

RESEARCH

Open Access



An extreme rainfall-induced landslide susceptibility assessment using autoencoder combined with random forest in Shimane Prefecture, Japan

Kounghoon Nam* and Fawu Wang

Abstract

Background: Landslide-affecting factors are uncorrelated or non-linearly correlated, limiting the predictive performance of traditional machine learning methods for landslide susceptibility assessment. Deep learning methods can take advantage of the high-level representation and reconstruction of information from landslide-affecting factors. In this paper, a novel deep learning-based algorithm that combine classifiers of both deep learning and machine learning is proposed for landslide susceptibility assessment. A stacked autoencoder (StAE) and a sparse autoencoder (SpAE) both consist of an input layer for raw data, hidden layer for feature extraction, and output layer for classification and prediction. As a study case, Oda City and Gotsu City in Shimane Prefecture, southwestern Japan, were used for susceptibility assessment and prediction of landslides triggered by extreme rainfall.

Results: The prediction performance was compared by analyzing real landslide and non-landslide data. The prediction performance of random forest (RF) was evaluated as better than that of a support vector machine (SVM) in traditional machine learning, so RF was combined with both StAE and SpAE. The results show that the prediction ratio of the combined classifiers was 93.2% for StAE combined with RF model and 92.5% for SpAE combined with RF model, which were higher than those of the SVM (87.4%), RF (89.7%), StAE (84.2%), and SpAE (88.2%).

Conclusions: This study provides an example of combined classifiers giving a better predictive ratio than a single classifier. The asymmetric and unsupervised autoencoder combined with RF can exploit optimal non-linear features from landslide-affecting factors successfully, outperforms some conventional machine learning methods, and is promising for landslide susceptibility assessment.

Keywords: Stacked autoencoder, Sparse autoencoder, Support vector machine, Random forest, Landslide susceptibility

Introduction

Landslide susceptibility assessment is a cogent research topic intended to determine the spatial probability of landslide occurrence since landslides continuously result in damages and casualties worldwide (Corominas et al. 2013). Spatial occurrence is called susceptibility, and landslide susceptibility maps generated from landslide-affecting factors using statistical approaches subdivide areas into different terrains that are likely to cause certain types of landslides (Segoni et al. 2018). Physical methods using GIS and

remote sensing are more accurate than statistical approaches (Alexakis et al. 2014; Ciampalini et al. 2015; Di Martire et al. 2016), while physical methods are not suitable for large areas (Tien Bui et al. 2016). Therefore, statistical approaches have received much attention because it is efficient for fast recognizing landslides in large areas (Chen et al. 2018b). It is necessary for decision makers to fast recognize large areas where landslides are expected to result in land use planning and disaster control. Landslide susceptibility prediction based on statistical approaches can achieve this goal efficiently (Borrelli et al. 2018; Huang et al. 2019). Most of the quantitative methods for producing landslide susceptibility maps refer to regression or

* Correspondence: geonamsoil@gmail.com

Department of Earth Science, Shimane University, 1060 Nishikawatsu-cho, Matsue, Shimane 690-8504, Japan

classification approaches between real landslide data and artificially created non-landslide data (Fell et al. 2008). The quantitative methods most widely used for landslide susceptibility mapping are such as logistic regression (Lee and Talib 2005; Ayalew and Yamagishi 2005; Bai et al. 2010; Aditian et al. 2018), naïve Bayes (Tien Bui et al. 2012; Tsangaratos and Ilia 2016), artificial neural networks (Pradhan et al. 2010; Arnone et al. 2016), support vector machines (Yao et al. 2008; Yilmaz 2010; Ballabio and Sterlacchini 2012; Xu et al. 2012), decision trees (Saito et al. 2009; Yeon et al. 2010), and random forest (Alessandro et al. 2015; Trigila et al. 2015; Hong et al. 2016; Chen et al. 2019b; Park et al. 2019) in machine learning techniques.

Recently, deep learning algorithms have made a series of revolutions in the field of machine learning (Huang et al. 2019) since the classification capability of a neural network to fit a decision boundary plane has become significantly more reliable (LeCun et al. 2015) which can successfully learn and extract patterns and unique features from big data (Ayinde et al. 2019). Deep learning also can effectively avoid local optimization and eliminates the need to set model parameters because of autonomous processes (Zhang et al. 2017). At the moment, the core techniques of deep learning are neural networks that have two or more hidden layers, including the following techniques: the adaptive neuro-fuzzy inference system (Park et al. 2012); recurrent neural networks (Chen et al. 2015); deep belief networks (Huang and Xiang 2018); long short-term memory (Xiao et al. 2018; Yang et al. 2019); and convolutional neural networks (Wang et al. 2019). Deep learning-based autoencoder is a semi-supervised learning method with no prior knowledge, such as landslide inventory, which means that landslide and non-landslide labels and linear and non-linear correlation assumptions are not needed (Huang et al. 2019). For landslide susceptibility assessment, traditional methods for de-correlation are based on the prior assumption that there are linear correlations between landslides and non-landslides. However, landslide-affecting factors are usually non-linear in practical applications. The autoencoder driven by data rather than prior knowledge can transform raw data into non-linear correlated features.

In this paper, novel deep learning algorithms, namely, both stacked autoencoder and sparse autoencoder combined with traditional machine learning, are proposed for landslide susceptibility prediction. StAE and SpAE are unsupervised learning as it does not require external labels on landslides information. The encoding and decoding process all happen in the dataset. The input and output data have the same number of dimensions, and the hidden layer has fewer dimensions. Autoencoders are learned automatically from dataset, which is easy to train specialized instances of the algorithm that will perform well on a specific type of landslide-affecting factors. The autoencoder technique takes advantage of dimension reduction by

stacked autoencoder and dropout by sparse autoencoder for non-linear correlations of the landslide-affecting factors and gives better feature descriptions than the original data. It does not require any additional methods which are required for appropriate training data. In summary, this study proposes the combined method of the advantage of deep learning and the benefits of machine learning for landslide susceptibility assessment. The landslides in Oda City and Gotsu City in Shimane Prefecture, southwestern Japan, are used as case study. A stacked autoencoder and sparse autoencoder are combined with random forest acquired from the results of a better predictive performance between support vector machine and random forest.

Study area

The study area is located in Oda City and Gotsu City, Shimane Prefecture, southwestern Japan (Fig. 1). The elevation varies from sea level to 1123 m (Table 1). The average annual rainfall recorded from the rainfall stations at Fukumitsu, Oda, and Sakurae are 1657 mm, 1786 mm, and 2011 mm from 2008 to 2018 (Fig. 2). The cumulative rainfall for 2013 recorded from the rainfall stations at Fukumitsu, Oda, and Sakurae are 2270 mm, 2102 mm, and 2656 mm, respectively (<http://www.jma.go.jp/jma/index.html>). In this study, a total of 90 landslides were caused by extreme rainfall from May to October 2013 (Table 2), and 69 of the landslides were triggered by extreme rainfall in August 2013. These landslides can be described as shallow landslides that were determined based on field investigation.

Spatial data setting

Landslide susceptibility prediction can be evaluated as a binary classification problem between landslides and non-landslides. A spatial database setting including landslide pixel grid, non-landslide pixel grid, and related landslide-affecting factors is needed for statistical analysis (Huang et al. 2019). This spatial database was divided into a training dataset and a validation dataset.

These real 90 of landslides and 90 of non-landslides artificially generated from ArcGIS software were randomly split into two parts with a ratio of 70% and 30%. Seventy percent of the landslide and non-landslide grid cells were selected for the training model, and the remaining 30% were used for the validation model. Furthermore, the landslide (event) and non-landslide (non - event) grid cells were set to 1 and 0, respectively, and the values of 1 and 0 were used for classification and prediction as the output variables of the landslide susceptibility prediction models. Thereafter, the calculated frequency ratio (FR) values were considered as numeric input variables of landslide susceptibility prediction models.

The landslide-affecting factors in study area are complex, and it is difficult to confirm which affecting factors

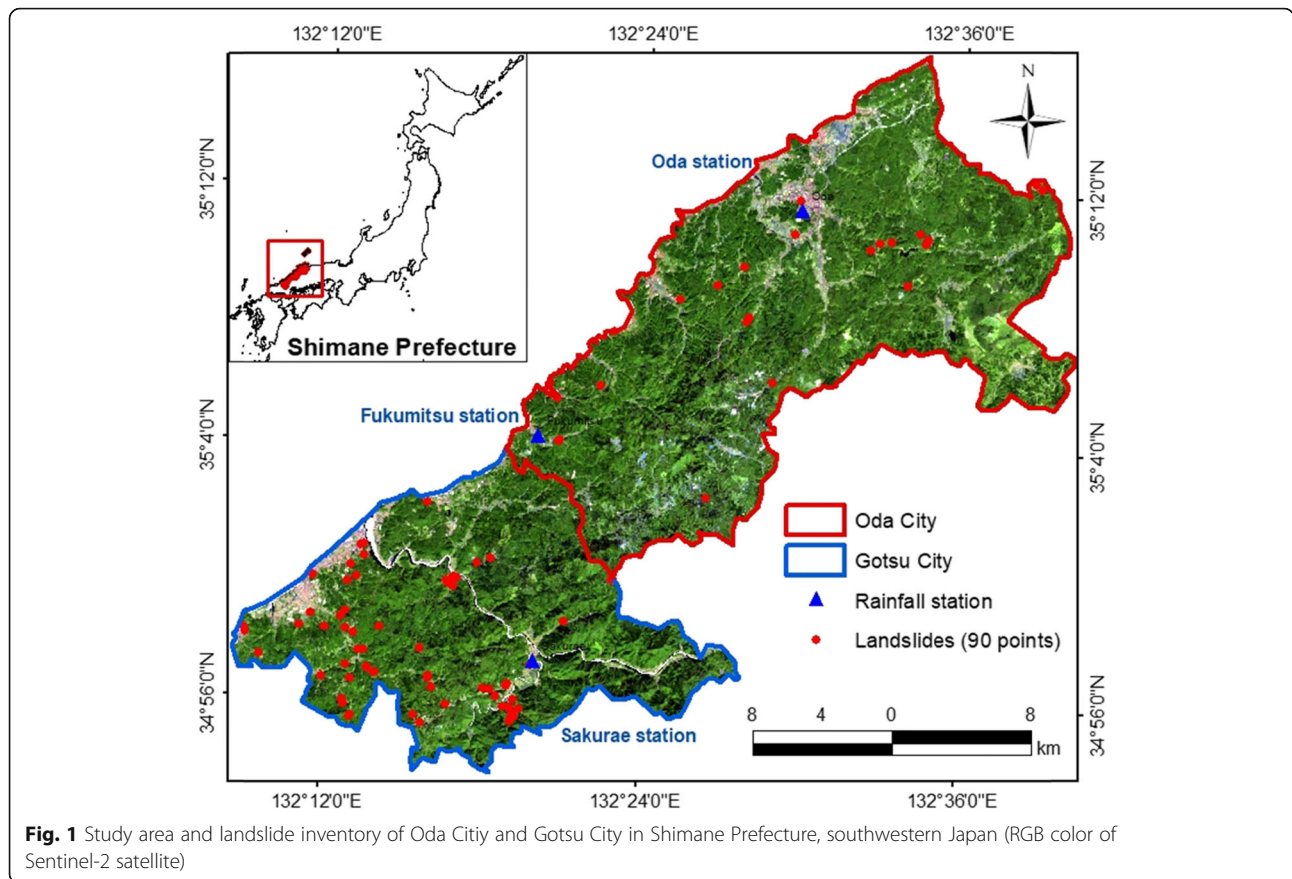


Table 1 Description and frequency ratio (FR) of topographical and distance to factors in the study area

Factors	Values	No. of Landslides	FR	Factors	Values	No. of Landslides	FR
Altitude (m)	0–105.77	59	2.15	Profile curvature	–37.55 - (– 3.71)	3	0.84
	105.78–215.95	21	0.80		–3.72 - (–1.18)	6	0.43
	215.96–339.36	10	0.43		–1.19 - 1.03	39	0.90
	339.37–550.90	0	0		1.04–3.87	36	1.53
	550.91–1123.84	0	0		3.88–43.08	6	1.06
Slope (degree)	0–9.50	2	0.78	Dis. to stream	< 101	58	1.64
	9.51–19.00	30	1.11		101–200	21	0.73
	19.01–28.21	29	1.08		201–300	8	0.44
	28.22–38.00	24	0.92		301–400	3	0.46
	38.01–73.40	5	0.71		> 401	0	0
Plan curvature	–49.05 - (–3.81)	2	0.78	Dis. to road	< 200	39	2.91
	–3.82 - (–1.11)	12	1.11		201–400	17	1.38
	–1.12 - 0.57	47	1.08		401–600	7	0.64
	0.58–2.60	24	0.92		601–800	5	0.52
	2.61–37.03	5	0.71		> 801	22	0.50

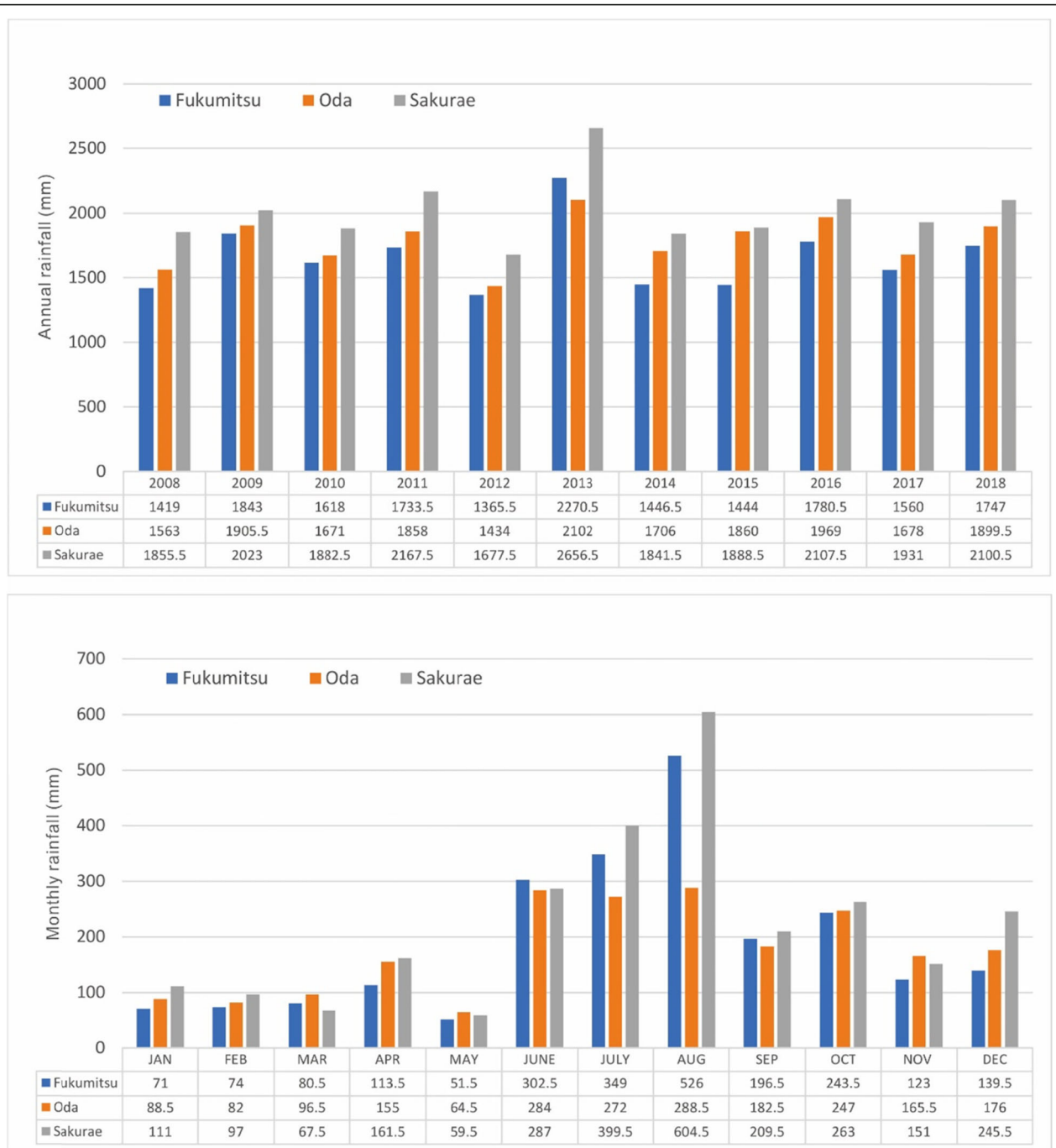


Fig. 2 Annual rainfall from 2008 to 2018, and monthly rainfall of 2013 in the study area

are the most important and necessary among the topographic, geological, hydrological, distance to stream and distance to road. In landslide susceptibility modeling, landslides may reoccur under conditions similar to those of past landslides (Westen et al. 2003; Lee and Talib 2005; Dagdelenler et al. 2016). A total of 14 affecting factors were acquired and chosen as input variables for landslide susceptibility models (Figs. 3, 4 and 5).

The topographic factors were acquired and calculated based on the digital elevation model (DEM), with a spatial resolution of 10 m, including altitude, slope angle, plan curvature, profile curvature (Yilmaz et al. 2012), distance to stream (Devkota et al. 2013; Guo et al. 2015), stream power index (SPI) (Park and Kim 2019), and topographic wetness index (TWI) (Althuwaynee et al. 2016; Colkesen et al. 2016). The distance to road (Alexakis et al. 2014; Roy

Table 2 Description and frequency ratio (FR) of remote sensing - derived index and hydrological factors in the study area

Factors	Values	No. of Landslides	FR	Factors	Values	No. of Landslides	FR
NDVI	-0.242 - 0.143	15	7.26	SPI	-13.816 - (-8.806)	3	0.65
	0.144-0.255	19	3.15		-8.805 - (-4.352)	14	0.83
	0.256-0.332	30	1.40		-4.351 - 0.101	20	0.76
	0.333-0.391	26	0.75		0.102-2.773	46	1.28
	0.392-0.650	0	0		2.774-14.574	7	1.04
NDWI	-0.324 - 0.046	1	0.22	TWI	-7.969 - (-2.240)	17	0.82
	0.047-0.125	22	2.19		-2.239 - 2.114	23	0.75
	0.126-0.179	36	1.75		2.115-4.520	38	1.43
	0.180-0.223	25	0.79		4.521-8.416	12	1.22
	0.224-0.483	6	0.26		8.417-21.250	0	0
BI	0.248-0.312	4	0.86	Rainfall (mm) (from May to Oct. in 2013)	1338.5-1424	13	0.54
	0.313-0.322	10	0.45		1424.1-1528.5	5	0.36
	0.323-0.329	33	0.97		1528.6-1633	2	0.17
	0.330-0.339	34	1.50		1633.1-1718.5	47	1.69
	0.340-0.436	9	1.41		1718.6-1823	23	1.82

and Saha 2019) was acquired from Geospatial Information Authority of Japan (<https://fgd.gsi.go.jp/download/menu.php>). Normalized difference vegetation index (NDVI) (Chen et al. 2019a), normalized difference water index (NDWI) (Luo et al. 2019), and bare soil index (BI) (Huang et al. 2019) were derived from the Landsat TM 8 image data, resampled with a 10 m resolution (Zhu et al. 2018). The geological factors were derived from the 1:200,000 scale geological map, which was obtained from the Geological Survey of Japan, AIST (<https://www.gsj.jp/en/>). These landslide-affecting factors were reflected using the raster format with a spatial resolution of 10 × 10 m, which results in raster format that has the advantages of regular shape, quick subdivision, and high modeling efficiency (Huang et al. 2019).

For continuous affecting factors, the Jenks natural break method was used to divide each continuous affecting factor into five classes. Then the frequency ratio of all subclasses of each landslide affecting factor was calculated as shown in Tables 1, 2 and 3. The frequency ratio allows that all 14 landslide-affecting factors have significant influences on landslide occurrence. Some studies have suggested that the correlations between affecting factors should be eliminated to reduce model noise for the landslide susceptibility assessment (Hong et al. 2017; Lin et al. 2017; Chen et al. 2018a). However, the number of input variables of the deep learning algorithm is generally hundreds or thousands due to their strong feature extraction ability, and 14 input variables will not result in information redundancy. On the other hand, some collinearity phenomena between landslide-affecting factors can be tolerated by the fast-developed machine learning models (Huang et al. 2019). These 14 landslide-affecting factors

provide valuable information for producing landslide susceptibility maps, as quantitative measurement determined by frequency ratio. Therefore, all 14 landslide-affecting factors are utilized as input variables in the model to evaluate their capabilities in performance and feature extraction for the landslide susceptibility assessment.

Methodology

This study was performed using the following main steps (Fig. 6): (1) correlation analysis between landslide inventory and landslide-affecting factors using frequency ratio, (2) landslide susceptibility prediction using SVM and RF models in machine learning, (3) landslide susceptibility prediction using StAE and SpAE employing back propagation neural network in deep learning, (4) evaluation of StAE and SpAE combined with machine learning acquired from a better prediction ratio between SVM and RF, and (5) validation and comparison of predictive performance from the area under the curves and landslide susceptibility maps produced by six models. The landslide samples were created after collecting and preparing the landslide inventory map, the DEM derived factors, and remote sensing and geological factors. The landslide inventory samples were counted and used to randomly generate non-landslide samples. The final data combined the landslides and non-landslides samples with a defined label (1 and 0, respectively) for each sample. Fourteen landslide-affecting factors were prepared from a spatial database. The values of the landslide-affecting factors at each sample location were utilized, and the derived information was prepared using RStudio. The dependent variable was converted with one-hot encoding. The data were then categorized into subsets: for

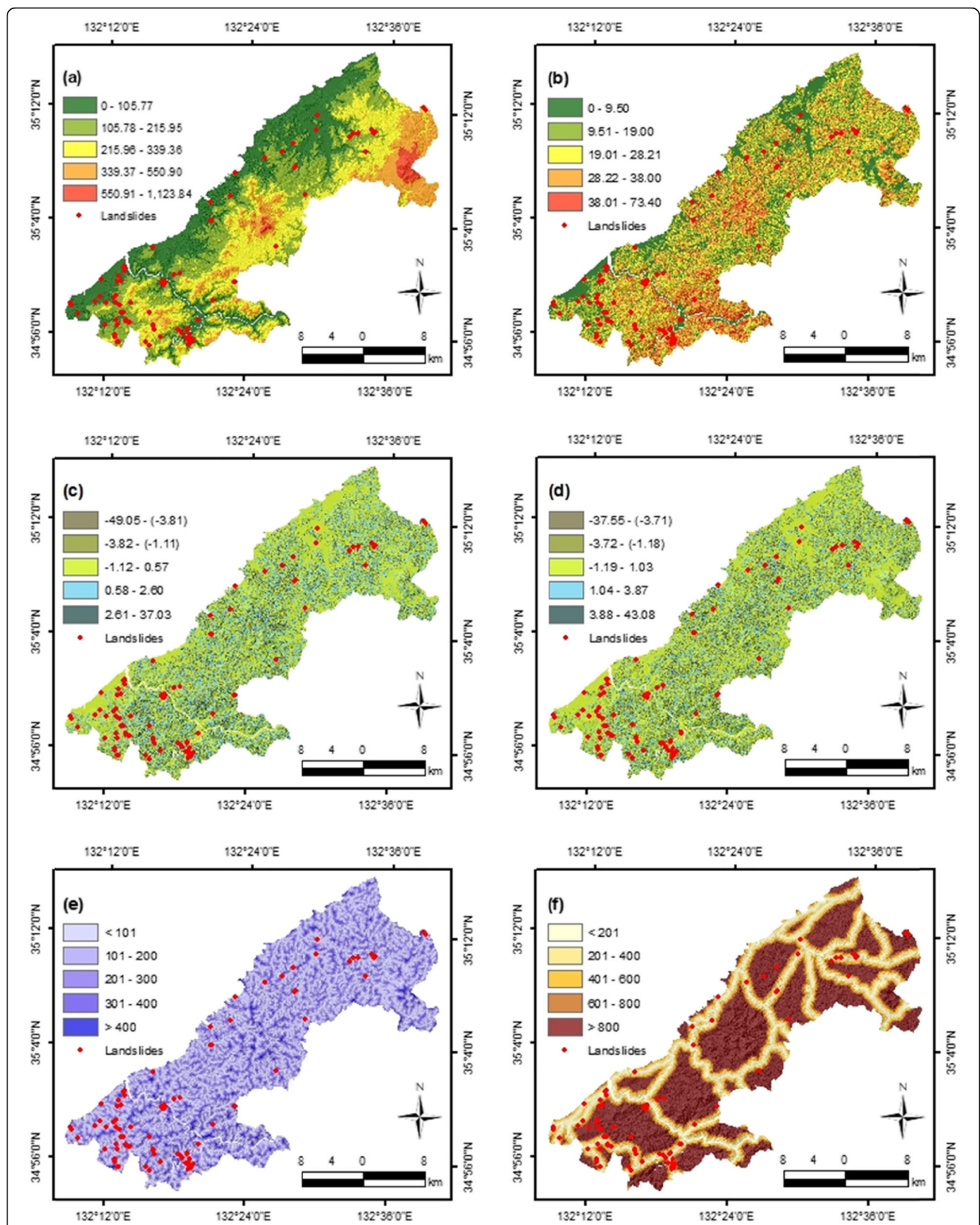


Fig. 3 Thematic maps of topographic factors (a-d) and distance to factors (e and f) considered in this study: **a** elevation (m), **b** slope angle (degree), **c** plan curvature, **d** profile curvature, **e** distance to stream (m), and **f** distance to road (m)

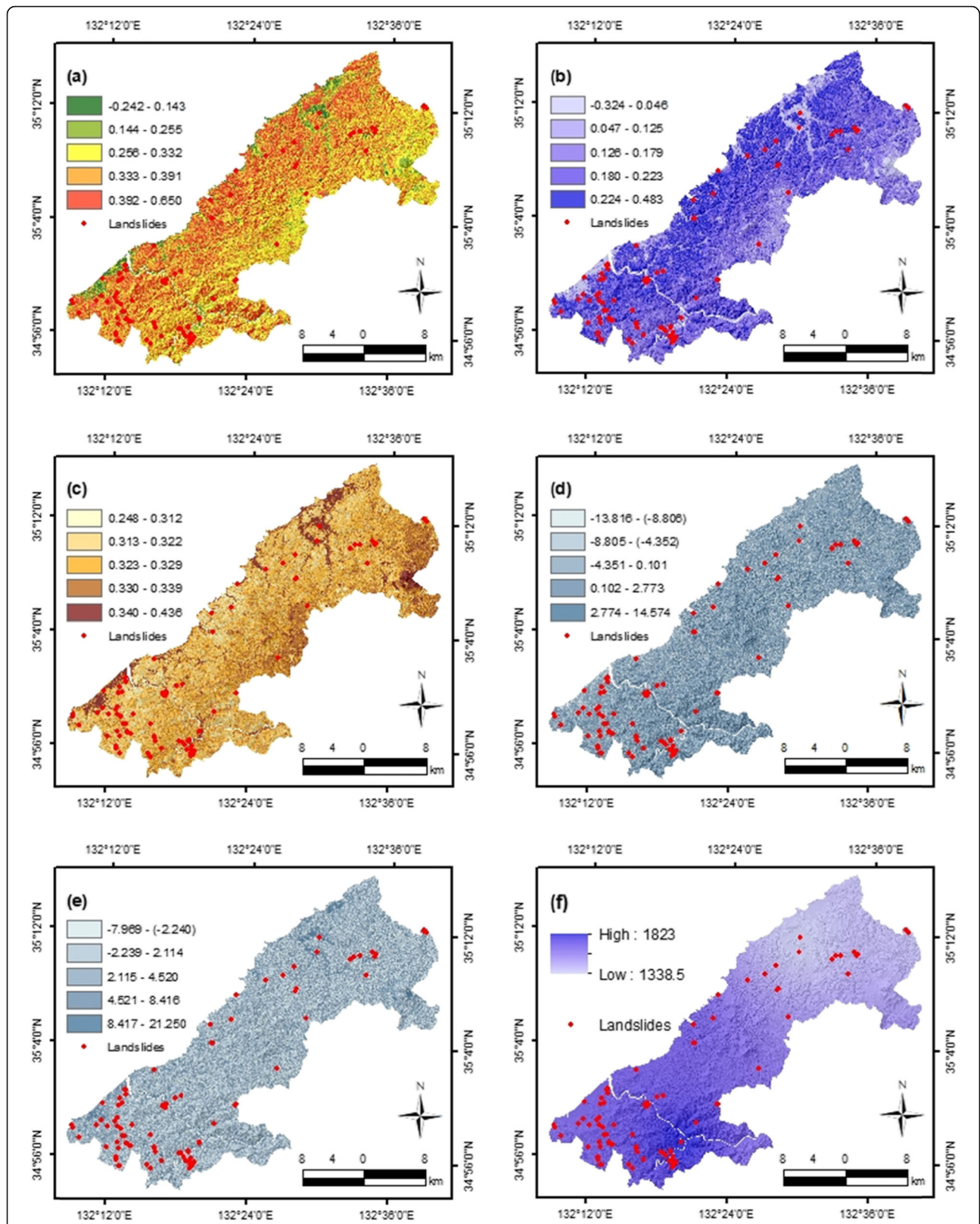


Fig. 4 Thematic maps of remote sensing - derived index (a-c) and hydrological factors (d-f) considered in this study: **a** NDVI, **b** NDWI, **c** BI, **d** SPI, **e** TWI, and **f** cumulative rainfall from March to October in 2013 (mm)

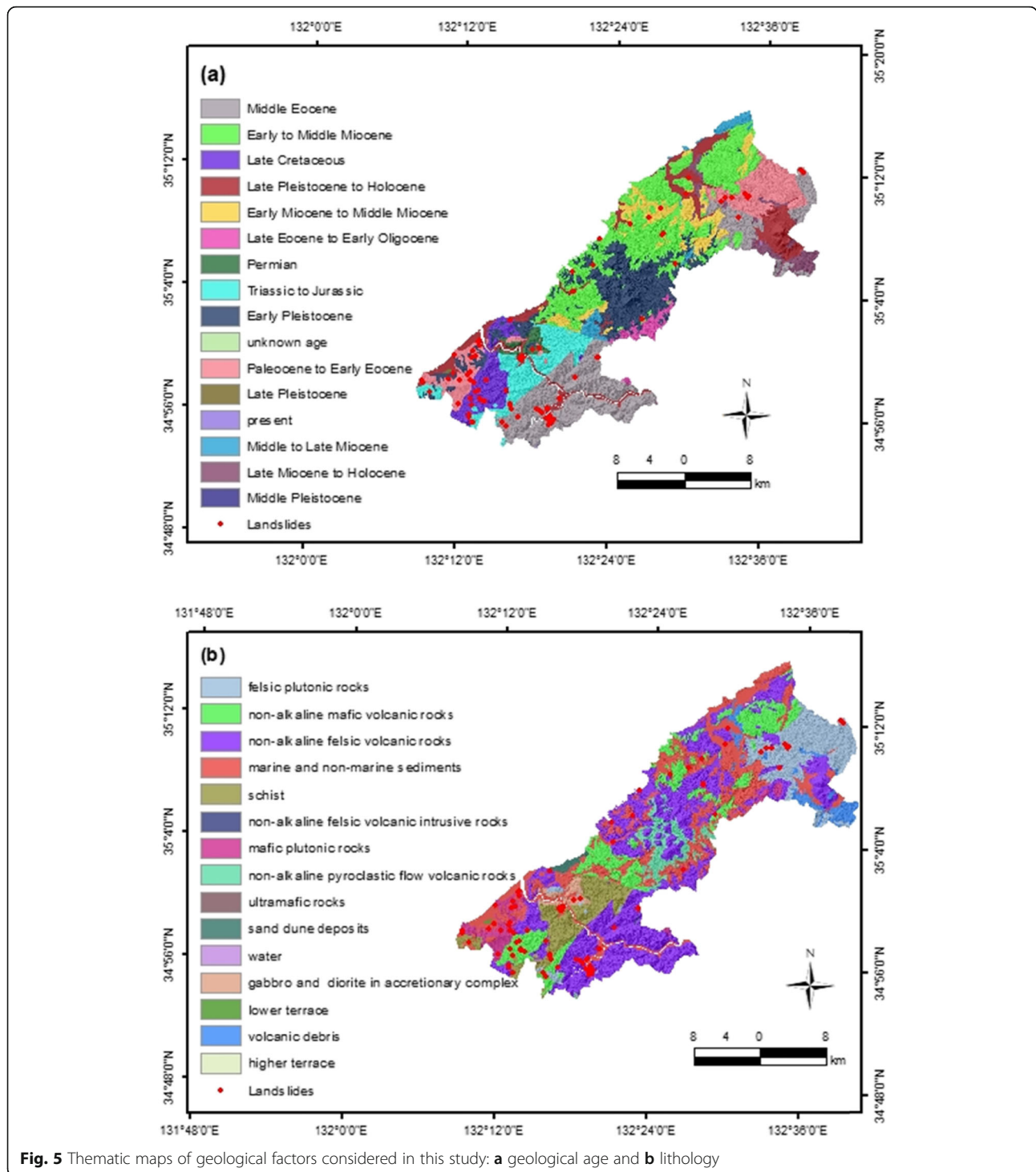


Fig. 5 Thematic maps of geological factors considered in this study: **a** geological age and **b** lithology

training (70%) and validation (30%). The StAE and SpAE model was trained in an unsupervised manner for feature extraction, and a set of new features was generated. These new features were used to train StAE-bpnn and SpAE-bpnn in deep learning, and anomaly detection based StAE with RF and SpAE with RF which is selected as better

prediction rate than SVM model. In this study, the validation of the proposed models was based on a well-known area under the receiver operating characteristic curve. Parameter tuning was also utilized to assess better accuracy. Finally, landslide susceptibility maps were generated using equal interval function in ArcGIS 10.6 software.

Table 3 Description and frequency ratio (FR) of geological factors in the study area

Factors	Values	No. of Ls	FR	Factors	Values	No. of Ls	FR
Geological age	Unknown age	22	1.31	Lithology	Felsic plutonic rocks	9	0.93
	Triassic to Jurassic	8	0.41		Gabbro and diorite in accretionary complex	13	1.12
	Present	14	3.09		Higher terrace	27	0.85
	Permian	13	1.53		Lower terrace	16	0.82
	Paleocene to Early Eocene	2	0.36		Mafic plutonic rocks	10	1.50
	Middle to Late Miocene	2	1.37		Marine and non-marine sediments	13	3.94
	Middle to Late Miocene	10	9.15		Non-alkaline felsic volcanic intrusive rocks	2	0.64
	Middle Eocene	1	0.15		Non-alkaline felsic volcanic rocks	0	0
	Late Pleistocene to Holocene	18	1.41		Non-alkaline mafic volcanic rocks	0	0
	Late Pleistocene	0	0		Non-alkaline pyroclastic flow volcanic rocks	0	0
	Late Miocene to Holocene	0	0		Sand dune deposits	0	0
	Late Eocene to Early Oligocene	0	0		Schist	0	0
	Late Cretaceous	0	0		Ultramafic rocks	0	0
	Early to Middle Miocene	0	0		Volcanic debris	0	0
	Early Pleistocene	0	0		Water	0	0
Early Miocene to Middle Miocene	0	0		0	0		

Frequency ratio (FR)

The number of landslide pixel grids in each class is evaluated, and the frequency ratio for each factor class is assigned by dividing the landslide ratio by the area ratio. The frequency ratio shows the correlation between landslides and affecting factors in a specific area. If this ratio is greater than 1, then the relationship between a landslide and the affecting factor’s class will be strong but if

the ratio is less than 1, then the relationship will be weak. If the value is 1, it means an average correlation (Meten et al. 2015).

Support vector machine (SVM)

Two main principles of SVM are the optimal classification hyperplane and the use of kernel features. The purpose of optimal sorting hyperplanes is to accurately distinguish

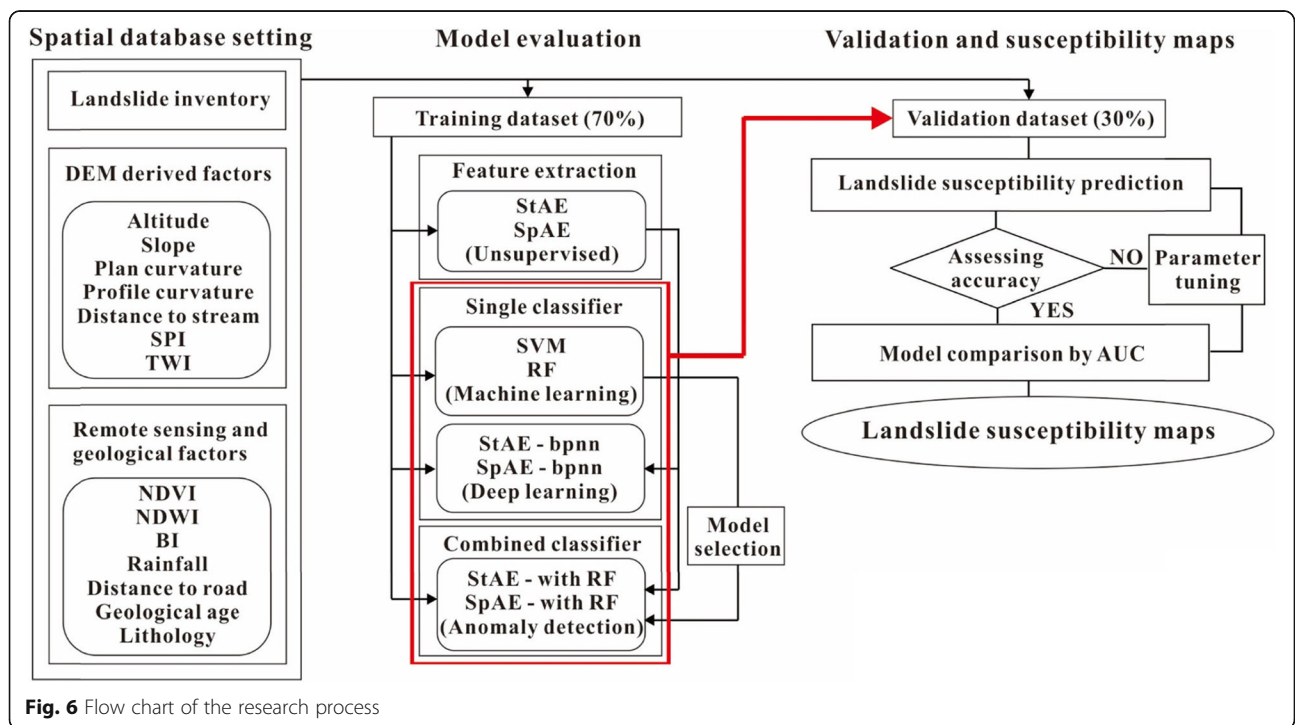


Fig. 6 Flow chart of the research process

the two types of samples between landslides and non-landslides while maximizing the sorting margin. Determining kernel function and optimal parameters are critical for evaluating landslide susceptibility using SVM. Polynomial kernels and radial basis function are the most commonly used kernels in the literature (Huang and Zhao 2018). To optimize two parameters, both penalty coefficient C and kernel function parameters are needed in the SVM model.

Random forest (RF)

The random forest, a classification tree algorithm with repeated dichotomy data, can significantly reduce the computations required for classification and regression. In RF algorithms, predictive models are established by utilizing many decision trees. Based on randomly selected variables and samples, these trees and their decisions are generated. Once the model is established, the samples are first sorted individually according to all decision trees in the model, and then by all trees (Huang and Zhao 2018). The proportion of decision tree estimates and generates landslide susceptibility indexes, which can predict landslide occurrence between all decision trees in the RF model (Goetz et al. 2015).

Stacked autoencoder (StAE)

The StAE is an artificial neural network, which is a special type of multi-layer perceptron. It is a type of unsupervised learning algorithm with an asymmetric structure, in which the middle layer represents the encoding of the input data in the bottleneck layer (Yu and Príncipe 2019). The bottleneck constrains the amount of information that can traverse the full network, forcing the learned compression of the input data. The StAE is trained to reconstruct the input of landslide-affecting factors onto the output layer for feature representation, which prevents the simple copying of the data and the network. The middle layer has a lower dimension to avoid overfitting, which can either select a subset of features with the highest importance or apply some dimension reduction techniques (Hinton and Salakhutdinov 2006; Charte et al. 2018). In this study, the StAE combined with back propagation neural network was processed for a lower dimension of features than the input data have, which can be used for learning the most important features of the data.

Sparse autoencoder (SpAE)

The SpAE consists of an input layer, hidden layers, and an output layer. Each layer in this neural network contains a sufficient number of neurons. Dropout can randomly classify the weight of some implicit layer nodes and reduce the mutual dependence between nodes to realize the normalization of neural networks. Additionally, dropout

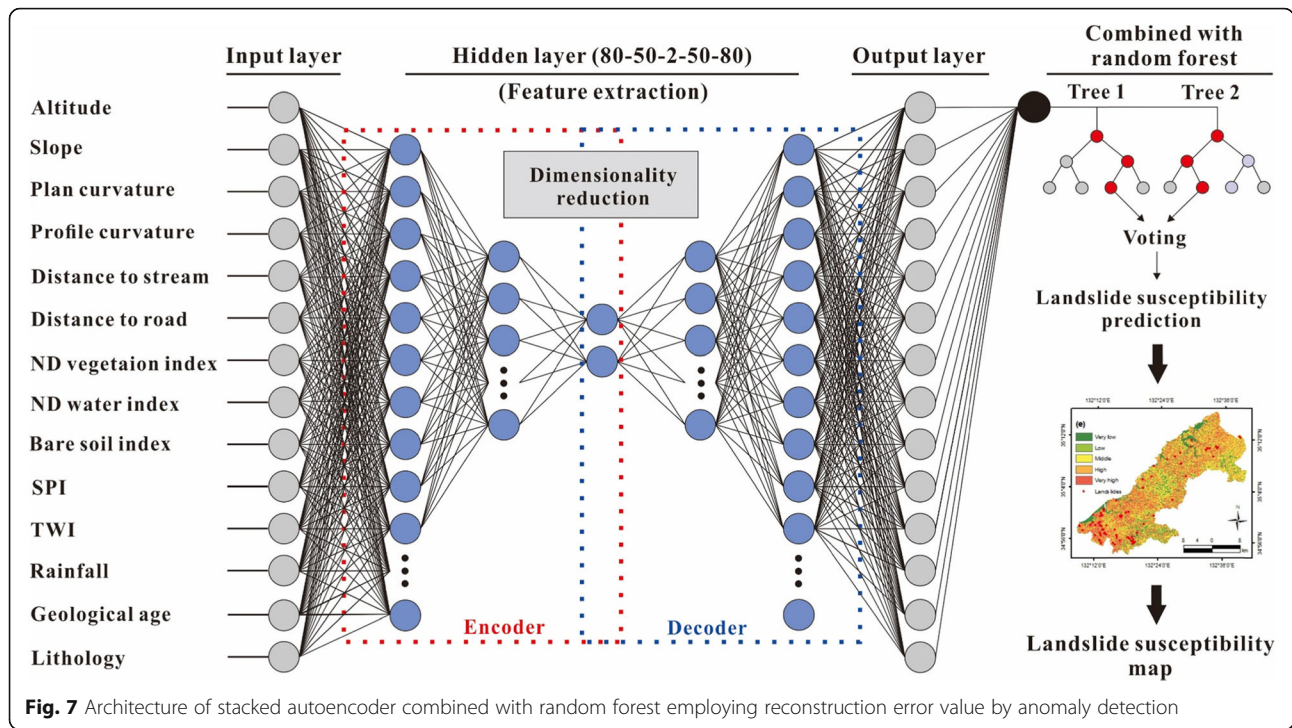
can effectively prevent overfitting and gradient disappearance (Huang et al. 2019). To initially achieve decorrelation among the 14 landslide-affecting factors, dropout was added to the input layer.

The process of StAE is as follows. First, some of the neurons in the network are randomly dropped in the mini-batch training samples and the remaining neurons are fed to the next layer. After obtaining this mini-batch training sample, the deleted neurons are recovered and some neurons in the network are randomly deleted once again. The corresponding parameters are updated based on the stochastic gradient descent method, performed on the neurons that have not been removed.

Results

Landslide susceptibility modelling using the six models

All models based on the deep learning and machine learning were coded in R language on RStudio. For the SVM model and RF model, parameters were determined using a 10-fold cross-validation approach. With radial basis function, SVM model was acquired from grid search for SVM parameter tuning. For RF model, it was composed of 'mtry' and 'tree', which were 3 and 300, respectively. The autoencoder models based on the deep neural network were coded in R language on RStudio using H2O packages. These algorithms were performed using hyperbolic tangent function (i.e., the tanh function) in every hidden layer which was used to encode and decode the input to the output in the undercomplete autoencoder. In the H2O library, five hidden layers with encoders and decoders were designed by using the tanh activation function in each layer. Stacked autoencoders (StAE) were constructed by organizing autoencoder on top of each other also known as deep autoencoder. StAE consists of multiple autoencoder stacked into multiple layers where the output of each layer was wired to the inputs of the successive layers, as seen in Fig. 7, which was composed of 80–50–2–50–80. To obtain good parameters, StAE employed greedy layer-wise training. The benefit of StAE was that it can evaluate the benefits of deep network, which has greater expressive power. Furthermore, it usually can capture useful hierarchical grouping of the input. Finally, model construction was determined by the majority vote among all trees using RF models. The aim of sparse autoencoder (SpAE) was to make a large number of neurons to have low average output so that neurons may be inactive most of the time. The limitation of autoencoders to have only small numbers of hidden units can be overcome by adding a sparsity constraint, where a large number of hidden units can be utilized usually more than one input. Three hidden layers with encoders and decoders were designed by using the tanh activation function in each layer in the H2O library. Sparsity can be achieved by introducing a



loss function during training or manually zeroing few strongest hidden unit activations, which was composed of 200–200–200 (Fig. 8). For classification, model class was constructed by RF model by means of the majority vote among all trees. Reconstruction error value employing mean square error was used by means of anomaly detection in both StAE and SpAE, which were 0.068 and 0.088, respectively.

Landslide susceptibility maps produced by the six models

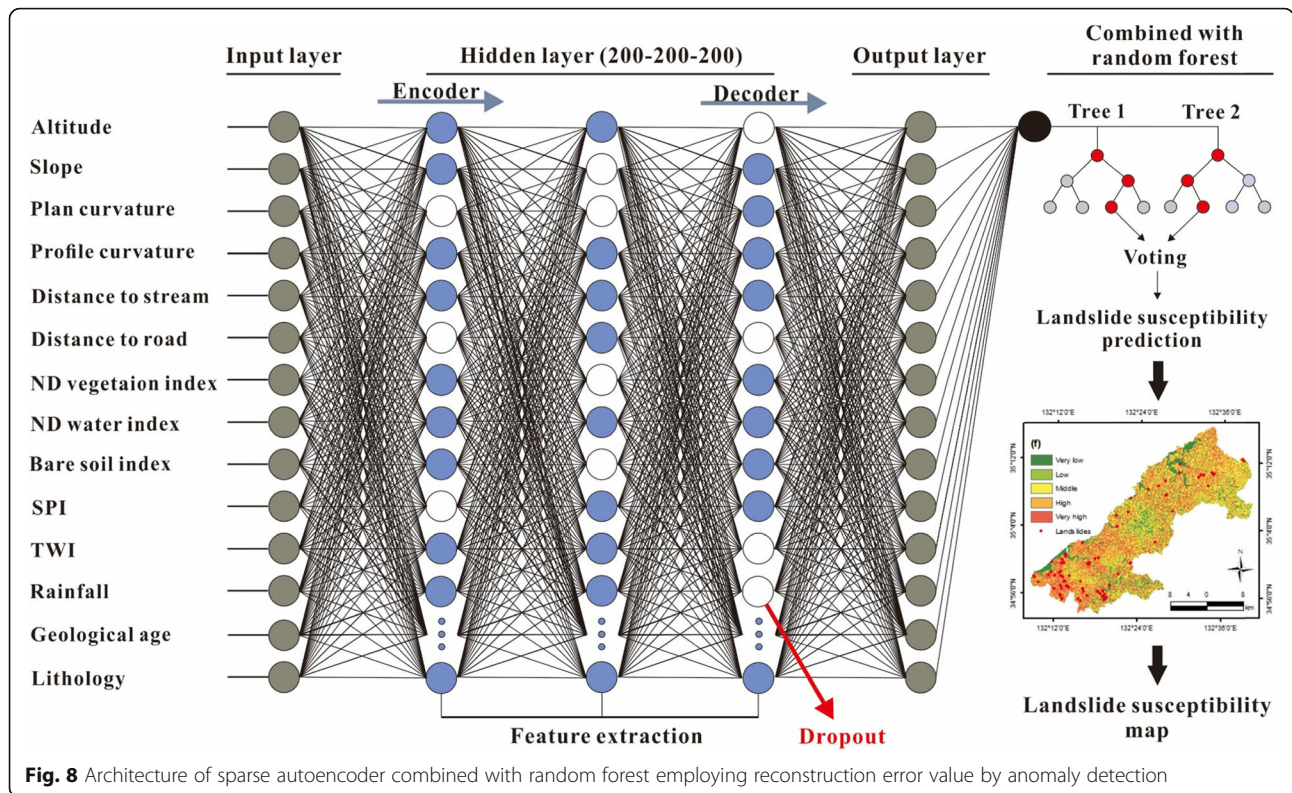
The landslide susceptibility maps were derived from SVM, RF, StAE, SpAE, StAE with RF, and SpAE with RF in the ArcGIS 10.6 software (Fig. 9). For better visualization and comparison, the indices were reclassified into five classes using the equal interval function: very low (0–0.2), low (0.2–0.4), moderate (0.4–0.6), high (0.6–0.8), and very high (0.8–1). The susceptibility class area of the StAE model as the best performance (Table 4) were 6.31%, 13.58%, 33.04%, 36.81%, and 10.26%, respectively. The susceptibility class area of the RF model (Fig. 9b) and StAE model (Fig. 9c) has very high value. The susceptibility index value of the SVM model (Fig. 9a) and StAE models (Fig. 9c) were prominent near the road (Fig. 3f). SpAE and SpAE with RF have lower values of class area percentage for a very high (0.8–1.0) index of the susceptibility map. RF and StAE have lower values of class area percentage for a moderate (0.4–0.6) index of the susceptibility map. StAE with RF and SpAE with RF have lower values of class area percentage for a very

low (0.0–0.2) index of the susceptibility map (Fig. 9d, e, f and Table 4).

Discussion

Validation of prediction performance

The landslide susceptibility assessment was verified using the area under the curve on the validation dataset for six models. The predictive ratio for landslide susceptibility assessment is mainly calculated by confusion matrix. The true positive rate (TPR) is defined as the ratio of true positive to the sum of true positive and false negative, and the false positive rate (FPR) is defined as the ratio of false positive to the sum of false positive and true negative to the number of validation samples (Zhang and Wang 2019). In general, the true positive defines the landslide grid cells that are predictive as landslides, true negative means non-landslide grid cells that are predictive as non-landslides, false-positive reflects non-landslide grid cells that are predictive as landslides, and false negative means landslide grid cells that are predictive as non-landslides (Huang et al. 2019). The area under the curve was applied to assess the prediction performance of landslide susceptibility index values on the validation dataset. The prediction rate values of SVM, RF, StAE, SpAE, StAE with RF and SpAE with RF model are obtained by calculating the area under the prediction rate curves. The StAE with RF and SpAE model of combined classifier have relatively higher prediction rates than using SVM, RF, StAE, and SpAE model of single classifier (Fig. 10). This means that the classifiers



combined with both autoencoder and traditional machine learning are better than using a single classifier. Autoencoder is unsupervised learning as it does not require external labels on landslide information. The encoding and decoding process all happen within the dataset. The input and output data have the same number of dimensions, and the hidden layer has fewer dimensions. Thus, it contains compressed information of the input layer, which is why it acts as a dimension reduction for the original input layer. From the hidden layer, the neural network is able to decode the information to its original dimensions. Autoencoders are learned automatically from data examples, which is a useful property. It means that it is easy to train specialized instances of the algorithm that will perform well on a specific type of input. It does not require any additional methods which are required for appropriate training data.

Sample size

One of the challenges for landslide susceptibility mapping is to suggest the sample size on the number of landslide inventories. Several articles have been reported to address adequate numbers of landslide inventories needed to make acceptable landslide susceptibility mapping where sample size varies from 0 to several thousand in different scales of study areas. The sample size affects

the result of the statistical analysis, as an increase in sample size, the result would be more acceptable. According to Demoulin and Chung (2007), in spite of the limited sample size using ten landslides in about 15 × 15 km scale, Bayesian method in machine learning delivered satisfying prediction rates. Heckmann et al. (2014) state that small samples result in large standard errors and wide confidence intervals for the population parameters. In the case of regression parameters, small samples cause the estimation to be uncertain, and there is a higher risk of coefficients being insignificant when the respective confidence interval includes zero. With respect to replicate sampling and model selection, it is expected that the diversity of models. However, increasing sample sizes causes standard errors and confidence intervals in parameter estimation to decrease. In a significance-based stepwise model selection, very large samples are expected to facilitate the inclusion of more and more variables. Reichenbach et al. (2018) present that some articles did not use any landslide inventory, which are based on the relative importance of the thematic maps as landslide-affecting factors (Adler and Huffman 2007). In this study, all models obtained from 84% to 93% prediction rate using 90 landslides (about 20 km square), which is similar to previous study (Sabokbar et al. 2014) of different study area where 82 landslides were used (about 24 km square).

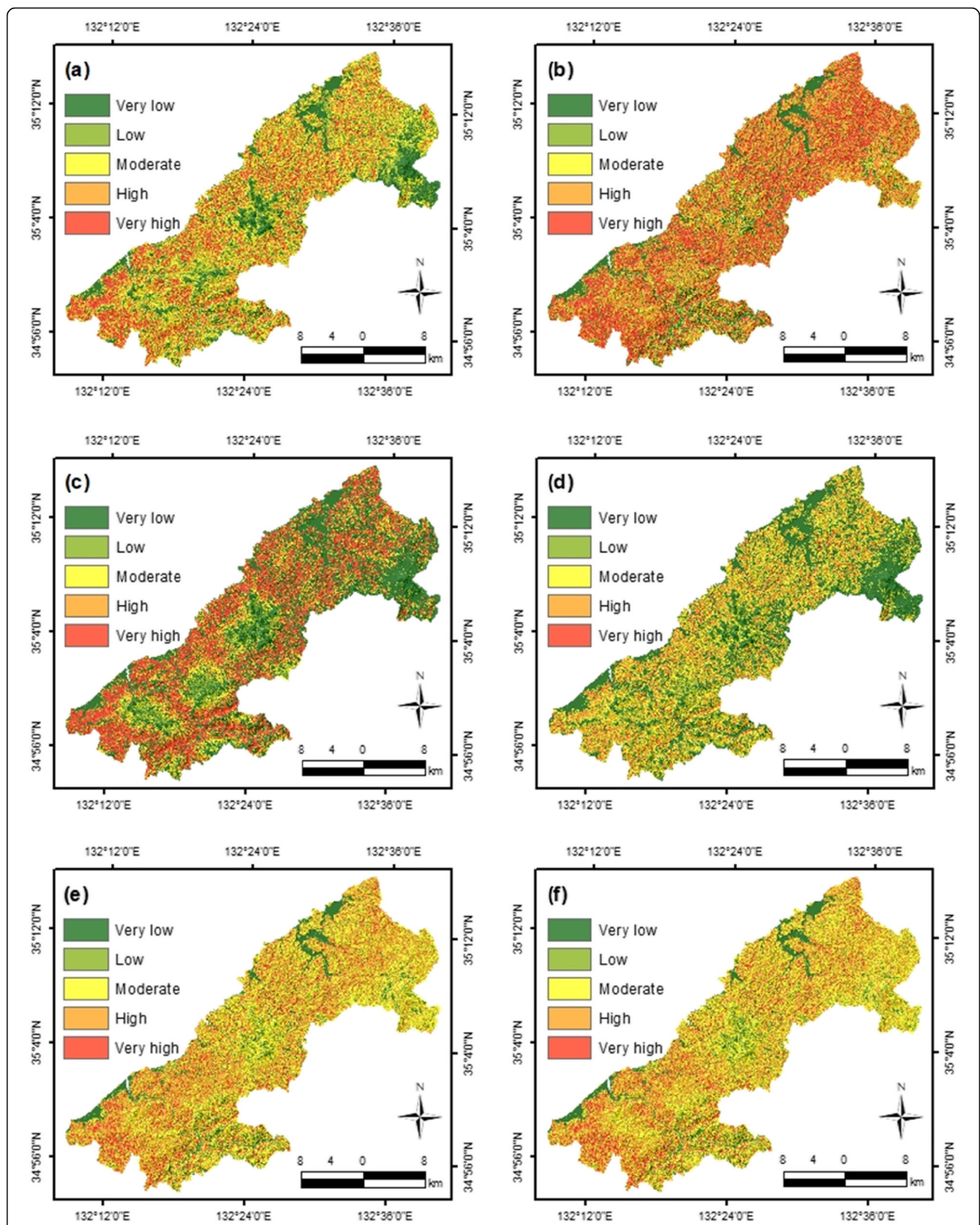


Fig. 9 Landslide susceptibility maps of Shimane prefecture using (a) support vector machine, (b) random forest, (c) stacked autoencoder, (d) sparse autoencoder, (e) stacked autoencoder combined with RF, and (f) sparse autoencoder combined with RF

Table 4 Number of landslides occurred and percentage of landslide susceptibility class area

Landslide susceptibility class(index value)	SVM		RF		StAE		SpAE		StAE with RF		SpAE with RF	
	No.	%	No.	%	No.	%	No.	%	No.	%	No.	%
Very low (0–0.2)	15	17.00	1	11.80	5	35.65	7	30.16	2	6.31	3	5.82
Low (0.2–0.4)	4	19.61	2	10.45	1	11.09	10	19.07	6	13.58	8	14.38
Moderate (0.4–0.6)	7	25.83	4	8.43	3	10.52	15	24.73	16	33.04	13	35.96
High (0.6–0.8)	17	24.94	13	33.43	4	13.07	44	22.60	66	36.81	66	33.25
Very high (0.8–1.0)	47	12.62	70	35.89	77	29.67	14	3.44	0	10.26	0	10.59
Sum	90	100	90	100	90	100	90	100	90	100	90	100

Study limitation

In this study, all landslide points were obtained through GPS by field investigation from May to October in 2013 without the aid of satellite imagery or unnamed aerial vehicle (UAV). As seen in Fig. 2f, most landslide points were in the vicinity of human activity near the roads in the mountains, not inside the mountainous area. The landslide inventory near the roads may affect landslide susceptibility maps (Fig. 9), which results in landslide susceptibility index value near the roads higher than in other areas.

Landslide susceptibility mapping is based on the probability of reoccurrence at the area where landslides already occurred, unlike mapping physically based on modeling, which relies on as follows: 1) the number of abundant landslide inventories for statistical analysis, 2) sampling

strategy to construct non-landslide for regression and classification, 3) scale of study area, 4) resolution of DEM, 5) relatively equal scatter distribution of landslide inventory in study area 6) setting boundary of study area to construct landslide-affecting factors, 7) reasonable selection of landslide-affecting factors. To construct distinct landslide inventory with distinguishing landslide triggering factors between rainfall and earthquake is considered the most important key step than using any advanced classifier for landslide susceptibility mapping.

Conclusion

In this study, the classifiers combined with both deep learning and traditional machine learning, StAE with RF and SpAE with RF models, are proposed for landslide

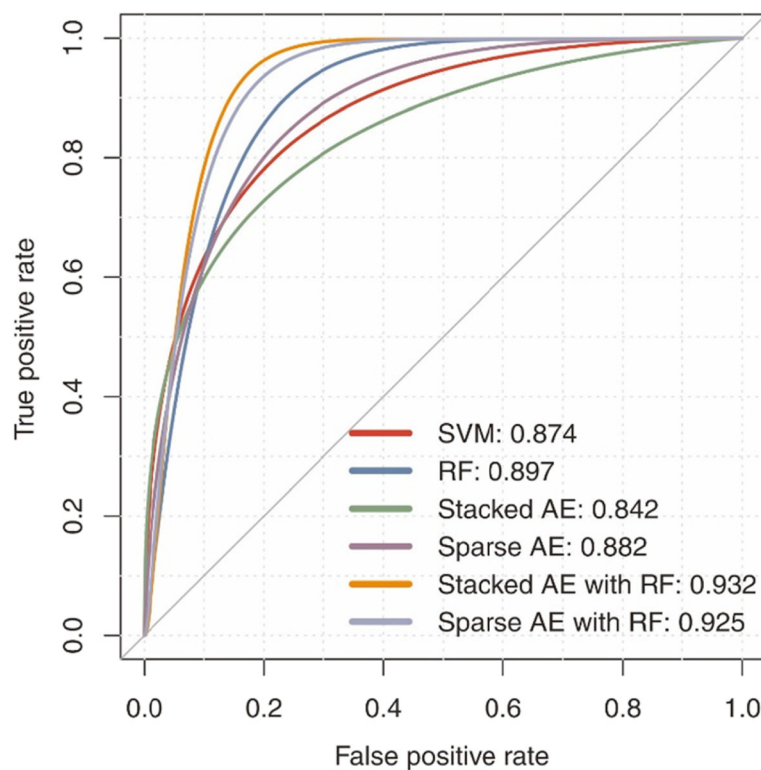


Fig. 10 The area under the curves for prediction ratio and validation of landslide susceptibility maps produced by the six models

susceptibility prediction. The autoencoder consists of input layers for raw data, hidden layers for feature extraction, and output layers for landslide susceptibility prediction. The combined classifiers have the advantage of both machine learning and deep learning, i.e., dimension reduction of the StAE model and dropout of the SpAE model for feature extraction.

The six models were applied in Oda City and Gotsu City, Shimane Prefecture, southwestern Japan. The correlation between landslides and landslide-affecting factors using frequency ratio was high in NDVI, distance to road, and altitude. Performance assessment was carried out with the SVM, RF, StAE, SpAE, StAE with RF, and SpAE with RF models. The results show that the proposed StAE with RF and SpAE with RF models have a relatively better prediction rate than a single classifier such as SVM, RF, StAE and SpAE models. In conclusion, the proposed combined classifier is promising for classification between landslide and non-landslide following landslide susceptibility prediction because it can overcome the limitations of conventional machine learning algorithms, extract features and pattern recognition, reduce computations, and improve performance.

Acknowledgements

The authors express their sincere gratitude to Zili DAI (Shimane University, Matsue, Japan), Prakash DHUNGANA (Shimane University, Matsue, Japan), Ran LI (Shimane University, Matsue, Japan), Rong Zhou (Shimane University, Matsue, Japan), and Akinori IIO (Shimane University, Matsue, Japan), for their kind assistance. The authors also acknowledge financial support from Shimane University, Japan.

Authors' contributions

FW conducted field investigation in 2013 and provided guidance in the study area of landslides triggered by extreme rainfall in Shimane Prefecture, southwestern Japan. KN carried out the landslide susceptibility assessment and produced landslide susceptibility maps using deep learning combined with machine learning. Both authors read and approved the final manuscript.

Funding

The study was financially supported by funding awarded under the study, "Initiation and motion mechanisms of long runout landslides due to rainfall and earthquake in the falling pyroclastic deposit slope area" (JSPS-B-19H01980, Principal Investigator: Fawu Wang).

Availability of data and materials

The DEM data utilized in this study are freely available from the Geospatial Information Authority of Japan (<https://fgd.gsi.go.jp/download/menu.php>). The Landsat 8 satellite image are acquired from the United States Geological Survey (USGS) (<https://earthexplorer.usgs.gov>). All six models for statistical computing and data visualization are coded by R (ver. 3.6.1) open-source software (<https://www.r-project.org/>).

Competing interests

The authors declare that they have no competing interests.

Received: 13 November 2019 Accepted: 21 January 2020

Published online: 30 January 2020

References

Adition 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
- Alessandro T, Carla I, Carlo E, Gabriele SM (2015) Comparison of logistic regression and random forests techniques for shallow landslide susceptibility assessment in Giampilieri (NE Sicily, Italy). *Geomorphology* 249:119–136
- Alexakis D, Agapiou A, Tzouvaras M, Themistocleous K, Neocleous K, Michaelides S, Hadjimitsis D (2014) Integrated use of GIS and remote sensing for monitoring landslides in transportation pavements: the case study of Paphos area in Cyprus. *Nat Hazards* 72:119–141
- Althuwaynee OF, Pradhan B, Lee S (2016) A novel integrated model for assessing landslide susceptibility mapping using CHAID and AHP pair-wise comparison. *Int J Remote Sens* 37(5):1190–1209
- Arnone E, Francipane A, Scarbaci A, Puglisi C, Noto LV (2016) Effect of raster resolution and polygon-conversion algorithm on landslide susceptibility mapping. *Environ Model Softw* 84:467–481
- Ayalew L, Yamagishi H (2005) The application of GIS-based logistic regression for landslide susceptibility mapping in the Kakuda-Yahiko Mountains, Central Japan. *Geomorphology* 65:15–31
- Ayinde BO, Inanc T, Zurada JM (2019) Regularizing deep neural networks by enhancing diversity in feature extraction. *IEEE Trans Neural Netw Learn Syst* 30(9):1–12
- Bai S, Wang J, Lü G, Zhou P, Hou S, Xu S (2010) GIS-based logistic regression for landslide susceptibility mapping of the Zhongxian segment in the three gorges area, China. *Geomorphology* 115:23–31
- Ballabio C, Sterlacchini S (2012) Support vector machines for landslide susceptibility mapping: the Staffora River basin case study, Italy. *Math Geosci* 44(1):47–70
- Borrelli L, Ciarleo M, Gullà G (2018) Shallow landslide susceptibility assessment in granitic rocks using GIS-based statistical methods: the contribution of the weathering grade map. *Landslides* 15(6):1127–1142
- Charte D, Charte F, García S, Jesus MJ, Herrera F (2018) A practical tutorial on autoencoders for nonlinear feature fusion: taxonomy, models, software and guidelines. *Inf Fusion* 44:78–96
- Chen H, Zeng Z, Tang H (2015) Landslide deformation prediction based on recurrent neural network. *Neural Process Lett* 41(2):169–178
- Chen W, Panahi M, Tsangaratos P, Shahabi H, Ilia I, Panahi S, Li S, Jaafari A, Ahmadv BB (2019a) Applying population-based evolutionary algorithms and a neuro-fuzzy system for modeling landslide susceptibility. *Catena* 172:212–231
- Chen W, Peng J, Hong H, Shahabi H, Pradhan B, Liu J, Zhu AX, Pei X, Duan Z (2018a) Landslide susceptibility modelling using GIS-based machine learning techniques for Chongren county, Jiangxi province, China. *Sci Total Environ* 626:1121–1135
- Chen W, Sun Z, Han J (2019b) Landslide susceptibility modeling using integrated ensemble weights of evidence with logistic regression and random forest models. *Appl Sci* 9(1):171
- Chen W, Zhang S, Li R, Shahabi H (2018b) Performance evaluation of the GIS-based data mining techniques of best-first decision tree, random forest, and naïve Bayes tree for landslide susceptibility modeling. *Sci Total Environ* 644:1006–1018
- Ciampalini A, Raspini F, Bianchini S, Frodella W, Bardi F, Lagomarsino D, Traglia F, Moretti S, Proietti C, Pagliara P, Onori R, Corazza A, Duro A, Basile G, Casagli N (2015) Remote sensing as tool for development of landslide databases: the case of the Messina Province (Italy) geodatabase. *Geomorphology* 249:103–118
- Colkesen I, Sahin EK, Kavzoglu T (2016) Susceptibility mapping of shallow landslides using kernel-based Gaussian process, support vector machines and logistic regression. *J Afr Earth Sci* 118:53–64
- Corominas J, Van Westen C, Frattini P, Cascini L, Malet J, Fotopoulou S, Catani F, Van Den Eeckhaut M, Mavrouli O, Agliardi F, Pitiakakis K, Winter M, Pastor M, Ferlisi S, Tofani V, Hervás J, Smith J (2013) Recommendations for the quantitative analysis of landslide risk. *Bull Eng Geol Environ* 73(2):209–263
- Dagdelenler G, Nefeslioglu HA, Gokceoglu C (2016) Modification of seed cell sampling strategy for landslide susceptibility mapping: an application from the eastern part of the Gallipoli peninsula (Canakkale, Turkey). *Bull Eng Geol Environ* 75(2):575–590
- Demoulin A, Chung C (2007) Mapping landslide susceptibility from small datasets: a case study in the pays de Herve (E Belgium). *Geomorphology* 89:391–404
- Devkota KC, Regmi AD, Pourghasemi HR, Yoshida K, Pradhan B, Ryu IC, Dhital MR, Althuwaynee OF (2013) Landslide susceptibility mapping using certainty factor, index of entropy and logistic regression models in GIS and their comparison at Mugling-Narayanghat road section in Nepal Himalaya. *Nat Hazards* 65(1):135–165

- Di Martire D, Tessitore S, Brancato D, Ciminelli MG, Costabile S, Costantini M, Graziano GV, Minati F, Ramondini M, Calcaterra D (2016) Landslide detection integrated system (LaDIS) based on in-situ and satellite SAR interferometry measurements. *Catena* 137:406–421
- Fell R, Corominas J, Bonnard C, Cascini L, Leroi E, Savage W (2008) Guidelines for landslide susceptibility, hazard and risk zoning for land use planning. *Eng Geol* 102(3–4):85–98
- Goetz JN, Brenning A, Petschko H, Leopold P (2015) Evaluating machine learning and statistical prediction techniques for landslide susceptibility modeling. *Comput Geosci* 81:1–11
- Guo C, David RM, Zhang Y, Wang K, Yang Z (2015) Quantitative assessment of landslide susceptibility along the Xianshuihe fault zone, Tibetan plateau, China. *Geomorphology* 248:93–110
- Heckmann T, Gegg K, Gegg A, Becht M (2014) Sample size matters: investigating the effect of sample size on a logistic regression susceptibility model for debris flows. *Nat Hazards Earth Syst Sci* 180:259–278
- Hinton GE, Salakhutdinov RR (2006) Reducing the dimensionality of data with neural networks. *Science* 313:504–507
- Hong H, Ilia I, Tsangaratos P, Chen W, Xu C (2017) A hybrid fuzzy weight of evidence method in landslide susceptibility analysis on the Wuyuan area, China. *Geomorphology* 290:1–16
- Hong H, Pourghasemi HR, Pourtaghi ZS (2016) Landslide susceptibility assessment in Lianhua County (China): a comparison between a random forest data mining technique and bivariate and multivariate statistical models. *Geomorphology* 259:105–118
- Hong Y, Adler RF, Huffman GJ (2007) Use of satellite remote sensing data in the mapping of global landslide susceptibility. *Nat Hazards* 43(2):245–256
- Huang F, Zhang J, Zhou C, Wang Y, Huang J, Zhu L (2019) A deep learning algorithm using a fully connected sparse autoencoder neural network for landslide susceptibility prediction. *Landslides* 17:217–229
- Huang L, Xiang LY (2018) Method for meteorological early warning of precipitation-induced landslides based on deep neural network. *Neural Process Lett* 48(2):1243–1260
- Huang Y, Zhao L (2018) Review on landslide susceptibility mapping using support vector machines. *Catena* 165:520–529
- LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521:436
- Lee S, Talib JA (2005) Probabilistic landslide susceptibility and factor effect analysis. *Environ Geol* 47(7):982–990
- Lin GF, Chang MJ, Huang YC, Ho JY (2017) Assessment of susceptibility to rainfall induced landslides using improved self-organizing linear output map, support vector machine, and logistic regression. *Eng Geol* 224:62–74
- Luo X, Lin F, Zhu S, Yu M, Zhang Z, Meng L, Peng J (2019) Mine landslide susceptibility assessment using IVM, ANN and SVM models considering the contribution of affecting factors. *Public Libr Sci* 14(4):1–18
- Meten M, Prakash B, Yatabe R (2015) Effect of landslide factor combinations on the prediction accuracy of landslide susceptibility maps in the Blue Nile gorge of Central Ethiopia. *Geoenvironmental Disasters* 2:1–17
- Park I, Choi J, Lee M, Lee S (2012) Application of an adaptive neuro-fuzzy inference system to ground subsidence hazard mapping. *Comput Geosci* 48:228–238
- Park S, Hamm S, Kim J (2019) Performance evaluation of the GIS-based data-mining techniques decision tree, random forest, and rotation forest for landslide susceptibility modeling. *Sustainability* 11(20):5659
- Park S, Kim J (2019) Landslide susceptibility mapping based on random forest and boosted regression tree models, and a comparison of their performance. *Appl Sci* 9(5):942
- Pradhan B, Lee S, Buchroithner MF (2010) GIS-based back-propagation neural network model and its cross-application and validation for landslide susceptibility analyses. *Comput Environ Urban Syst* 34:216–235
- Reichenbach P, Rossi M, Malamud BD, Mihir M, Guzzetti F (2018) A review of statistically-based landslide susceptibility models. *Earth Sci Rev* 146:60–91
- Roy J, Saha S (2019) Landslide susceptibility mapping using knowledge driven statistical models in Darjeeling District, West Bengal, India. *Geoenvironmental Disasters* 6:1–18
- Sabokbar HF, Roodposhti MS, Tazik E (2014) Landslide susceptibility mapping using geographically-weighted principal component analysis. *Geomorphology* 226:15–24
- Saito H, Nakayama D, Matsuyama H (2009) Comparison of landslide susceptibility based on a decision-tree model and actual landslide occurrence: the Akaishi Mountains, Japan. *Geomorphology* 109(3):108–121
- Segoni S, Tofani V, Rosi A, Catani F, Casagli N (2018) Combination of rainfall thresholds and susceptibility maps for dynamic landslide hazard assessment at regional scale. *Front Earth Sci* 6(85):1–11
- Tien Bui D, Ho TC, Pradhan B, Pham BT, Nhu VH, Revhaug I (2016) GIS-based modeling of rainfall-induced landslides using data mining-based functional trees classifier with AdaBoost, bagging, and MultiBoost ensemble frameworks. *Environ Earth Sci* 75:1–22
- Tien Bui D, Pradhan B, Lofman O, Revhaug I (2012) Landslide susceptibility assessment in Vietnam using support vector machines, decision tree, and Naïve Bayes models. *Math Probl Eng* 2012:1–26
- Tsangaratos P, Ilia I (2016) Comparison of a logistic regression and naïve Bayes classifier in landslide susceptibility assessments: the influence of models complexity and training dataset size. *Catena* 145:164–179
- Wang Y, Fang Z, Hong H (2019) Comparison of convolutional neural networks for landslide susceptibility mapping in Yanshan County, China. *Sci Total Environ* 666:975–993
- Westen CJV, Rengers N, Soeters R (2003) Use of geomorphological information in indirect landslide susceptibility assessment. *Nat Hazards* 30(3):399–419
- Xiao L, Zhang Y, Peng G (2018) Landslide susceptibility assessment using integrated deep learning algorithm along the China-Nepal highway. *Sensors* 18:1–13
- Xu C, Dai F, Xu X, Lee YH (2012) GIS-based support vector machine modeling of earthquake-triggered landslide susceptibility in the Jianjiang River watershed, China. *Geomorphology* 145–146:70–80
- Yang BB, Yin KL, Lacasse S, Liu ZQ (2019) Time series analysis and long short-term memory neural network to predict landslide displacement. *Landslides* 16(4):677–694
- Yao X, Tham LG, Dai FC (2008) Landslide susceptibility mapping based on support vector machine: a case study on natural slopes of Hong Kong, China. *Geomorphology* 101:572–582
- Yeon Y, Han J, Ryu K (2010) Landslide susceptibility mapping in Injae, Korea, using a decision tree. *Eng Geol* 116:274–283
- Yilmaz C, Topal T, Süzen ML (2012) GIS-based landslide susceptibility mapping using bivariate statistical analysis in Devrek (Zonguldak-Turkey). *Environ Earth Sci* 65(7):2161–2178
- Yilmaz I (2010) Comparison of landslide susceptibility mapping methodologies for Koyulhisar, Turkey: conditional probability, logistic regression, artificial neural networks, and support vector machine. *Environ Earth Sci* 61:821–836
- Yu S, Principe JC (2019) Understanding autoencoders with information theoretic concepts. *Neural Netw* 117:104–123
- Zhang R, Isola P, Efros AA (2017) Split-brain autoencoders: unsupervised learning by cross-channel prediction. Paper presented at the Computer Vision & Pattern Recognition
- Zhang S, Wang FW (2019) Three-dimensional seismic slope stability assessment with the application of Scoops3D and GIS: a case study in Atsuma, Hokkaido. *Geoenvironmental Disasters* 6:1–14
- Zhu X, Miao Y, Yang L, Bai S, Liu J, Hong H (2018) Comparison of the presence-only method and presence-absence method in landslide susceptibility mapping. *Catena* 171:222–233

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen® journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)