# Geoscience Frontiers

## Spatial Landslide Susceptibility Assessment Using Machine Learning Techniques Assisted by Additional Data Created with Generative Adversarial Networks
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GSF-D-20-00191R2 |
| Article Type: | Research Paper |
| Keywords: | Landslide susceptibility;  Machine learning;  Generative adversarial networks;  Convolutional neural networks;  Geographic information systems |
| Corresponding Author: | Biswajeet Pradhan, PhD<br>University Technology Sydney<br>Sydney, AUSTRALIA |
| First Author: | Husam A. H. Al-Najjar |
| Order of Authors: | Husam A. H. Al-Najjar |
| | Biswajeet Pradhan, PhD |
| Abstract: | In recent years, landslide susceptibility mapping has substantially improved with advances in machine learning. However, there are still challenges remain in landslide mapping due to the availability of limited inventory data. In this paper, a novel method that improves the performance of machine learning techniques is presented. The proposed method creates synthetic inventory data using generative adversarial networks (GANs) for improving the prediction of landslides. In this research, landslide inventory data of 156 landslide locations were identified in Cameron Highlands, Malaysia, taken from previous projects the authors worked on. Altitude, slope, aspect, plan curvature, profile curvature, total curvature, lithology, land use and land cover, distance from the road, distance from the stream, stream power index, sediment transport index, terrain roughness index, topographic wetness index and vegetation density are geo-environmental factors considered in this study based on suggestions from previous works on Cameron Highlands. To show the capability of GANs in improving landslide prediction models, this study tests the proposed GAN model with artificial neural network (ANN), support vector machine (SVM), decision trees (DT), random forest (RF) and bagging ensemble models with ANN and SVM models. These models were validated using the area under the receiver operating characteristic curve (AUROC). The DT, RF, SVM, ANN and Bagging ensemble could achieve the AUROC values of (0.90, 0.94, 0.86, 0.69 and 0.82) for the training; and the AUROC of (0.76, 0.81, 0.85, 0.72 and 0.75) for the test, subsequently. When using additional samples, the models achieved the AUROC values of (0.92, 0.94, 0.88, 0.75 and 0.84) for the training and (0.78, 0.82, 0.82, 0.78 and 0.80) for the test, respectively. Without the use of additional samples created by the GAN model, SVM achieved the highest AUROC of 0.85, whereas ANN had the lowest AUROC of 0.72. RF and SVM achieved AUROC of 0.82 when the additional samples were used for training these models. Using the additional samples improved the test accuracy of all the models except SVM. As a result, in data-scarce environments, this research showed that utilizing GANs to generate supplementary samples is promising because it can improve the predictive capability of common landslide prediction models. |

1 # Spatial Landslide Susceptibility Assessment Using Machine
2 # Learning Techniques Assisted by Additional Data Created with
3 # Generative Adversarial Networks

4

5 Husam A. H. Al-Najjar [a] and Biswajeet Pradhan [a,b]*

6 [a] Centre for Advanced Modelling and Geospatial Information Systems (CAMGIS), Faculty of Engineering and
7 IT, University of Technology Sydney, 2007 NSW, Australia; Husam.al-najjar@student.uts.edu.au);
8 Biswajeet.Pradhan@uts.edu.au

9 [b] Department of Energy and Mineral Resources Engineering, Sejong University, Choongmu-gwan, 209,
10 Neungdong-ro, Gwangin-gu, Seoul 05006, Korea

11 *Corresponding author Email: biswajeet24@gmail.com or Biswajeet.Pradhan@uts.edu.au

12
13 **Abstract**

14 In recent years, landslide susceptibility mapping has substantially improved with advances in

15 machine learning. However, there are still challenges remain in landslide mapping due to the

16 availability of limited inventory data. In this paper, a novel method that improves the

17 performance of machine learning techniques is presented. The proposed method creates

18 synthetic inventory data using Generative Adversarial Networks (GANs) for improving the

19 prediction of landslides. In this research, landslide inventory data of 156 landslide locations

20 were identified in Cameron Highlands, Malaysia, taken from previous projects the authors

21 worked on. Elevation, slope, aspect, plan curvature, profile curvature, total curvature, lithology,

22 land use and land cover (LULC), distance to the road, distance to the river, stream power index

23 (SPI), sediment transport index (STI), terrain roughness index (TRI), topographic wetness

24 index (TWI) and vegetation density are geo-environmental factors considered in this study

25 based on suggestions from previous works on Cameron Highlands. To show the capability of

26 GANs in improving landslide prediction models, this study tests the proposed GAN model with

27 benchmark models namely Artificial Neural Network (ANN), Support Vector Machine (SVM),

28 Decision Trees (DT), Random Forest (RF) and Bagging ensemble models with ANN and SVM

29 models. These models were validated using the area under the receiver operating characteristic

curve (AUROC). The DT, RF, SVM, ANN and Bagging ensemble could achieve the AUROC

values of (0.90, 0.94, 0.86, 0.69 and 0.82) for the training; and the AUROC of (0.76, 0.81, 0.85,

0.72 and 0.75) for the test, subsequently. When using additional samples, the same models

achieved the AUROC values of (0.92, 0.94, 0.88, 0.75 and 0.84) for the training and (0.78,

0.82, 0.82, 0.78 and 0.80) for the test, respectively. Using the additional samples improved the

test accuracy of all the models except SVM. As a result, in data-scarce environments, this

research showed that utilizing GANs to generate supplementary samples is promising because

it can improve the predictive capability of common landslide prediction models.

**Keywords:** Landslide susceptibility, Inventory, Machine learning, Generative adversarial

network, Convolutional neural network, Geographic information system

**1. Introduction**

Natural hazards are major challenges worldwide, and many countries are spending a significant

amount of their yearly budget to control and prevent them. Landslides pose a serious risk to

human habitats. The risk of landslides is a major barrier to agricultural and urban development

practices. In addition, ongoing urbanization is placing vast demands on infrastructure and

escalating the threat to property and human lives. As a result, landslide hazard assessment has

become a major step in planning the most suitable for risk mitigation measures. Experts

frequently use the maps generated from this assessment to identify regions where thorough in-

situ studies should be conducted. Landslide hazard assessment is a complex task that includes

comprehension of the science of geotechnics, geomorphology, hydrology and statistics (Glade

et al., 2012). This objective has motivated computational modeling studies, particularly the

evaluation of landslide susceptibility. Statistical and physical models are often used to

accomplish this task (Formetta et al., 2014).

Physical-based models combine susceptibility analysis with soil and rock mechanics, creating a physical basis for this method (Wang et al., 2019). They are appropriate at a local scale such as single slope, basin/ catchment and requires site-specific geotechnical data (Park et al., 2019). Generally, the infinite slope model is used in the analysis of slope stability with hydrological or earthquake models. Although reliable geotechnical parameters are essential for such models, lack of geotechnical data throughout a large scale area and the expensiveness remain the main obstacles in the physical-based models (Lee et al., 2014). Various landslides studies and assessments were carried out to develop landslide-prone areas in Malaysia (Fanos and Pradhan, 2019; Mezaal and Pradhan, 2018; Pradhan and Lee, 2010; Sameen et al., 2020).

Statistical models, also known as empirical models, use landslide inventories and other conditioning factors (e.g. terrain and land use), which can be extracted at large scales using remote sensing data and Geographical Information Systems (GIS). Such techniques have gained popularity in the field of landslide susceptibility assessment, especially when addressing the challenge of landslide mapping of prone areas at large scales, where enough geotechnics information is not available to perform physical-based models (Goetz et al., 2011).These models have also been supported by the latest progress in the availability and accessibility of remote sensing-based derived information, such as topography, land cover and precipitation products, thereby improving the application of the method at large scales.

Several scholars have evaluated various statistical models to assess landslide susceptibility (Akbar and Chen, 2018; Braun et al., 2018; Ciurleo et al., 2017; Goetz et al., 2015; Huang and Zhao, 2018; Kavzoglu et al., 2019; Süzen and Doyuran, 2004; Xiao et al., 2019; Zêzere et al., 2017). Early approaches to modeling landslide susceptibility are based on field investigations. Such techniques, however, are costly and site-specific, and they heavily involve extensive expertise in geology and geomorphology. Statistical approaches of landslide susceptibility modeling have become very popular during the last two decades. Recently, several scholars

78    including (Kawabata and Bandibas, 2009; Lee and Sambath, 2006; Mandal and Mondal, 2019;

79    Pradhan, 2013) evaluated several statistical models, such as frequency ratio (FR), logistic

80    regression (LR), artificial neural network (ANN), certainty factor (CF), analytical hierarchy

81    process (AHP) and fuzzy logic (FL). They suggested that ANN-, CF and FR-based FL are the

82    most reliable techniques in assessing and predicting landslide susceptibility, at least for their

83    case study. Regardless of the type of models and where they belong (statistical or machine

84    learning), they are good for landslide susceptibility assessment of large areas. Statistical models

85    can also be evaluated quantitatively at lower costs than evaluating a physical model. In

86    addition, these models are computationally more efficient than physical models because the

87    latter require simulations with numerous iterations to determine some geotechnical parameters

88    that are used to prepare the susceptibility products. However, they have certain limitations,

89    which include difficulties in explaining the results of the black box models and over-fitting in

90    the case of limited training samples.

91    **2. Related Works**

92    Landslide susceptibility mapping has improved substantially during the last decade because of

93    new data processing techniques such as sampling methods, machine learning models, and

94    validation measures. Some studies have focused on sampling strategies, selection of training

95    samples and addressing the effects of incomplete inventory datasets. In landslide susceptibility

96    mapping, training data play a critical role in determining the accuracy and generalization of the

97    model. The size of the training data has a significant effect on the accuracy of the susceptibility

98    model. For example, in the training data under some sample threshold limits, Hussin et al.

99    (2016) showed that model performance was very low, while the use of a large number of

100   landslides above the threshold created a plateau effect, with no increase in model performances.

101   Tsangaratos and Ilia (2016) also reported that the size of training data influences the prediction

102   accuracy when using models such as LR and Naive Bayes.

103    Several studies have attempted to improve landslide susceptibility models by proposing new

104    factors into the process including conditioning factors optimization (Al-Najjar et al., 2019;

105    Canoglu et al., 2019; Dou et al., 2015; Kavzoglu et al., 2015; Kornejady et al., 2018; Samia et

106    al., 2018; Soma et al., 2019). Moreover, model parameterization and integration methods have

107    been studied to improve landslide susceptibility mapping. Statistical and machine learning

108    models are often affected by the selection of proper hyper-parameters for a specific case study

109    (Can et al., 2019; Feizizadeh et al., 2017). Moreover, the model's integration has also been

110    active research for improving the landslide susceptibility in the last few years (Kalantar et al.,

111    2018). Examples of model integration studies include ensemble models (Bragagnolo et al.,

112    2020; Kadavi and Lee, 2018) and integration of data-driven and knowledge-based models

113    (Ashournejad et al., 2019; Yan et al., 2019; Zhang et al., 2019).

114    Studies on sampling strategies for landslide susceptibility mapping have been active in recent

115    years. Hussin et al. (2016) assessed different landslide sampling strategies (scarp centroid,

116    points populating the scarp and entire scarp polygon) in a grid-based statistical model. These

117    strategies achieve the highest performance when sampling shallow landslides as grid points

118    and debris flow scarps as polygons. Yilmaz and Ercanoglu (2019) discussed the necessity of

119    studying the selection of data mining techniques; they emphasized that sampling methods such

120    as polygon features or seed cells representative pre-failure settings appear to be more genuine

121    in obtaining truthful maps than other methods. Lai et al. (2019) also explored the influence of

122    sampling strategies for improving landslide susceptibility mapping.

123    In addition to the size of training data and the sampling strategy, studies have investigated

124    various ways of selecting training samples. Conoscenti et al. (2016) performed landslide

125    susceptibility mapping through investigating the impact of landslide absence (negative

126    samples) on the models; they extracted the landslide absence using randomly distributed circles

127    that have a diameter equivalent to the mean width of the landslide source areas. Moreover, the

128    individual grid cells were randomly distributed to distinguish the non-landslide zones (absence

129    selection). Experiments from this study based on multivariate adaptive regression splines

130    showed that absences selection using random circles are significantly better than the other

131    method when learning and validation samples were extracted from the same area, and no

132    significant difference was observed when testing the models outside the training area. Kalantar

133    et al. (2018) evaluated the impact of landslide samples varieties on the SVM, LR and ANN

134    methods; their investigation demonstrated that randomness in the training sample selection has

135    a significant effect on the susceptibility models. The outcome showed that, in the section of

136    training samples, the LR model is less sensitive than the SVM and ANN models. Zhu et al.

137    (2019) proposed a method based on similarity sampling for absence selection; their

138    experiments on a common machine learning models showed that this new method

139    outperformed the existing methods, such as buffer control and target space exteriorization.

140    Hong et al. (2019) assessed the impact of absence data selection on the RF model. Aktas and

141    San (2019) developed a new automatic sampling method based on a two-level random

142    sampling.

143    The impact of landslide inventory incompleteness on susceptibility mapping was also carried

144    out in recent studies. Du et al. (2020) assessed landslide susceptibility in Tibet Chinese

145    Himalayas, with a multinomial logistic regression model with reported average AUC of 0.867;

146    however, there were some uncertainties in the landslide-prone areas defined by their AHP

147    model. Steger et al. (2016)  assessed the impact of spatially heterogeneous completeness of

148    landslide information on statistical landslide susceptibility models (e.g. logistic regression) by

149    artificially introducing two different mapping biases into available landslides and synthetically

150    generated landslides. Although they reported AUROCs greater than 0.85, they suggested the

151    method needed to be evaluated with other different models. In another study, Lee et al. (2018)

152    employed optimized data mining and statistical methods for various scenarios considering

153  limited inventories. In their model, SVM achieved the AUC of 0.85 when either the full or

154  limited landslide inventories were used; however, generating additional inventories was not

155  considered in their study. Steger et al. (2016) suggested that models directly associated with

156  inventory-based incompleteness should be rejected regardless of their performance.

157  Furthermore, they proposed using mixed-effects modeling if systematically missing landslide

158  information can be attributed to a spatial variable (Steger et al., 2018).

159  The aforementioned studies indicate several ways to improve landslide susceptibility models,

160  such as data-related methods and others that target the model construction and training process.

161  This study aims to develop a new method for additional landslide sample creation with

162  generative adversarial networks (GANs) which could be useful in the inventory-scarce

163  environment. Several machine learning models, such as ANN, SVM, DT, RF and Bagging

164  ensemble, with ANN and SVM as base classifiers, are used to evaluate the new method of

165  landslide susceptibility mapping. These methods are compared in a case study in Cameron

166  Highlands, Malaysia.

167  **3. Study area and materials**

168  **3.1. Study area**

169  The Cameron Highlands district, located in the state of Pahang, Malaysia (Fig. 1), was selected

170  as a study area because it often experiences landslides and flash floods. These events are caused

171  by heavy and prolonged rainfall causing significant damages to properties. In this tropical

172  mountainous area, landslides are common as shown by government reports and past studies by

173  (Matori and Basith, 2012; Pradhan and Lee, 2010).

174  From the geomorphology aspect, the region is characterized as hilly, and altitudes are in the

175  range of 840–2110 meters (Sameen and Pradhan, 2019). The primary drainage characteristics

176  of the area consist of two rivers, namely, the Bertam and the Telom. Considerable types of

vegetation in Cameron Highlands include tropical forest and tea plantations, flower fields and temperate crops. Concerning lithology, the greater part of the region contains mega crystal biotite granites and phyllite as well as some schists layers (Pradhan and Lee, 2010). The area has a fair climate with an average annual rainfall starting from March to May and from November to December. The average nightly temperature of the study area is 14 °C, whereas the daily temperature reaches 24 °C. Approximately 8.0% (55 km²) of the area is classified as cropland, 86% (600 km²) is categorized as cultivated area, and 4.0% (27.5 km²) represents as residential areas.

**3.2. Landslide inventory map**

Data-driven landslide susceptibility assessment requires landslide inventories for model training and validation. Landslide inventory can be prepared using field investigations, historical landslide events from news and government reports and remote sensing data analysis. In this investigation, landslide inventories were taken from the study compiled by (Mezaal and Pradhan, 2018; Pradhan and Lee, 2010; Sameen et al., 2020). Overall, 156 landslides were identified and verified in the study area.

**3.3. Landslide conditioning factors**

Fifteen conditioning factors including elevation, slope, aspect, plan curvature, profile curvature, total curvature (Fig. 2a-f), lithology, LULC, distance to road, distance to river, SPI, STI (Fig. 2g-l), TRI, TWI and vegetation density (Fig. 2m-o) were selected as geo-environmental factors because they have been widely used in landslide susceptibility studies (Al-Najjar et al., 2019; Can et al., 2019; Canoglu et al., 2019; Huang and Zhao, 2018; Lee and Sambath, 2006). The related data were obtained over the study area on 15 January 2015 by utilizing a light detection and ranging (LiDAR) airborne system with a specification of 25,000 HZ pulse frequency rate and a density of 8 points/m². Then, a one-meter spatial resolution of

201   the digital elevation model was generated after removing non-ground points. Non-ground point

202   removal was performed utilizing multi-scale curvature and inverse distance weighted

203   interpolation approaches via ArcGIS Pro 2.4 software.

204   This study used six geomorphological factors, i.e. total curvature, plan curvature, profile

205   curvature, slope, elevation and slope aspect in the susceptibility mapping given that landslides

206   are influenced by terrain type. The elevation was included because it affects the extent of rock

207   weathering and is used by many scholars for landslide susceptibility assessment (Ayalew and

208   Yamagishi, 2005). The elevation of the investigation region was in the range from 690 to 1487

209   meters. The slope is another important factor, often included in landslide susceptibility studies

210   (Kamp et al., 2008). The slope values ranged from 0° to 78.88°. We also included the slope

211   direction (also known as slope aspect) because its task is to control concentrations of

212   topographic wetness affected by precipitation and solar radiation. In addition, plan, profile and

213   total curvature were also used (Ozdemir and Altural, 2013). In general, curvature affects slope

214   instability. Plan curvature represents the curvature when it is vertical to the path of the highest

215   slope. Profile curvature is parallel to the slope and designates the maximum slope orientation.

216   It affects the speeding up and slowing down of stream movement (Lee et al., 2004). The total

217   curvature is formed by combining the plane and profile curvatures (Romer and Ferentinou,

218   2016). If the surface is convex, the curvature is considered as positive; if it is concave, then it

219   is considered as negative. The value of zero reveals a linear surface (Al-Najjar et al., 2019).

220   Lithology and LULC were also used as conditioning factors for the preparation of landslide

221   susceptibility mapping. Lithology is important for landslide susceptibility assessment studies

222   because it affects the nature and system of landslides as rocks vary in form of mineral structure

223   besides internal formation (Kornejady et al., 2017). The lithology types in the study area are

224   mostly granite. The study area also contains schist, phyllite and slate types of lithology

225   (Pradhan and Lee, 2010). Human activities are also considered influential to landslides because

9

226  they affect patterns of land use and land cover. The LULC map of the study area obtained from

227  the Department of Survey and Mapping, Malaysia which shows that the area contains forest,

228  agricultural areas, urban areas, water bodies, transportation, barren lands and others (industrial,

229  infrastructure and utilities, institutions and community facilities). Also, the distance to the road

230  and river were included in our analysis.

231  Moreover, four hydrological factors were used in this study. These factors are topographic

232  wetness index (TWI), sediment transport index (STI), stream power index (SPI) and terrain

233  roughness index (TRI). SPI represents the movement of solid particles when gravity plays its

234  role on deposits (Rotigliano et al., 2012). STI represents slope failure and deposition. TRI

235  describes the coarseness of the local terrain which affects the topographic and hydrological

236  processes in the development of landslide occurrence. TWI reflects the direction and slope of

237  the flow, which is considered as a measurement for mastering the hydrological processes.

238  These factors were calculated using the following formulas (Yilmaz, 2009). Finally, vegetation

239  density was also used as a landslide conditioning factor. The vegetation density was calculated

240  using the normalized difference vegetation index variable (Pradhan, 2013) extracted from

241  Landsat 8 images. A vegetation density map was classified under four types, i.e. high-density

242  vegetation, medium density, poor density and non-vegetation.

243 $$\text{SPI} = A_s \times \tan\beta \qquad (1)$$

244 $$\text{STI} = \left(\frac{A_s}{22.13}\right)^{0.6} \times \left(\frac{\sin\beta}{0.0896}\right)^{1.3} \quad (2)$$

245 $$\text{TRI} = \sqrt{Abs\,(\max^2 - \min^2)} \qquad (3)$$

246 $$\text{TWI} = \ln\left(\frac{A_s}{\tan\beta}\right) \qquad (4)$$

247 where, *As* is defined as a specific area of the catchment (m²/m); ($\beta$) in radian, is a slope gradient

248 (in °); min and max values represent the highest and lowest number of rectangular cells within

249 nine DTM windows, respectively. The definition of the specific catchment is the area of the

250 slope in the upper slide per unit of the length of a contour, which is the area of cells divided by

251 the size of the cell (Kalantar et al., 2018).

252 **4. Methodology**

253 **4.1. Overview**

254 The proposed method creates synthetic inventory data using GANs for improving the

255 prediction of landslides. Fig. 3 illustrates the overall workflow of the current study. First, a

256 landslide inventory of 156 landslide locations and 15 conditioning factors were set as inputs

257 for the models. The inventory dataset was split into 70% of training and 30% of testing samples.

258 Then, five machine learning models (e.g. DT, RF, SVM, RF and Bagging ensemble) utilized

259 to evaluate the landslide susceptibility without additional samples. Thereafter, the GAN

260 method was used to create additional training samples with the existing inventory dataset; these

261 new samples were combined with the original training dataset and used to train the same

262 machine learning models again. Once the models were trained, they were tested with the same

263 test dataset used in the first case (without additional samples). Finally, the landslide

264 susceptibility maps were produced by the proposed models. Each map was classified into five

265 susceptibility categorical classes. These models were then validated and assessed using the area

266 under the receiver operating characteristic curve (AUROC).

267 **4.2. Description of machine learning techniques**

268 The following subsections describe the machine learning models used in this study.

269 **4.2.1. ANN model**

270 ANNs exhibit advantages over traditional computational methods (e.g. rule-based) because the
271 model does not require a straightforward practice to estimate desired yields (Jain et al.,
272 1996).After deciding on the number of hidden layers and the number of processing units in an
273 individual layer, the ANN starts learning from the training samples (Aditian et al., 2018).

### 274 4.2.2. SVM model

275 The goal of SVM models is to find the widest margin between two classes in feature space, by
276 a hyperplane (Vapnik, 1995). In landslide susceptibility, the aim is to discriminate between
277 susceptible (1) and not susceptible (−1) pixels. Its main advantages include mapping the data
278 to a high dimensional space where it is easier to classify with linear decision surfaces, also
279 reformulating problems so that data is mapped implicitly into this space.

### 280 4.2.3. Decision tree (DT) model

281 The DT model is a supervised and nonparametric machine learning technique that is operable
282 without prior knowledge about data distribution, with easy interpretation and capability to
283 model as well as it handles the reduction of data complexity and the relationships between
284 variables. Compared to other models, it is a flexible, fast, and robust algorithm that can be used
285 to control the nonlinearity between the input features and discrete classes so that nonlinear
286 relationships between parameters do not affect tree performance. Moreover, DT models are
287 simple to construct and clarify for decision-makers (Kadavi et al., 2019; Saito et al., 2009;
288 Yeon et al., 2010).

### 289 4.2.4. Random forest (RF) model

290 RF is a group of DTs that form an ensemble learning model used for classification and
291 regression problems (Liaw and Wiener, 2002). These models are effective for prediction
292 because they utilize the strength of each tree and their correlations and less sensitive to over-
293 fitting problems. The difference between RF and DT is that a decision tree is built on a whole

294  dataset, utilizing all the variables of interest, while a random forest randomly adopts

295  observations and specific variables to construct multiple decision trees from, and then averages

296  the results. In the present study, samples for landslide and non-landslide events were selected

297  to construct the classification tree (30% of the samples were kept aside from the training and

298  500 nodes were set as a favorite value).

**4.2.5. Bagging ensemble model**

299

300  In machine learning, several classifiers sometimes are combined and trained to boost the

301  prediction competence of a model (Polikar, 2012). Several combination methods, such as

302  Bagging, AdaBoost, multi boost and stacking can be used such as averaging or majority voting

303  (Breiman, 1996; Freund and Schapire, 1995; Kadavi and Lee, 2018; Webb, 2000). In landslide

304  susceptibility, Bagging has shown superiority over the other methods. Bagging, which is also

305  known as bootstrap aggregating, is a method of sub-dataset generation and combining learners.

306  In this study, the bootstrap samples were employed to build base learners utilizing similar

307  classification approaches, such as SVM and ANN. These based learners were then united by

308  the dominant voting technique.

**4.3. Additional data creation with GANs**

309

310  The GAN which was introduced by Goodfellow et al., (2014) is a type of neural network that

311  trained in an adversarial pattern to produce novel data mimicking specific divisions or

312  distributions. Since their invention, numerous upgraded versions of GANs (concerning

313  firmness of training and perceptual quality) have been developed, including Wasserstein,

314  conditional, Laplacian pyramid and deep convolutional GANs. GANs have been applied for

315  the generation of images, image in-painting, semi-supervised learning and image super-

316  resolution in various domains.

317    The general design of a GAN consists of two functions (Goodfellow et al., 2014), i.e. a

318    generator (G) and a discriminator (D) which its functionality is demonstrated in (Fig. 4). In

319    consideration of a random uniform distribution, the G maps a sample from the data distribution.

320    Meanwhile, the D is trained to discriminate whether the generated sample has a place in the

321    genuine distribution of the data. The G and D are generally learned together following game

322    theory, although they can be learned through other approaches and techniques.

323    For each duty, a sample from arbitrary noise $z$ is created by the $G$ to mislead $D$. Then, the real

324    samples are presented by the $D$, as well as the samples created by the $G$, to categorize the

325    samples as fake or real. By producing samples that can fool the $D$, the $G$ is rewarded. By

326    generating correct classification, $D$ is also rewarded. Both tasks are continuously revised until

327    a Nash equilibrium is obtained. Then, the repetition is paused. More particularly, let $D$ (s) be

328    the likelihood that $s$ originates from genuine information (real data) rather than the generator.

329    $G$ and $D$ play a minimax game with the following value function (Goodfellow et al., 2014).

330    $$\min_{G} \max_{D} V(D, G) = E_{s \sim p_{data}(s)}[\log D(s)] + E_{z \sim p_{z}(z)}[\log(1 - D(G(z)))] \quad (5)$$

331    **4.4. Validation of susceptibility maps**

332    For a given set of models, the validation was tested by calculating the area under the receiver

333    operating characteristic curve (AUROC). The inventory dataset was split into 70% of training

334    and 30% of testing samples. The ROC was created by plotting the sensitivity of the model

335    versus 1-specificity. The values of AUROC ranged from 0.5 to 1.0, where a high value

336    indicates the superiority of a model.

337

338    **5. Results**

339    **5.1. Application of RFs in selecting factors for modeling**

340 This study applied RFs to remove irrelevant factors from the analysis. The model was used

341 with 180 base estimators and the entire inventory dataset. After the model was trained, the

342 importance values of the 15 factors along with the standard deviation values were computed.

343 Table 1 shows the results of this analysis. The results indicate that the slope factor has the

344 greatest importance value (0.178), followed by LULC (0.171) and aspect (0.125). Most of the

345 landslides have occurred in moderate to high steep areas (slope > 18°). This characteristic

346 allowed the model to distinguish slides from non-slide pixels easily. Similarly, past landslides

347 have occurred in certain land use areas, such as forest, agriculture and barren lands. Schist

348 bedrock is more frequently exposed to slopes facing north through the southwest. The

349 remaining factors, except SPI and STI, also have significant contributions to landslide

350 occurrence. Thus, only SPI and STI were removed from the analysis in this study.

**5.2. Evaluation of five applied models (without additional data)**

352 The five models were evaluated by the most commonly used statistical measure, AUROC,

353 where 70% and 30% of the inventory samples were used as training and test data, respectively.

354 In all five models, the best values of the hyper-parameters as computed by the grid search over

355 a specific search space were used, which is shown in (Table 2). Table 3 shows the results

356 obtained for the studied models. The highest AUROC values for the training and test datasets

357 were achieved by the RF (0.94) and SVM (0.85) models, respectively. Using either the training

358 or test dataset, the ANN model has the lowest AUROC value compared with the other models.

359 The Bagging ensemble model was disadvantageous in the current study when SVM was used

360 as a base learner. The training and test accuracy of the SVM model was decreased by 0.04 and

361 0.1, respectively after the Bagging ensemble model was used. Therefore, the SVM model was

362 a good choice for the study area. However, SVM still faces challenges. For example, it slows

363 down with additional factors, its predictive capability can be degraded with a smaller training

364     sample size and it requires careful optimization of the penalty parameter and the kernel

365     function.

366     **5.3. Evaluation of applied models (with additional data)**

367     Additional training samples were generated by the proposed GAN model. These new samples

368     were combined with the original training dataset and used to train the same models again. Once

369     the models were trained, they were tested with the same test dataset used in the previous section

370     (Section 5.2). Thus, a fair comparison was conducted to evaluate the proposed GAN model.

371     Table 4 shows the AUROC values obtained for the five models using the training (with

372     additional samples) and test datasets. The highest training accuracy was achieved by the RF

373     model (0.94). The RF and SVM models achieved the same accuracy (0.82) using the test

374     dataset. ANN has the lowest training accuracy of 0.75. However, ANN is as accurate as of the

375     DT model on the test dataset.

376     The additional samples created by the GAN model contributed to increasing the training

377     accuracy of the five models, except that the RF model that achieved the same accuracy in both

378     cases. The ANN model gained the greatest benefit from the additional samples as its training

379     accuracy increased by 0.06. Using the test dataset, the additional samples improved the

380     predictive capability of the models, except that of the SVM model whose test accuracy was

381     decreased by 0.03.

382     By employing the proposed models, five landslide susceptibility maps were generated from the

383     study area using natural break methods (Fig. 5). Each map was classified into five categorical

384     classes, i.e. very low, low, moderate, high and very high. The blue indicates a low susceptible

385     area, whereas red indicates a highly susceptible area.

386     **5.4. Influence of additional samples created by GANs on model performance**

387    Various numbers of additional samples (5, 10, 20, 30, 40, 50, 100 and 500) were tested to

388    analyze the influence of the number of generated samples on the performance of the models'

389    prediction (Fig. 6). The analysis showed that the DT model performed the best using 10

390    additional samples on the training dataset, but performed worse using more than 50 additional

391    samples. On the test dataset, the DT model performed the best with 40 additional samples.

392    Meanwhile, the SVM model suffered from over-fitting on the training dataset using additional

393    samples. With 500 additional samples, the SVM model achieved 0.97 AUROC on the training

394    dataset, but it achieved only 0.72 AUROC on the test dataset. Similar results were observed for

395    the Bagging ensemble model. With 500 additional samples, the model achieved 0.94 AUROC

396    on the training dataset and 0.65 AUROC on the test dataset, thereby indicating over-fitting.

397    Similarly, the ANN model also suffered from over-fitting on the training dataset. It achieves

398    0.75 AUROC with 5 additional samples and 0.91 with 500 additional samples. Among the

399    models, the RF model was less sensitive to the number of additional samples. The best accuracy

400    remained with the 50 additional samples on both datasets. The generation of samples with

401    GANs does not always guarantee to improve model accuracy. Various tests should be evaluated

402    before deciding on the final susceptibility models.

403    **6. Discussion**

404    Machine learning has been an effective landslide susceptibility mapping method. However,

405    with insufficient data, these machine learning models often suffer from generalizing to areas

406    other than the training area. Especially in landslide susceptibility mapping, gathering inventory

407    data is expensive, and some areas have not experienced a large number of landslides.

408    Nevertheless, many studies have attempted to develop models that work with insufficient data.

409    For example, sampling strategy and validation methods have been validated to address the

410    challenges of modeling with limited data effectively. Given that randomness of the training,

411    data selection influences the model performance (Kalantar et al., 2018), sampling strategies

412   that avoid model over-fitting to the training data have been proposed (Aktas and San, 2019;

413   Conoscenti et al., 2016). More often than not, landslide inventory data are incomplete. Such

414   incomplete data affect the selection of the absence samples. For this problem, Steger et al.

415   (2016) suggested that models can correlate with landslide inventory incompleteness, and thus,

416   they should be rejected regardless of their performance. Techniques such as factor

417   optimization, development of new factors and model ensembling have also been extensively

418   discussed in the recent literature.

419   Removing insignificant factors was useful to decrease the impact of model over-fitting due to

420   the limited training. The RF model showed that SPI and STI were not influential and thus were

421   removed from the analysis. Estimation of the factors also plays an important role in obtaining

422   insights into the factors included in the model. Similar to previous studies, the present study

423   found the slope to be a significant factor. The landslide inventory dataset showed that most of

424   the landslides have occurred in moderate to high steep areas. A significant number of past

425   landslides have occurred in certain land use areas, such as forest, agriculture and barren lands.

426   The results of the RF model were also consistent with the inventory data, where LULC and

427   aspect were found to be significant.

428   The evaluation of the models with and without additional samples showed that the proposed

429   GAN can improve the performance of the susceptibility model. When the training data were

430   used, the GAN model improved the accuracy of all the models except RF. Some models, such

431   as ANN, performed better than others. Using the test data contributed to increasing the

432   accuracy of all the models except SVM. Moreover, the number of additional samples

433   significantly affected the modeling performance. The DT, SVM and ANN models over-fitted

434   the training data when a large number of additional samples were included in the training set.

435   The RF model was less sensitive to the number of additional samples than other models. Thus,

436   adding newly generated samples to the training set may not always lead to an increase in model

18

437   accuracy, especially on the test data. Therefore, the number of additional samples should be

438   considered as a parameter and fine-tuned before training any machine learning model.

439   **7. Conclusions**

440   This study addressed the aforementioned problem with a GAN-based method. This model was

441   used to create an additional training sample with the existing inventory dataset. The proposed

442   method was evaluated on a dataset taken from Cameron Highlands, Malaysia. Five machine

443   learning and statistical models were implemented to assess the proposed GAN model. The

444   outcomes revealed that using additional samples created by the proposed GAN model can

445   improve the predictive capability of the studied models, except SVM.

446   Generative models, such as GANs, can be useful for landslide susceptibility mapping,

447   especially when the training data for the area under study are inadequate. However, the used

448   models should be carefully analyzed to avoid over-fitting to the training samples. In addition,

449   the hyper-parameters of the used models can be optimized to improve the overall performance

450   of the landslide susceptibility models when samples created by generative models are used.

451   Improvements in landslide susceptibility maps can help in the implementation of land use

452   planning and the design of landslide mitigation strategies. Improvements in landslide

453   susceptibility models also contribute towards improving landslide hazard and risk assessment.

454   The proposed method, therefore, can be a useful tool for engineers, geoscientists and planners.

during this research. We are also thankful to the Department of Mineral and Geosciences, the Department of Surveying Malaysia, and the Federal Department of Town and Country Planning.

**References**

Aditian, A., Kubota, T., Shinohara, Y., 2018. Comparison of GIS-based landslide susceptibility models using frequency ratio, logistic regression, and artificial neural network in a tertiary region of Ambon, Indonesia. Geomorphology 318, 101–111. https://doi.org/10.1016/j.geomorph.2018.06.006

Akbar, A.Q., Chen, G., 2018. Comparison of major statistical methods and their combination using matrix validation for landslide susceptibility mapping. Lowland Technology International 20, 401–412.

Aktas, H., San, B.T., 2019. Landslide susceptibility mapping using an automatic sampling algorithm based on two level random sampling. Computers and Geosciences 133, 104329. https://doi.org/10.1016/j.cageo.2019.104329

Al-Najjar, H.A.H., Kalantar, B., Pradhan, B., Saeidi, V., 2019. Conditioning factor determination for mapping and prediction of landslide susceptibility using machine learning algorithms, in: In Earth Resources and Environmental Remote Sensing/GIS Applications. International Society for Optics and Photonics. p. 19. https://doi.org/10.1117/12.2532687

Ashournejad, Q., Hosseini, A., Pradhan, B., Hosseini, S.J., 2019. Hazard zoning for spatial planning using GIS-based landslide susceptibility assessment: a new hybrid integrated data-driven and knowledge-based model. Arabian Journal of Geosciences 12. https://doi.org/10.1007/s12517-019-4236-0

484    Ayalew, L., Yamagishi, H., 2005. The application of GIS-based logistic regression for

485        landslide susceptibility mapping in the Kakuda-Yahiko Mountains, Central Japan.

486        Geomorphology 65, 15–31. https://doi.org/10.1016/j.geomorph.2004.06.010

487    Bragagnolo, L., Silva, R.V. d., Grzybowski, J.M.V., 2020. Artificial neural network ensembles

488        applied to the mapping of landslide susceptibility. Catena 184, 104240.

489        https://doi.org/10.1016/j.catena.2019.104240

490    Braun, A., Urquia, Elias Leonardo Garcia Lopez, Rigoberto Moncada Yamagishi, H., 2018.

491        Landslide Susceptibility Mapping in Tegucigalpa, Honduras, Using Data Mining

492        Methods, in: Slope Stability: Case Histories, Landslide Mapping, Emerging

493        Technologies. pp. 207–215.

494    Breiman,    L.,    1996.    Bagging    Predictors.    Machine    Learning    24,    123–140.

495        https://doi.org/10.1007/BF00058655

496    Can, A., Dagdelenler, G., Ercanoglu, M., Sonmez, H., 2019. Landslide susceptibility mapping

497        at Ovacık-Karabük (Turkey) using different artificial neural network models:

498        comparison of training algorithms. Bulletin of Engineering Geology and the

499        Environment 78, 89–102. https://doi.org/10.1007/s10064-017-1034-3

500    Canoglu, M.C., Aksoy, H., Ercanoglu, M., 2019. Integrated approach for determining spatio-

501        temporal variations in the hydrodynamic factors as a contributing parameter in landslide

502        susceptibility assessments. Bulletin of Engineering Geology and the Environment 78,

503        3159–3174. https://doi.org/10.1007/s10064-018-1337-z

504    Ciurleo, M., Cascini, L., Calvello, M., 2017. A comparison of statistical and deterministic

505        methods for shallow landslide susceptibility zoning in clayey soils. Engineering

506        Geology 223, 71–81. https://doi.org/10.1016/j.enggeo.2017.04.023

507 Conoscenti, C., Rotigliano, E., Cama, M., Caraballo-Arias, N.A., Lombardo, L., Agnesi, V.,

508       2016. Exploring the effect of absence selection on landslide susceptibility models: A

509       case study in Sicily, Italy. Geomorphology 261, 222–235.

510       https://doi.org/10.1016/j.geomorph.2016.03.006

511 Dou, J., Bui, D.T., Yunus, A.P., Jia, K., Song, X., Revhaug, I., Xia, H., Zhu, Z., 2015.

512       Optimization of causative factors for landslide susceptibility evaluation using remote

513       sensing and GIS data in parts of Niigata, Japan. PLoS ONE 10.

514       https://doi.org/10.1371/journal.pone.0133262

515 Du, J., Glade, T., Woldai, T., Chai, B., Zeng, B., 2020. Landslide susceptibility assessment

516       based on an incomplete landslide inventory in the Jilong Valley, Tibet, Chinese

517       Himalayas. Engineering Geology 270, 105572.

518       https://doi.org/10.1016/j.enggeo.2020.105572

519 Fanos, A.M., Pradhan, B., 2019. A novel hybrid machine learning-based model for rockfall

520       source identification in presence of other landslide types using LiDAR and GIS. Earth

521       Systems and Environment 3, 491–506. https://doi.org/10.1007/s41748-019-00114-z

522 Feizizadeh, B., Roodposhti, M.S., Blaschke, T., Aryal, J., 2017. Comparing GIS-based support

523       vector machine kernel functions for landslide susceptibility mapping. Arabian Journal

524       of Geosciences 10. https://doi.org/10.1007/s12517-017-2918-z

525 Formetta, G., Rago, V., Capparelli, G., Rigon, R., Muto, F., Versace, P., 2014. Integrated

526       Physically based System for Modeling Landslide Susceptibility. Procedia Earth and

527       Planetary Science 9, 74–82. https://doi.org/10.1016/j.proeps.2014.06.006

528 Freund, Y., Schapire, R.E., 1995. A desicion-theoretic generalization of on-line learning and

529       an application to boosting., in: In European Conference on Computational Learning

530       Theory. pp. 23–37. https://doi.org/10.1007/3-540-59119-2

531    Glade, T., Anderson, M., Crozier, M.J., 2012. Landslide hazard and risk.
532        https://doi.org/10.1002/9780470012659

533    Goetz, J.N., Brenning, A., Petschko, H., Leopold, P., 2015. Evaluating machine learning and
534        statistical prediction techniques for landslide susceptibility modeling. Computers and
535        Geosciences 81, 1–11. https://doi.org/10.1016/j.cageo.2015.04.007

536    Goetz, J.N., Guthrie, R.H., Brenning, A., 2011. Integrating physical and empirical landslide
537        susceptibility models using generalized additive models. Geomorphology 129, 376–
538        386. https://doi.org/10.1016/j.geomorph.2011.03.001

539    Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville,
540        A., Bengio, Y., 2014. Generative Adversarial Nets, in: In Advances in Neural
541        Information Processing Systems (Pp. 2672-2680). pp. 2672–2680.

542    Hong, H., Miao, Y., Liu, J., Zhu, A.X., 2019. Exploring the effects of the design and quantity
543        of absence data on the performance of random forest-based landslide susceptibility
544        mapping. Catena 176, 45–64. https://doi.org/10.1016/j.catena.2018.12.035

545    Huang, Y., Zhao, L., 2018. Review on landslide susceptibility mapping using support vector
546        machines. Catena 165, 520–529. https://doi.org/10.1016/j.catena.2018.03.003

547    Hussin, H.Y., Zumpano, V., Reichenbach, P., Sterlacchini, S., Micu, M., van Westen, C.,
548        Bălteanu, D., 2016. Different landslide sampling strategies in a grid-based bi-variate
549        statistical susceptibility model. Geomorphology 253, 508–523.
550        https://doi.org/10.1016/j.geomorph.2015.10.030

551    Jain, A.K., Mao, J. and, Mohiuddin, K.M., 1996. Artificial neural networks: A tutorial.
552        Computer 29, 31–44.

553  Kadavi, P.R., Lee, C., 2018. Application of Ensemble-Based Machine Learning Models to
554      Landslide    Susceptibility    Mapping.    Remote    Sensing    10,    1–18.
555      https://doi.org/10.3390/rs10081252

556  Kadavi, P.R., Lee, C.W., Lee, S., 2019. Landslide-susceptibility mapping in Gangwon-do,
557      South Korea, using logistic regression and decision tree models. Environmental Earth
558      Sciences 78, 1–17. https://doi.org/10.1007/s12665-019-8119-1

559  Kalantar, B., Pradhan, B., Amir Naghibi, S., Motevalli, A., Mansor, S., 2018. Assessment of
560      the effects of training data selection on the landslide susceptibility mapping: a
561      comparison between support vector machine (SVM), logistic regression (LR) and
562      artificial neural networks (ANN). Geomatics, Natural Hazards and Risk 9, 49–69.
563      https://doi.org/10.1080/19475705.2017.1407368

564  Kamp, U., Growley, B.J., Khattak, G.A., Owen, L.A., 2008. GIS-based landslide susceptibility
565      mapping for the 2005 Kashmir earthquake region. Geomorphology 101, 631–642.
566      https://doi.org/10.1016/j.geomorph.2008.03.003

567  Kavzoglu, T., Colkesen, I., Sahin, E.K., 2019. Landslides: Theory, Practice and Modelling 50,
568      283–301. https://doi.org/10.1007/978-3-319-77377-3

569  Kavzoglu, T., Kutlug Sahin, E., Colkesen, I., 2015. Selecting optimal conditioning factors in
570      shallow translational landslide susceptibility mapping using genetic algorithm.
571      Engineering Geology 192, 101–112. https://doi.org/10.1016/j.enggeo.2015.04.004

572  Kawabata, D., Bandibas, J., 2009. Landslide susceptibility mapping using geological data, a
573      DEM from ASTER images and an Artificial Neural Network (ANN). Geomorphology
574      113, 97–109. https://doi.org/10.1016/j.geomorph.2009.06.006

575     Kornejady, A., Ownegh, M., Bahremand, A., 2017. Landslide susceptibility assessment using

576          maximum entropy model with two different data sampling methods. Catena 152, 144–

577          162. https://doi.org/10.1016/j.catena.2017.01.010

578     Kornejady, A., Ownegh, M., Rahmati, O., Bahremand, A., 2018. Landslide susceptibility

579          assessment using three bivariate models considering the new topo-hydrological factor:

580          HAND.          Geocarto          International          33,          1155–1185.

581          https://doi.org/10.1080/10106049.2017.1334832

582     Lai, J.S., Chiang, S.H., Tsai, F., 2019. Exploring influence of sampling strategies on event-

583          based landslide susceptibility modeling. ISPRS International Journal of Geo-

584          Information 8. https://doi.org/10.3390/ijgi8090397

585     Lee, J.H., Sameen, M.I., Pradhan, B., Park, H.J., 2018. Modeling landslide susceptibility in

586          data-scarce environments using optimized data mining and statistical methods.

587          Geomorphology 303, 284–298. https://doi.org/10.1016/j.geomorph.2017.12.007

588     Lee, M.L., Ng, K.Y., Huang, Y.F., Li, W.C., 2014. Rainfall-induced landslides in Hulu Kelang

589          area, Malaysia. Natural Hazards 70, 353–375. https://doi.org/10.1007/s11069-013-

590          0814-8

591     Lee, S., Ryu, J.H., Won, J.S., Park, H.J., 2004. Determination and application of the weights

592          for landslide susceptibility mapping using an artificial neural network. Engineering

593          Geology 71, 289–302. https://doi.org/10.1016/S0013-7952(03)00142-X

594     Lee, S., Sambath, T., 2006. Landslide susceptibility mapping in the Damrei Romel area,

595          Cambodia using frequency ratio and logistic regression models. Environmental

596          Geology 50, 847–855. https://doi.org/10.1007/s00254-006-0256-7

597     Liaw, A., Wiener, M., 2002. Classification and regression by randomForest. R News 2, 18–22.

598   Mandal, S., Mondal, S., 2019. Statistical approaches for landslide susceptibility assessment

599         and prediction, in: Statistical Approaches for Landslide Susceptibility Assessment and

600         Prediction. pp. 181–189. https://doi.org/10.1007/978-3-319-93897-4

601   Matori, A.N., Basith, A., 2012. Evaluation of landslide causative factors towards efficient

602         landslide susceptibility modelling in the Cameron Highlands, Malaysia. WIT

603         Transactions     on     Engineering     Sciences     73,     207–218.

604         https://doi.org/10.2495/DEB120181

605   Mezaal, M.R., Pradhan, B., 2018. An improved algorithm for identifying shallow and deep-

606         seated landslides in dense tropical forest from airborne laser scanning data. Catena 167,

607         147–159. https://doi.org/10.1016/j.catena.2018.04.038

608   Ozdemir, A., Altural, T., 2013. A comparative study of frequency ratio, weights of evidence

609         and logistic regression methods for landslide susceptibility mapping: Sultan mountains,

610         SW     Turkey.     Journal     of     Asian     Earth     Sciences     64,     180–197.

611         https://doi.org/10.1016/j.jseaes.2012.12.014

612   Park, J.Y., Lee, S.R., Lee, D.H., Kim, Y.T., Lee, J.S., 2019. A regional-scale landslide early

613         warning methodology applying statistical and physically based approaches in sequence.

614         Engineering Geology 260, 105193. https://doi.org/10.1016/j.enggeo.2019.105193

615   Polikar,     R.,     2012.     Ensemble     Learning,     In     Ensemble     machine     learning.

616         https://doi.org/10.1007/978-1-4419-9326-7

617   Pradhan, B., 2013. A comparative study on the predictive ability of the decision tree, support

618         vector machine and neuro-fuzzy models in landslide susceptibility mapping using GIS.

619         Computers and Geosciences 51, 350–365. https://doi.org/10.1016/j.cageo.2012.08.023

620    Pradhan, B., Lee, S., 2010. Regional landslide susceptibility analysis using back-propagation

621        neural network model at Cameron Highland, Malaysia. Landslides 7, 13–30.

622        https://doi.org/10.1007/s10346-009-0183-2

623    Romer, C., Ferentinou, M., 2016. Shallow landslide susceptibility assessment in a semiarid

624        environment - A Quaternary catchment of KwaZulu-Natal, South Africa. Engineering

625        Geology 201, 29–44. https://doi.org/10.1016/j.enggeo.2015.12.013

626    Rotigliano, E., Cappadonia, C., Conoscenti, C., Costanzo, D., Agnesi, V., 2012. Slope units-

627        based flow susceptibility model: Using validation tests to select controlling factors.

628        Natural Hazards 61, 143–153. https://doi.org/10.1007/s11069-011-9846-0

629    Saito, H., Nakayama, D., Matsuyama, H., 2009. Comparison of landslide susceptibility based

630        on a decision-tree model and actual landslide occurrence: The Akaishi Mountains,

631        Japan. Geomorphology 109, 108–121. https://doi.org/10.1016/j.geomorph.2009.02.026

632    Sameen, M.I., Pradhan, B., 2019. Landslide detection using residual networks and the fusion

633        of spectral and topographic information. IEEE Access 7, 114363–114373.

634        https://doi.org/10.1109/access.2019.2935761

635    Sameen, M.I., Pradhan, B., Bui, D.T., Alamri, A.M., 2020. Systematic sample subdividing

636        strategy for training landslide susceptibility models. Catena 187, 104358.

637        https://doi.org/10.1016/j.catena.2019.104358

638    Samia, J., Temme, A., Bregt, A.K., Wallinga, J., Stuiver, J., Guzzetti, F., Ardizzone, F., Rossi,

639        M., 2018. Implementing landslide path dependency in landslide susceptibility

640        modelling. Landslides 15, 2129–2144. https://doi.org/10.1007/s10346-018-1024-y

641    Soma, A.S., Kubota, T., Mizuno, H., 2019. Optimization of causative factors using logistic

642        regression and artificial neural network models for landslide susceptibility assessment

643        in Ujung Loe Watershed, South Sulawesi Indonesia. Journal of Mountain Science 16,

644        383–401. https://doi.org/10.1007/s11629-018-4884-7

645    Steger, S., Brenning, A., Bell, R., Glade, T., 2018. Incompleteness matters − An approach to

646        counteract inventory-based biases in statistical landslide susceptibility modelling 20,

647        8551.

648    Steger, S., Brenning, A., Bell, R., Glade, T., 2016. The influence of systematically incomplete

649        shallow landslide inventories on statistical susceptibility models and suggestions for

650        improvements. Landslides 14, 1767–1781. https://doi.org/10.1007/s10346-017-0820-0

651    Süzen, M.L., Doyuran, V., 2004. A comparison of the GIS based landslide susceptibility

652        assessment methods: Multivariate versus bivariate. Environmental Geology 45, 665–

653        679. https://doi.org/10.1007/s00254-003-0917-8

654    Tsangaratos, P., Ilia, I., 2016. Comparison of a logistic regression and Naïve Bayes classifier

655        in landslide susceptibility assessments: The influence of models complexity and

656        training dataset size. Catena 145, 164–179.

657        https://doi.org/10.1016/j.catena.2016.06.004

658    Vapnik, V.N., 1995. Constructing learning algorithms. In The nature of statistical learning

659        theory.

660    Wang, H.J., Xiao, T., Li, X.Y., Zhang, L.L., Zhang, L.M., 2019. A novel physically-based

661        model for updating landslide susceptibility. Engineering Geology 251, 71–80.

662        https://doi.org/10.1016/j.enggeo.2019.02.004

663    Webb, G.I., 2000. MultiBoosting: a technique for combining boosting and wagging. Machine

664        Learning 40, 159–196. https://doi.org/10.1023/A:1007659514849

665  Xiao, T., Yin, K., Yao, T., Liu, S., 2019. Spatial prediction of landslide susceptibility using

666      GIS-based statistical and machine learning models in Wanzhou County, Three Gorges

667      Reservoir, China. Acta Geochimica 38, 654–669. https://doi.org/10.1007/s11631-019-

668      00341-1

669  Yan, F., Zhang, Q., Ye, S., Ren, B., 2019. A novel hybrid approach for landslide susceptibility

670      mapping integrating analytical hierarchy process and normalized frequency ratio

671      methods with the cloud model. Geomorphology 327, 170–187.

672      https://doi.org/10.1016/j.geomorph.2018.10.024

673  Yeon, Y.K., Han, J.G., Ryu, K.H., 2010. Landslide susceptibility mapping in Injae, Korea,

674      using a decision tree. Engineering Geology 116, 274–283.

675      https://doi.org/10.1016/j.enggeo.2010.09.009

676  Yilmaz, I., 2009. Landslide susceptibility mapping using frequency ratio, logistic regression,

677      artificial neural networks and their comparison: A case study from Kat landslides

678      (Tokat-Turkey). Computers and Geosciences 35, 1125–1138.

679      https://doi.org/10.1016/j.cageo.2008.08.007

680  Yilmaz, I., Ercanoglu, M., 2019. Natural Hazards GIS-Based Spatial Modeling Using Data

681      Mining Techniques, in: Landslide Inventory, Sampling and Effect of Sampling

682      Strategies on Landslide Susceptibility/Hazard Modelling at a Glance. pp. 205–224.

683  Zêzere, J.L., Pereira, S., Melo, R., Oliveira, S.C., Garcia, R.A.C., 2017. Mapping landslide

684      susceptibility using data-driven methods. Science of the Total Environment 589, 250–

685      267. https://doi.org/10.1016/j.scitotenv.2017.02.188

686  Zhang, T. yu, Han, L., Zhang, H., Zhao, Y. hua, Li, X. an, Zhao, L., 2019. GIS-based landslide

687      susceptibility mapping using hybrid integration approaches of fractal dimension with

688       index of entropy and support vector machine. Journal of Mountain Science 16, 1275–

689       1288. https://doi.org/10.1007/s11629-018-5337-z

690   Zhu, A.X., Miao, Y., Liu, J., Bai, S., Zeng, C., Ma, T., Hong, H., 2019. A similarity-based

691       approach to sampling absence data for landslide susceptibility mapping using data-

692       driven methods. Catena 183, 104188. https://doi.org/10.1016/j.catena.2019.104188

693

694   **Figure caption**

695   **Fig. 1.** Location of the study area and landslide inventory map

696   **Fig. 2a-f.** Maps of landslide conditioning factors: (a) Elevation, (b) Slope, (c) Aspect, (d) Plan

697   curvature, (e) Profile curvature, and (f) Total curvature.

698   **Fig. 2g-l.** Maps of landslide conditioning factors: (g) Lithology, (h) LULC, (i) Distance to road, (j)

699   Distance to river, (k) SPI, and (l) STI.

700   **Fig. 2m-o.** Maps of landslide conditioning factors: (m) TRI, (n) TWI, and (o) Vegetation density.

701   **Fig. 3.** Overall workflow used in this study.

702   **Fig. 4.** The general architecture of GANs.

703   **Fig. 5.** Landslide susceptibility maps produced by proposed (a) DT, (b) RF, (c) SVM, (d) Bagging

704   ensemble, and (e) ANN models.

705   **Fig. 6.** Training and test AUROC values calculated for the five models trained with original training

706   dataset and additional samples created by GANs.

707

708   **Table caption**

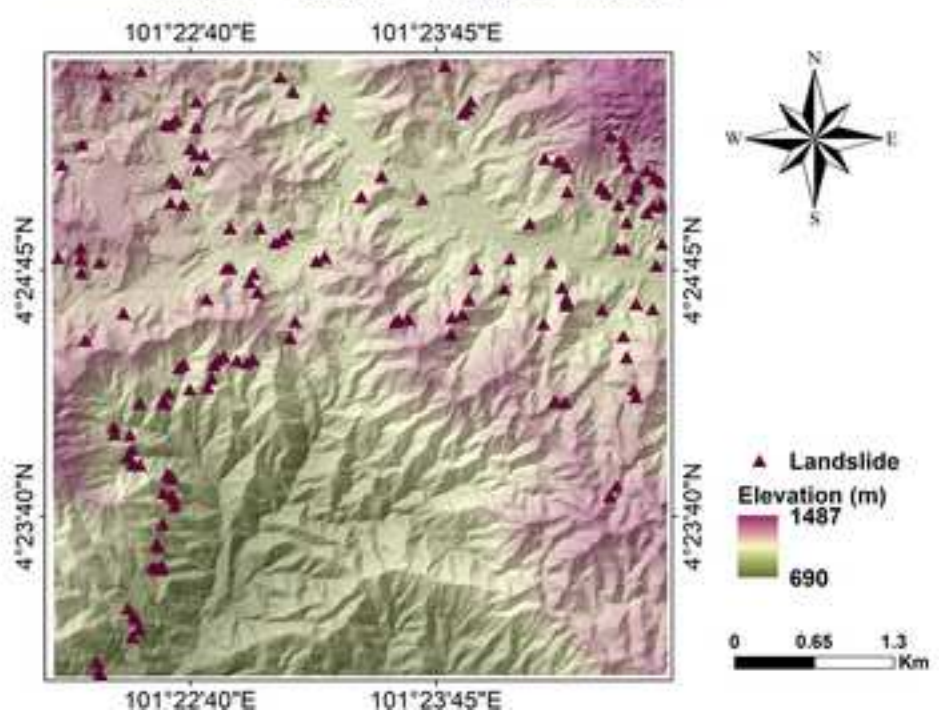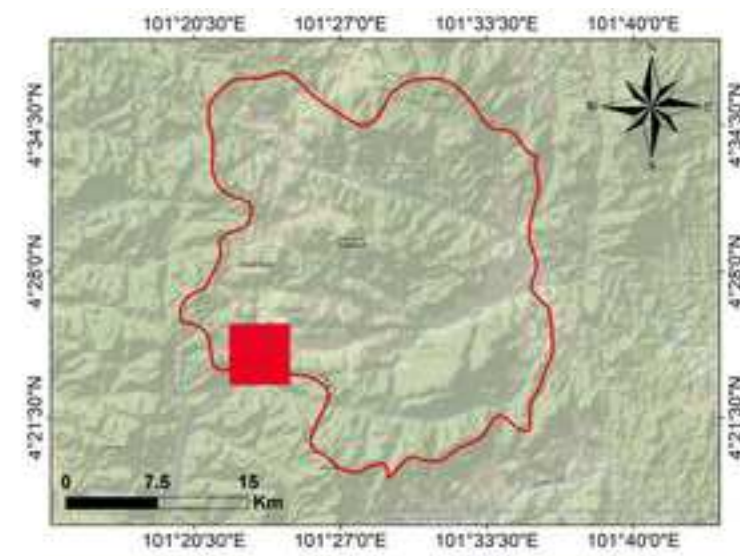709   **Table 1.** Importance of affecting factors.

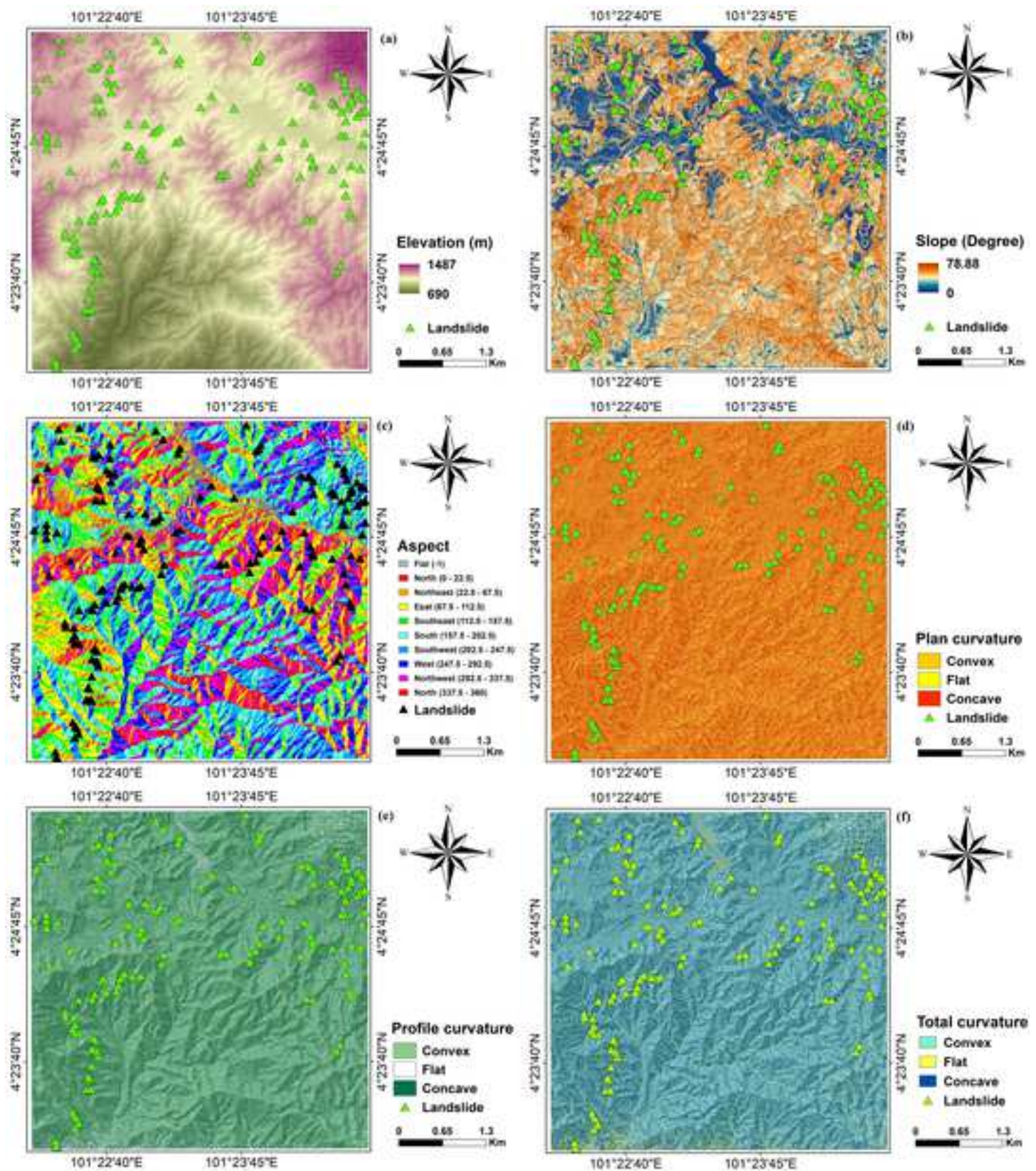710   **Table 2.** Optimised parameters of five models and search spaces.

711    **Table 3.** AUROC values of five models using training and test datasets.
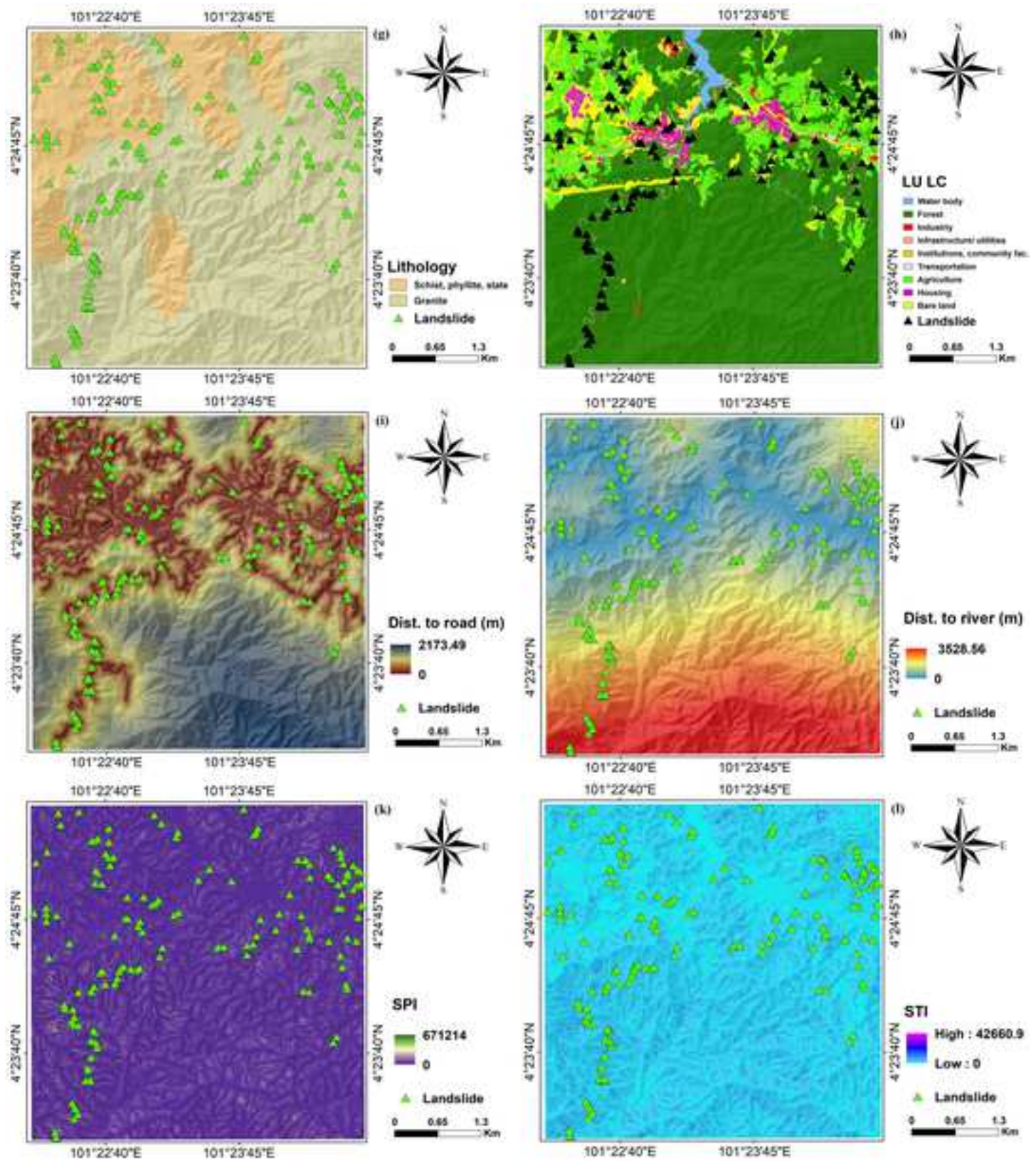
712    **Table 4.** AUROC values of models using training (with additional samples) and test datasets.
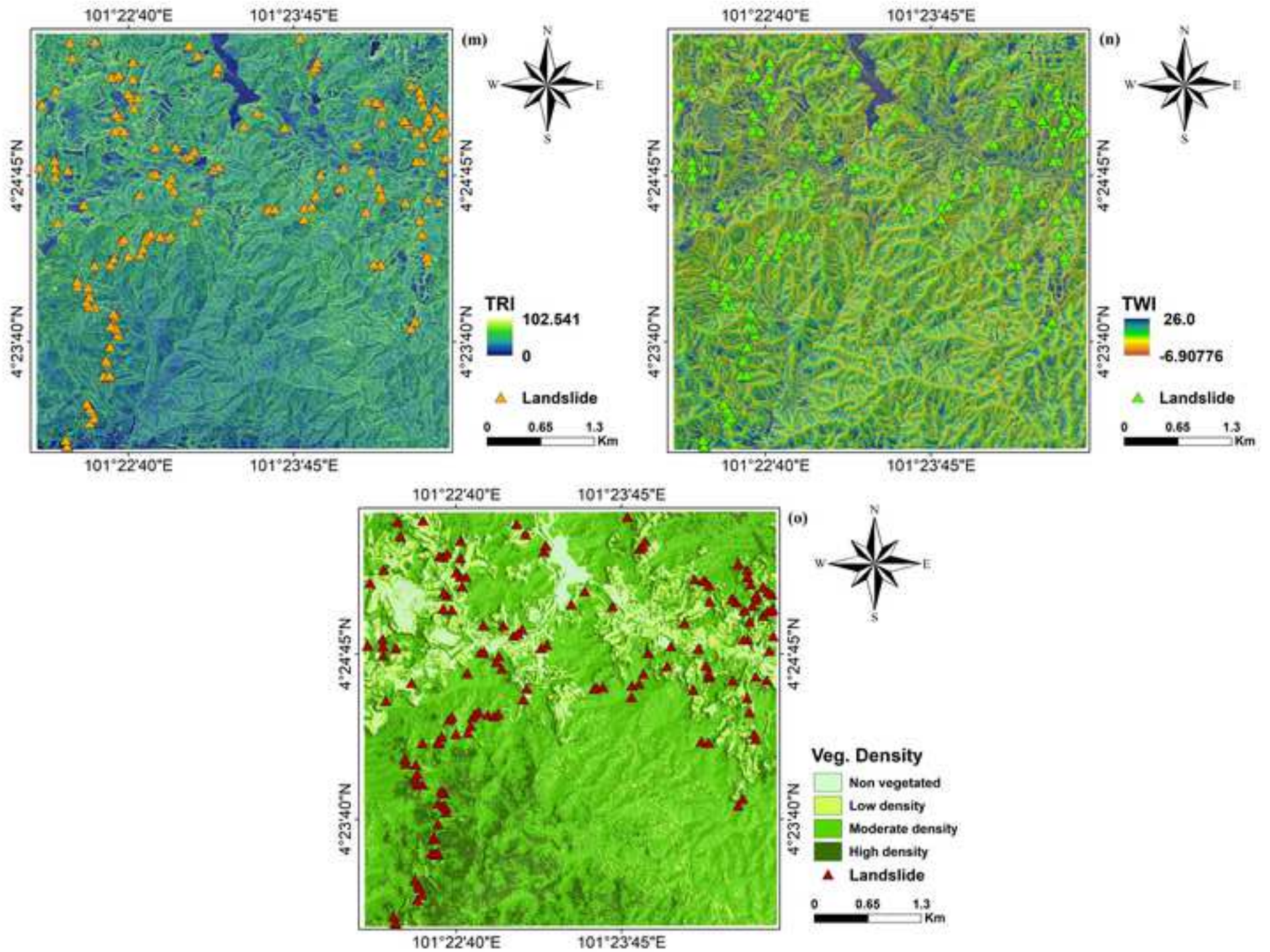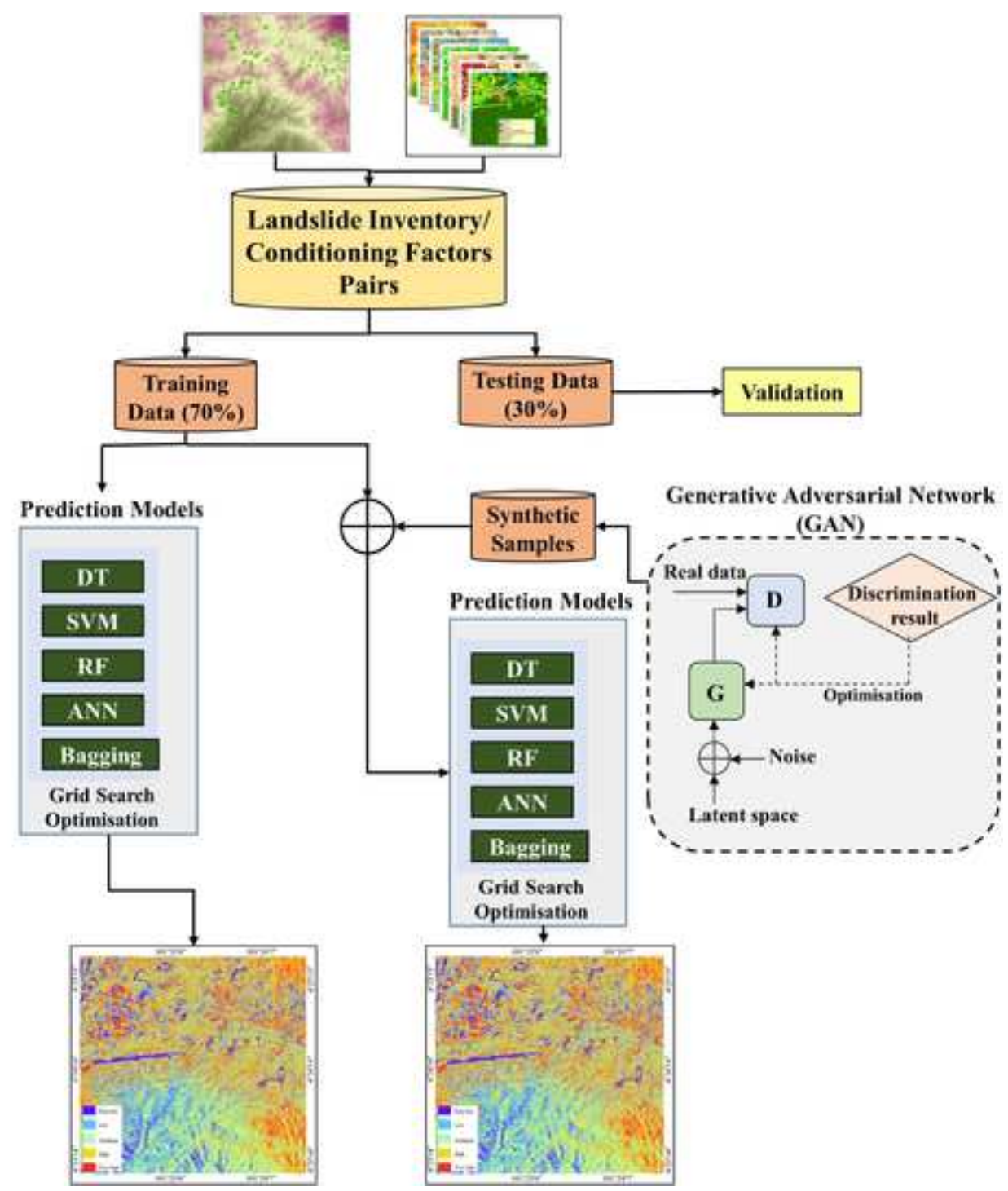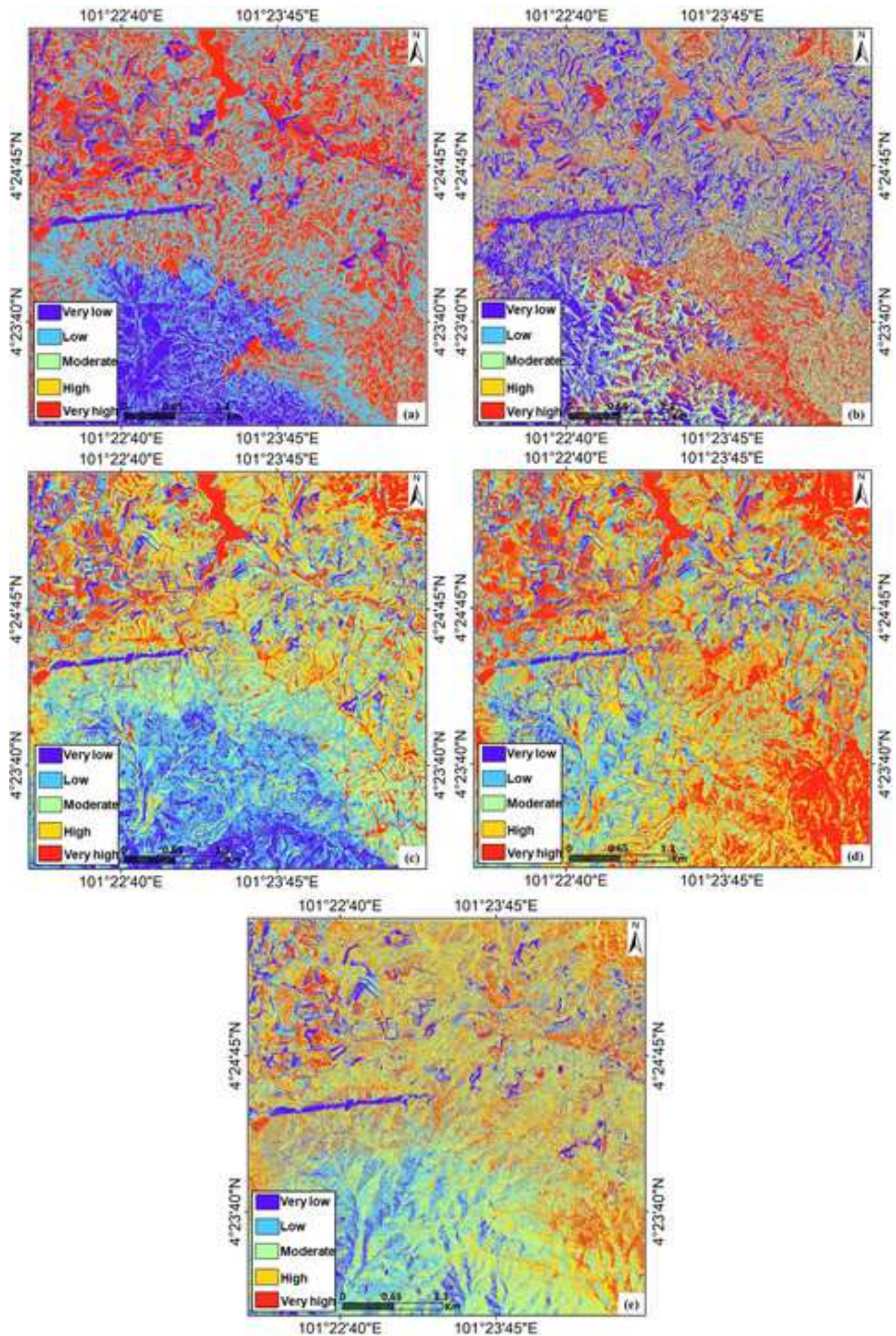
713

Figure 1

Click here to access/download;Figure;Fig.1.jpg ±

Figure 3

Figure 4

Click here to access/download;Figure;Fig. 4.tif ⬇

Figure 5                                    Click here to access/download;Figure;Fig. 5 .jpg ⬇

Figure 6     Click here to access/download;Figure;Fig. 6 .tif ⬇

Table 1

| Factor | RF importance value | Standard deviation |
|---|---|---|
| Slope | 0.178 | 0.022 |
| LULC | 0.171 | 0.012 |
| Aspect | 0.125 | 0.013 |
| Elevation | 0.09 | 0.026 |
| Vegetation density | 0.073 | 0.017 |
| TWI | 0.063 | 0.023 |
| Distance to river | 0.061 | 0 |
| TRI | 0.06 | 0.018 |
| Lithology | 0.041 | 0.019 |
| Total curvature | 0.038 | 0.014 |
| Plane curvature | 0.034 | 0.039 |
| Profile curvature | 0.032 | 0.015 |
| Distance to road | 0.029 | 0.045 |
| SPI | 0 | 0 |
| STI | 0 | 0.015 |

Table 2

| Model | Parameters | Search space | Best value (grid search with 10-fold cross validation) |
|---|---|---|---|
| DT | Maximum tree depth | [2–13] | 5 |
| RF | Number of base estimators | [10–1000] | 180 |
| SVM | C | [1–1000] | 5 |
| | Kernel function | [Linear, RBF, Sigmoid] | RBF |
| | Learning rate | $[10^{-5}–1.0]$ | 0.01 |
| | Activation function | [ReLU, Tanh, Sigmoid, Linear] | ReLU |
| ANN | Number of hidden layers | [1–12] | 1 |
| | Number of hidden units in a hidden layer | [2–1024] | 62 |
| Bagging ensemble | Base learner | [DT, RF, SVM, ANN] | SVM |

Table 3

| Model | Training AUROC | Test AUROC |
|---|---|---|
| DT | 0.9 | 0.76 |
| RF | 0.94 | 0.81 |
| SVM | 0.86 | 0.85 |
| ANN | 0.69 | 0.72 |
| Bagging ensemble | 0.82 | 0.75 |

Table 4

| Model | Training AUROC | Test AUROC |
|---|---|---|
| DT | 0.92 | 0.78 |
| RF | 0.94 | 0.82 |
| SVM | 0.88 | 0.82 |
| ANN | 0.75 | 0.78 |
| Bagging ensemble | 0.84 | 0.8 |

**Declaration of interests**

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

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