on the slope. The slope structure is divided according to Table 3 (Figure 7(h)). In this study area, landslides mainly occurred in the B2 region. 35.06% of the total pixels are distributed in this category, whose information value is 0.42.

**Table 4.** Lithological classification in this study area.

Category	Geologic group	Main lithology
L1	$\delta_{2-1}$ , Pt	Granite and diorite
L2	Z, $\epsilon_1$ , $\epsilon_{2+3}$ , O, T <sub>1j</sub> , T <sub>2b3</sub>	Limestone, Shale, Malmstone
L3	T <sub>1d</sub> , T <sub>2b4+5</sub> , J <sub>1x</sub> , J <sub>2s</sub> , J <sub>3s</sub>	Marl mudstone
L4	S, J <sub>2x</sub>	Shale, Mudstone and Shi Ying Sandstone, Muddy Siltstone, etc.
L5	T <sub>3s</sub> , J <sub>1-2n</sub> , J <sub>3p</sub>	Malmstone (Feldspar sandstone, Shi Ying Sandstone, etc.) with coal seam

**4.2.9. Lithology**

We divide the lithology in the study area into five classes (Table 4 and Figure 7(i)). In the stratiform structure containing weak strata, especially in the stratified clastic rocks and the carbonate rocks developed on the weak bedrock, the large and medium-sized landslides more formed, and its information value is 0.42. On the other hand, few landslides developed in the hard rocks, such as granite and diorite, with the information value being the lowest of -3.94.

**4.2.10. Topographic wetness index**

Topographic Wetness Index (TWI) reflects topography’s influence on soil water saturation. It can be calculated using the following formula:

$$TWI = \ln\left(\frac{A_s}{\tan\beta}\right) \tag{6}$$

where,  $A_s$  is the upstream gathering area and  $\beta$  is the slope. The TWI is divided into six classes: (1.37, 3), (3, 4.5), (4.5, 6), (6, 7.5), (7.5, 9), and (9,  $-\infty$ ] (Figure 7(j)). When the TWI value is within the range (7.5~9), it shows the most substantial positive influence on landslide occurrence, whose information value is the highest of 0.72.

**4.2.11. Distance to rivers**

This study area is obviously affected by the hydrogeological environment, whose main river system is the Yangtze River and its tributaries (Figure 1). After the impoundment of the TGRA, the stability of the bank slopes is influenced by the periodical fluctuation of reservoir water level, river erosion, and softening effect. The factor of distance to rivers represents the intensity of its influence. We divide the distance to rivers into four classes, namely (0~300 m), (300 m~900 m), (900 m~1,200 m), and  $\geq 1,200$  m (Figure 7(k)). The maximum information value is 0.92 within the distance range of 300 m. With the distance increasing, the influence of the river system on landslides gradually weakened, and the information value decreased.

**4.2.12. Distance to faults**

Due to the severely broken rock mass and intense tectonic movement in the area, minor faults develop and serve as triggers for landslides. The distance of the faults to the landslide prone area control the intensity of landslides. The distance to faults is classified into five categories, namely (0~500 m), (500~1,000 m), (1,000~1,500 m), (1,500~2,000 m), and  $\geq 2,000$  m (Figure 7(l)). Their information values are -0.04, 0.08, 0.24, 0.23, and -0.22, respectively.

**Table 5.** Multi-collinearity analysis of the causal factors.

Influencing factors	T	VIF
Elevation	0.43	2.30
Slope	0.26	3.80
Aspect	0.96	1.04
Terrain Roughness Index	0.31	3.18
Topographic relief	0.32	3.12
Slope shape	0.64	1.52
Landuse	0.81	1.22
Slope structure	0.26	3.80
Stratigraphic lithology	0.94	1.06
Topographic Wetness Index	0.75	1.34
Distance to rivers	0.45	2.25
Distance to faults	0.95	1.05

### 4.3. Landslide susceptibility modelling

#### 4.3.1. Multi-collinearity analysis

The collinearity factors will affect the performance of the evaluation model. Therefore, it is necessary to carry out a collinearity analysis prior to the LSM, to ensure that the factors are independent. The multi-collinear analysis is performed using Tolerance (T) and Variance Inflation Factor (VIF). When T is greater than 0.2 or VIF is less than 5, it is considered that there is no multi-collinearity between the factors. Both indices of T and VIF are calculated using *SPSS Statistics 26.0*, and the results are shown in [Table 5](#), showing that all the twelve factors are independent with no collinearity.

#### 4.3.2. Non-landslide sampling using LR

In this study, we randomly selected 70% of landslide pixels for training and the remaining 30% for validation, and the method of k-fold cross validation is applied for training. Simultaneously, the same number of non-landslide samples are selected for model training. We propose a non-landslide sampling method to extract high-quality samples by the LR algorithm. At first, we produce a preliminary landslide susceptibility map by randomly selecting non-landslide as an example in landslide-free areas, using the following equations:

$$\text{Logit}(P) = -10.87 + 2.175 \cdot x_1 + 1.17 \cdot x_2 + 6.028 \cdot x_3 + 1.079 \cdot x_4 + 0.750 \cdot x_5 + 1.071 \cdot x_6 + 0.987 \cdot x_7 + 0.600 \cdot x_8 - 1.263 \cdot x_9 + 0.672 \cdot x_{10} + 0.559 \cdot x_{11} + 0.814 \cdot x_{12} \quad (7)$$

where  $x_1, x_2, \dots, x_{12}$  are independent variables, which indicate the factor values of the slope, aspect, altitude, slope shape, landuse, topographic relief, TRI, TWI, slope structure, lithology, distance to rivers, distance to faults;  $P$  is the probability of landslide susceptibility. The landslide susceptibility index is divided into five levels: very high (5%), high (10%), medium (15%), low (20%), and very low (50%). The non-landslide samples for training are randomly selected only in the Very Low area. The preliminary susceptibility map ([Figure 8](#)) shows different classes of susceptibility (very high, high, moderate, low, very low and non-landslide) zone along both sides of the Yangtze River. The non-landslide samples are distributed throughout the study area, concentrated in the areas with high altitudes, steep slopes and few human engineering activities. Due to topographical and lithological constraints, landslides rarely develop in these areas. Therefore, the engineering geological conditions of the selected samples are quite different from those of the landslide, and they are more representative of landslide-free areas.

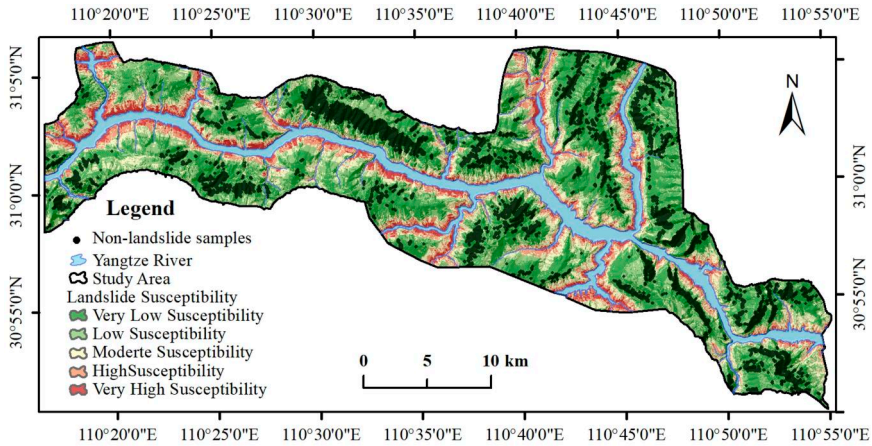


Figure 8. Preliminary susceptibility map and the distribution of non-landslide samples.

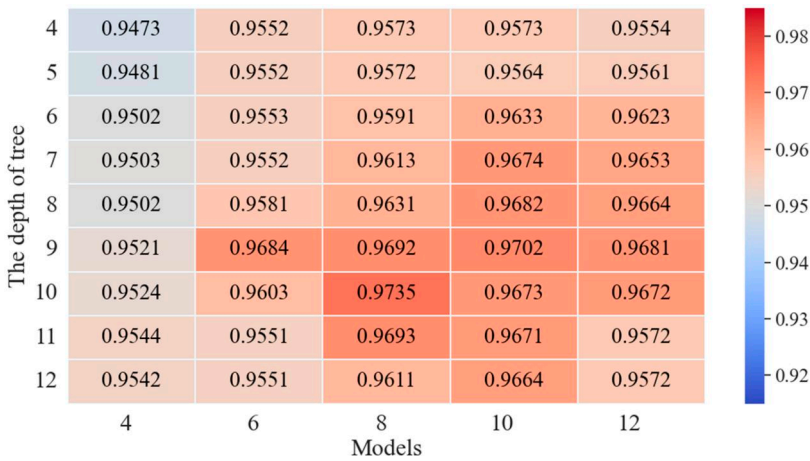


Figure 9. Accuracy statistics of CART-Boosting with various parameters.

### 4.3.3. Parameter setting

**4.3.3.1. Single models of CART and MLP.** Maximum tree depth is the crucial parameter for CART Modelling. In this article, a trial-and-error method is used to determine the maximum tree depth of the CART as 8. Regarding MLP, the number of hidden layers and neurons affects its modelling accuracy. After multiple sets of tests, we find that a model structure of MLP with two hidden layers is suitable. The neuron number of the first and second layers are set to 8 and 25, respectively.

**4.3.3.2. CART-Boosting.** In the same way, we obtain the relationship between the parameters and model accuracy of CART-Boosting through multiple sets of trials (Figure 9). Similarly with CART-Bagging, the accuracy of CART-Boosting first increases and then decreases with the classifier number when the maximum tree depth is constant. In the case of fewer than 10 classifiers, the model accuracy increases with the maximum tree depth; while the classifier number is larger than 10, the model accuracy showed a



**Figure 10.** Accuracy statistics of CART-Bagging with various parameters.

**Table 6.** Statistics of the classifier number and accuracy.

No. of classifier	6	8	10	12	14	16	18	20
MLP-Bagging	0.9682	0.9740	0.9756	0.9764	0.9780	0.9767	0.9765	0.9761
MLP-Boosting	0.9823	0.9829	0.9848	0.9852	0.9857	0.9855	0.9853	0.9852

downward trend after the growth. Therefore, we set the maximum tree depth and the classifier number of CART-Boosting to 10 and 14, respectively.

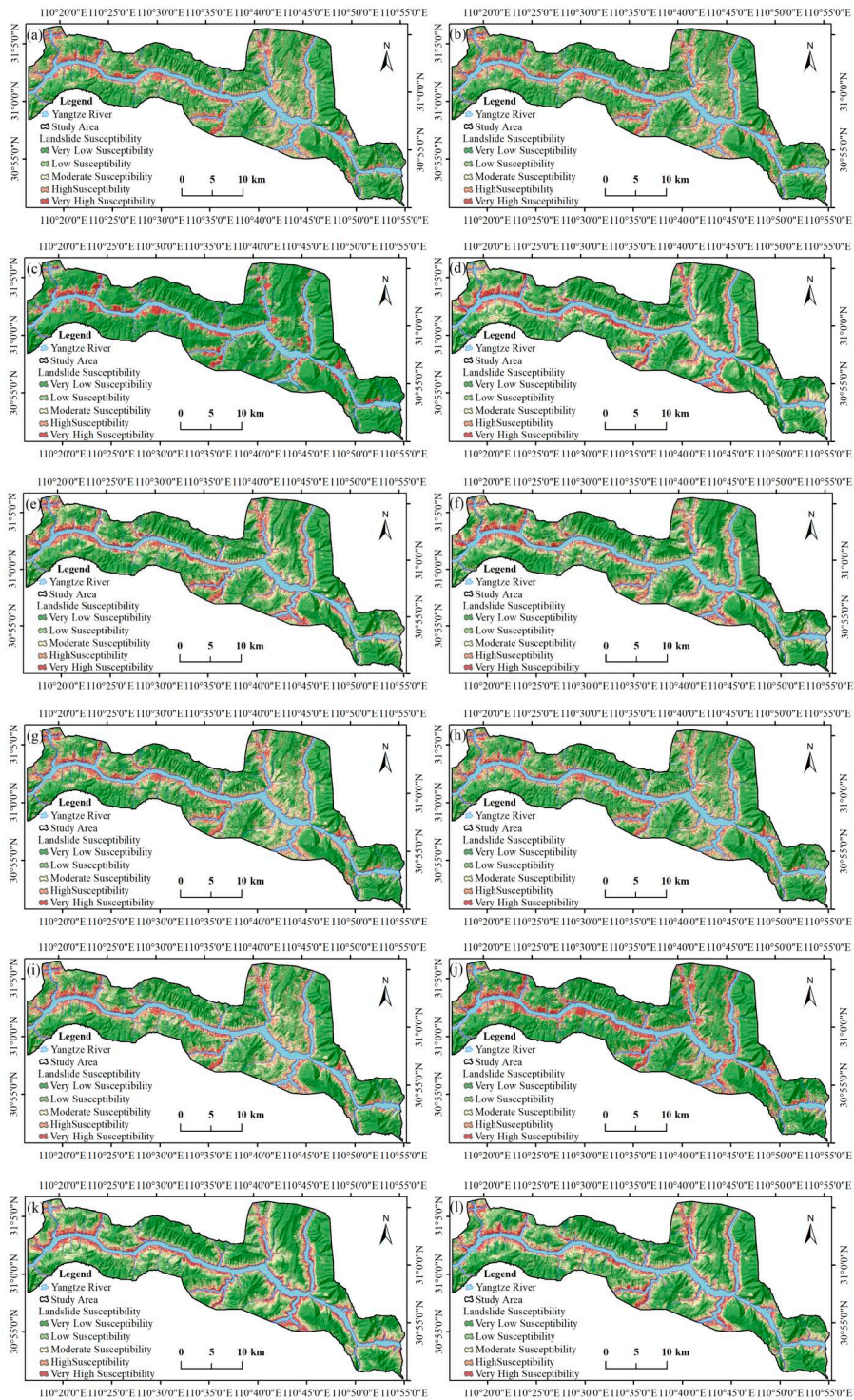
**4.3.3.3. CART-Bagging.** The number of classifier and the maximum tree depth are the significant parameters for the CART-based ensemble learning models. To determine the optimal parameters, we obtain the relationship between the two parameters and the model accuracy of CART-bagging through multiple trials. As shown in Figure 10, when the maximum tree depth is determined, the modelling accuracy increases with the classifier number within a specific range. Similarly, when the classifier number is constant, and the maximum tree depth is less than 10, the modelling accuracy increases with the increase of maximum tree depth. Therefore, we set the maximum tree depth and the classifier number of CART-bagging as 8 and 10, respectively.

**4.3.3.4. MLP-Bagging and MLP-Boosting.** The relationship between classifier number and the accuracy of MLP-based ensemble learning models is shown in Table 6. For MLP-Bagging and MLP-Boosting, the model accuracy first increases and then decreases with the classifier number. For example, both models achieve the highest accuracy when the classifier number is 14. With this, we set the classifier number of MLP-Bagging and MLP-Boosting models to 14.

#### 4.4. Landslide susceptibility mapping

The probability of landslide susceptibility is calculated by applying CART-Bagging, CART-Boosting, MLP-Bagging, MLP-Boosting, single CART, and single MLP models, respectively, which are implemented through *SPSS Modeler*. According to the ratio of 0.5: 1: 1.5: 2: 5, the probability value of landslide susceptibility is divided into five levels, namely Very High, High, Moderate, Low, and Very Low. The produced landslide





**Figure 11.** Landslide susceptibility maps using different methods (a) LR-MLP-Boosting, (b) LR-MLP-Bagging, (c) LR-CART-Boosting, (d) LR-CART-Bagging, (e) LR-MLP, (f) LR-CART, (g) No-MLP-Boosting, (h) No-MLP-Bagging, (i) No-CART-Boosting, (j) No-CART-Bagging, (k) No-MLP, and (l) No-CART.



susceptibility maps are presented in [Figures 11\(a–f\)](#). In addition, to verify the quality of the non-landslide samples selected by the LR constraint method, we used random sampling method to choose a set of non-landslide samples under no constraint condition (the whole landslide-free areas) for comparison. Similarly, these six models are used for LSM as well. The produced landslide susceptibility maps are shown in [Figures 11\(g–l\)](#).

## 5. Discussion

### 5.1. The relationship between landslide development and the main factors

The statistics of information value ([Table 1](#)) and susceptibility maps ([Figure 11](#)) indicate that the spatial development of landslides in the study area is mainly controlled by altitude, lithology, and distance to rivers. The widely distributed mudstone, marlstone, and weak strata, as well as the layered clastic rock strata containing weak interlayers, significantly reducing the sliding mass strength and making the slope vulnerable to instability ([Tang et al. 2019](#)). In the study area, the majority of landslides manifest at low altitudes. Altitudes lower than 240 meters exert the most pronounced influence on landslide development, with a maximum information value of 1.49. This is because many human engineering activities occur in this area, where the thick loose deposits provide the material basis for landslide occurrence. The periodic fluctuation of reservoir water level significantly changes the hydrogeological conditions of bank slopes, and plenty of seepage-driven and buoyancy-driven landslides are triggered ([Zhou et al. 2022a](#)). In the mountainous regions along the road, excessive rainfall enhances flow in the drainage, leading to the erosion of rock masses and the initiation of landslides. The statistics in [Table 1](#) suggest that the closer the slope to the river, the more it is affected. When the distance to rivers is less than 300 m, its information value is high at 0.92. It is to be noted that the information value method is a typical statistical method whose reliability depends

**Table 7.** Statistical results of susceptibility zoning.

Levels	Landslide pixles		Domain pixles		Ratio (a/b)
	No.	% (a)	No.	% (b)	
<b>No-CART</b>					
Very Low	861	3.56	353,136	50.55	0.07
Low	1,506	6.22	138,069	19.76	0.31
Moderate	4,882	20.18	101,150	14.48	1.39
High	8,393	34.69	69,174	9.90	3.50
Very High	8,553	35.35	37,117	5.31	6.65
<b>No-MLP</b>					
Very Low	846	3.35	351,927	50.37	0.07
Low	1,666	6.89	134,798	19.29	0.36
Moderate	5,174	21.38	108,194	15.49	1.38
High	8,250	34.09	68,527	9.81	3.48
Very High	8,259	34.14	35,200	5.04	6.77
<b>LR-CART</b>					
Very Low	645	2.67	341,817	48.93	0.05
Low	1,543	6.38	145,048	20.76	0.30
Moderate	4,674	19.32	107,442	15.38	1.26
High	8,726	36.07	68,256	9.77	3.69
Very High	8,607	35.57	36,083	5.16	6.89
<b>LR-MLP</b>					
Very Low	652	2.69	352,147	50.4	0.05
Low	1,408	5.82	136,945	19.6	0.30
Moderate	4,869	20.12	107,515	15.39	1.31
High	8,854	36.59	67,333	9.64	3.80
Very High	8,412	34.77	34,706	4.97	7.00

on the number of samples. Additionally, we conducted importance analysis of factors using the random forest method. Our findings reveal that altitude and lithology are the two factors with the highest contribution values, 0.15 and 0.11, respectively. These results align closely with those obtained from the statistics of information value.

The purpose of LSM is to quantitatively and accurately establish the non-linear relationship between influencing factors and the probability of landslide occurrence. Earlier studies have suggested that the selection of appropriate input factors is essential for effectively LSM using ML techniques (Zhou et al. 2018; Gentilucci et al. 2023; Qazi et al. 2023). If essential factors are not included as ML inputs, the model is at risk of underfitting. At the same time, selecting too many factors may lead to noise introduced by unimportant indicators, thereby decreasing the accuracy of LSM. Therefore, we believe that prior to conducting ML-based LSM, it is essential to understand the impact of each causal factor on landslide occurrence and select appropriate causal factor inputs.

## **5.2. Performance comparison of the used algorithms**

### **5.2.1. Machine learning algorithms**

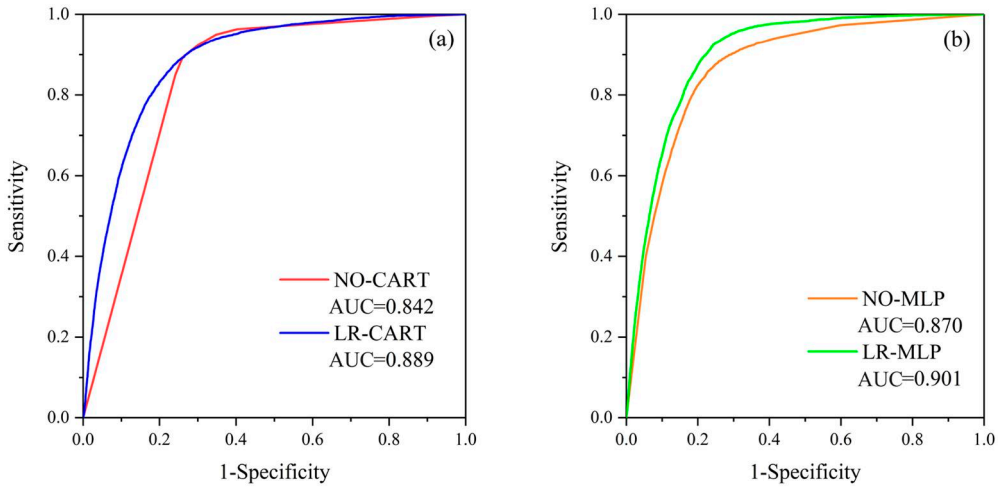
To verify the performance of the machine learning algorithms, we count the pixel distribution of four landslide susceptibility maps produced by LR-CART, LR-MLP, No-CART, and No-MLP. The statistics (Table 7) indicate that the maps produced by these four models are similar. The landslide development law is consistent with the mapping results, which indicates that the LSM is reliable. Regarding LR-MLP, the higher the landslide susceptibility level, the higher the landslide ratio. 34.77% of the landslide pixels are located in the Very High susceptibility area, with the highest landslide ratio of 7.00.

Conversely, only 2.67% of the landslide pixels lie in the Very Low susceptibility area, and the landslide ratio is 0.05. In the results of LR-MLP, the landslide ratio in the Very High susceptibility area is the highest at 6.89, where 35.57% of landslide pixels are distributed in this area. The same characteristics are presented in the results of the No-MLP and No-CART models. The statistics suggest that MLP is slightly performed better than CART.

The receiver operating characteristic (ROC) curve is a commonly used performance evaluation method in LSM. The area under the ROC curve (AUC) is used to assess model performance, and the model with a larger AUC is considered better. As shown in Figure 12, LR-MLP outperforms LR-CART, and their AUCs are 0.901 and 0.889, respectively. No-MLP achieves better accuracy than No-CART as well. The ROC curves suggest that both algorithms of MLP and CART perform excellently in LSM. We can infer MLP algorithm can more accurately establish the nonlinear relationship between landslide occurrence and its influencing factors than CART.

### **5.2.2. Ensemble learning algorithms**

The results of the eight coupling models are shown in Table 8. In the landslide susceptibility map generated by LR-MLP-Boosting, 43.75% of the landslide pixels are concentrated in the Very High susceptibility area, exhibiting the highest landslide ratio at 8.594. In contrast, the Very Low susceptibility area registers the lowest ratio at 0.028. LR-MLP-Boosting not only has the highest prediction accuracy but also has the lowest false negative error which may lead to catastrophic losses. Highest landslide ratios in Very High susceptibility areas and lowest landslide ratios in Very Low susceptibility areas suggest that this method achieves the most excellent performance. The results indicate that LR-MLP-Boosting performs better compared to the LR-CART-Boosting, and



**Figure 12.** ROC curves of single models: (a) LR- and No-CART, and (b) LR- and No-MLP.

LR-MLP-Bagging outperforms LR-CART-Bagging. The same comparison results are presented in the results of the No-MLP-Boosting, No-MLP-Bagging, No-CART-Boosting, and No-CART-Bagging models.

We have also utilized the ROC curves to quantify the performance of the coupling models. The same conclusion about the performance ranking of the coupled models can be drawn from the ROC Curves (Figure 13). LR-MLP-Boosting achieved the best prediction accuracy with the highest AUC of 0.985. Apparently, the coupling models outperformed the single machine learning models. Boosting and Bagging improved the accuracy of LR-MLP by 0.085 and 0.077, respectively, while improving the accuracy of LR-CART by 0.092 and 0.084, respectively (Table 9). Accuracy improvement from Boosting was more significant than the improvement from Bagging. In the Boosting method, the prediction of all the classifiers was sequentially integrated into the training process to achieve the final results. However, it was parallel to the Bagging method. The Boosting method was more effective at reducing the deviation and variance, which enhances the prediction ability of the coupling model.

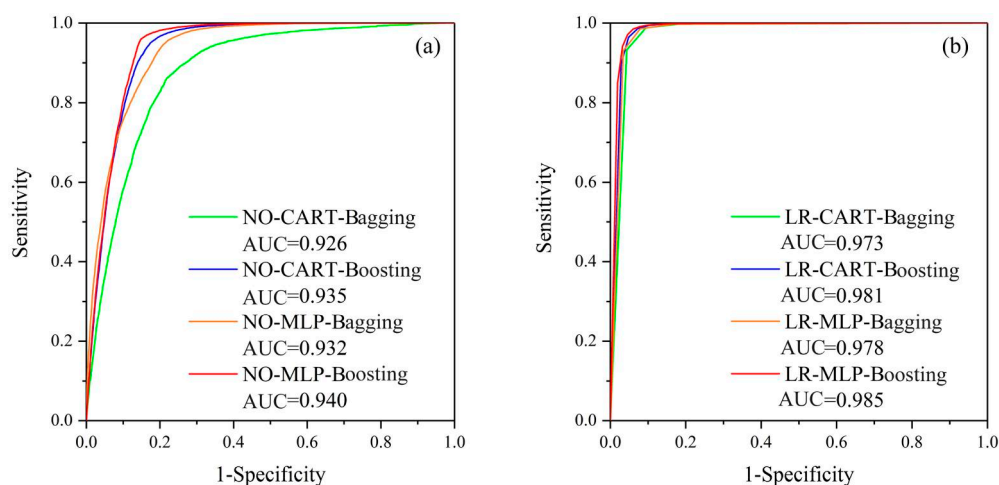
### 5.3. The advantages of the proposed method for non-landslide sampling

High-quality non-landslide samples are critical to the performance improvement of LSM. Numerous methods have been proposed to control the sampling range of non-landslide areas, aiming to acquire high-quality samples. Fu et al. (2023) investigated the impact of four non-landslide sampling methods on LSM and proposed a novel integrative sampling approach that incorporates multiple existing methods, which achieve the best performance. Theoretically, controlling the range of sampling can enhance the quality of non-landslide samples, and this notion supported by many studies (Dou et al. 2023; Hong et al. 2024). However, an improper control range can result in a decrease in the quality of non-landslide sampling. Wang et al. (2022) designed two distinct non-landslide sampling areas, covering both the entire study region and the mountainous area of Anhui Province, China, to explore the impacts of different spatial extents on LSM. The results indicate that sampling from the mountainous area yields inferior performance of LSM.

**Table 8.** Statistical results of each landslide susceptibility zoning.

Levels	Landslide pixels		Domain pixels		Ratio (a/b)
	No	% (a)	No	% (b)	
<b>LR-MLP-Boosting</b>					
Very Low	346	1.43	352,156	50.41	0.028
Low	706	2.92	135,816	19.44	0.150
Moderate	3,222	13.32	105,327	15.08	0.883
High	9,336	38.59	69,792	9.99	3.863
Very high	10,585	43.75	35,555	5.09	8.597
<b>LR-MLP-Bagging</b>					
Very Low	423	1.75	359,283	51.43	0.034
Low	933	3.86	131,112	18.77	0.205
Moderate	4,185	17.30	105,520	15.10	1.145
High	8,998	37.19	68,320	9.78	3.803
Very High	9,656	39.91	34,411	4.93	8.103
<b>LR-CART-Boosting</b>					
Very Low	394	1.63	364,801	52.22	0.031
Low	770	3.18	122,823	17.58	0.181
Moderate	3,829	15.83	107,890	15.44	1.025
High	9,006	37.22	67,320	9.64	3.863
Very High	10,196	42.14	35,812	5.13	8.221
<b>LR-CART-Bagging</b>					
Very Low	413	1.71	342,258	48.99	0.035
Low	1,105	4.57	145,233	20.79	0.220
Moderate	3,786	15.65	107,928	15.45	1.013
High	8,883	36.71	67,105	9.61	3.820
Very High	10,008	41.36	36,122	5.17	8.000
<b>No-MLP-Boosting</b>					
Very Low	595	2.46	342,258	51.43	0.048
Low	929	3.84	145,233	18.77	0.205
Moderate	4,124	17.04	107,928	15.07	1.131
High	8,834	36.51	67,105	9.78	3.735
Very High	9,713	40.14	36,122	4.96	8.094
<b>No-MLP-Bagging</b>					
Very Low	601	2.48	353,804	50.64	0.049
Low	1,089	4.50	132,710	18.99	0.237
Moderate	4,034	16.67	108,589	15.54	1.073
High	8,707	35.99	67,991	9.73	3.698
Very High	9,764	40.36	35,552	5.09	7.930
<b>No-CART-Boosting</b>					
Very Low	623	2.57	352,231	50.16	0.051
Low	1,106	4.57	135,474	19.39	0.236
Moderate	4,121	17.03	107,995	15.46	1.102
High	8,786	36.31	67,976	9.73	3.732
Very High	9,559	39.51	34,970	5.01	7.893
<b>No-CART-Bagging</b>					
Very Low	622	2.57	342,258	48.99	0.050
Low	1,340	5.54	145,233	20.79	0.260
Moderate	3,750	15.50	107,928	15.45	1.000
High	8,688	35.91	67,105	9.61	3.740
Very High	9,795	40.48	36,122	5.17	7.830

In this study, we use two methods for non-landslide sampling, and two single models and four coupling models for susceptibility mapping are established using selected samples. As shown in the ROC curves and statistics (Figures 12 and 13, Tables 8 and 9), all the LR- models achieve a better performance compared to the corresponding No- models. It indicates that the application of the LR model to constrain the selection range of non-landslide samples can effectively improve sample quality. Many machine learning methods can produce more accurate initial susceptibility maps to constrain the range of non-landslide sampling. However, as reported in earlier studies (Zhou et al. 2018; Sun et al. 2022),



**Figure 13.** ROC curves of coupling models: (a) no constraint sampling and (b) LR constrained sampling.

**Table 9.** Statistics of modelling accuracy.

Single model	Original AUC	Bagging		Boosting	
		AUC	Improvement	AUC	Improvement
<b>No sampling</b>					
CART	0.842	0.926	0.084	0.935	0.093
MLP	0.870	0.932	0.062	0.940	0.070
<b>LR sampling</b>					
CART	0.889	0.973	0.084	0.981	0.092
MLP	0.901	0.978	0.077	0.986	0.085

the performance of machine learning methods varies in regions, and high-quality data is required. As a result, a poor prediction may occur in some landslide-prone regions. Due to the simplicity of operation, high accuracy and stable performance, LR is widely used and consistently achieves acceptable results. In comparison, LR is a better choice to ensure the generalization of the non-landslide sampling method.

Random sampling from the whole landslide-free areas is the most utilized method for non-landslide sampling. The non-landslide pixels obtained through random sampling without any constrained conditions may have engineering geological conditions that are susceptible to landslides. The mixing of these pixels will reduce the quality of non-landslide samples. Furthermore, various engineering geological conditions contribute to the absence of landslides. Some non-landslide sampling range constraint methods, such as the low-slope method, may fail to capture non-landslide samples with diverse geological conditions. This can easily exert a negative impact on LSM. LR model produces an initial susceptibility map, and non-landslide samples are only selected from the Very Low susceptibility areas. This method can effectively avoid the mis-selection of samples in landslide-prone areas and keep the diversity of non-landslide sample characteristics. In general, our proposed non-landslide sampling method is conducive to improving LSM performance and can be applied worldwide.

## 6. Conclusion

In this study, we considered twelve causal factors as inputs for LSM after multi-collinearity analysis. The relationship between landslide spatial development and each causal factor

is quantitatively analyzed. We found that the altitude (<240 m) and distance to rivers (<300 m) emerged as important factors for landslide occurrence in the study area. Their information values were the highest at 1.49 and 0.92, respectively. LR-MLP-Boosting achieved the highest prediction accuracy with an AUC of 0.985. The accuracy comparison indicates that MLP performs better than CART. The ensemble methods outperformed the corresponding single classifiers, and Boosting algorithm slightly performed better than Bagging algorithm. LR is a reliable method to generate a preliminary susceptibility map to determine the Very Low susceptibility area. The non-landslide samples selected from the low susceptibility area are of higher quality than those selected from the entire landslide-free areas. High-quality non-landslide samples, which can be effectively obtained by using the LR model to constrain its sampling range, is able to enhance the performance of LSM. These results will be of great help to the community and to the scientists to monitor landslide susceptible locations and to get early information about the occurrence of landslide events to minimize loss and damage.

## Acknowledgements

We are grateful to the anonymous reviewers for providing useful comments/suggestions that have helped us to improve an earlier version of the manuscript. The first author would like to thank the China Scholarship Council for funding his research at the German Research Centre for Geosciences.

## Disclosure statement

No potential conflict of interest was reported by the author(s).

## Funding

This research is funded by the National Natural Science Foundation of China (No. 42371094 and No. 41702330) and the Key Research and Development Program of Hubei Province (No. 2021BCA219).

## Data availability statement

Data used in this study are available from the corresponding author by request.

## References

- Akinci H, Zeybek M. 2021. Comparing classical statistic and machine learning models in landslide susceptibility mapping in Ardanuc (Artvin), Turkey. *Nat Hazards*. 108(2):1515–1543. doi:10.1007/s11069-021-04743-4.
- Cao Y, Yin K, Alexandre DE, et al. 2016. Using an extreme learning machine to predict the displacement of step-like landslides in relation to controlling factors. *Landslides*. 13:725–736.
- Dou H, He J, Huang S, Jian W, Guo C. 2023. Influences of non-landslide sample selection strategies on landslide susceptibility mapping by machine learning. *Geomat Nat Hazard Risk*. 14(1):2285719. doi:10.1080/19475705.2023.2285719.
- Fu Z, Wang F, Dou J, Nam K, Ma H. 2023. Enhanced absence sampling technique for data-driven landslide susceptibility mapping: a case study in Songyang County, China. *Remote Sens*. 15(13):3345. doi:10.3390/rs15133345.
- Gameiro S, Riffel ES, de Oliveira GG, Guasselli LA. 2021. Artificial neural networks applied to landslide susceptibility: the effect of sampling areas on model capacity for generalization and extrapolation. *Appl Geogr*. 137:102598. doi:10.1016/j.apgeog.2021.102598.



- Gardner MW, Dorling SR. 1998. Artificial neural networks (the multilayer perceptron)—a review of applications in the atmospheric sciences. *Atmos Environ.* 32(14-15):2627–2636. doi:10.1016/S1352-2310(97)00447-0.
- Gentilucci M, Pelagagge N, Rossi A, Domenico A, Pambianchi G. 2023. Landslide susceptibility using climatic–environmental factors using the weight-of-evidence method—a study area in Central Italy. *Appl Sci.* 13(15):8617. doi:10.3390/app13158617.
- Hong H. 2023. Assessing landslide susceptibility based on hybrid Best-first decision tree with ensemble learning model. *Ecol Indic.* 147:109968. doi:10.1016/j.ecolind.2023.109968.
- Hong H. 2023. Assessing landslide susceptibility based on hybrid multilayer perceptron with ensemble learning. *Bull Eng Geol Environ.* 82(10):382. doi:10.1007/s10064-023-03409-8.
- Hong H, Wang D, Zhu AX, Wang Y. 2024. Landslide susceptibility mapping based on the reliability of landslide and non-landslide sample. *Expert Syst Appl.* 243:122933. doi:10.1016/j.eswa.2023.122933.
- Hu X, Mei H, Zhang H, Li Y, Li M. 2021. Performance evaluation of ensemble learning techniques for landslide susceptibility mapping at the Jinping County, Southwest China. *Nat Hazards.* 105(2):1663–1689. doi:10.1007/s11069-020-04371-4.
- Hu J, Motagh M, Wang J, Qin F, Zhang J, Wu W, Han Y. 2021. Karst collapse risk zonation and evaluation in Wuhan, China based on analytic hierarchy process, logistic regression, and InSAR angular distortion approaches. *Remote Sensing.* 13(24):5063. doi:10.3390/rs13245063.
- Kavzoglu T, Sahin EK, Colkesen I. 2014. Landslide susceptibility mapping using GIS-based multi-criteria decision analysis, support vector machines, and logistic regression. *Landslides.* 11(3):425–439. doi:10.1007/s10346-013-0391-7.
- Liu B, Guo H, Li J, Ke X, He X. 2024. Application and interpretability of ensemble learning for landslide susceptibility mapping along the Three Gorges Reservoir area, China. *Nat Hazards.* 1–32. doi:10.1007/s11069-023-06374-3.
- Long J, Liu Y, Li C, Fu Z, Zhang H., 2021. A novel model for regional susceptibility mapping of rainfall-reservoir induced landslides in Jurassic slide-prone strata of western Hubei Province, Three Gorges Reservoir area. *Stoch Environ Res Risk Assess.* 35(7):1403–1426. doi:10.1007/s00477-020-01892-z.
- Lucchese LV, de Oliveira GG, Pedrollo OC. 2021. Investigation of the influence of nonoccurrence sampling on landslide susceptibility assessment using Artificial Neural Networks. *Catena.* 198:105067. doi:10.1016/j.catena.2020.105067.
- Meena SR, Puliero S, Bhuyan K, Floris M, Catani F. 2022. Assessing the importance of conditioning factor selection in landslide susceptibility for the province of Belluno (region of Veneto, northeastern Italy). *Nat Hazards Earth Syst Sci.* 22(4):1395–1417. doi:10.5194/nhess-22-1395-2022.
- Mohammed A, Kora R. 2023. A comprehensive review on ensemble deep learning: opportunities and challenges. *J King Saud Univ Comput Inform Sci.* 35(2):757–774.
- Pham BT, Jaafari A, Nguyen-Thoi T, Van Phong T, Nguyen HD, Satyam N, Masroor M, Rehman S, Sajjad H, Sahana M, et al. 2021. Ensemble machine learning models based on reduced error pruning tree for prediction of rainfall-induced landslides. *Int J Digital Earth.* 14(5):575–596. doi:10.1080/17538947.2020.1860145.
- Qazi A, Singh K, Vishwakarma DK, Abdo HG. 2023. GIS based landslide susceptibility zonation mapping using frequency ratio, information value and weight of evidence: a case study in Kinnaur District HP India. *Bull Eng Geol Environ.* 82(8):332. doi:10.1007/s10064-023-03344-8.
- Roy P, Martha TR, Vinod Kumar K, Chauhan P, Rao VV. 2023. Cluster landslides and associated damage in the Dima Hasao district of Assam, India due to heavy rainfall in May 2022. *Landslides.* 20(1):97–109. doi:10.1007/s10346-022-01977-6.
- Sabokbar HF, Roodposhti MS, Tazik E. 2014. Landslide susceptibility mapping using geographically-weighted principal component analysis. *Geomorphology.* 226:15–24. doi:10.1016/j.geomorph.2014.07.026.
- Sharma N, Saharia M, Ramana GV. 2024. High resolution landslide susceptibility mapping using ensemble machine learning and geospatial big data. *Catena.* 235:107653. doi:10.1016/j.catena.2023.107653.
- Song J, Wang Y, Fang Z, Peng L, Hong H. 2020. Potential of ensemble learning to improve tree-based classifiers for landslide susceptibility mapping. *IEEE J Sel Top Appl Earth Observations Remote Sensing.* 13:4642–4662. doi:10.1109/JSTARS.2020.3014143.
- Sun D, Gu Q, Wen H, Xu J, Zhang Y, Shi S, Xue M, Zhou X. 2022. Assessment of landslide susceptibility along mountain highways based on different machine learning algorithms and mapping units by hybrid factors screening and sample optimization. *Gondwana Res.* 123:89–106. doi:10.1016/j.gr.2022.07.013.
- Tang M, Xu Q, Yang H, Li S, Iqbal J, Fu X, Huang X, Cheng W. 2019. Activity law and hydraulics mechanism of landslides with different sliding surface and permeability in the Three Gorges Reservoir Area, China. *Eng Geol.* 260:105212. doi:10.1016/j.enggeo.2019.105212.

- Tanyu BF, Abbaspour A, Alimohammadlou Y, Tecuci G. 2021. Landslide susceptibility analyses using Random Forest, C4. 5, and C5. 0 with balanced and unbalanced datasets. *Catena*. 203:105355. doi:10.1016/j.catena.2021.105355.
- Wang C, Lin Q, Wang L, Jiang T, Su B, Wang Y, Mondal SK, Huang J, Wang Y. 2022. The influences of the spatial extent selection for non-landslide samples on statistical-based landslide susceptibility modeling: a case study of Anhui Province in China. *Nat Hazards*. 112(3):1967–1988. doi:10.1007/s11069-022-05252-8.
- Yavuz Ozalp A, Akinci H, Zeybek M. 2023. Comparative analysis of tree-based ensemble learning algorithms for landslide susceptibility mapping: a case study in Rize, Turkey. *Water*. 15(14):2661. doi:10.3390/w15142661.
- Youssef AM, Pourghasemi HR, Pourtaghi ZS, Al-Katheeri MM. 2016. Landslide susceptibility mapping using random forest, boosted regression tree, classification and regression tree, and general linear models and comparison of their performance at Wadi Tayyah Basin, Asir Region, Saudi Arabia. *Landslides*. 13(5):839–856. doi:10.1007/s10346-015-0614-1.
- Yu L, Cao Y, Zhou C, Wang Y, Huo Z. 2019. Landslide susceptibility mapping combining information gain ratio and support vector machines: a case study from Wushan Segment in the Three Gorges Reservoir Area, China. *Appl Sci*. 9(22):4756. doi:10.3390/app9224756.
- Yu L, Zhou C, Wang Y, Cao Y, Peres DJ. 2022. Coupling data-and knowledge-driven methods for landslide susceptibility mapping in human-modified environments: a case study from Wanzhou County, Three Gorges Reservoir Area, China. *Remote Sensing*. 14(3):774. doi:10.3390/rs14030774.
- Zhou C, Cao Y, Hu X, Yin K, Wang Y, Catani F. 2022. Enhanced dynamic landslide hazard mapping using MT-InSAR method in the Three Gorges Reservoir Area. *Landslides*. 19(7):1585–1597. doi:10.1007/s10346-021-01796-1.
- Zhou C, Cao Y, Yin K, Intrieri E, Catani F, Wu L. 2022. Characteristic comparison of seepage-driven and buoyancy-driven landslides in Three Gorges Reservoir area, China. *Eng Geol*. 301:106590. doi:10.1016/j.enggeo.2022.106590.
- Zhou C, Yin KL, Cao Y, Li Y. 2020a. Landslide susceptibility assessment by applying the coupling method of radial basis neural network and adaboost: A case study from the Three Gorges Reservoir area. *Earth Sci*. 45(6):1865–1876. doi:10.3390/rs14030774
- Zhou C, Cao Y, Yin K, Wang Y, Shi X, Catani F, Ahmed B. 2020b. Landslide characterization applying sentinel-1 images and InSAR technique: the Muyubao landslide in the three Gorges Reservoir Area, China. *Remote Sensing*. 12(20):3385. doi:10.3390/rs12203385.
- Zhou C, Yin K, Cao Y, Ahmed B. 2016. Application of time series analysis and PSO–SVM model in predicting the Bazimen landslide in the Three Gorges Reservoir, China. *Eng Geol*. 204:108–120. doi:10.1016/j.enggeo.2016.02.009.
- Zhou C, Yin K, Cao Y, Ahmed B, Li Y, Catani F, Pourghasemi HR., 2018. Landslide susceptibility modeling applying machine learning methods: a case study from Longju in the Three Gorges Reservoir area, China. *Comput Geosci*. 112:23–37. doi:10.1016/j.cageo.2017.11.019.
- Zhu AX, Miao Y, Liu J, Bai S, Zeng C, Ma T, Hong H. 2019. A similarity-based approach to sampling absence data for landslide susceptibility mapping using data-driven methods. *Catena*. 183:104188. doi:10.1016/j.catena.2019.104188.