# Using Statistical Learning Algorithms in Regional Landslide Susceptibility Zonation with Limited Landslide Field Data

WANG Yi-ting[1,2]  http://orcid.org/0000-0002-2498-3123;  e-mail: wangyiting01@gmail.com

SEIJMONSBERGEN Arie Christoffel[3]  http://orcid.org/0000-0002-7454-7637;
e-mail: A.C.Seijmonsbergen@uva.nl

BOUTEN Willem[3]  http://orcid.org/0000-0002-5250-8872; e-mail: W.Bouten@uva.nl

CHEN Qing-tao[4]  http://orcid.org/0000-0003-3628-1441; e-mail: cqt@cdut.edu.cn

1 State Key Laboratory of Remote Sensing Science, School of Geography, Beijing Normal University, Beijing 100875, China

2 National Marine Data & Information Service, Tianjin 300171, China

3 Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, the Netherlands

4 Institute of Remote Sensing and GIS, Chengdu University of Technology, Chengdu 610059, China

**Abstract:** Regional Landslide Susceptibility Zonation (LSZ) is always challenged by the available amount of field data, especially in southwestern China where large mountainous areas and limited field information coincide. Statistical learning algorithms are believed to be superior to traditional statistical algorithms for their data adaptability. The aim of the paper is to evaluate how statistical learning algorithms perform on regional LSZ with limited field data. The focus is on three statistical learning algorithms, Logistic Regression (LR), Artificial Neural Networks (ANN) and Support Vector Machine (SVM). Hanzhong city, a landslide prone area in southwestern China is taken as a study case. Nine environmental factors are selected as inputs. The accuracies of the resulting LSZ maps are evaluated through landslide density analysis (LDA), receiver operating characteristic (ROC) curves and *Kappa* index statistics. The dependence of the algorithm on the size of field samples is examined by varying the sizes of the training set. The SVM has proven to be the most accurate and the most stable algorithm at small training set sizes and on all known landslide sizes. The accuracy of SVM shows a steadily increasing trend and reaches a high level at a small size of the training set, while accuracies of LR and ANN algorithms show distinct fluctuations. The geomorphological interpretations confirm the strength of SVM on all landslide sizes. Our results show that the strengths of SVM in generalization capability and model robustness make it an appropriate and efficient tool for regional LSZ with limited landslide field samples.

**Keywords:** Landslide Susceptibility Zonation (LSZ); Logistic Regression (LR); Artificial Neural Network (ANN); Support Vector Machine (SVM); Regional scale; Southwest China

## Introduction

Landslide Susceptibility Zonation (LSZ) is an inventory method that considers relevant environmental attributes and identifies landslide-susceptible areas (Fell et al. 2008). The result is a zonation map that visualizes the terrain's susceptibility to potential landsliding. LSZ has

been widely used as a tool for landslide investigation and mitigation at different scales: local (Ercanoglu et al. 2004; Bai et al. 2007); provincial (Yao et al. 2008; Wang et al. 2009); and country-wide (Bălteanu et al. 2010).

Methods for LSZ can be heuristic, deterministic or statistical. For a more detailed review of these algorithms, please refer to Aleotti and Chowdhurry (1999), Dai et al. (2002) and Wang et al. (2005). Among these, statistical algorithms have been the most widely used in regional studies because of their objectivity, reproducibility and rapid assessment. They identify landslide-susceptible areas from the relationship of causative factors and occurrence of landslide events established on records of existing landslide events in the area (Carrara et al. 1991; Bălteanu et al. 2010). A series of statistical algorithms have been proposed to establish such relationship, including bivariate (Brabb et al. 1972; DeGraff and Romesburg 1980; Jade and Sarkar 1993; Yilmaz and Keskin 2009) and multivariate analyses (Carrara et al. 1991; Nefesliglu et al. 2008; Yilmaz 2010 a, b; Pradhan 2010; Bui et al. 2012; Holec et al. 2013). In bivariate analysis landslide maps are compared with each parametric map (e.g. lithology, slope angle and aspect) separately, while multivariate analysis integrates all relevant parametric maps to create a single susceptibility zonation map.

Based on the analysis of past landslide events, a reliable landslide inventory presenting locations and outlines of landslides is mandatory for regional LSZ analysis. For preparing a landslide inventory, researchers use different sampling strategies that record landslides such as point, scarp and seed cell (Yilmaz 2010b). Sampling strategy used in LSZ analysis should reflect the environmental conditions prior to landsliding. Yilmaz (2010b) has compared the effects of sampling strategies on LSZ analyses and recommended that scarp strategy that distinguishes main scarps from the accumulation/depletion zone as an optimal choice. In reality, the availability of existing landslide records is always a major concern and often determines the sampling strategy in regional LSZ practice.

However, records of existing landslides are generally insufficient in terms of length, details and spatial coverage due to rugged inaccessible terrain and the lack of instrumentation. Especially in southwestern China, where landslides are particularly prevalent south of the Qinling fold system and east of the Tibetan Plateau. These areas are lacking quantitative predictions of landslide occurrences using detailed LSZ studies. The rugged, inaccessible terrain hampers the collection of detailed, field-based landslide data. This situation demands an LSZ to be developed, tested and evaluated in order to overcome the scarcity of reliable landslide information in this area.

As field collection of landslide records is time and labor consuming, a possible solution is to introduce automated analysis techniques. Traditional statistical algorithms, such as multi-linear regression, are well-known for their weaknesses in modeling complex nonlinear relationships and dependence on the size and the distribution of the training samples. In LSZ, a sample records the location, occurrence and ambient environmental attributes of landslide events. Statistical learning algorithms have been introduced to LSZ studies with enhanced generalization capability and data adaptability, such as: logistic regression (LR) (Yilmaz 2009; Bai et al. 2010; Pradhan 2010; Yilmaz 2010a; Yalcin et al. 2011); artificial neural network (ANN) (Lee et al. 2003; Ermini et al. 2005; Chauhan et al. 2010; Pradhan et al. 2010); support vector machine (SVM) (Yao et al. 2008; Yilmaz 2009; Samui and Kothari 2010; Yilmaz 2010a; Chong et al. 2012).

In spite of statistical learning algorithms being used by a number of regional LSZ practices, few contributions addressed how they could work with limited field-based landslide data. For example, LSZ analysis are mostly performed in small study areas with a wealth of landslide records, such as the research conducted by Yao et al. (2008), Rossi et al. (2010) and Pradhan (2013).

In contrast to the demand for advanced analysis techniques in regional LSZ, little emphasis has been put on investigating potentials of different statistical learning algorithms by comparing the performance of algorithms over large areas. Comparative studies are commonly limited to small areas with abundant landslide records, which raises the question of whether the findings are valid for large areas with few landslide records. For example, Yesilnacara and Topal (2005) have compared LR and ANN algorithms in an area of

290 km². Kanungo et al. (2006) and Gupta et al. (2008) both have presented insights into conventional, ANN, fuzzy set and neuro-fuzzy weighting procedures, in 254 km². Yilmaz (2009, 2010a) has compared frequency ratio, LR, ANN and SVM algorithms on areas of 25 km² and 131.6 km² respectively. Marjanović et al. (2011) have investigated the performances of SVM, LR and decision tree algorithms with varying sizes of the training data set of 100 km². Yalcin et al. (2011) have compared frequency ratio, analytical hierarchy process, bivariate statistics and LR algorithms on an area of 4660 km² using the largest extent of 4660 km². Therefore, it remains unclear, how statistical learning algorithms work in regional LSZ with limited landslide samples.

Hanzhong City, a landslide-prone area in southwestern China, occupying an area of 27,246 km², is investigated for regional LSZ mapping. Three statistical learning algorithms, LR, ANN and SVM, are selected for comparison of their performance in regional LSZ analysis.

ANN has been regarded as capable of approximating any given nonlinear function to any degree of accuracy as one can adapt the number of neurons or hidden layers to the complexity of the target problem (Haykin 1998). This background makes the ANN algorithm among the most effective and popular methods in LSZ studies. However, a number of difficulties with the ANN algorithm have been reported, such as: i) the ANN training (model building) process is operator-dependent and cannot provide objective and steady output (Ermini et al. 2005; Kanungo et al. 2006; Yao et al. 2008); ii) ANN requires a large amount of samples for effective training; iii) the final solution, or ANN weights, that derives the best results is not unique, as many networks with different sets of parameters can derive very similar results (Haykin 1998; Balabin and Lomakina 2011).

The LR algorithm is claimed to overcome the operator-dependence problem of ANN, as it can derive an objective and steady output from a least squares algorithm (Yao et al. 2008). The least squares algorithm requires a large dataset and uniform data distribution to achieve good results. Similarly, ANN also requires abundant samples, as based on the principle of Empirical Risk Minimization (ERM), which tries to minimize the difference between model output values and true values. As a result, both algorithms might suffer from weak generalization capability with limited training samples. This suggests that the algorithms can work well on known samples but may fail in real data, as existing samples are always finite and cannot cover infinite cases in the real world. This confirms the need to investigate the performance of the widely used ANN and LR algorithms with limited samples in real regional LSZ practice.

The SVM algorithm has the potential to be a suitable candidate for LSZ on limited samples. The SVM algorithm has many advantages (Vapnik 1995), those we are mostly interested in include: i) that it overcomes the operator-dependence problem of ANN, as established on a solid foundation of mathematic theories; ii) that it tries to minimize the structural risk, instead of the empirical risk. The structural risk is composed of empirical risk and confidence risk, and the latter is a function of sample amount and nonlinearity (quantified as Vapnik-Chervonenkis dimensions). In this way, SVM achieves strong generalization ability. Training on a small set of samples can still work well on independent samples.

The aim of this study is to assess the performances of these three algorithms in regional LSZ analysis with limited landslide field samples. Three LSZ maps are derived from the LR, ANN and SVM algorithms respectively. Landslide density analysis (LDA), receiver operating characteristic curves (ROC) and Kappa index are applied to evaluate the performances. Furthermore, the dependence of the algorithms on sample sizes are investigated by varying the sizes of training sets to analyze how they work compared to known landslides in Hanzhong City. The results will provide an improved understanding of algorithm suitability for regional LSZ, which in turn will provide reference for future detailed landslide studies.

# 1 Methodology

## 1.1 Background information

The basic assumption in the use of statistical algorithms for LSZ is that a number of environmental factors must transpire to bring about the occurrence of a landslide (Fell et al.

2008). The relationship between landslide occurrence probability and the causative factors can be expressed as:

$$P = f(x_1, x_2, \ldots, x_n) \qquad (1)$$

where $x_1$, $x_2$, ..., $x_n$ are the causative factors; $P$ represents the probability of landslide occurrence, thus indicating the degree of landslide susceptibility (in this paper used as the landslide susceptibility index (LSI) value). Therefore using statistical algorithms for LSZ, our task can be formulated as: given $k$ landslide instances ($p_1$, $p_2$, ... , $p_k$) and their associated environmental parameters ($x_{1,1}, x_{2,1}, \ldots, x_{n,1}$), ($x_{1,2}, x_{2,2}, \ldots, x_{n,2}$), ... , ($x_{1,k}, x_{2,k}, \ldots, x_{n,k}$), we train the ANN/SVM/LR models and decide the function $f(x)$ that approximates the actual condition best from these known samples. Then the determined function $f(x)$ is ready for prediction of landslide susceptibility on the basis of environmental information.

The overview work flow is presented in Figure 1 and comprises five steps: (1) preparation of input data; (2) selection and parameterization of environmental factors; (3) training of the three machine learning algorithms; (4) calculation and classification of LSI values from the trained algorithms; (5) model evaluations.

## 1.2 Statistical learning algorithms

LR is a probabilistic model that is suitable for dealing with dichotomous dependent variables (George and Mallery 2000). Given the probability of landslide occurrence $P$ and its absence (1-$P$), the LR model between $P$ and independent variables is established through the logistic transformation of P, as:

$$Logit(P) = \ln(\frac{P}{1-P}) = b_0 + b_1 \cdot x_1 + b_2 \cdot x_2 + \ldots + b_n x_n \qquad (2)$$

where $b_0 \sim b_n$ are partial LR coefficients to be determined and $x_1 \sim x_n$ are the causative environmental parameters.

An ANN is defined as a collection of basic units (neurons), interconnected with one another (Chauhan et al. 2010). It 'learns' through the training process, in which weights between neurons are adjusted repetitively in response to the errors between target and predicted output values until the targeted minimal error is achieved. For technical details of the ANN is referred to as Bishop (1995). The most popular ANN for LSZ is the 'Multi-Layer Perceptron' with a 'Back-error Propagation' learning algorithm (MLP-BP) (Pradhan et al. 2010), which is also applied here.

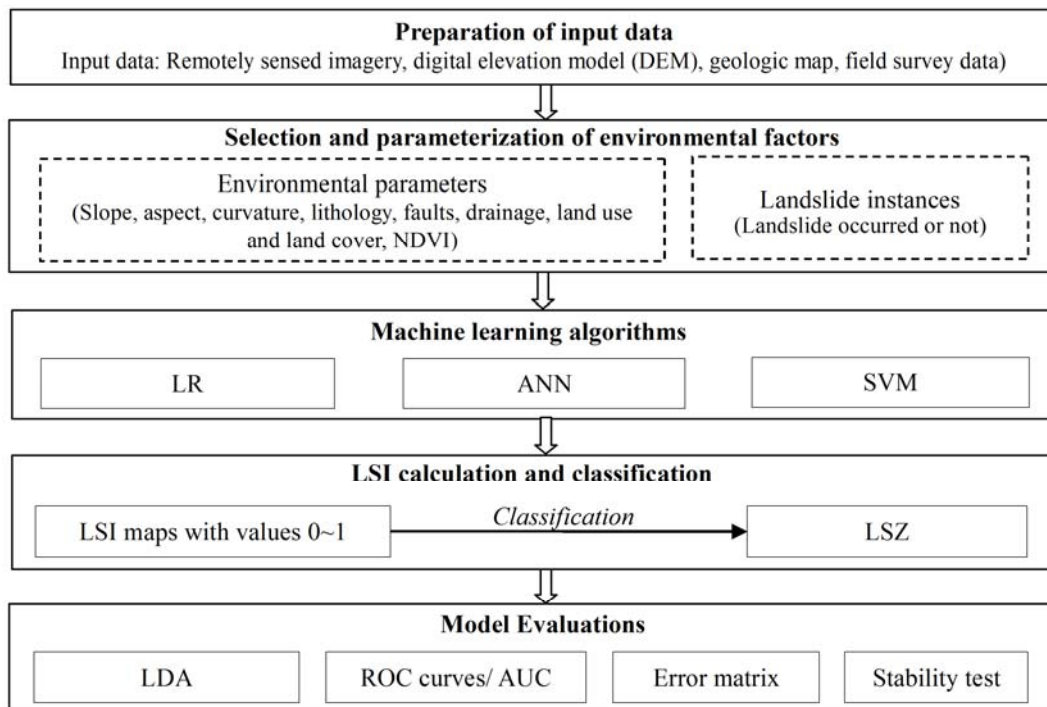SVMs were initially developed for



**Figure 1** Overview workflow for comparing model performance in regional LSZ (landslide susceptibility zonation) assessment.

classification purpose and were then extended to regression tasks (Vapnik 1995). Two main concepts form the basis for the SVM algorithm: i) an optimal hyper plane that separates two classes in the feature space; ii) a kernel function that transfers a non-linear inseparable data pattern into a format that is linearly separable in a high-dimensional feature space (Vapnik 1995). The technical details of SVM are described by Vapnik (1995), Yao et al. (2008) and Dai et al. (2012).

Most studies apply SVM as a one-class or two-class classifier to landslide susceptibility mapping (Yao et al. 2008; Marjanović et al. 2011), in which a sign function is applied on a target function to label output to a corresponding class (usually landslide or non-landslide area). This study uses a support vector regression (SVR) algorithm, as the focus is to derive continuous landslide susceptibility index values to compare with LR and ANN algorithms on the same basis. The only difference is that the output of target function is derived directly without using a sign function. The target function $f(x)$ is expressed as:

$$f(x) = (w \cdot \phi(x) + b) = \sum_{i=1}^{L} (a - a_i^*)K(x, x_i) + b \quad (3)$$

where $w$ is the vector of weights; $x$ is the input variable; and $b$ is a scalar; $(\cdot)$ denotes the scalar product operation; $\phi(x)$ denotes the nonlinear mapping of $x$; $L$ is the number of support vectors; $a_i$, $a_i^*$ and $b$ are parameters to determine the optimal hyper plane; $K(x, x_i)$ is the kernel function.

The LSI values are calculated from the output prediction values of $f(x)$.

## 2 Study Area and Data

### 2.1 Description of the study area

Hanzhong is a prefectural-level city located in the southwest of Shaanxi Province, China (Figure 2), most threatened city by geo-hazards in the province. The official geo-hazard report (DLR(Shaanxi) 2003) reveals that geo-hazards had caused a direct economic loss of nearly 100 million U.S. Dollars and 403 casualties between 1981 and 2001. Hanzhong lies within an east-west fault controlled subsiding basin, which is delimited by the Qinling fold system to the north and the Yangtze table land on the south. The complex geological structure in combination with prolonged tectonic movements have created a series of faults, which determine the present topography, the location of rivers and the activity of landslides. The topographical units distinguished in the area by Liu (1997) can be subdivided into three areas: the Hanzhong Basin and Hills (elevations below 600 m); the Low to Medium High Mountains (elevations between 600 m and 1000 m); and the Medium to High Mountains (elevations above 1000 m). The Hanzhong City area is underlain by sedimentary, metamorphic and igneous rock
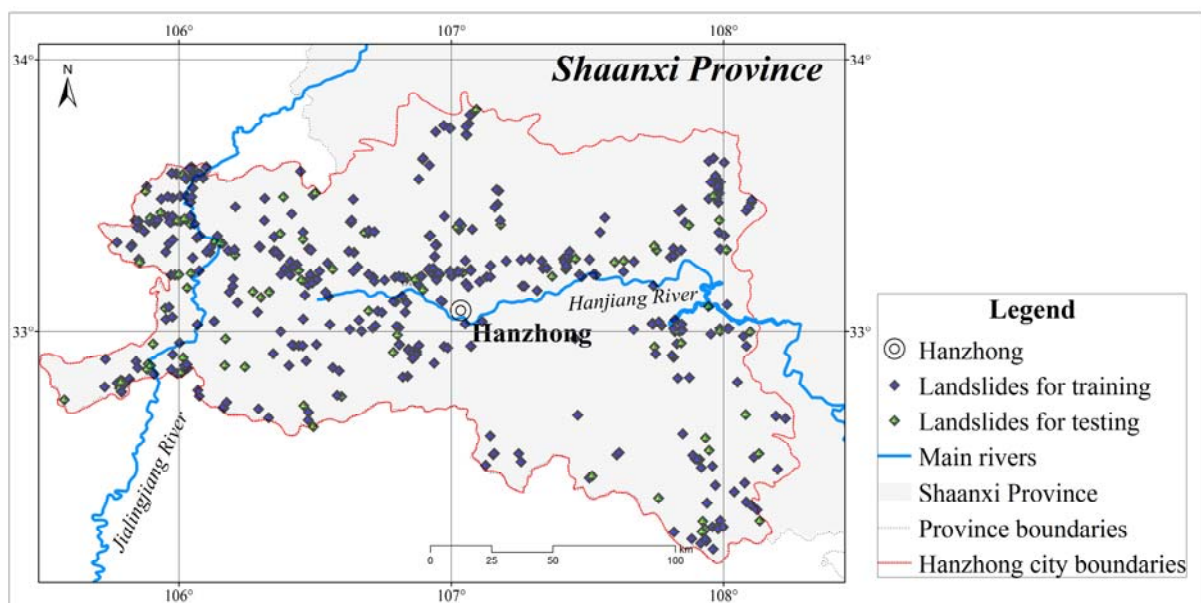


**Figure 2** Selected geological locations for database validation during field surveys.

**Table 1 Data used in the study**

| Data | Description | Specific use |
|---|---|---|
| Multispectral data | Landsat TM (path: 127-129; row: 36-38), band 1-5 & 7, 30m cell size, 2005 | Land use/land cover, NDVI |
| Topographic data | ASTER Global DEM (path 32-34, row 105-108), 30m cell size, 2009 | Terrain attributes (slope angle, aspect, curvature), surface drainage |
| Geologic map | Shaanxi Province, 1:1,750,000, 2002  Hanzhong City, 1:750,000, 2003 | Lithology, faults  Existing landslides, faults |
| Field data | Field surveys, 2001-2003, 2011 | Existing landslides, land use/land cover |

formations, ranging in age from the Proterozoic to the Quaternary. The lithological variations consist of a mainly sedimentary rocks (sandstone, siltstone, mudstone and limestone), some metamorphic rocks (mainly phyllite) and only occasionally some igneous rocks are exposed.

Rivers belong to the Yangtze River watershed, including the Jialingjiang River and Hanjiang River and many tributary streams, which combines a total annual runoff amount of 14.1 billion $m^3$ (Zhang and Jin 1995; Liu 1997). The climate belongs to the humid subtropical zone (Cwa), with average annual temperatures of 13.5°C. Annual precipitation is between 800-1000 mm, of which 70% concentrates from May to October bringing frequent rainstorms to the area.

## 2.2 Characteristics of landslides

Various types of geo-hazards, ranging from landslides, debris flows, avalanches, and ground subsidences occur in the study area. This paper focuses only on landslides as recorded in the official Atlas of Geological hazard Shaanxi Province (DLR(Shaanxi) 2003). According to this official geo-hazard report, there are 990 sites of existing landslides in the study area, categorized into four classes: 7 'very large landslides' (>10 million $m^3$); 79 'large landslides' (1-10 million $m^3$); 264 'medium landslides' (0.1-1 million $m^3$); and 640 'small landslides' (<0.1 million $m^3$). The most wide-spread type of landslide is the translational shallow landslide, whereas deep-seated landslides occasionally developed in the area. In terms of the source materials of landslides, many small landslides developed in alluvial deposits, whereas landslides in loess deposits and in swelling clay soils also occurred.

## 2.3 Data used

A map of 1:250,000 scale was selected (Cascini 2008), which is assumed to fit the regional analysis of Hanzhong City as well as Landsat imagery and ASTER DEM. Data collected includes nine scenes of Landsat TM imagery (http://glcf.umiacs. umd.edu) and ASTER DEM imagery (http://asterweb.jpl.nasa.gov), a geologic map of Shaanxi Province at the scale of 750,000 (Ma 2002), a geo-hazard map in Hanzhong City (DLR(Shaanxi) 2003) and field observation data surveyed during 2001-2003 and 2011 (Table 1). The Landsat TM pre-processed 2005 satellite imagery were downloaded from the Global land Cover Facility (GLCF), of which the three visible bands 1-3, and the infrared band 4 were used in the analyses. The ASTER DEM, derived from ASTER's along-track stereo images, is advantageous for its high spatial resolution (30 m) and high vertical accuracy (20 m). The lithological boundaries depicted on the geological maps were validated during field surveys, to ensure that all existing landslides fall in the correct lithology.

Records of landslide events were derived from the geo-hazard map in Hanzhong City (DLR(Shaanxi) 2003), which was compiled by a group of experts under the organization of local administration. All the geo-hazard locations in the map have been validated in the field investigation conducted during 2001-2003. The investigation has been performed under a series of rules and standards that regulates the behavior in field surveys and the criteria in map compiling (DLR(Shaanxi) 2003). Therefore, the map is regarded as official and reliable, forming the basis of the analysis in our work.

However, available landslide information is very limited compared to other LSZ studies. Firstly,

although 990 landslides have been reported, the publicized geo-hazard map only records 310 of them, while many small landslides are excluded from the map. Therefore, the spatial coverage of recorded landslides is limited. Secondly, the geo-hazard map records the location and the scale of landslides without documenting their areal distribution. Although more detailed landslide information would be preferred in the local area, in reality there is limited coverage of accurate landslide field data. Therefore, the focus is whether a quick and reliable LSZ analysis can be achieved with the existing limited landslide field data.

## 3 Experimental Setup

### 3.1 Preparation of environmental parameters

Landslides are related to both permanent conditioning and triggering factors. Triggering factors that involve a time frame (e.g. the recurrence interval of a rain storm) are normally used for frequency estimates in landslide hazard assessment (Fell et al. 2008), but are beyond the scope of this LSZ analysis. This study considers the empirical permanent environmental factors, although the LSZ model is flexible enough to incorporate such triggering factors (Bălteanu et al. 2010). Based on a thorough examination of the past landslide events in the study area, the following eight potential landslide factors are selected: slope angle; aspect; curvature; lithology; distance from faults; distance from drainage; land use/land cover (LULC); and the Normalized Differential Vegetation Index (NDVI). These permanent conditioning factors were frequently mentioned and used in previously conducted research (Ermini et al. 2005; Fell et al. 2008; Chauhan et al. 2010; Pradhan et al. 2010l; Marjanović et al. 2011) within the context of regional landslide mapping and on a regional scale in the Hanzhong City area.

These selected environmental permanent conditioning factors are calculated using GIS (Geographical Information System) and remote sensing (RS) techniques from the data (Table 1), as shown in Figure 3. Three geomorphometric derivatives, slope, aspect and curvature were

calculated from the 30 m resolution GDEM. The spatial distribution of lithology and fault locations is extracted from the existing geological map. The Landsat TM images are classified into four land cover classes and were used to calculate the NDVI values. The analysis cell size is set to 30 m, which matches the cell sizes of the GDEM and Landsat data. During the field survey landslide locations were checked against the existing landslide map of Shaanxi province.

Parameterization is a preliminary step to integrate environmental causative factors in the LSZ analysis, as the values of the factors are either categorical (e.g. lithology) or continuous (e.g. slope angle). A common method is to group the attribute values into categories for each environmental variable according to its contribution to landslide occurrence (Fell et al. 2008), which is also employed here. To further eliminate the influence of varying data scales, the attribute values of each causative factor are normalized to the range 0~1 using the minimum-maximum method.

To fully appreciate the parameterization, three issues should be taken into account: (1) this work only considers the occurrence of landslide events and did not distinguish different types of landslides. A landslide event that was labeled as 'occurred' refers to 'high landslide susceptibility'; (2) the influence of environmental parameters is, in general, different among types of landslides, therefore the parameterization and the LSZ analysis is based on the main types of landslides in the study area; (3) as each environmental variable is categorized into ordinal values based on prior knowledge, this may potentially influence the ranking of attribute values. Theoretically, such influence is limited, as statistical learning algorithms are supposed to learn from nonlinear data patterns. In practice, appropriate parameterization would certainly help reduce the non-linearity and the complexity of the problem. Therefore, the parameterization is carefully controlled by abundant prior knowledge and thus kept the same among the three algorithms. This means that the learning capability of three statistical learning algorithms can be compared on the same basis.

The three geomorphological derivatives, slope angle, aspect and curvature are calculated from the ASTER DEM using standard GIS tools included in

ESRI's ArcGIS 10.1. The slope angle is divided into five 10° interval classes and a >50° class (Anbalagan 1992; Gong 1996). The factor of slope aspect is considered as two variables of north-exposedness and west-exposedness, as was suggested by Brenning and Trombotto (2006). The north-exposedness is divided into flat south and north categories, while the west-exposedness is divided into flat east and west categories. The curvature is divided into negative, zero and positive values for concave, straight and convex slopes, respectively.

The bedrock geology is grouped into four categories according to their variation in rock strength properties and permeability (Bălteanu et al. 2010): (a) stable: igneous rocks; (b) moderately



**Figure 3** The resulting thematic layers in the study area of (a) slope angle, (b) North-exposedness, (c) west-exposedness, (d) slope curvature, (e) lithology, (f) distance from fault, (g) distance from drainage, (h) land cover, (i) NDVI  (-To be Continued-)
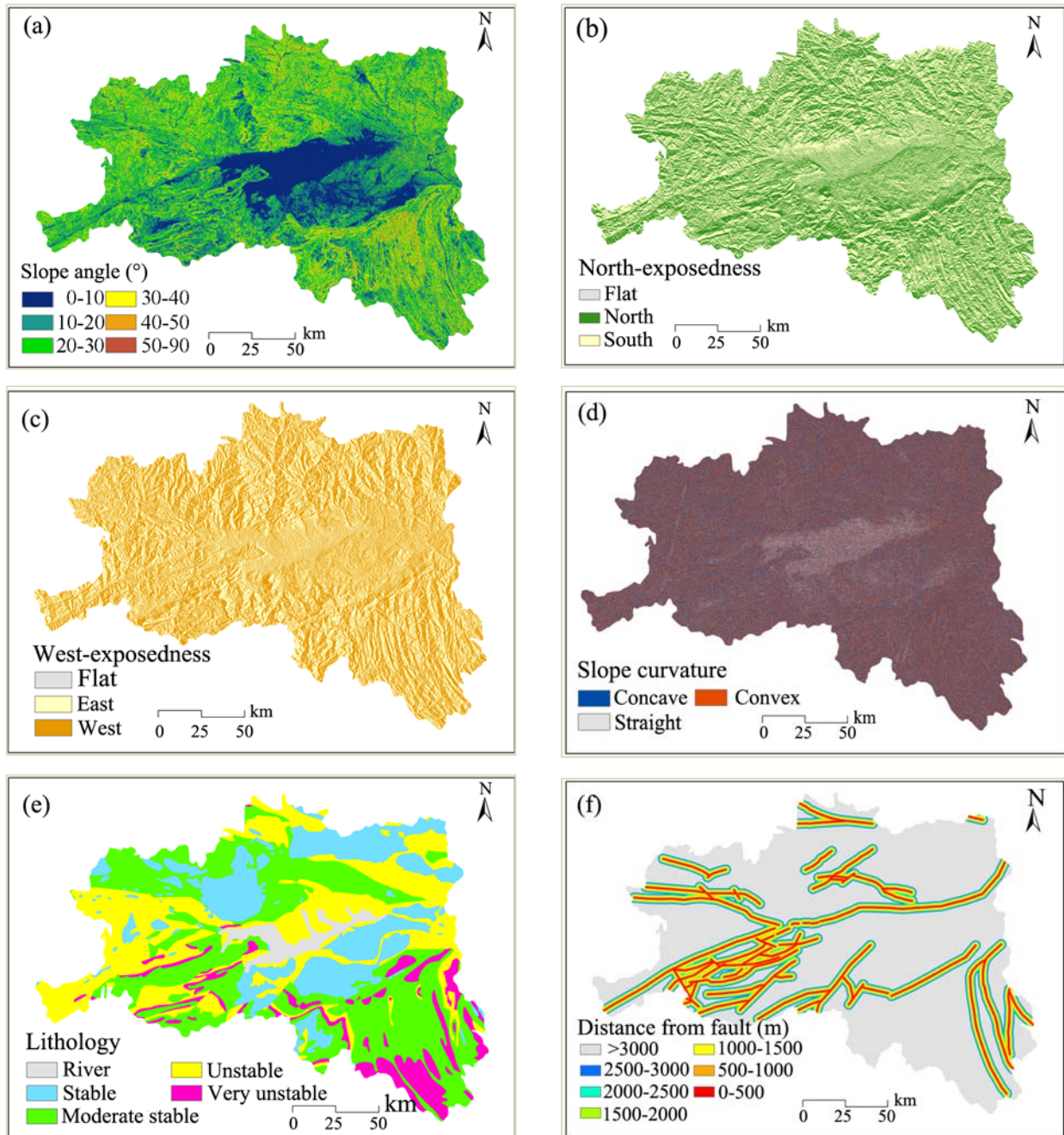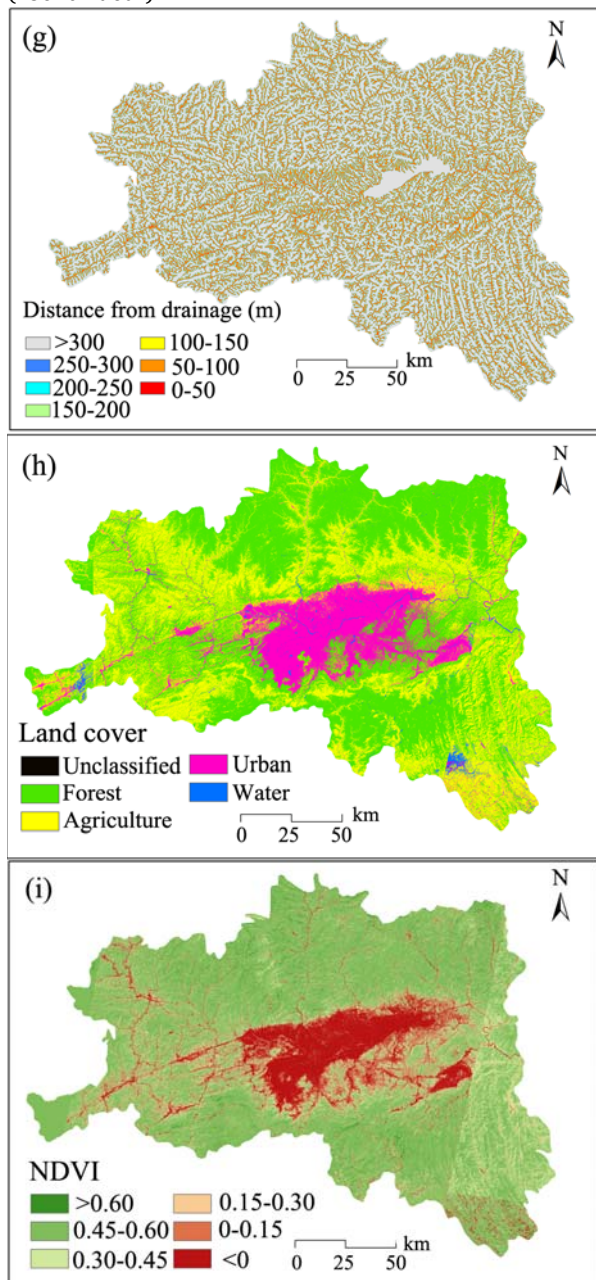
(-Continued-)



**Figure 3** The resulting thematic layers in the study area of (a) slope angle, (b) North-exposedness, (c) west-exposedness, (d) slope curvature, (e) lithology, (f) distance from fault, (g) distance from drainage, (h) land cover, (i) NDVI.

stable: limestones, dolomitic, marl and volcanic formations; (c) unstable: phyllite and shale; (d) very unstable: mudstone formations.

A buffer zone is created around the digitized faults in intervals of 500 m in line with the procedure of Abdallah et al. (2005) and Chauhan et al. (2010). The upper limit of the buffer distance to faults was set to 3000 m.

The surface drainage is extracted from ASTER DEM using the maximum slope gradient rule. The distances of a raster cell to the nearest drainage section are calculated and classified into seven categories at 50 m intervals, as was suggested by Pradhan and Lee (2007) and Pradhan et al. (2010).

Land use/land cover is derived from the mosaicked TM imagery by applying a supervised maximum likelihood classifier. Four land cover classes are used in training: forest, agriculture, urban/infrastructure and water. The overall accuracy and the *Kappa* index were calculated following the method described by Congalton (1991) as 86.0% and 0.789 respectively, indicating high accuracy and good agreement with ground truth of the classification result.

NDVI is calculated from the red and infrared bands of the TM images using the ENVI 4.6 software. The non-vegetation surfaces coincide with negative to close to zero NDVI values, whereas high positive values indicate a dense vegetation cover. Frequency analysis revealed that 83.81% of the cells distribute from 0 to 0.6. Therefore, the NDVI values are categorized into four classes of 0.15 intervals from 0 to 0.6 and two additional categories representing values below 0 and above 0.6. It is assumed that non-vegetated areas have a relation with landslide occurrence, particularly where the scars of landslides are related to the removal of vegetation in the lower depositional zones. Table 2 summarizes the parameterization and landslide density analysis of the eight environmental factors.

### 3.2 Preparation of the training set and input data

This study applies a point-based sampling strategy (Yilmaz 2010b), which was selected with respect to available landslide data source. The analysis is based on raster cells, where each landslide sample corresponds to a pixel that indicates landslide instances. Training samples have been selected from observed landslide areas and non-landslide areas. Although non-landslide areas are generally much larger than landslide affected areas, non-landslide samples have been carefully selected based on field surveys. In this way, it avoids the many samples on flat areas that

**Table 2 Parameterization and analysis of the environmental factors**

| Environmental factors | Categories | Attribute values | Standardized values | Area a (%) | Past landslides b (%) | Landslide density (b/a) |
|---|---|---|---|---|---|---|
| Slope angle | 0-10° | 1 | 0.167 | 20.331 | 31.290 | 1.539 |
| | 10-20° | 2 | 0.333 | 28.652 | 31.290 | 1.092 |
| | 20-30° | 3 | 0.500 | 29.749 | 23.226 | 0.781 |
| | 30-40° | 4 | 0.667 | 16.003 | 10.000 | 0.625 |
| | 40-50° | 5 | 0.833 | 4.438 | 3.548 | 0.800 |
| | 50-90° | 6 | 1.000 | 0.828 | 0.645 | 0.780 |
| North-exposedness | Flat | 1 | 0.333 | 0.620 | 0.323 | 0.520 |
| | North | 2 | 0.667 | 48.441 | 49.032 | 1.012 |
| | South | 3 | 1.000 | 50.939 | 50.645 | 0.994 |
| West-exposedness | Flat | 1 | 0.333 | 0.620 | 0.323 | 0.52 |
| | East | 2 | 0.667 | 50.489 | 49.677 | 0.984 |
| | West | 3 | 1.000 | 48.891 | 50.000 | 1.023 |
| Curvature | Concave | 1 | 0.333 | 46.794 | 47.097 | 1.006 |
| | Straight | 2 | 0.667 | 8.291 | 7.419 | 0.895 |
| | Convex | 3 | 1.000 | 44.915 | 45.484 | 1.013 |
| Lithology | Stable | 1 | 0.250 | 22.753 | 7.765 | 0.336 |
| | Moderate stable | 2 | 0.500 | 35.388 | 39.192 | 1.098 |
| | Unstable | 3 | 0.750 | 33.487 | 44.261 | 1.324 |
| | Very unstable | 4 | 1.000 | 8.172 | 8.786 | 1.082 |
| Distance from fault | >3000 m | 1 | 0.143 | 82.417 | 51.290 | 0.622 |
| | 2500-3000 m | 2 | 0.286 | 2.468 | 6.129 | 2.483 |
| | 2000-2500 m | 3 | 0.429 | 2.647 | 6.774 | 2.559 |
| | 1500-2000 m | 4 | 0.571 | 2.870 | 6.774 | 2.360 |
| | 1000-1500 m | 5 | 0.714 | 3.072 | 6.774 | 2.205 |
| | 500-1000 m | 6 | 0.857 | 3.232 | 10.323 | 3.194 |
| | 0-500 m | 7 | 1.000 | 3.294 | 11.935 | 3.624 |
| Distance from drainage | >300 m | 1 | 0.143 | 71.222 | 28.710 | 0.403 |
| | 250-300 m | 2 | 0.286 | 4.363 | 7.742 | 1.775 |
| | 200-250 m | 3 | 0.429 | 4.623 | 11.935 | 2.582 |
| | 150-200 m | 4 | 0.571 | 3.699 | 9.355 | 2.529 |
| | 100-150 m | 5 | 0.714 | 5.237 | 10.645 | 2.033 |
| | 50-100 m | 6 | 0.857 | 5.171 | 12.581 | 2.433 |
| | 0-50 m | 7 | 1.000 | 5.685 | 19.032 | 3.348 |
| Land use/land cover | Forest | 1 | 0.250 | 52.125 | 29.355 | 0.563 |
| | Agriculture | 2 | 0.500 | 34.118 | 44.839 | 1.314 |
| | Urban | 3 | 0.750 | 12.326 | 21.290 | 1.727 |
| | Water | 4 | 1.000 | 1.432 | 4.516 | 3.154 |
| NDVI | 0.6-1.0 | 1 | 0.167 | 1.978 | 0.645 | 0.326 |
| | 0.45-0.6 | 2 | 0.333 | 59.121 | 34.839 | 0.589 |
| | 0.3-0.45 | 3 | 0.500 | 18.974 | 22.903 | 1.207 |
| | 0.15-0.3 | 4 | 0.667 | 4.124 | 6.452 | 1.564 |
| | 0-0.15 | 5 | 0.833 | 1.587 | 3.548 | 2.235 |
| | <0 | 6 | 1.000 | 14.215 | 31.613 | 2.224 |

are trivial in the algorithm learning. The final sample dataset contains 310 landslide occurrence and 310 non-landslide samples, in which landslide and non-landslide areas are kept at a balanced ratio of 1:1. Each of the samples are connected to eight environmental parameter values (independent variables) and one output value (dependent variable), in which landslide occurrence is indicated as '1' and landslide absence is indicated as '0'.

As suggested by Basheer and Hajmeer (2000) and Chauhan et al. (2010), the total dataset is split into a randomly selected training set, 80% (496 samples), and a testing set, 20% (124 samples). To make sure that training and testing set have similar statistical distribution as the whole dataset, the

whole dataset is sorted in ascending order according to the attribute values of the environmental parameters. Then the testing set is selected as every one out of five (20%) records, while the remaining dataset is taken as the training set. The training and testing sets are kept the same in the three LR, ANN and SVM algorithms. After the training phase, the models can predict LSZ for the entire area.

### 3.3 LSI calculation and classification

The maps of LSI are calculated by applying the three algorithms to the entire study area. The final LSZ map is a reclassification of LSI values into five zones, categorized as: very low susceptibility (VLS), low susceptibility (LS), moderate susceptibility (MS), high susceptibility (HS) and very high susceptibilities (VHS), following the classifications used by Clerici et al. (2002), Bălteanu et al. (2010), Chauhan et al. (2010), and Pradhan and Lee (2010). Success rate curves (Saha et al. 2005) are calculated to define the thresholds for zoning the LSI values based on the mean ($\mu_o$) and standard deviation ($\sigma_o$) values, defined as ($\mu_o - 1.5m\sigma_o$), ($\mu_o - 0.5m\sigma_o$), ($\mu_o + 0.5m\sigma_o$) and ($\mu_o + 1.5m\sigma_o$), where $m$ is a positive value. The value of $m$ is usually decided by trial and error in the vicinity of $m$=1 (Chauhan et al. 2010).

### 3.4 Stability test

To evaluate the influence of limited landslide samples, the stabilities of the three statistical learning algorithms with varying sample size is further examined. As the learning algorithms are essentially data-driven, the sizes of the training set might influence the model performances. As the LR, ANN and SVM algorithms are all free of data distribution (Yao et al. 2008; Bai et al. 2010; Chauhan et al. 2010), it is investigated whether and how the sizes of the training set influence the performances of the models. To guarantee the same basis for comparison, at each implementation the training set is kept identical among the algorithms.

While the training set is increased from 10% to 100% of the entire training set by 10% steps, repetitive model runs are performed and the result of each implementation is evaluated. In detail, the

training set is first randomly divided into ten equal subsets; a training set is actually a random combination from these ten subsets, as used in the SVM evaluation in Yao et al. (2008). Then for each size of the training set (e.g. 50%), the experiments are performed repetitively until each subset has been used at least once. The average accuracy of the repetitive experiments is taken as the indicator of model performance.

### 3.5 Model evaluations

The performance of the algorithms in predicting landslide susceptibilities is evaluated in three metrics: 1) landslide density; 2) the receiver operating curves (ROC) and the area under curves (AUC); and 3) error matrix. In addition to these criteria, the dependence of the three statistical learning algorithms on sample size is determined by varying the sizes of the training set.

Landslide density is defined as the ratio of the percentage of existing landslide area to the percentage of zonation area. It evaluates the density of existing landslides in each susceptibility zone. The basic assumption is that areas where landslides have occurred are more susceptible to landslides, so landslide density is expected to be higher in VHS and HS zones than it is in VLS and LS zones. According to Chauhan et al. (2010), landslide density is supposed to increase gradually from VLS to VHS zones.

ROC curves and AUC are commonly used as accuracy criteria for prediction models in natural hazard assessments (Chauhan et al. 2010; Pradhan et al. 2010). The curves are obtained by plotting cumulative rates of true positives (correctly classified landslide cells) versus false positives (landslide absence, misclassified landslide cells). An AUC value of 1 indicates perfect classification performance, whereas a value of 0.5 is considered equal to a random model. The details of ROC and AUC can be found in Frattini et al. (2010).

Error matrix analysis is used to evaluate the correctly predicted cases versus the failed cases. *Kappa* index is employed as a quantitative indicator. It is advantageous by removing effects of random agreement between predictions and ground truths by considering all elements in the error matrix (Lillesand et al. 2008; Chauhan et al. 2010). Unlike ROC curve, *Kappa* index is cutoff-

dependent, which means that a cutoff should be decided for classifying LSI values into two classes, corresponding to landslide occurrence or absence in ground truth data. In this study, the cutoff is determined as the value that has the highest *Kappa* index by iterating cutoff values from 0 to 1 (Fieldings and Bell 1997). Following Monserud and Leemans (1992), accuracy levels evaluated by *Kappa* index are: (1) <0.4, poor; (2) 0.4~0.55, moderate; (3) 0.55~0.70, good; (4) 0.70~0.85, very good; and (5) >0.85, nearly perfect.

# 4 Implementation

## 4.1 LSZ using LR

The LR is stepwise, backwardly implemented. In each step, Wald statistics are calculated to test the significance of each of the nine environmental variables, with a threshold set to 0.05 (George and Mallery 2000). The variable with Wald statistical value below 0.05 is kept in the model, otherwise the variable fails the significance test and is removed from the equation. One initially included variable (land cover types), failed the significance test and was removed from the equation, hence the established LR model is:

$$Logit(P) = -3.834 - 1.271x_1 - 1.101x_2 + 2.084x_3 \\ + 0.683x_4 + 2.058x_5 + 1.017x_6 + 0.895x_7 + 1.910x_9 \quad (4)$$

where $P$ is the estimated probability of landslide occurrence; $x_1$, $x_2$, $x_3$, $x_4$, $x_5$, $x_6$, $x_7$ and $x_9$ are slope angle, north-exposedness, west-exposedness, curvature, lithology, distance from fault, distance from drainage and NDVI; the coefficient value of variable $x$ is the weight attributed to each input variable. Through equation (4), the LSI values of the entire study area are derived with the mean value $\mu_0 = 0.41$ and standard deviation value $\sigma_0 = 0.21$.

To determine the boundaries of LSZ based on LSI values, three representative success rate curves corresponding to m=1, 1.1 and 1.2 are plotted and analyzed following the method used by Chauhan et al. (2010). The curve corresponding to *m*=1 is determined as the optimal success rate. Accordingly, the boundaries for LSZ are fixed at LSI values of 0.10, 0.31, 0.52 and 0.72, deriving: very low susceptibility (VLS) category with LSI

values lower than 0.10; low susceptibility (LS) category with LSI values between 0.10 and 0.31; moderate susceptibility (MS) category with LSI values between 0.31 and 0.52; high susceptibility (HS) category with LSI values between 0.52 and 0.72; and very high susceptibility (VHS) category with LSI values higher than 0.72. The derived LSI map is shown in Figure 4.

## 4.2 LSZ using ANN

The LSZ using ANN is implemented in the MATLAB interface. The MLP-BP ANN consists of nine input neurons corresponding to selected environmental parameters and one output neuron that indicates landslide occurrence. The one hidden layer structure is assumed to be sufficient to handle nine input neurons. The number of neurons in the hidden layer is usually determined through trial and error (Ermini et al. 2005; Chauhan et al. 2010). Thus ANN with varying numbers of neurons in hidden layers are tested to determine the optimal structure. The learning rate is 0.01 and the initial weights are selected randomly from -1 to 1. The maximum epochs are 2000 and the root mean square error (RMSE) goal is set to 0.01, either of which criterion will stop the training once met.

The difference in accuracies between predictions on training and testing sets is taken as the criterion to determine the optimal ANN. An ANN with good generalization capability is supposed to derive predictions with nearly identical high accuracies on both known and unknown samples. By varying the number of neurons in the hidden layer from five to twenty-five, repetitive experiments are performed. Both the predictions on the training and testing sets are derived and evaluated. Finally, the ANN with 9×22×1 architecture that derives testing and training accuracies as 0.70 and 0.72 respectively is determined as the optimal one. The LSI map based on the MLP-BP ANN algorithm is shown in Figure 4b.

An ANN learns from known samples by automatically adjusting weights between neurons. Although adjustments of weights are kept in a black box, we can capture the updated weight matrices between layers, e.g. a 9×22 weight matrix for input-hidden connections and a 22×1 weight matrix for hidden-output connections in the case of
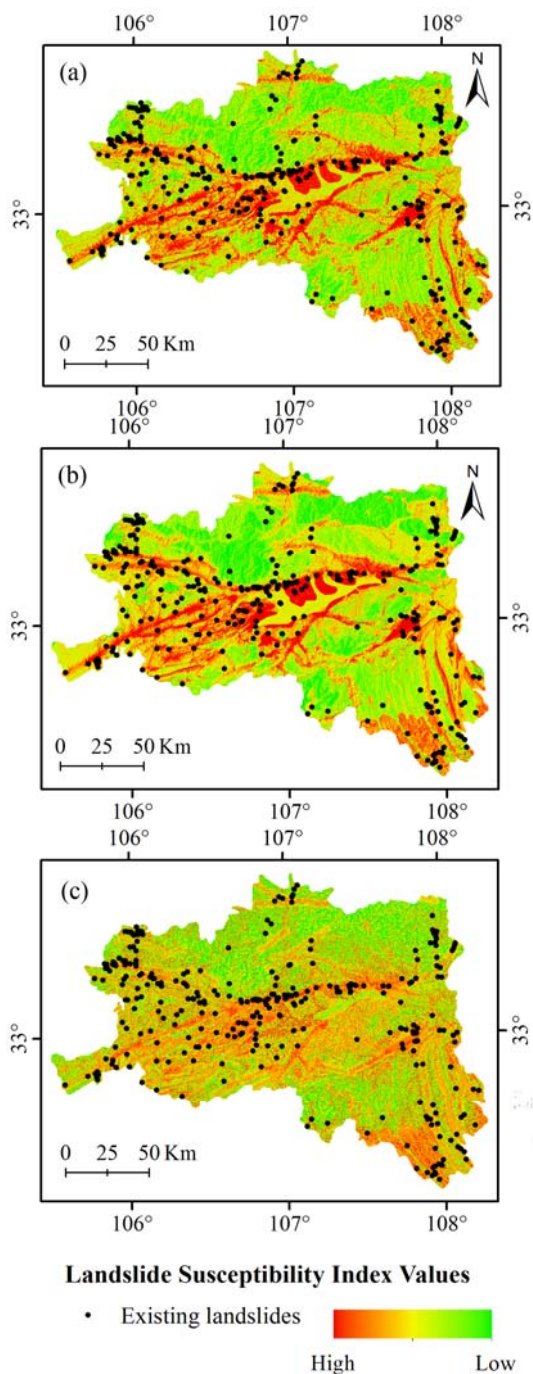
**Figure 4** The three LSZ (landslide susceptibility zonation) maps derived from (a) Logistic Regression (LR), (b) Artificial Neural Networks (ANN) and (c) Support Vector Machine (SVM).

this study. Then matrix multiplication is performed on each of the matrices in a sequential manner to obtain a 9×1 weight matrix, similar to Chauhan et al. (2010). These weights are assumed to reflect the contribution of the environmental variables on slope failure.

### 4.3 LSZ using SVM

The LSZ using SVM is implemented in the MATLAB interface of the software LibSVM 2.8.9 (Chang and Lin 2011). The environmental parameters are taken as input variables and the output is the probability of landslide occurrence. The most widely used radial basis function (RBF) is employed as the kernel function, which is insensitive to outliers (Tax and Duin 1999). The kernel function parameter $\gamma$ and penalty parameter $C$ of SVM require careful determination to ensure the model performance. Thus, a cross validation method is used to determine optimal values for $C$ and $\gamma$. For each pair of $(C, \gamma)$, the training dataset is randomly divided into five equal folds, in which four folds are for training and one fold for testing. By iterating each fold as a test set and the other four folds as a training set, five mean squared error (MSE) values are calculated. The best values of $C$ and $\gamma$ are determined as the pair that produces the minimal average MSE value. A rough parameter search is performed by increasing $C$ and $\gamma$ from $2^{-8}$ to $2^8$ at an increasing step of $\log_2 C=1$ and $\log_2 \gamma=1$. Based on the rough search result, a precise parameter search is performed by increasing $C$ and $\gamma$ from $2^{-4}$ to $2^4$ at an increasing step of $\log_2 C=0.5$ and $\log_2 \gamma=0.5$. The determined optimal values of $C$ and $\gamma$ are 5.657 and 11.314 respectively, deriving the final LSI map (Figure 4c).

## 5 Results

### 5.1 Model comparisons

The resulting LSI maps derived from the three algorithms are presented and compared to existing landslides in Figure 4. The resulting maps show distinct spatial intensity patterns of landslide susceptibility. In all maps, high values of LSI maps generally coincide with the location of existing landslides. The SVM map however, shows less extreme variations of the LSI. In addition, the LSI map from SVM shows more detailed spatial patterned distribution of landslides. This is attributed to its capability of dealing with high-dimensional nonlinear data. For example, the spatial pattern of the SVM result shows more influence of drainage patterns, while LR and ANN consider drainage as a less influential factor, thus

**Table 3 Landslide data in the five landslide susceptible zones for the three algorithms used**

| Landslide susceptibility zones | | VLS | LS | MS | HS | VHS |
|---|---|---|---|---|---|---|
| LR | % of the area of LS zones (a) | 2.38 | 34.48 | 33.07 | 20.28 | 9.78 |
| | % of landslides of LS zones (b) | 0.65 | 15.16 | 27.42 | 29.35 | 27.42 |
| | Landslide density (b/a) | 0.27 | 0.44 | 0.83 | 1.45 | 2.80 |
| ANN | % of the area of LS zones (a) | 4.03 | 29.31 | 38.94 | 19.07 | 8.65 |
| | % of landslides of LS zones (b) | 0.65 | 12.90 | 34.52 | 28.06 | 23.87 |
| | Landslide density (b/a) | 0.16 | 0.44 | 0.89 | 1.47 | 2.76 |
| SVM | % of the area of LS zones (a) | 2.30 | 26.78 | 37.20 | 27.32 | 6.41 |
| | % of landslides of LS zones (b) | 0.65 | 9.35 | 24.52 | 35.16 | 30.32 |
| | Landslide density (b/a) | 0.28 | 0.35 | 0.66 | 1.29 | 4.73 |

**Notes:** LR for logistic regression; ANN for Artificial Neural Network; SVM for Support Vector Machine; VLS for very low susceptibility; LS for low susceptibility; MS for moderate susceptibility; HS for high susceptibility; VHS for very high susceptibility.

suppressing such potential patterns.

The predicted susceptibility zones and landslide density are calculated based on raster cells on LSZ maps, as shown in Table 3, in order to evaluate whether the zonation areas coincide with existing landslides. It is generally assumed that highest landslide frequencies are expected to occur in VHS and HS zones, and low landslide frequencies are expected in VLS and LS zones. The results presented in Table 3 confirm that SVM predicts the susceptibility zones best, as most (65.48%) landslides occur in VHS and HS zones and only 10.00% of existing landslides fall in VLS and LS zones. LR predicts 56.77% of landslides in VHS and HS zones, while 51.93% is predicted by ANN. ANN predicts 13.35% of landslides in VLS and LS zones, while 15.81% is predicted by LR.

Landslide density shows the degree of concentration of existing landslide in the zonal area. Generally a gradually increasing trend of landslide density values is expected from VLS to VHS zones. The VHS zone in a LSZ map is supposed to derive the highest landslide density values (Chauhan et al. 2010). Table 3 shows that all the three algorithms show increasing trends in landslide density from VLS to VHS zones. SVM derives the highest landslide density in the VHS zone. The landslide density differences from the VLS to HS zones are less prominent among the three approaches. This suggests that SVM is very reliable in predicting highly concentrated distribution of existing landslides in zones of very high susceptibility.

The ROC curves are plotted using the testing set of landslide samples (Figure 5). The largest difference of the three curves lies in the range

between 0.05 and 0.45 of false positive rates. Here, the curve of SVM shows far better performance than ANN and LR. The calculated AUC values reveal that SVM derives the highest accuracy (AUC=0.853), whereas LR (AUC=0.718) and ANN (AUC=0.763) derives lower accuracies. Calculated *Kappa* indices of 0.355 (LR), 0.403 (ANN) and 0.613 (SVM) respectively, agree well with the



**Figure 5** The receiver operating characteristic (ROC) curves of landslide susceptibility index (LSI) plotted for the three algorithms.

model performances evaluated by the ROC analysis.

The quantitative evaluations demonstrate that, given the entire training set, SVM performs significantly better than LR and ANN, and that ANN outweighs LR in LSZ studies. In addition, the evaluation of coincidence of existing landslides with susceptibility zones also shows that SVM performs best, while the performances of LR and

ANN do not show prominent differences.

## 5.2 Model stability under varying sizes of training set

Model stability is investigated as the dependence of model performance on sample sizes. A selection of test results for LR coefficients under five sizes of the training set are presented in Table 4. The regression coefficients derived from LR with varying sizes of training set highly fluctuate. It is observed that eligible variables under varying sizes of training set are different, and the same variable may have different coefficient values during different trials. This suggests that the established LR equation (equation 2) is easily affected by the size of the training samples, which might raise questions about the reliability of the derived coefficients in describing the relative contributions of the environmental factors.

In the ANN algorithm, the network architecture, including both the number of neurons in each layer and the layers themselves, are kept the same among the implementations. We have investigated the changes in the updated weight

matrix to determine the influences of the size of the training set on the trained ANN model, with examples shown in Table 5. Assumed as indicating the contributions of the environmental variables on landslide, the weights derived from the ANN model trained with varying sizes of training set are found to fluctuate. This indicates that the ANN model is easily influenced by the size of training set.

In the SVM algorithm, the optimal values of parameters $C$ and $\gamma$ are also different among varying sizes of training sets (Table 6). This indicates that the size of training set does affect the established SVM model. SVM is essentially a nonlinear model, the changes in the parameters ($C$, $\gamma$) of the kernel function indicate the adaption of the SVM model to the varying training sets, rather than the changes in contribution of the environmental variables on landslides.

All the three data-driven algorithms are affected by the size of training set. Besides the sample size, random sampling of a subset from the entire training set also casts potential effects on model performance. For example, even under the same size of training set, LR or ANN could also derive different coefficient values for the same

**Table 4 Logistic Regression (LR) coefficients under five sizes of training samples**

| Samples | | 10% | | 30% | | 50% | | 70% | | 90% | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Trial | | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd |
| Coefficients | $b_1$ | - | - | - | - | - | - | - | - | - | - |
| | $b_2$ | -5.94 | - | - | -2.73 | -2.91 | -1.99 | -2.22 | -2.59 | -2.12 | -2.23 |
| | $b_3$ | - | 8.21 | 2.50 | 3.63 | 3.06 | 3.38 | 2.66 | 3.80 | -3.13 | 3.38 |
| | $b_4$ | - | - | - | - | - | 1.10 | - | 0.95 | - | - |
| | $b_5$ | - | - | 1.71 | 1.80 | 1.68 | 1.39 | 1.69 | 1.50 | 1.63 | 1.48 |
| | $b_6$ | - | - | - | - | - | - | - | - | - | - |
| | $b_7$ | - | - | 1.37 | - | - | - | 0.95 | - | 0.91 | 0.83 |
| | $b_8$ | - | - | - | - | - | - | - | - | - | - |
| | $b_9$ | 6.04 | - | 2.74 | 2.22 | 2.63 | 1.81 | 2.44 | 1.86 | 2.15 | 1.93 |
| | $b_{10}$ | - | -14.01 | - | - | - | -4.52 | -2.71 | -3.78 | -3.37 | -3.28 |

**Note:** "-" indicates the variable fails the significance test and is removed

**Table 5 Artificial Neural Networks (ANN) derived weights under five sizes of training set**

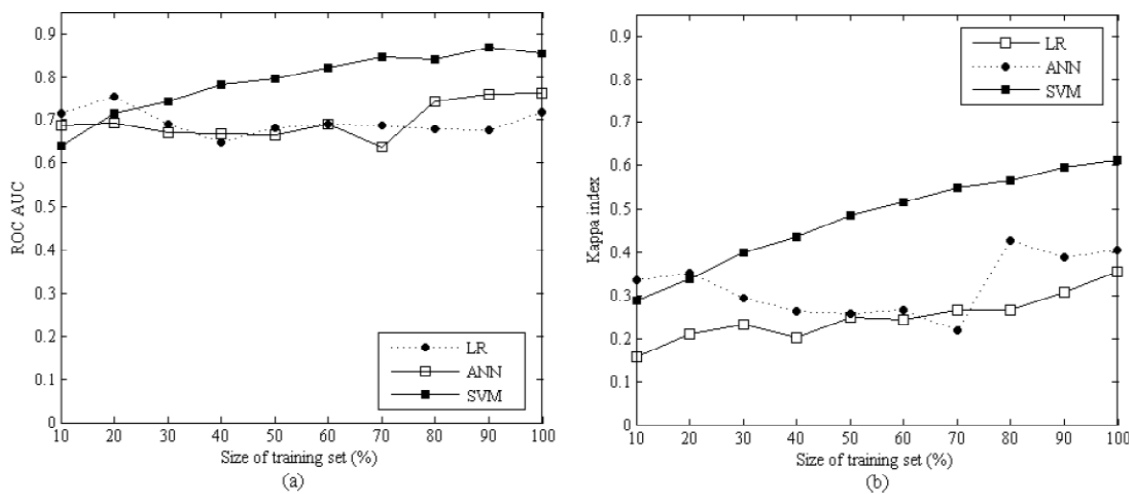| Samples | | 10% | | 30% | | 50% | | 70% | | 90% | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Trial | | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd |
| Derived weights | $w_1$ | -1.44 | -2.10 | -5.92 | -5.16 | 5.38 | 4.68 | 1.44 | 0.86 | 1.18 | -0.11 |
| | $w_2$ | -0.48 | -0.29 | 2.85 | 0.87 | -1.32 | -2.12 | -1.93 | -2.39 | -0.33 | -2.23 |
| | $w_3$ | 1.87 | 3.69 | 4.32 | 4.92 | 2.37 | 3.28 | 5.60 | 6.19 | 3.13 | 5.63 |
| | $w_4$ | 3.30 | 4.67 | 3.01 | 3.16 | 1.78 | 2.63 | 1.13 | 1.47 | -0.24 | 0.50 |
| | $w_5$ | 5.00 | 5.48 | 0.19 | 1.19 | -0.79 | 0.06 | 0.29 | 1.15 | 2.78 | 3.16 |
| | $w_6$ | -1.42 | 3.47 | 1.35 | 2.23 | 0.04 | 0.79 | 0.23 | 0.58 | 1.10 | 1.95 |
| | $w_7$ | 2.64 | 2.93 | 4.27 | 4.75 | 2.48 | 2.04 | 2.52 | 2.01 | 5.23 | 4.90 |
| | $w_8$ | 0.02 | -0.06 | 1.52 | 1.30 | 1.68 | 1.97 | 0.75 | 0.28 | 3.14 | 2.62 |
| | $w_9$ | 3.72 | 2.26 | 2.16 | 1.14 | 1.88 | 1.86 | 5.68 | 4.68 | 2.81 | 2.69 |

**Figure 6** Values of (a) the area under curves (AUC) and (b) Kappa index of the three models with varying sizes of training set.

**Table 6 Support Vector Machine (SVM) optimal parameters under five sizes of training set**

| Samples | | 10% | | 30% | | 50% | | 70% | | 90% | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Trial | | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd | 1st | 2nd |
| Optimal parameters | $C$ | 1.41 | 2.83 | 0.35 | 0.35 | 1.41 | 4 | 5.66 | 11.31 | 4 | 16 |
| | $\gamma$ | 5.66 | 2 | 5.66 | 4 | 4 | 8 | 0.16 | 0.16 | 8 | 8 |

variable. Therefore, the average accuracies of repetitive implementations of a model under a given size of the training set, are calculated and compared, as shown in Figure 6. In the LR algorithm, a fluctuating trend is observed for AUC values, whereas the curve of *Kappa* index shows a steadily increasing trend. The LR algorithm has the highest AUC values using small sizes of training set (10%-20%), although differences with ANN and SVM are rather small. The SVM algorithm has the highest AUC values between 30% and 100% of the training set size. The values of AUC of LR are close to those of the ANN between 10% and 60% of the training set size, but are lowest between 80% and 100% of the training set. The Kappa index curve for LR remains below 0.3 between 10% and 100% of the training set, which seems not sufficient to achieve an acceptable accuracy.

In the ANN algorithm, the curves of both AUC and *Kappa* index exhibit marked fluctuations in the entire range. The ANN algorithm derives fair accuracies on a small size (10%-20%) of training set. Accuracies decrease steadily until a drop at the 70% initiates a sharp increase from 70% to 80% of the training set. The increase from 90% to 100%

indicates that the ANN performance would be more predictive with increased size of the training set.

In the SVM algorithm, a steady increase in both AUC and *Kappa* index values are observed from 10% to 100% in the training set. On a small size of the training set (10%-20%), the SVM obtains moderate accuracy. The values of AUC and *Kappa* index steadily increase and outweigh those of the LR and ANN algorithms above 30% of the training set size. Between 60% and 100% of the training set size, the values of AUC and *Kappa* index maintain highest levels above 0.8 and 0.5 respectively, indicating accurate results.

**5.3 Geomorphological interpretations**

The LSZ results and validation of the LR, ANN and SVM models have some consequences for geomorphological interpretations. The focus is to predict landslides, which, in terms of the area and volumes involved, have been reported in four categories. In Table 7 the model outcome is listed according to the size of the reported landslides. Most of the landslide training samples fall within the small and medium size category. The percentages of

**Table 7 Landslide data within the five susceptibility zones for the three models**

| Size of landslide | Area (%) | LSZ zones | Area of each LSZ zone (%) | | |
| --- | --- | --- | --- | --- | --- |
| | | | LR | ANN | SVM |
| Small (<0.1 million m³) | 41.94 | VLS | 0.77 | 0.77 | 0.00 |
| | | LS | 10.77 | 10.77 | 9.23 |
| | | MS | 33.08 | 35.38 | 26.15 |
| | | HS | 31.54 | 30.77 | 38.46 |
| | | VHS | 23.85 | 22.31 | 26.15 |
| Medium (0.1-1 million m³) | 42.58 | VLS | 0.75 | 0.75 | 1.50 |
| | | LS | 17.29 | 12.03 | 9.02 |
| | | MS | 23.31 | 34.59 | 21.05 |
| | | HS | 29.32 | 28.57 | 36.84 |
| | | VHS | 29.32 | 24.06 | 31.58 |
| Large (1-10 million m³) | 13.23 | VLS | 0.00 | 0.00 | 0.00 |
| | | LS | 19.51 | 19.51 | 12.20 |
| | | MS | 24.39 | 34.15 | 24.39 |
| | | HS | 24.39 | 17.07 | 21.95 |
| | | VHS | 31.71 | 29.27 | 41.46 |
| Very large (>10 million m³) | 1.94 | VLS | 0.00 | 0.00 | 0.00 |
| | | LS | 33.33 | 33.33 | 0.00 |
| | | MS | 16.67 | 16.67 | 66.67 |
| | | HS | 16.67 | 33.33 | 16.67 |
| | | VHS | 33.33 | 16.67 | 16.67 |

**Notes:** LSZ means landslide susceptibility zonation; LR means Logistic Regression; ANN means Artificial Neural Networks; SVM means Support Vector Machine.

each size of landslides falling into five susceptibility zones are calculated to test how the algorithms work on various sizes of landslide. Generally a higher percentage of landslides should be classified in the LSZ results with susceptibility classes of very high and high. This is based on the assumption that the area where landslides have occurred previously is highly susceptible to landslides. As shown in Table 7, the SVM algorithm performs best on all sizes of landslide, as it classifies the least existing landslides as VLS and LS categories and the most existing landslides as VHS and HS zones. None of the very large landslides have been classified by the SVM as VLS and LS zones, while both the LR and ANN algorithms classifies 33.33% of them as LS zones. On all sizes of landslide, the LR algorithm classifies higher percentages of landslides as HS and VHS zones than the ANN. The ANN algorithm classifies a lower percentage of landslides as VLS and LS zones than the LR does on medium-sized landslides.

Furthermore, the LSZ maps are analyzed in comparison with basic terrain attributes of lithology (Table 8), to emphasize the degree of susceptibility of various lithologies to landslides in the area. As expected (Table 8), mudstone formations are most susceptible to landslides, where 100% of the area falls in the MS, HS and

VHS zones in all the three algorithms. Phyllite and shale formations rank second for their susceptibility to landslide occurrence, with an average of 91.61% of the area classified as medium to very high susceptibility. Igneous rocks show less landslide susceptibility than limestone, dolomite, marlite and volcanic formations. These findings confirm that lithology type is a major factor in controlling the susceptibility of the terrain to landslide and is in lines with earlier prior knowledge about lithological stabilities (Bălteanu et al. 2010).

## 6    Discussion

In our LSZ case study, the SVM algorithm outperforms both the ANN and LR algorithms in all evaluation metrics. In the majority of LSZ related studies, there seems to be a trend to regard SVM as superior to most other algorithms (Bui et al. 2012), such as in analytical chemistry (Balabin and Lomakina 2011). For example, our results agree well with the good performance of SVM algorithm used by Marjanović et al.(2011), which outperforms both LR and decision tree algorithms. The success of SVM can be attributed to its strong

generalization capacity and high robustness.

The ANN algorithm outperforms LR algorithm based on evaluations of the AUC and *Kappa* index. This is in agreement with the findings in the comparative studies by Yilmaz (2009) and

sizes of training set. In both the LR and ANN algorithms, fluctuations in accuracies are mostly observed from 10% to 80% of the training set. This might infer that both algorithms require most landslide samples to derive accurate results.

**Table 8 The areal distribution of LSZ zones in each lithological unit**

| Lithology | Area (%) | LSZ zones | Area of each LSZ zone (%) | | | |
|---|---|---|---|---|---|---|
| | | | LR | ANN | SVM | Mean |
| River | 3.31 | - | - | - | - | -- |
| Igneous rocks | 22.00 | VLS | 0.00 | 0.00 | 0.00 | 0.00 |
| | | LS | 27.27 | 31.82 | 0.00 | 19.70 |
| | | MS | 31.82 | 36.36 | 50.00 | 39.39 |
| | | HS | 27.27 | 27.27 | 22.73 | 25.76 |
| | | VHS | 13.64 | 4.55 | 27.27 | 15.15 |
| Limestones, dolomite, marlite and volcanic formations | 36.62 | VLS | 0.88 | 0.88 | 0.88 | 0.88 |
| | | LS | 23.89 | 23.01 | 13.27 | 20.06 |
| | | MS | 38.05 | 37.17 | 16.81 | 30.68 |
| | | HS | 21.24 | 23.01 | 36.28 | 26.84 |
| | | VHS | 15.93 | 15.93 | 32.74 | 21.53 |
| Phyllite and shale | 29.86 | VLS | 0.00 | 0.00 | 0.76 | 0.25 |
| | | LS | 9.16 | 4.58 | 10.69 | 8.14 |
| | | MS | 17.56 | 32.06 | 22.90 | 24.17 |
| | | HS | 34.35 | 29.01 | 36.64 | 33.33 |
| | | VHS | 38.93 | 34.35 | 29.01 | 34.10 |
| Mudstone formations | 8.21 | VLS | 0.00 | 0.00 | 0.00 | 0.00 |
| | | LS | 0.00 | 0.00 | 0.00 | 0.00 |
| | | MS | 20.69 | 17.24 | 37.93 | 25.29 |
| | | HS | 37.93 | 48.28 | 41.38 | 42.53 |
| | | VHS | 41.38 | 34.48 | 20.69 | 32.18 |

**Notes:** LSZ means landslide susceptibility zonation; LR means Logistic Regression; ANN means Artificial Neural Networks; SVM means Support Vector Machine.

Yesilnacara and Topal (2005). However, in this study case and in Yilmaz (2009) and Yesilnacara and Topal (2005), the accuracy differences between ANN and LR algorithms are not prominent. The qualitative interpretations also do not show distinct differences between ANN and LR. The operator-dependence of ANN is another reason that this study would rather avoid drawing any confirmative conclusions regarding the performance of LR and ANN.

The model stability analyses shows that all the three data-driven algorithms can be affected by varying sizes of training set. Amongst the three algorithms, SVM is the most accurate and stable algorithm at 30% or more of the training set size. Therefore, it is a very reliable and practical method. The ANN algorithm is found to be the least steady algorithm, as it suffers from pronounced fluctuations in the accuracy curves with varying

The superiority of SVM in stability tests can be attributed to its structural-risk minimum principle in a mathematical context. This principle takes the empirical risk, the size of the training set and the complexity of the problem into consideration, while trying to minimize them as a whole. The lowest accuracy of the LR result can be mainly attributed to its less strong capability of dealing with high-dimensional nonlinear data. The insufficiency of landslide samples is another potential cause. As the differences of accuracies and stabilities between LR and ANN algorithms are not prominent, it might be too opportunistic to conclude whether LR or ANN performs better. The differences in physiography of the study area, the qualities of environmental variables and potential data uncertainty will all have an effect on the differences in algorithm performance. The operator-dependence of ANN is hard to capture,

which naturally raises doubts about whether the network output is correct. However, in this study case, the SVM algorithm achieves the best results with limited samples and is less sensitive to the size of the training set. For large, landslide-prone areas in the southwest of China, and in the absence of sufficient landslide samples, the SVM algorithm is the most suitable option.

The application and comparison of the three statistical learning algorithms in the LSZ is conducted in a large area with limited landslide data, which reflects the real situation in the southwestern mountainous less-developed areas. As a result, some simplifications had to be made, for example, only to consider landslide occurrence. This was done, because the scale-dependent relationship between landslide types and environmental variables would complicate the parameterization and would decrease the number of available samples for each type of landslide in such a large study area. The dependence of conditioning factors on the nature of landslides has recently been approached in the study conducted by Micheletti et al. (2013).

Another issue is the absence of a detailed landslide inventory, therefore a point-based sampling strategy is applied, conforming to the location records of landslides. As a result of this point-based approach, each landslide location corresponds to a cell recorded as a landslide. The auto-correlation of spatial data is reduced, as there are no two samples located on the same landslide. Still, the total dataset is very limited if compared to the large study area. A consequence is the occurrence of heterogeneous and scattered points in the LSI map, which can of course be filtered out. To reduce this effect, sampling strategy adaption and influence of spatial autocorrelation of sample data can be reduced (Micheletti et al. 2013; Yilmaz 2010b; Nefeslioglu et al. 2008).

The performance of statistical learning algorithms under varying sizes of training set is crucial to obtain correct answers to the problem domain; in this case landslide susceptibility. Usually, more landslides samples contain more information, and a learning algorithm might theoretically work better with enhanced information. In contrast, learning algorithms might perform worse, due to increased complexity of the problem. Random sampling of the entire training set into subsets also influences model performance. The effects of varying sizes of the training set are probably a combination of the size and the composition of the data. Although it is difficult to discriminate between these two effects, the fluctuations in LR and ANN demonstrate that the increase in sample sizes does not necessarily promote better performance. Steadily increasing performance of the SVM affirms its strong learning capability.

## 7 Conclusion

Three LSZ maps derived from LR, ANN and SVM algorithms on a regional scale, using nine environmental input parameters are evaluated based on LDA, ROC/AUC and *Kappa* index and tested with varying sizes of the training set. We conclude that the SVM algorithm is the best performing, followed by the ANN and LR algorithms, for preparing LSZ maps in large area such as Hanzhong City, with limited field-based landslide information.

We further conclude that SVM is the most stable algorithm, followed by LR and ANN under varying sizes of the training set. The accuracy of SVM shows a steadily increasing trend and reaches a high level at a small size of the training set (30%), while accuracies of LR and ANN algorithms show distinct fluctuations.

From a geomorphological perspective, SVM is the best performing model on all sizes of landslides, as it classifies the lowest landslide densities into the VLS and LS categories and the highest landslide densities into VHS and HS zones. The LR classifies higher percentages of landslides as HS and VHS zones than the ANN, across all sizes of landslide. Whereas, the ANN algorithm classifies a lower percentage of landslides as VLS and LS zones than the LR for medium-sized landslides.

Overall, SVM is regarded to be the best performing model for LSZ analysis on a regional scale, given the same environmental parameters and landslide samples. For very large areas where only limited landslide samples are available and various sizes of landslides occur, SVM algorithm is potentially the most optimal choice for its strong learning capability, compared to ANN and LR algorithms.

## Acknowledgement

## References

Abdallah C, Chorowicz J, Bou Kheir R, et al. (2005) Detecting major terrain parameters relating to mass movements' occurrence using GIS, remote sensing and statistical correlations, case study Lebanon. Remote Sensing of Environment 99: 448-461.

Aleotti P, Chowdhurry R (1999) Landslide hazard assessment: summary review and new perspectives. Bulletin of Engineering Geology and the Environment 58: 21-44.

Anbalagan R (1992) Landslide susceptibility evaluation and zonation mapping in mountainous terrain. Engineering Geology 32: 269-277.

Bai S, Wang J, Lv G, et al. (2007) GIS-based and data drive bivariate landslide susceptibility mapping in the Three Gorge Area, China. Journal of Mountain Science 25(1): 85-92. (In Chinese).

Bai SB, Wang J, Lü G-N, et al. (2010) GIS-based logistic regression for landslide susceptibility mapping of the Zhongxian segment in the Three Gorges area, China. Geomorphology 115: 23-31.

Balabin RM and Lomakina EI (2011) Support vector machine regression (SVR/LS-SVM)—an alternative to neural networks (ANN) for analytical chemistry? Comparison of nonlinear methods on near infrared (NIR) spectroscopy data. Analyst 136: 1703-1712.

Bălteanu D, Chendes V, Sima M, et al. (2010) A country-wide spatial assessment of landslide susceptibility in Romania (in press). Geomorphology.

Basheer IA, Hajmeer M (2000) Artificial neural networks: fundamentals, computing, design and application. Journal of Microbiological Methods 43: 3-31.

Bishop CM (1995). Neural Networks for Pattern Recognition. Oxford University Press, Oxford, UK.

Brabb EE, Pampeyan EH, Bonilla M (1972) Landslide susceptibility in the San Mateo County, California, scale 1:62,500. USGS Miscellaneous Field Studies Map, MF344.

Brenning A, Trombotto D (2006) Logistic regression modeling of rock glacier and glacier distribution: topographic and climatic controls in the semi-arid Andes. Geomorphology 81: 141-154.

Bui DT, Pradhan B, Lofman O, et al. (2012) Landslide susceptibility assessment in Vietnam using suport vector machines, decision tree, and naïve bayes models. Mathematical problems in Engineering: 1-26.

Carrara A, Cardinali M, Detti R, et al. (1991) GIS techniques and statistical models in evaluating landslide hazard. Earth Surface Processes and Landforms 16: 427-445.

Cascini L (2008) Applicability of landslide susceptibility and hazard zoning at different scales. Engineering Geology 102: 164-177.

Chang CC, Lin CJ (2011) LIBSVM: a library for support vector machines. ACM transactions on Intelligent Systems and Technology 2(27): 1-27.

Chauhan S, Sharma M, Arora MK, et al. (2010) Landslide susceptibility zonation through ratings derived from Artificial Neural Network. International Journal of Applied Earth Observation and Geoinformation 12: 340-350.

Chong X, Dai FC, Xu XW, et al. (2012) GIS-based support vector machine modelling of earthquake-triggered landslide suscpetibility in the Jianjiang River watershed, China. Geomorphology 145: 70-80.

Clerici A, Perego S, Tellini C, et al. (2002) A procedure for landslide susceptibility zonation by the conditional analysis method. Geomorphology 48: 349-364.

Congalton RG (1991) A review of assessing the accuracy of classifications of remotely sensed data. Remote Sensing of Environment 37: 35-46.

Dai FC, Lee CF, Nagi Y (2002) Landslide risk assessment and management: an overview. Engineering Geology 64: 65-87.

DeGraff J, Romesburg H (1980) Regional landslide-susceptibility assessment for wildland management: a matrix approach. In: Coates D, Vitek J (Eds.) Thresholds in geomorphology. George Allen and Unwin, London, UK. pp 401-414.

DLR (Shaanxi) (2003) Atlas of geological hazard in Shaanxi Province (Hanzhong Fascicule). Department of Land and Resources of Shaanxi Province, Xi'an, China.

Ercanoglu M, Gokceoglu C, Van Asch TWJ (2004) Landslide susceptibility zoning north of Yenice (NW Turkey) by multivariate statistical techniques. Natural Hazards 32: 1-23.

Ermini L, Catani F, Casagli N (2005) Artificial neural networks applied to landslide susceptibility assessment. Geomorphology 66: 327-343.

Fell R, Corominas J, Bonnard C, et al. (2008) Guidlines for landslide susceptibility, hazard, risk zoning for land-use planning. Engineering Geology 102: 99-111.

Fieldings AH, Bell JF (1997) A review of methods for the assessment of prediction errors in the conservation presence/absence methods. Environmental Conservation 24(1): 38-49.

Frattini P, Crosta G, Carrara A (2010) Techniques for evaluating the performance of landslide susceptibility models. Engineering Geology 111: 62-72.

George D, Mallery P (2000). SPSS Windows Step by Step: A Simple Guide and Reference. Allyn and Bacon, New York, NY, USA.

Gong P (1996) Integrated analysis of spatial data for multiple sources: using evidential reasoning and artificial neural network techniques for geological mapping. Photogrammetric Engineering and Remote Sensing 62(5): 513-523.

Gupta RP, Kanungo DP, Arora MK et al. (2008) Approaches for comparative evaluation of raster GIS-based landslide susceptibility zonation maps. International Journal of Applied Earth Observation and Geoinformation (10): 330-341.

Haykin S (1998) Neural Networks: A Comprehensive Foundation. Prentice Hall, London, UK.

Holec J, Bednarik M, Sabo M, et al. (2013) A small-scale landslide susceptibility assessment for the territory of Western Carpathians. Natural Hazards 69 (1): 1081-1107.

Jade S, Sarkar S (1993) Statistical models for slope instability classification. Engineering Geology 36: 91-98.

Kanungo DP, Arora MK, Starker S, et al. (2006) A comparative study of conventional, ANN black box, fuzzy and combined neural and fuzzy weighting procedures for landslide susceptibility zonation in Darjeeling Himalayas. Engineering Geology 85(3-4): 347-366.

Lan HX, Zhou CH, Wang LJ, et al. (2004) Landslide hazard spatial analysis and prediction using GIS in the Xiaojiang watershed, Yunnan, China. Engineering Geology 76: 109-128.

Lee S, Ryu J, Min K, et al. (2003) Landslide susceptibility analysis using GIS and artificial neural network. Earth Surface Processes and Landforms 28: 1361-1376.

Lillesand TM, Kiefer RW, Chipman JR (2008) Remote Sensing and Image Interpretation, Sixth Edition. John Wiley & Sons Inc., Beijing, China.

Liu MG, Ed. (1997) Atlas of Physical Geography in China. SinoMaps Press, Beijing, China.

Ma LF, Ed. (2002) Atlas of China Geology. Geology Publishing House, Beijing, China.

Marjanović M, Kovačević M, Bajat B, et al. (2011) Landslide susceptibility assessment using SVM machine learning algorithm. Engineering Geology 123: 225-234.

Micheletti N, Foresti L, Robert S, et al. (2013) Machine learning feature selection methods for landslide susceptibility mapping. Mathematical Geosciences 46: 33-57.

Monserud RA, Leemans R (1992) Comparing global vegetation maps with the Kappa statistics. Ecological Modelling 62: 275-293.

Nefeslioglu HA, Gokcegolu C, Sonmes H (2008) An assessment on the use of logistic regression and artificial neural networks with different sampling strategies for the preparation of landslide susceptibility maps. Engineering Geology 97: 171-191.

Pradhan B (2010) Remote sensing and GIS-based landslide hazard analysis and cross-validation using multivariate logistic regression model on three test areas in Malaysia. Advances in Space Research 45: 1244-1256.

Pradhan B (2013) A comparative study on the predictive ability of the decision tree, support vector machine and neuro-fuzzy models in landslide susceptibility mapping using GIS. Computers & Geosciences 51: 350-365.

Pradhan B, Lee S (2007) Utilization of optical remote sensing data and GIS tools for regional landslide hazard analysis using an artificial neural network model. Earth Science Frontiers 14(6): 143-152.

Pradhan B, Lee S (2010) Landslide susceptibility assessment and factor effect analysis: backpropagation artificial neural networks and their comparison with frequency ratio and bivariate logistic regression modelling. Environmental Modelling & Software 25: 747-759.

Pradhan B, Lee S, Buchroithner MF (2010) A GIS-based back-propagation neural network model and its cross-application and validation for landslide susceptibility analysis. Computers, Environment and Urban Systems (34): 216-235.

Rossi M, Guzzetti F, Reichenbach P, et al. (2010) Optimal landslide susceptibility zonation based on multiple forecasts. Geomorphology 114(3): 129-142.

Saha AK, Gupta RP, Starker I, et al. ( 2005) An approach for GIS based statistical landslide susceptibility zonation—with a case study in the Himalayas. Landslides 2: 61-69.

Samui P, Kothari DP (2010) Utilization of a least square support vector machine (LSSVM) for slope stability analysis. Scientia Iranica 18(1): 53-58.

Tax D, Duin E (1999) Support vector domain description. Pattern Recognition Letter 20: 1191-1199.

Vapnik VN (1995). The Nature of Statistical Learning Theory. Springer, New York, NY, USA.

Wang H, Liu G, Xu W, et al. (2005) GIS-based landslide hazard assessment: an overview. Progress in Physical Geography 29: 548-567.

Wang WD, Cui CM, Du XG (2009) Landslides susceptibility mapping in Guizhou province based on fuzzy theory. Mining Science and Technology (19): 399-404.

Yalcin A, ReiS S, Aydinoglu AC, et al. (2011) A GIS-based comparative study of frequency ratio, analytical hierarchy process, bivariate statistics and logistics regression methods for landslide susceptibility mapping in Trabzon, NE Turkey. Catena 85: 274-287.

Yao X, Tham LG, Dai FC (2008) Landslide susceptibility mapping based on Support Vector Machine: a case study on natural slopes for Hongkong, China. Geomorphology 101: 572-582.

Yesilnacara E, Topal T (2005) Landslide susceptibility mapping: A comparison of logistic regression and neural networks methods in a medium scale study, Hendek region (Turkey). Engineering Geology 79: 251-266.

Yilmaz I (2009) Landslide susceptibility mapping using frequency ratio, logistic regression, artificial neural networks and their comparison: A case study from Kat landslides (Tokat—Turkey). Computers & Geosciences 35: 1125-1138.

Yilmaz I (2010a) Comparison of landslide susceptibility mapping methodologies for Koyulhisar, Turkey: conditional probability, logistic regression, artificial neural networks, and support vector machine. Environmental Earth Sciences 61: 821-836.

Yilmaz I (2010b) The effect of the sampling strategies on the landslide susceptibility mapping by conditional probability and artificial neural networks, Environmental Earth Science 60(3): 505-519.

Yilmaz I, Keskin I (2009) GIS based statistical and physical approaches to landslide susceptibility mapping (Sebinkarahisar, Turkey). Bulletin of Engineering Geology and the Environment 68 (4): 459-471.

Zhang WB, Jin SL, Eds. (1995). Atlas of China. SinoMaps Press, Beijing, China.