

2.6 Statistical Methods for Geomorphic Distribution Modeling

J Hjort, University of Oulu, Oulu, Finland

M Luoto, University of Helsinki, Helsinki, Finland

© 2013 Elsevier Inc. All rights reserved.

2.6.1	Introduction	60
2.6.2	Modeling Steps	62
2.6.3	Review of Statistical Methods	63
2.6.3.1	Generalized Linear Model	63
2.6.3.1.1	Technical background	63
2.6.3.1.2	GLM in geomorphology	64
2.6.3.1.3	Strengths and weaknesses of GLM	64
2.6.3.2	Generalized Additive Model	65
2.6.3.2.1	Technical background	65
2.6.3.2.2	GAM in geomorphology	65
2.6.3.2.3	Strengths and weaknesses of GAM	65
2.6.3.3	Artificial Neural Network	66
2.6.3.3.1	Technical background	66
2.6.3.3.2	ANN in geomorphology	66
2.6.3.3.3	Strengths and weaknesses of ANN	67
2.6.3.4	Statistical Boosting and BRT	67
2.6.3.4.1	Technical background	67
2.6.3.4.2	Boosting in geomorphology	68
2.6.3.4.3	Strengths and weaknesses of statistical boosting	68
2.6.4	SWOT Analysis of Statistical Modeling in Geomorphology	68
2.6.4.1	Strengths and Opportunities	68
2.6.4.2	Weaknesses and Threats	69
2.6.5	Future Challenges	70
References		71

Glossary

- Artificial neural network** A computational (or mathematical) model that tries to simulate the structure and/or functional aspects of biological neural networks.
- Boosted regression tree** A statistical ensemble method that combines machine learning and regression tree approaches.
- Calibration, statistical model** A selection of predictor variables and a construction (e.g., parameter estimation) of a statistical model.
- Evaluation, statistical model** Assessment of the realism of fitted response functions and predictor variables, model's fit to data, characteristics of residuals, and predictive performance on test data.
- Formulation, statistical model** A choice of a proper statistical approach and a suitable algorithm for modeling a particular type of response variable.
- Generalized additive model** A semiparametric extension of generalized linear model; the only underlying

assumption made is that the functions are additive and that the components are smooth.

Generalized linear model An extension of ordinary least squares regression model that allow for nonlinearity and nonconstant variance structures in the data.

Geomorphic distribution model Empirical/numerical model relating geomorphic (field) observations to predictor variables (i.e., environmental variables).

Multicollinearity A case of multiple regression in which the predictor variables are themselves highly correlated.

Over-fit, statistical model Over-fitted models include too many predictors, are exceedingly complex, and may begin to fit random noise in the data.

Spatial autocorrelation Spatial autocorrelation occurs when values of a variable sampled at nearby locations are more similar than those sampled at locations more distant from each other.

Hjort, J., Luoto, M., 2013. Statistical methods for geomorphic distribution modeling. In: Shroder, J. (Editor in Chief), Baas, A.C.W. (Ed.), Treatise on Geomorphology. Academic Press, San Diego, CA, vol. 2, Quantitative Modeling of Geomorphology, pp. 59–73.

Abstract

Statistically based geomorphic distribution modeling (GDM) has become popular among geoscientists as an efficient approach for analysis and prediction. Here, we provide a cross section of the concept of GDM. First, we introduce the main steps in the GDM process. Second, we provide an overview of statistical techniques, which have shown to be promising in geomorphic modeling. Third, we draw attention to important advantages and pitfalls of GDM. Finally, we highlight some future challenges in the application of the GDM approach. The general aim is to aid the geomorphic community to gain novel insights into Earth surface process–environment relationships using the concept of GDM.

2.6.1 Introduction

Determination of the environmental factors controlling Earth surface processes and landform patterns is one of the central themes in physical geography (e.g., Goudie, 1995). However, the identification of the main drivers of geomorphic processes is characteristically challenging, particularly if complex, multivariate systems are under investigation. In recent years, statistically based geomorphic distribution models (GDMs) with geographic information (GI) and remote-sensing (RS) data have become more popular among geoscientists as an efficient approach for analysis and prediction (Carrara and Pike, 2008; Harris et al., 2009; Remondo and Oguchi, 2009).

GDMs are empirical models relating field observations to explanatory variables (i.e., predictor variables and environmental variables), based on statistically or theoretically derived response surfaces. Geomorphic data can be simple presence, presence-absence, or abundance (e.g., cover and activity) observations on landforms, processes, or feature assemblages (Guzzetti et al., 1999; Hjort, 2006). Environmental variables can be acquired from various sources, commonly from GI and RS data (Moore et al., 1991; Etzelmüller et al., 2001). Environmental variables can exert direct (causal variable) or indirect (noncausal variable) effects on geomorphic features (e.g., Ayalew and Yamagishi, 2005; cf. Austin, 2002).

GDMs can be used to simplify complex systems (model reduction), to provide understanding of process–environment

relationships (explanatory models), and to predict distributions not only across space, but also in time (predictive models). Model simplifications utilize variable reduction techniques in the analytical phase and have as their goal a model that explains and/or predicts the occurrences of geomorphic features with a restricted number of explanatory variables (Figure 1). The concept of parsimony, that the simplest explanation is best, is intrinsic in such a modeling approach. Explanatory models seek to provide insights into the geomorphic processes and physical conditions that determine the distribution of landforms and processes (Figure 2). By contrast, predictive models typically seek to provide the user with a statistical relationship between the response and a series of explanatory variables for use in predicting the feature occurrence or estimating abundance of geomorphic features at new, previously unmapped areas (Luoto and Hjort, 2005; Figure 3).

GDMs are important tools in mapping remote regions (Guzzetti et al., 1999; Etzelmüller et al., 2006), analyzing processes across scales (Luoto and Hjort, 2006), predicting hazards (McKillop and Clague, 2007; Carrara and Pike, 2008), and exploring the potential consequences of climate change on Earth surface processes and landforms (Guzzetti et al., 2005; Fronzek et al., 2006, 2010). Moreover, they can be used to model feature assemblages (e.g., geodiversity and geomorphic process units) and to identify the shapes of the geomorphic process–environment relationships. In general, GDMs provide a mathematical basis for the interpretation of relationships between response and explanatory variables and

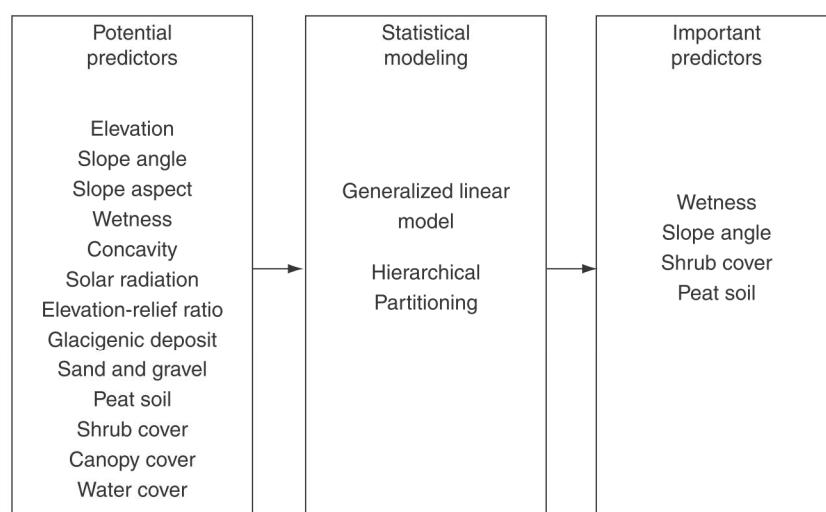


Figure 1 Important predictors (i.e., explanatory variables) for the distribution of cryoturbation features were drawn from a set of potentially important predictors (Hjort, 2006). The analyses were conducted in northern Finland at a mesoscale resolution (grid cell size = 500 × 500 m).

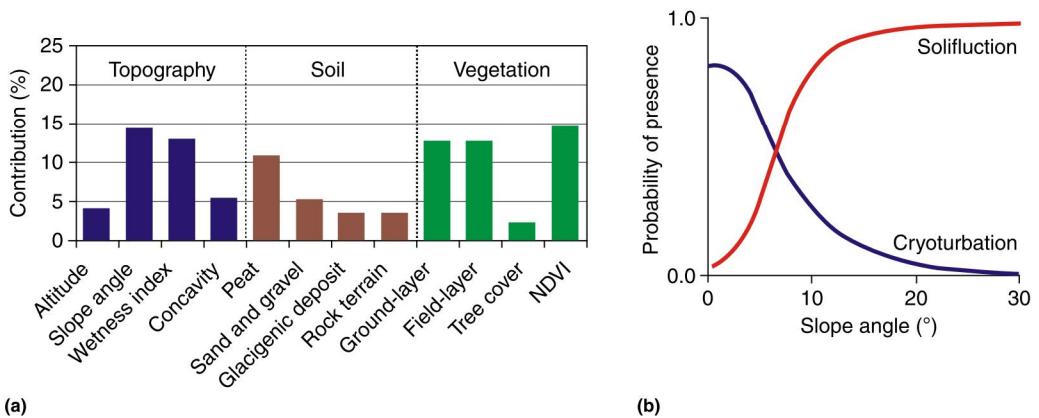


Figure 2 (a) Contribution of various explanatory variables for the distribution of active cryoturbation features in alpine areas in northernmost Finnish Lapland. (b) Simulated generalized additive modeling (GAM) based response curves for the distribution of solifluction and cryoturbation in northernmost Finnish Lapland. (a) Modified from Hjort, J., Luoto, M., 2009. Interaction of geomorphic and ecologic features across altitudinal zones in a subarctic landscape. *Geomorphology* 112, 324–333. (b) Modified from Hjort, J., Luoto, M., 2011. Novel theoretical insights into geomorphic process-environment relationships using simulated response curves. *Earth Surface Processes and Landforms* 36, 363–371.

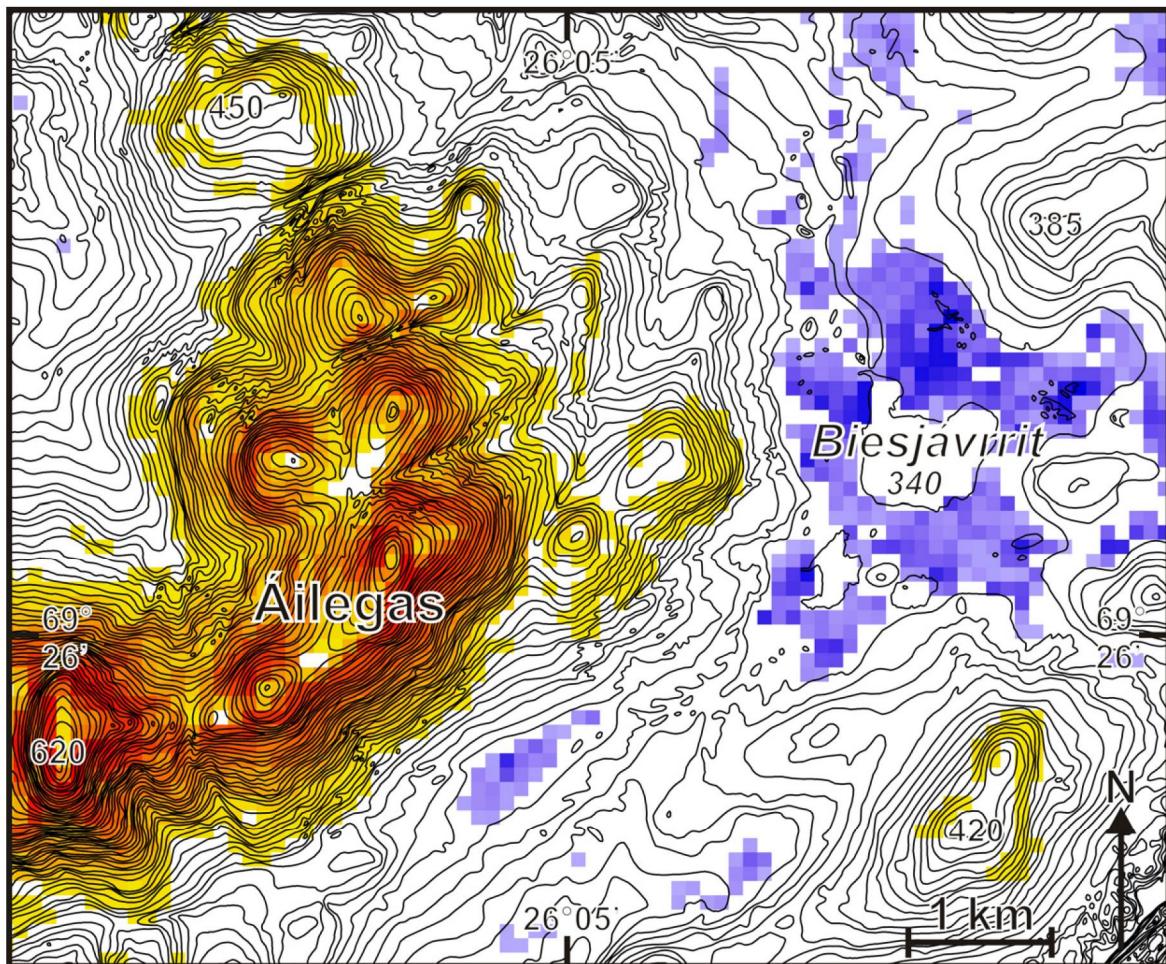


Figure 3 A prediction map based on a generalized linear model (GLM) on the occurrence of solifluction (from yellow to red) and permafrost features, palsas, (blue) in northern Finland. The darker the color within a grid square, the higher the likelihood of an occurrence. The vertical interval of the contours is 5 m (copyright National Land Survey of Finland 2010).

they provide a valuable source of information for generating and testing scientific hypotheses (Hjort and Luoto, 2011).

Here, we provide a cross section of the concept of statistically based distribution modeling in geomorphology. First, we briefly introduce the modeling steps in GDM. Second, we provide an overview of four statistical techniques, which includes two extensions of linear regression (generalized linear model, GLM, and generalized additive model, GAM), one machine learning technique (artificial neural networks, ANNs), and one method combining both machine learning and regression approaches (boosted regression tree, BRT). Third, we draw attention to the most important advantages and pitfalls of GDM through a SWOT analysis where the strengths, weaknesses, opportunities, and threats are discussed briefly. Finally, we highlight some future challenges in the use and application of GDM.

Statistical modeling of Earth surface processes and landforms can be approached in various ways (e.g., prediction or explanation, univariate or multivariate analysis, and spatial or nonspatial modeling). Moreover, the utilized techniques vary depending on the focus of the study (e.g., geostatistics, least square regression, and machine learning). It is important to note that our aim is not to cover all the statistical concepts and techniques for geomorphic modeling. For example, geostatistical techniques such as kriging and variogram modeling are beyond the scope of this chapter (see Cressie, 1993; Goovaerts, 1999; Bivand et al., 2008). Moreover, this is not a comprehensive literature review of statistical geomorphic studies. Consequently, we focus on issues and methods that have and hopefully will aid the geomorphic community to gain novel insights into the process–environment relationships and landscape development on the Earth but potentially on other planets as well. In future, the focus of statistically based modeling should shift from description and spatial prediction to an emphasis on explanation and hypothesis testing.

2.6.2 Modeling Steps

Key steps in sound GDM practice include the following (Figure 4; Guisan and Zimmermann, 2000; Hjort and Mar-mion, 2008; Elith and Leathwick, 2009): (1) establishment of a conceptual model and setting relevant study questions, (2) gathering geomorphic (response) and environmental (explanatory variable) data, (3) data exploration, (4) statistical formulation, (5) model calibration, (6) model evaluation, and (7) prediction and/or interpretation of the results.

First, a conceptual model based on solid geomorphic theory should be proposed before a statistical model is even considered. This is extremely important because in addition to the study problem, the conceptual framework outlines the subsequent steps in data collection and modeling.

Second, the compilation of response data and selection of appropriate explanatory variables for the statistical modeling can be a complicated and difficult task without a firm conceptual model. There are neither universal criteria nor widely accepted guidelines for the selection of explanatory variables and hence the study aims guide the procedure (Ayalew and Yamagishi, 2005). Commonly, the variables are gathered from field work, digital and paper maps, RS, maps obtained from

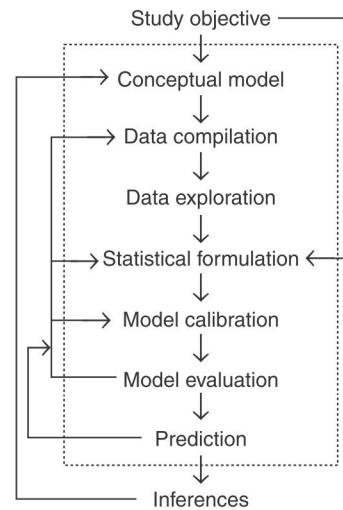


Figure 4 Schematic presentation of the key steps of statistically based geomorphic distribution modeling (GDM).

GI-based modeling, and other digital databases (Moore et al., 1991; Etzelmüller et al., 2001). A spatial grid-based approach is commonly used as the data compilation system (e.g., Guzzetti et al., 1999). An advantage of the grid approach is the possibility to convert commonly fuzzy spatial variables to numeric form enabling numerical analysis and the possibility to utilize GI and RS data as a source of explanatory variables (e.g., Figure 3).

Third, assessing the adequacy of the data in explorative analysis should not be overlooked (Ott and Longnecker, 2010; Section 2.6.4.2). For example, frequency-size distribution (e.g., normality), prevalence, abundance, and spatial properties (e.g., scale, autocorrelation, and trends) of the response data should be considered (e.g., Bivand et al., 2008). Scatter plots, correlation analysis, and geographical plots are useful in the exploration of the relationships between environmental variables as well as responses and environmental variables (e.g., Reimann et al., 2008).

Fourth, statistical formulation means the choice of a proper statistical approach with regard to the modeling context and a suitable algorithm for modeling a particular type of response variable and estimating the model coefficients. In addition to the explorative analysis, previous studies are used to guide this stage (e.g., Sokal and Rohlf, 1995; Crawley, 2007).

Fifth, the environmental variables are selected to the final model and the statistical model is constructed (e.g., estimation of model parameters) in model calibration. Traditionally, the model selection has been based on *p*-values, but a recent shift has seen much greater emphasis on Akaike's information criterion (or related information theories) and multimodel inference (Burnham and Anderson, 2002). This shift is seen to be useful for reducing reliance on models selected by stepwise approaches and for understanding the error tendencies of conventional selection approaches (Whittingham et al., 2006; Elith and Leathwick, 2009).

Sixth, evaluation of the generated model is a vital step in the model building process (Oreskes et al., 1994). Evaluating the model includes the assessment of the realism of fitted

response functions and explanatory variables, the model's fit to data, characteristics of residuals, and predictive performance on test data (e.g., Sokal and Rohlf, 1995; Ott and Longnecker, 2010). For predictive purposes, it is advisable to assess the model performance using spatially independent evaluation data. However, this generally is not possible due to the data constraints. Thus, cross-validation and split-sample approaches (data are split to separate calibration and evaluations sets) are often used (Venables and Ripley, 2002; Crawley, 2007). The final stage includes mapping predictions to geographical space and/or iterating the process to improve the model in light of knowledge gained throughout the process, or the modeling outcomes can directly be used to draw conclusions. All the above stages are interconnected and ultimately controlled by the objectives of the study.

2.6.3 Review of Statistical Methods

Statistical modeling increased its popularity among earth scientists when: (1) the techniques permitted more liberties related to the data used; and (2) new methods and systems (e.g., geographic information system, GIS) allowed robust and detailed preparation of digital models of the Earth's surface properties, interpolation of climate parameters, and RS of surface conditions (e.g., Guzzetti et al., 1999). More precisely, multivariate statistical modeling of geomorphic features gained attention from the late 1970s (Carrara, 1983 and references therein). At the beginning, regression and discriminant analysis were the most common approaches (e.g., Neuland, 1976; Carrara, 1983; Carrara et al., 1991). In the late 1980s and early 1990s, GLM increased its popularity as a statistical technique (e.g., Atkinson et al., 1998; Guzzetti et al., 1999 and references therein). After the mid-1990s, machine-learning techniques (e.g., ANNs) were introduced in geomorphology (e.g., Lees, 1996; Aleotti and Chowdhury, 1999; Guzzetti et al., 1999) and opened up new possibilities for modeling complex and multivariate features.

In the twenty-first century, the diversity of statistical techniques used to study Earth surface processes and landforms has exploded. The list of techniques range from traditional least square (LS) regression methods to highly advanced machine-learning techniques (Brenning, 2005; Melchiorre et al., 2008;

Marmion et al., 2009; Rossi et al., 2010). Here, we consider four methods, namely GLMs (McCullagh and Nelder, 1989), GAMs (Hastie and Tibshirani, 1990), statistical boosting and especially BRTs (Friedman et al., 2000), and ANNs (Ripley, 1996) (Table 1). These methods were selected because they have shown to be highly promising methods in various fields of physical geography (Guisan and Zimmermann, 2000; Guisan and Thuiller, 2005; Luoto and Hjort, 2005; Elith et al., 2006; Heikkilä et al., 2006; Marmion et al., 2008, 2009).

2.6.3.1 Generalized Linear Model

2.6.3.1.1 Technical background

An important statistical development of the last decades has been the advance in regression analyses provided by GLMs (e.g., Nelder and Wedderburn, 1972). GLM is more flexible and better suited for analyzing geomorphic relationships than the linear LS regression method that has implicit statistical assumptions (Sokal and Rohlf, 1995). Technically, GLMs are close to linear regressions and thus relatively easy to utilize.

GLMs are mathematical extensions of linear models that allow for nonlinearity and nonconstant variance (heteroscedasticity) structures in the data (McCullagh and Nelder, 1989). GLMs have three components: (1) the response variables Y_1, Y_2, \dots, Y_n , which are assumed to share the same distribution from the exponential family; (2) a set of parameters α and β and explanatory variables; and (3) a link function g , which allows transformation to linearity and the predictions to be maintained within the range of coherent values for the response variable (McCullagh and Nelder, 1989).

For GLMs, we have data (Y_i, x_i) ($i = 1, 2, \dots, n$) where n is the number of observations and $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$ is a vector of p explanatory variables. The mean of the response variable at $X = x$, namely, $\mu_i = \mu_i(x) = E(Y_i)$, is related to the covariate information by

$$g(\mu) = \alpha + \beta^T x = \alpha + \sum_{j=1}^p \beta_j x_j \quad [1]$$

where α is the constant (i.e., intercept) and $\beta = (\beta_1, \beta_2, \dots, \beta_p)^T$ is a vector of regression coefficients.

In GLMs, the model is formulated through deviance reduction that is comparable to LS model's variance reduction. However, the regression coefficients of the model cannot be

Table 1 Summary of the strengths and weaknesses of generalized linear models (GLMs), generalized additive models (GAMs), artificial neural networks (ANNs), and statistical boosting in geomorphic modeling

	GLMs	GAMs	ANNs	Boosting
Flexibility	Moderate	Rather high	High	High
Data requirements				
- Number of observations	Rather low	Moderate	High	Rather high
- Statistical assumptions	Rather high	Moderate	Low	Low
Expert knowledge	Rather low	Moderate	High	High
Over-fitting risk	Rather low	Moderate	High	Rather high
Model interpretability	High	Rather high	Low	Rather low
Usability in explanation	High	Rather high	Low	Low
Usability in prediction	Rather high	High	High	High

estimated with the ordinary LS method. Instead, maximum likelihood techniques, where the estimation method maximizes the log-likelihood function, are used to calculate these parameters. Further information about GLMs can be found in McCullagh and Nelder (1989) and Dobson (2002).

2.6.3.1.2 GLM in geomorphology

GLMs have been utilized in various fields of geomorphology, although most of the examples are from either slope hazard and landslide studies or from periglacial geomorphology. Especially, logistic-regression models have shown to be very useful in geomorphology. This is because response variables are generally in a binary form (geomorphic feature present/geomorphic feature absent). In slope hazard and landslide analysis, the number of studies that have applied GLMs is huge. For example, Atkinson and Massari (1998), Rowbotham and Dudycha (1998), Dai and Lee (2002), Ayalew and Yamagishi (2005), Van Den Eeckhaut et al. (2006), and Das et al. (2010) have used GLMs in slope process studies. Examples in periglacial geomorphology include Luoto and Seppälä (2002), Lewkowicz and Ednie (2004), Luoto and Hjort (2004), Brenning and Trombotto (2006), Hjort et al. (2007), and Brenning and Azócar (2010). Moreover, GLMs have been used in glacial (Atkinson et al., 1998), fluvial (e.g., Bledsoe and Watson, 2001; McKillop and Clague, 2007), and karst (Lamelas et al., 2008) geomorphology. However, GLMs are rarely used in distribution modeling context in other fields of geomorphology beyond landslide and periglacial research.

2.6.3.1.3 Strengths and weaknesses of GLM

GLMs constitute a more flexible family of methods than traditional LS regression techniques. GLMs handle nonlinear relationships and different types of statistical distributions of

geographical data, such as discrete, categorical, ordinal, and continuous data. Therefore, GLMs provide a useful modeling framework for testing the shapes of the response functions and significance of variables describing environmental gradients. However, the technique and data-related constraints should also be considered. For example, GLMs, which are generalizations of LS regression, assume that all explanatory variables are measured without error. GLMs may also distort inferences about the relative importance of explanatory variables. This is because these approaches do not take into account the intercorrelation of the variables (i.e., the multicollinearity problem, Section 2.6.4.2). Furthermore, spatial autocorrelation can hamper the detection of causal correlates because the presence of autocorrelation may inflate the degrees of freedom in the test of significance (Section 2.6.4.2).

Sometimes GLMs are not flexible enough to capture the shape of the relationships between environmental variables and responses. Even by adding higher polynomial terms (e.g., a cubic term), the approximation may still be inadequate (Figure 5(a)). The solution to the detection of more complex responses can be the utilization of nonparametric methods that allow a wider range of response curves to be modeled (e.g., Luoto and Hjort, 2005). However, nonparametric techniques may have little statistical theory to support them and it is easy to over-fit and over-explain features of the data. In addition, nonparametric methods have been criticized because they can produce very complex model outputs that are difficult to interpret (e.g., Venables and Ripley, 2002). Thus, corresponding parametric functions such as GLMs may capture most of the same variation and have a more realistic (a mathematical formula) explanation (Figure 5(b)). In summary, the strengths of GLMs in geomorphic modeling are related to the resistance to over-fitting and interpretability of the

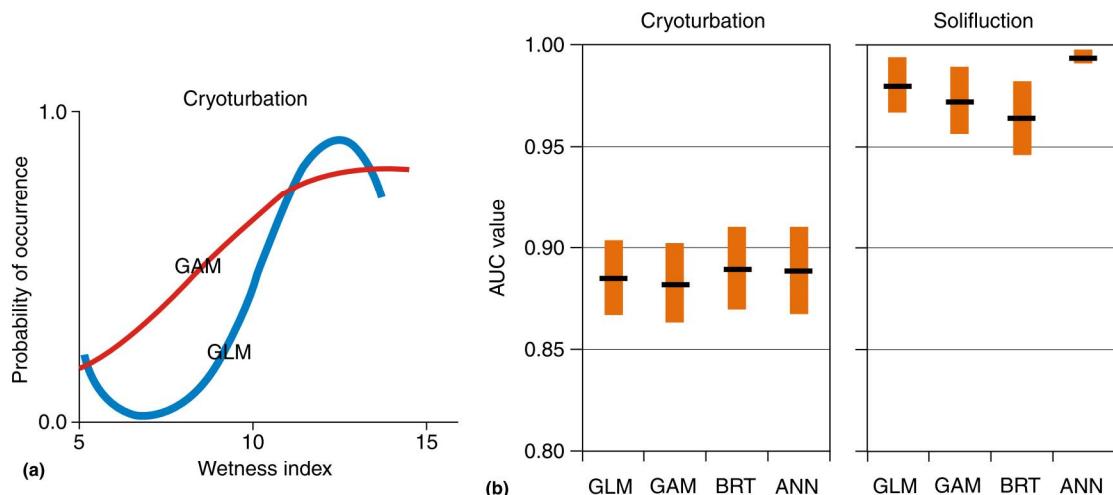


Figure 5 (a) An example where generalized linear modeling (GLM) is not flexible enough to capture the leveling-off effect of the shape of the relationship between wetness index and cryoturbation, whereas generalized additive modeling (GAM) provides a more realistic response shape. (b) The prediction ability of GLM is approximately at the same level as more flexible methods (BRT = boosted regression tree, ANN = artificial neural network) (Hjort and Marmion, 2009). The area under the curve (AUC) values were computed using evaluation data (error bars with 95% confidence intervals). Based on the AUC values, the models can be classified as excellent (0.90–1.00), good (0.80–0.90), fair (0.70–0.80), poor (0.60–0.70), and failed (0.50–0.60) (Swets, 1988). (a) Modified from Hjort, J., Luoto, M., 2011. Novel theoretical insights into geomorphic process-environment relationships using simulated response curves. *Earth Surface Processes and Landforms* 36, 363–371.

results, whereas the inflexibility and lower predictive ability compared to nonparametric techniques may be considered as weaknesses (**Table 1**).

2.6.3.2 Generalized Additive Model

2.6.3.2.1 Technical background

GAMs are semi-parametric extensions of GLMs in which the linear predictor variable is substituted with a smoothing function that can take various forms (Hastie and Tibshirani, 1990). GAMs are parametrized just like GLMs, except that some explanatory variables can be modeled nonparametrically in addition to linear and polynomial terms for other variables. In general, GAMs are designed to capitalize on the strengths of GLMs without requiring the problematic steps of postulating a response curve shape or specific parametric response function (Austin, 2002; Venables and Ripley, 2002). GLMs relate the mean response to the explanatory variables via eqn [1], but GAMs relax this to

$$g(\mu) = \alpha + \sum_{j=1}^p f_j(x_j) \quad [2]$$

where f_j are unspecified smooth functions. In practice, the f_j are estimated from the data by using techniques developed for smoothing scatterplots (e.g., cubic smoothing splines and local polynomial regression) (Hastie and Tibshirani, 1990). Consequently, GAMs are more data driven than their parametric GLM counterparts. This is because the data determine the nature of the relationship between the response and the set of explanatory variables rather than assuming some form of parametric relationship. GAMs can handle any of the data types that GLMs are used for (e.g., Gaussian, binomial, multinomial, and Poisson data) as well as certain types of survival data. The only underlying assumption made is that the functions are additive and that the components are smooth (Hastie and Tibshirani, 1990).

A crucial step in the use of GAMs is the selection of an appropriate level of smoothing for an explanatory variable. The level of smoothing depends on the size of the neighborhood that is used to calculate the smoothed value at a particular point. Small neighborhoods mean that there is little smoothing, whereas large neighborhoods result in (very) smooth curves. Hastie and Tibshirani (1990) and Venables and Ripley (2002) examined various general scatterplot smoothers that can be applied to the explanatory variable values, with the target criterion to maximize the quality of prediction of the response variable values.

A commonly used scatterplot smoother is the cubic spline smoother that minimizes the penalized residual sum of squares. The degree of smoothing is defined by the number of degrees of freedom (d.f.). High number of d.f. means that there is not much smoothing, but the response tracks closely to the data points. A low number of d.f. means much smoothing, at the extreme, one d.f. defining a linear fit. In the geomorphic context, a cubic spline smoother with a maximum of four d.f. is a good starting point (e.g., Hjort and Luoto, 2006; Brenning et al., 2007). Basically, this means that the complexity of the response curve is about the same as a

polynomial regression of degrees 4. However, the cubic spline smoother is much more flexible than a polynomial regression (Hastie and Tibshirani, 1990).

Many of the standard result statistics computed by GAMs are similar to those customarily reported by linear or non-linear model fitting procedures. For example, predicted and residual values for the final model can be computed, and various graphs of the residuals can be displayed to help the user identify, for example, possible outliers. Further information about scatterplot smoothers and fitting GAMs can be obtained from Hastie and Tibshirani (1986, 1990) and Wood (2006).

2.6.3.2.2 GAM in geomorphology

In geomorphology, GAMs are clearly less commonly used when compared with GLMs. For example, Brenning (2008) and Park and Chi (2008) used GAMs in landslide studies and Fronzek et al. (2006), Brenning (2009), and Brenning and Azócar (2010) used GAMs in periglacial geomorphology. In addition, López-Moreno and Nogués-Bravo (2005) applied additive models in snowpack modeling and López-Moreno et al. (2006) used GAMs in glacier studies.

2.6.3.2.3 Strengths and weaknesses of GAM

GAMs are useful in exploratory analysis or when analysts have weak *a priori* ideas as to the functional form relating explanatory variables to response variables (**Table 1**). GAMs are particularly useful to study the shape of the response function (i.e., the relationship between geomorphic feature and environmental variable) (cf. Austin et al., 2006). For example, in the variable selection, GAMs have an advantage over GLMs in that the smoother automatically takes into account the shape of the curve for that variable. Consequently, it is not necessary to choose whether a higher order term should be included, a decision that needs to be made for each case when using a GLM (Yee and Mitchell, 1991). Thus, GAMs may offer certain benefits over GLMs due to their greater flexibility and capacity to reveal more complicated relationships between dependent and environmental variables (Figure 5(a); Austin et al., 2006; Brenning et al., 2007). Still, like in GLMs, the relationships between dependent and explanatory variables are explicit and interpretable. However, in some particular situations (e.g., when there exist sharp discontinuities), the flexibility of GAMs may be inadequate (Elith et al., 2008).

GAMs are more complicated to fit and require greater judgment, and it is possible to over-fit features in the data when compared with GLMs. Over-fitted models include too many predictors, are exceedingly complex, and may begin to fit random noise in the data. Thus, the predictive abilities of over-fitted models are often poor, especially if the models are extrapolated to new data or areas. In addition, the interpretation of the GAMs can be challenging, particularly when they involve complex nonlinear effects of some or all of the explanatory variables. Thus, it is advisable to compare the quality of the fit obtained from GAMs to the fit obtained via GLMs. In other words, evaluate whether the added complexity of GAMs is necessary in order to obtain a satisfactory fit to the data. If the fits are comparable, the simpler GLM is preferable to the more complex GAM. Moreover, models using interaction terms are difficult to build when utilizing basic GAMs.

In addition, GAMs are based on standard regression theory (e.g., Sokal and Rohlf, 1995) and, for example, the effects of measurement error and intercorrelation of the explanatory variables as well as autocorrelation and nonstationarity of the responses should be considered in depth. In the end, even with spatially independent evaluation data, the prediction ability of GAMs is generally higher when compared with parametric techniques (e.g., Marmion et al., 2008, 2009).

2.6.3.3 Artificial Neural Network

2.6.3.3.1 Technical background

An ANN, usually simply called a ‘neural network’, is a computational model that tries to simulate the structure and/or functional aspects of the human brain (Bishop, 1995). The key element of this method is the novel structure of the information processing system. ANNs consist of an interconnected group of artificial neurons, and process information using a connectionist approach to computation (Crawley, 2007). In most cases, an ANN is an adaptive system that changes its structure based on external or internal information that flows through the network during the learning phase (Venables and Ripley, 2002).

An ANN is usually used to model complex relationships between inputs and outputs or to find patterns in data. ANNs can process problems involving very nonlinear and complex data even if the data are imprecise and noisy. They have been shown to be universal and highly flexible function approximators for any data (Smith, 1993). ANNs are adaptive models that can learn from the data and generalize things learned (Melchiorre et al., 2008). They extract the essential characteristics from the numerical data as opposed to memorizing all of it. This offers a convenient way to reduce the amount of data as well as to form an implicit model without having to form a traditional, physical model of the underlying phenomenon (Bishop, 1995; Crawley, 2007).

In contrast to traditional models, which are theory rich and data poor, the ANN is data rich and theory poor in a way that little or no *a priori* knowledge of the problem is present. This makes ANN a powerful tool for modeling, especially when the underlying data relationships are unknown (Lek and Guégan, 1999). ANNs have recently become the focus of much attention, largely because of their wide range of applicability and the ease with which they can treat complicated problems (Ermini et al., 2005; Melchiorre et al., 2008; Oh et al., 2010; Rossi et al., 2010).

Research into ANNs has led to the development of various types of algorithms, suitable to solve different kinds of problems: auto-associative memory, generalization, optimization, data reduction, and prediction tasks in various scenarios (Lek and Guégan, 1999). The descriptions of these methods can be found in various books such as Bishop (1995) and Venables and Ripley (2002). The choice of the type of network depends on the nature of the problem to be solved. At present, two popular ANNs are: (1) multi-layer, feed-forward neural networks trained by a back-propagation algorithm, that is, back-propagation network (BPN); and (2) Kohonen self-organizing mapping (SOM) (Kohonen, 1984).

The BPN, also called multi-layer, feed-forward neural network or multi-layer perceptron, is popular and is used more than other neural network types for a wide variety of tasks (Bishop, 1995; Venables and Ripley, 2002). The BPN is based on the supervised procedure, that is, the network constructs a model based on examples of data with known outputs. It has to build the model up solely from the examples presented, which are together assumed to implicitly contain the information necessary to establish the relation. A BPN is a powerful system, commonly capable of modeling complex relationships between variables (Lek and Guégan, 1999). It allows prediction of an output object for a given input object. The architecture of the BPN is a layered feed-forward neural network, in which the nonlinear elements (neurons) are arranged in successive layers, and the information flows unidirectionally, from input layer to output layer, through the hidden layer(s) (Lek and Guégan, 1999).

SOM falls into the category of unsupervised learning methodology, in which the relevant multivariate algorithms seek clusters in the data to produce a low-dimensional, discretized representation of the input space of the training samples (Kohonen, 1984). The SOM is an algorithm used to visualize and interpret large high-dimensional data sets. Self-organizing maps are different from other ANNs in the sense that they use a neighborhood function to preserve the topological properties of the input space.

The most important part of ANN modeling is the generalization, the development of a model that is reliable in geomorphic modeling. Over-fitting (i.e., a model describes random error or noise instead of the underlying relationship) can be minimized by having two validation samples in addition to the training sample. In the generalization, the data are divided typically into three subsets: for example, 40% for training, 30% to prevent over-fitting, and 30% for testing (Smith, 1993). Training on the training set should stop at the epoch when the average error term computed on the second set begins to rise (the second set is not used for training but merely to decide when to stop training). Then, the third set is used to examine how well the model performs (Bishop, 1995).

2.6.3.3.2 ANN in geomorphology

ANNS have been applied in various fields of geomorphology, especially examples in hillslope and fluvial geomorphology are numerous. For example, Lee et al. (2003, 2004), Lee (2007), Nefeslioglu et al. (2008), and Falaschi et al. (2009) applied ANNs in landslide studies and Campolo et al. (1999), Gautam et al. (2000), and Sarangi and Bhattacharya (2005) in studying fluvial systems. Moreover, ANNs have been used, for example, in periglacial (e.g., Leverington and Duguay, 1997; Luoto and Hjort, 2005), aeolian (e.g., Ehsani and Quiel, 2008), volcanic (Ibanez et al., 2009), and karst (Wu et al., 2008) geomorphology.

Recently, ANNs have been used for various hazard assessments and geo-engineering applications (Lee, 2007; Lee et al., 2004; Ermini et al., 2005; Melchiorre et al., 2008; Oh et al., 2010; Rossi et al., 2010) because they allow the modeling of a process, which starts from the database containing the variables that describe that particular process.

They have already been applied in multiple landslide studies, in particular, to the indirect determination of triggering parameters and also to landslide susceptibility mapping with physical terrain factors (Lee, 2007; Melchiorre et al., 2008).

2.6.3.3.3 Strengths and weaknesses of ANN

ANNs offer a number of advantages, including requiring less formal statistical training, the ability to implicitly detect complex nonlinear relationships between dependent and independent variables, the ability to detect efficient interactions between predictor variables, and the availability of multiple training algorithms (Tu, 1996). ANNs are computationally intensive methods for finding patterns in data sets that are so large, and contain so many predictors, that standard methods such as multiple regression are impractical (Crawley, 2007). ANNs have been highly efficient in offering solutions to problems, where traditional models have failed or are very complicated to build. Due to the nonlinear nature of the ANNs, they are able to express much more complex phenomena than some linear modeling techniques. Additionally, the transformations of the variables are generally automated in the computational process. ANNs can identify and learn correlated patterns between input data sets and corresponding target values (Lek and Guégan, 1999) and can be used to predict the output of new independent input data (cf. **Figure 5(b)**). Thus, they are ideally suited for the modeling of geomorphic data which are known to be very complex and often nonlinear (Phillips, 2003, 2009).

Disadvantages of the ANN include its black box nature, greater computational burden, proneness to over-fitting, and the empirical nature of model development (Ripley, 1996; Tu, 1996; **Table 1**). The individual relations between the input variables and the output variables of ANN are not developed by theoretical judgment so that the model tends to be an input–output table without solid analytical basis. Moreover, in applications where the goal is to create a system that generalizes well in unseen examples (e.g., spatial prediction), the problem of overtraining has emerged. This arises in over-complex or over-specified systems when the capacity of the network significantly exceeds the needed free parameters (Bishop, 1995; Heikkinen et al., 2006). There are two schools of thought for avoiding this problem. The first is to use cross-validation and similar techniques to check for the presence of overtraining. The cross-validation helps to optimize the fit in three ways: (1) by limiting the number of hidden units; (2) by limiting the number of iterations; and (3) by inhibiting network use of large weights (Bishop, 1995; Tu, 1996). The second is to use some form of regularization. This is a concept that emerges naturally in a probabilistic (Bayesian) framework, where the regularization can be performed by selecting a larger prior probability over simpler models; but also in statistical learning theory, where the goal is to minimize over two quantities: the empirical risk and the structural risk, which roughly correspond to the error over the training set and the predicted error in unseen data due to over-fitting (Venables and Ripley, 2002).

In addition, drawbacks of ANNs include the requirement of large quantities of data to train, validate, and test the

network, and the limited insights into the contributions of the predictors in the modeling process (but see Olden and Jackson, 2002). Moreover, an ANN does not allow comprehensive examination of the response curves of features against environmental gradients (Manel et al., 1999; Pearson et al., 2002). One of the most critical aspects of the use of ANN as a modeling tool is the level of knowledge needed. In general, limited expertise exists in modeling with ANN for practical applications. ANN has a multipart model structure and the skill levels required to achieve reasonable results are higher than when using other modeling approaches (Ermini et al., 2005).

2.6.3.4 Statistical Boosting and BRT

2.6.3.4.1 Technical background

In geomorphology, ANNs have been utilized clearly more frequently when compared with other machine-learning techniques. Recently, other learning algorithms such as random forests, bagging, and boosting have received attention (Hastie et al., 2001). Of these, boosting is seen to be one of the major improvements in statistical modeling (Friedman et al., 2000). Boosting was developed by Freund and Schapire (1996) but was not fully understood until examined in depth by Friedman et al. (2000). In general, boosting is used, first, to improve the performance of models calibrated using traditional statistical methods and, second, to overcome problems related to more conventional modeling techniques. Boosting typically occurs in iteration by incrementally combining single models into a final complex model.

Here, we focus briefly on a BRT method. BRT combines the strengths of two commonly used techniques: regression trees and boosting (Friedman et al., 2000; Elith et al., 2008). Similar to GLM, BRT models can be fitted to a variety of response types (e.g., Gaussian, Poisson, and binomial) by specifying the error distribution and the link.

BRT is a model-averaging (ensemble) method that differs fundamentally from more often used statistical techniques (e.g., GLM). In BRT, each of the individual models consists of a simple classification or regression tree (Hastie et al., 2001). The boosting algorithm uses an iterative method for developing a final model in a forward stage-wise fashion, progressively adding trees to the model, while re-weighting the data to emphasize cases poorly predicted by the previous trees (Friedman et al., 2000).

BRT utilizes a numerical optimization technique for minimizing a loss function (like deviance) by adding a new tree at each step. Predictor variables are input into a first regression tree, which reduces the loss function to a minimum. It should be noted that each consecutive tree is built for the prediction residuals of an independently drawn random sample. The introduction of a certain degree of randomness into the boosted model usually improves accuracy and speed and reduces over-fitting (Friedman, 2002). Thus, a second tree is fitted to the residuals of the first and the model is updated to contain two trees, and the residuals from these are then calculated. This residual is then input into another tree to improve the classification. The sequence is then repeated for as long as necessary. The process is stagewise, not stepwise,

because existing trees are left unchanged as the model is enlarged. The final BRT model is a linear combination of many trees (often hundreds to thousands) that can be thought of as a regression model where each term is a tree. Further information about the boosting and BRT method can be found in Ridgeway (1999), Friedman et al. (2000), Friedman (2001, 2002), Hastie et al. (2001), and Elith et al. (2008).

2.6.3.4.2 Boosting in geomorphology

Boosting methods have been applied, for example, in hydrology (Snelder et al., 2009), soil science (Brown et al., 2006; Brown, 2007), and ecology (e.g. Elith et al., 2006). However, there is a paucity of examples in geomorphology. To our knowledge, the only examples are in modeling the distributions of periglacial landforms and processes (Hjort and Marmion, 2008, 2009; Marmion et al., 2008, 2009; Luoto et al., 2010).

2.6.3.4.3 Strengths and weaknesses of statistical boosting

Boosting methods have several strengths that encourage their utilization in modeling complex geomorphic features (**Table 1**). First, boosting provides an opportunity to capture complex phenomena-environment relationships by taking into account nonlinearities and interactions in the data (Friedman et al., 2000). For example, geomorphic processes are commonly linked to the interaction between two or more environmental factors. Moreover, important interactions can be identified (Elith et al., 2008). Second, boosting methods are less affected by outliers (Friedman et al., 2000). For example, this has significance in GDM at medium- and coarse-scale resolutions (scales with a grain size over 1 ha). Third, boosting is relatively immune to over-fitting, a rather uncommon problem for machine-learning techniques (Friedman et al., 2000; Friedman, 2002). Still, the over-fitting and poor extrapolation ability may be a problem if compared with parametric techniques (e.g., GLM). Thus, it is important to evaluate the models using (independent) evaluation data sets.

Fourth, scalability to large data sets is a desired property in modeling extensive areas because massive datasets can be collected cost efficiently using RS techniques and GI data banks. Fifth, the relative influence of predictors on the response can be estimated (Friedman, 2001; Friedman and Meulman, 2003). This is an advantage in explorative analysis and when the variables are ranked according to their contribution. Sixth, in the exploration of response shapes, boosting enables the detection of sharp discontinuities (Friedman, 2001; Friedman and Meulman, 2003). This has relevance when modeling the distributions of landforms and processes that occur over only a small proportion of the sampled environmental space.

Finally, the prediction ability of the boosted models has shown to be very high (Elith et al., 2006). Several ecological studies have suggested that boosting methods outperform conventional modeling techniques (e.g., Brown et al., 2006; Leathwick et al., 2006). In the geomorphic context, the differences between BRT and the other techniques presented in this chapter (GLM, GAM, and ANN) have not been especially clear (Hjort and Marmion, 2008, 2009; Marmion et al., 2008, 2009;

Figure 5(b)). In terms of disadvantages, a potential weakness of boosting is that insufficient or noisy data may result in an inconsistent model (Bauer and Kohavi, 1999; Hjort and Marmion, 2008). Moreover, the computation time may be excessive with large data sets (more than thousands to tens of thousands observations). The presence of spatial autocorrelation in the response data may also be problematic, resulting in inconsistency in the models, but this weakness is common for all statistical techniques (Diniz-Filho et al., 2003; De'ath, 2007). In general, the utilization of boosting may be a complicated task. In addition, for those seeking a single best model, the boosting techniques may be an unsuitable approach.

2.6.4 SWOT Analysis of Statistical Modeling in Geomorphology

Below, we address the strengths, weaknesses, opportunities, and threats of statistically based distribution modeling in geomorphology. Aspects can often be both a strength and an opportunity or a weakness and a threat. Thus, to avoid overlap, we treat strengths and opportunities together as well as weaknesses and threats. For both groups, we highlight eight different issues. Naturally, many of the arguments presented would deserve more explanation, but to keep the results of the SWOT analysis relatively simple, we avoid excessive argumentation. Moreover, some common issues in traditional statistical analysis (e.g., sample design, unit of observation, and observational vs. experimental data) and spatial modeling (e.g., scale related) are not considered here (Cochran, 1977; Sokal and Rohlf, 1995; Bivand et al., 2008; Ott and Longnecker, 2010). A summary of the SWOT analysis is shown in **Figure 6**.

2.6.4.1 Strengths and Opportunities

The possibility to: (1) increase the objectivity of interpretations; (2) simplify complex geomorphic systems; (3) predict the occurrences of landforms and processes in changing environmental conditions; (4) explore remote areas; (5) analyze

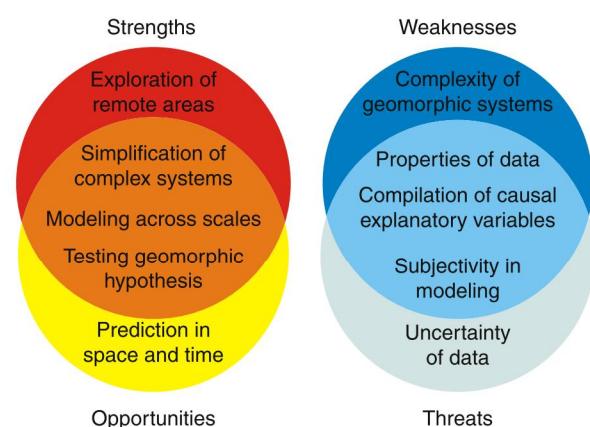


Figure 6 A summary of the SWOT analysis where five strengths/opportunities and weaknesses/threats are highlighted.

and predict geomorphic processes across scales; (6) identify the shapes of responses of environmental drivers and geomorphic processes; (7) develop a controlled study setting to test specific hypothesis; and (8) utilize various statistical and GIS software and working packages are considered to be the most important strengths and opportunities of the GDM approach.

First, many of the findings in geomorphology are traditionally based on the experience and skills of a researcher (or interpreter). There is a risk that the hypothesis/expectations affect the interpretations and you see what you want to see. Thus, the level of objectivity can be increased by formulating and testing hypotheses with statistical techniques (e.g., Ayalew and Yamagishi, 2005). Naturally, the study setting has to be sound and support the quantitative test of hypothesis (Section 2.6.2).

Second, many geomorphic systems are multivariate in nature and the relationships are complex with nonlinear effects and feedback mechanisms (Phillips, 2003, 2009). Multivariate statistical techniques help to simplify complex relationships and identify the key drivers controlling geomorphic processes. Third, a GDM approach enables the researcher to pose what if questions and predict, for example, effects of climate change on the distribution or activity of geomorphic features (Fronzek et al., 2006, 2010). Fourth, efficient data-acquisition techniques (e.g., laser scanning) and data management systems (e.g., GIS) with multivariate statistical techniques enable the exploration of extensive and remote regions in a way that fundamentally differs from traditional geomorphic mapping (Guzzetti et al., 1999; Etzelmüller et al., 2006).

Fifth, the analyses and predictions can be performed across scales and the effects of scale can be addressed from local to global scales. Although scale issues have received considerable attention in biogeographical distribution modeling (e.g., Menke et al., 2009), surprisingly few have focused on these issues in geomorphology (Luoto and Hjort, 2006). For example, it should be noted that those environmental factors important at one scale may not be important at another. This scale dependency has been revealed in previous modeling studies, in which both the important explanatory variables and model performances showed changes according to the modeling scale (Luoto and Hjort, 2006; cf. Menke et al., 2009). Moreover, GDM enables up- and downscaling of geomorphic occurrences (e.g., Luoto and Hjort, 2008).

Sixth, statistical modeling enables a detailed study of the shapes of the response functions (Brenning and Trombotto, 2006; Hjort and Luoto, 2011). A thorough understanding of the shapes of responses is crucial to improve our understanding of the determinants of Earth surface processes. For example, in multivariate analysis, the assumption of the shape of the response function may be incorrect if none of the effects of the other explanatory variables are taken into consideration (Brenning and Trombotto, 2006; Hjort and Luoto, 2011). Thus, detailed study of these response shapes introduces new aspects and opens new theoretical discussions. Seventh, the possibility to test geomorphic hypotheses of potentially important environmental factors in a setting where the other affecting factors (e.g., sample size, data distribution, and mapping intensity) are controlled is a significant advantage. Moreover, statistical methods enable the use of artificial data

that open new possibilities for methodological developments in GDM.

Finally, there exist numerous software solutions and working packages that can be used to analyze and predict geomorphic features. For example, the R statistical software (R Foundation for Statistical Computing, Vienna, Austria) (<http://www.R-project.org>) and S-PLUS (Insightful Corporation, Seattle, WA, USA) include numerous useful applications as well as the statistical working packages BIOMOD (Thuiller, 2003; Thuiller et al., 2009), GARP (Stockwell and Peters, 1999), GRASP (Lehmann et al., 2002), and MAXENT (Phillips et al., 2006). GIS software often utilized are ArcGIS (Environmental Systems Research Institute Inc., Redlands, CA, USA), GRASS (Geographic Resources Analysis Support System (GRASS GIS) Software, ITC-irst, Trento, Italy), and SAGA (System for Automated Geoscientific Analyses, Göttingen, Germany). The choice of software and working package depends on the scope of the study, although the R statistical software is recommended in many modeling settings (Crawley, 2007). Altogether, the methodological developments in statistics and Earth observation techniques provide a completely new approach to analyze and predict geomorphic landforms and processes in a spatial context.

2.6.4.2 Weaknesses and Threats

The issues that may hinder the utilization of GDM in exploration, explanation, and prediction in geomorphology include: (1) complexity of geomorphic processes and non-linearity of the responses; (2) difficulties in compilation of causal spatial variables; (3) interpretability of the results; (4) data orientation of statistical techniques; (5) postulation of rather static conditions; (6) geographic properties of data in combination with strict assumptions of several statistical techniques; (7) uncertainty and error of input data; and (8) subjectivity during various stages of modeling, especially during model calibration (Section 2.6.2).

First, although GDMs are useful in shedding light on complex relationships, the complexity of geomorphic systems (e.g., nonlinearity, interaction, multi-scale, and feedback mechanisms) can be unmanageably high (cf. Phillips, 2003; Murray et al., 2009). Consequently, statistical techniques can be incapable of capturing the true relationship between geomorphic process and environmental variable. Second, it is desirable to model geomorphic processes based on causal parameters. However, to include such factors in GDM studies is often problematic (e.g., Ayalew and Yamagishi, 2005). In practice, a spatially coherent study design based on causal variables requires costly and laborious field measurements. Thus, variables are often drawn from accessible GI and RS data layers. This can lead to another problem, namely the scale difference between the process modeled and the data used.

Third, complex processes are most efficiently modeled using flexible methods such as non-parametric statistical techniques. However, the modeling results may be difficult or impossible to interpret especially if the relationships between variables are explored (e.g., Ermini et al., 2005). Fourth, one common problem in GDM studies is that novel statistical modeling techniques are often data driven (e.g., Luoto and Hjort, 2005).

Consequently, the observed shapes of the response functions may be misleading and prediction abilities overestimated. This may cause serious problems when the models are transferred to another area (spatial extrapolation) or time (temporal extrapolation, e.g., studies of climate change effects). Fifth, statistical models often postulate static conditions, but processes in nature are dynamic and may not fulfill the equilibrium assumption (cf. Guisan and Zimmermann, 2000). However, most of the geomorphic features are rather stable over short time periods.

Sixth, typical properties (e.g., spatial autocorrelation, confounding factors and multicollinearity, outliers, closure, and nonstationary) of modeling data may bias the results due to the assumptions related to input data in several statistical techniques (especially regression methods) (Sokal and Rohlf, 1995; Ott and Longnecker, 2010; cf. Guzzetti et al., 1999). For example, spatial autocorrelation is a very general statistical property of geomorphic variables observed across geographic space. A variable is spatially autocorrelated if a measure made at one location can be used to estimate a measure made at another location, and autocorrelation is positive when subjects close to each other are more alike than distant things (Goodchild, 1986). Spatial autocorrelation can hamper attempts to identify plausible relationships between geomorphic phenomena and environmental correlates, because the use of statistical tests may be invalidated by a strong spatial structure (e.g., Diniz-Filho et al., 2003). Multi-collinearity (i.e., explanatory variables are highly correlated) may result in excluding more causal factors from multivariate models. Moreover, confounding factors (i.e., uncontrolled variables that correlate with response and explanatory variables) can hamper the detection of the actual key environmental factors underlying process–environment relationships (e.g., Ott and Longnecker, 2010).

Seventh, uncertainty and errors of the explanatory variables (especially GIS and RS) are often ignored in statistical analysis, although these issues can have profound effects on the outcomes of models. For example, explanatory variables are commonly derived from a digital elevation model (DEM) in geomorphic analyses. DEMs characteristically contain systematic and nonsystematic errors that are amplified when first- (e.g., slope angle) and second-order (e.g., topographical wetness index) derivatives are calculated (e.g., Moore et al., 1991). Therefore, the quality of the DEMs should be assessed in detail.

Finally, several subjective choices are made during the modeling process (Section 2.6.2). The investigator has to make decisions on data (e.g., distribution, source, amount, and scale), modeling technique (e.g., parametric or nonparametric), model selection approach (e.g., *p*-value vs. information theory and stepwise vs. model averaging), and model evaluation methods (e.g., split-sample approach or independent validation). Moreover, no clear guidelines exist for measuring model performance and the assessment of the goodness of the model can be rather subjective.

As presented above, several critical issues may affect the usability of GDM in geomorphology. Commonly, it is impossible to consider all the potential problems although various ways to cope with data and technique-related problems exist (cf. Guisan and Thuiller, 2005; Heikkilä et al., 2006;

Elith and Leathwick, 2009). Thus, the key is to focus on those issues critical to the study problem.

2.6.5 Future Challenges

In future, one of the main tasks in GDM is the generation of more robust models. Robust geomorphic models are better transferred in space and time but, more importantly, would improve our understanding of geomorphic systems. In this context, we highlight the incorporation of solid geomorphic theory into the modeling process. For example, models should be calibrated using causal explanatory variables instead of surrogates of environmental determinants. Models become increasingly robust and less location specific as the environmental variables become more process oriented and relevant to geomorphic processes. Moreover, integration of statistical and mechanistic models could increase the robustness of the models and provide new insights into geomorphic systems (e.g., Frattini et al., 2008).

The traditional problems of data quality and nature should be considered more seriously in statistically based modeling. Data uncertainties are seldom addressed although errors in GI and RS data may cause flawed results. The effects of spatial autocorrelation on the reliability of the results could be studied using, for example, autoregressive models (see Dormann et al., 2007). However, the true harmfulness of autocorrelation in data is still under discussion (e.g., Bini et al., 2009). Multi-collinearity issues have been addressed in some studies but deserve more attention (e.g., Luoto, 2007; Hjort and Luoto, 2009). For example, hierarchical partitioning (HP) (Chevan and Sutherland, 1991) and variation partitioning (VP) (Borcard et al., 1992) are efficient approaches for tackling multi-collinearity problems. HP and VP are quantitative statistical methods, which could be useful to study Earth surface process–environment relationships by decomposing the variation of response variables into independent and joint components.

Several modeling techniques have been used in GDM (Luoto and Hjort, 2005; Marmion et al., 2008). However, there exist different untested and underused approaches such as Bayesian (e.g., Ellison, 2004), presence only (e.g., Phillips et al., 2006), statistical consensus (e.g., Marmion et al., 2009), support vector machine (e.g., Brenning, 2005), and quantile regression (e.g., Cade et al., 2005) methods. For example, presence-only methods would enable the use of geomorphic maps and data sets with insufficient mapping intensity (i.e., data consisting of records describing known occurrences but lacking information about known absences). In predictive settings, consensus approaches could be very useful because of their high prediction ability when compared with single method approaches (e.g., Marmion et al., 2009).

Finally, statistically based modeling in geomorphology has a rather short history when compared with biological and ecological applications. Consequently, we highly recommend interdisciplinary cooperation between geoscientists and ecologists as well as statisticians to improve the usability of the GDM approach to gain novel insights into the drivers and processes shaping the Earth's surface.

References

- Aleotti, P., Chowdhury, R., 1999. Landslide hazard assessment: summary review and new perspectives. *Bulletin of Engineering Geology and the Environment* 58, 21–44.
- Atkinson, P.M., Massari, R., 1998. Generalized linear modeling of susceptibility to landsliding in the central Apennines, Italy. *Computers and Geosciences* 24, 373–385.
- Atkinson, P., Jiskoot, H., Massari, R., Murray, T., 1998. Generalized modelling in geomorphology. *Earth Surface Processes and Landforms* 23, 1185–1195.
- Austin, M.P., 2002. Spatial prediction of species distribution: an interface between ecological theory and statistical modelling. *Ecological Modelling* 157, 101–118.
- Austin, M.P., Belbin, L., Meyers, J.A., Doherty, M.D., Luoto, M., 2006. Evaluation of statistical models used for predicting plant species distributions: role of artificial data and theory. *Ecological Modelling* 199, 197–216.
- 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.
- Bauer, E., Kohavi, R., 1999. An empirical comparison of voting classification algorithms: bagging, boosting, and variants. *Machine Learning* 36, 105–139.
- Bini, L.M., Diniz, J.A.F., Rangel, T.F.L.V., et al., 2009. Coefficient shifts in geographical ecology: an empirical evaluation of spatial and non-spatial regression. *Ecohydrology* 32, 193–204.
- Bishop, C., 1995. *Neural Networks for Pattern Recognition*. Oxford University Press, Oxford, 504 pp.
- Bivand, R.S., Pebesma, E.J., Gomez-Rubio, V., 2008. *Applied Spatial Data Analysis with R*. Springer, New York, NY, 378 pp.
- Bledsoe, B.P., Watson, C.C., 2001. Logistic analysis of channel pattern thresholds: meandering, braiding, and incising. *Geomorphology* 38, 281–300.
- Borcard, D., Legendre, P., Drapeau, P., 1992. Partialling out the spatial component of ecological variation. *Ecology* 73, 1045–1055.
- Brenning, A., 2005. Spatial prediction models for landslide hazards: review, comparison and evaluation. *Natural Hazards and Earth System Sciences* 5, 853–862.
- Brenning, A., 2008. Statistical geocomputing combining R and SAGA: the example of landslide susceptibility analysis with generalized additive models. In: Böhner, J., Blaschke, T., Montanarella, L. (Eds.), *SAGA – Seconds Out. Hamburger Beiträge zur Physischen Geographie und Landschaftsökologie*, vol. 19, pp. 23–32.
- Brenning, A., 2009. Benchmarking classifiers to optimally integrate terrain analysis and multispectral remote sensing in automatic rock glacier detection. *Remote Sensing of Environment* 113, 239–247.
- Brenning, A., Azócar, G.F., 2010. Statistical analysis of topographic and climatic controls and multispectral signatures of rock glaciers in the dry Andes, Chile (27°–33° S). *Permafrost and Periglacial Processes* 21, 54–66.
- 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.
- Brenning, A., Grasser, M., Friend, D.A., 2007. Statistical estimation and generalized additive modeling of rock glacier distribution in the San Juan Mountains, Colorado, United States. *Journal of Geophysical Research* 112, F02S15.
- Brown, D.J., 2007. Using a global VNIR soil-spectral library for local soil characterization and landscape modeling in a 2nd-order Uganda watershed. *Geoderma* 140, 444–453.
- Brown, D.J., Shepherd, K.D., Walsh, M.G., Mays, M.D., Reinsch, T.G., 2006. Global soil characterization with VNIR diffuse reflectance spectroscopy. *Geoderma* 132, 273–290.
- Burnham, K.P., Anderson, D.R., 2002. *Model Selection and Multimodel Inference: A Practical Information-theoretic Approach*, Second ed. Springer, New York, 488 pp.
- Cade, B.S., Noon, B.R., Flather, C.H., 2005. Quantile regression reveals hidden bias and uncertainty in habitat models. *Ecology* 86, 786–800.
- Campolo, M., Andreussi, P., Soldati, A., 1999. River flood forecasting with a neural network model. *Water Resources Research* 35, 1191–1197.
- Carrara, A., 1983. Multivariate methods for landslide hazard evaluation. *Mathematical Geology* 15, 403–426.
- Carrara, A., Pike, R.J., 2008. GIS technology and models for assessing landslide hazard and risk. *Geomorphology* 94, 257–260.
- Carrara, A., Cardinali, M., Detti, R., et al., 1991. GIS techniques and statistical models in evaluation landslide hazard. *Earth Surface Processes and Landforms* 16, 427–445.
- Chevan, A., Sutherland, M., 1991. Hierarchical partitioning. *American Statistician* 45, 90–96.
- Cochran, W.G., 1977. *Sampling Techniques*, Third ed. Wiley, New York, NY, 428 pp.
- Crawley, M.J., 2007. *The R Book*. Wiley, Chichester, 942 pp.
- Cressie, N.A.C., 1993. *Statistics for Spatial Data*. Wiley, New York, NY, 900 pp.
- Dai, F.C., Lee, C.F., 2002. Landslide characteristics and slope instability modeling using GIS, Lanatau Island, Hong Kong. *Geomorphology* 42, 213–338.
- Das, I., Sahoo, S., van Westen, C., Stein, A., Hack, R., 2010. Landslide susceptibility assessment using logistic regression and its comparison with a rock mass classification system, along a road section in the northern Himalayas (India). *Geomorphology* 114, 627–637.
- De'ath, G., 2007. Boosted trees for ecological modeling and prediction. *Ecology* 88, 243–251.
- Diniz-Filho, J.A.F., Bini, L.M., Hawkins, B., 2003. Spatial autocorrelation and red herrings in geographical ecology. *Global Ecology and Biogeography* 12, 53–64.
- Dobson, A.J., 2002. *An Introduction to Generalized Linear Models*, Second ed. Chapman and Hall/CRC, Boca Raton, FL, 225 pp.
- Dormann, C.F., McPherson, J.M., Araújo, M.B., et al., 2007. Methods to account for spatial autocorrelation in the analysis of species distributional data: a review. *Ecography* 30, 609–628.
- Ehsani, A.H., Quiel, F., 2008. Application of self organizing map and SRTM data to characterize yardangs in the Lut desert, Iran. *Remote Sensing of Environment* 112, 3284–3294.
- Elith, J., Leathwick, J., 2009. Species distribution models: ecological explanation and prediction across space and time. *Annual Review of Ecology, Evolution and Systematics* 40, 677–697.
- Elith, J., Leathwick, J.R., Hastie, T., 2008. A working guide to boosted regression trees. *Journal of Animal Ecology* 77, 802–813.
- Elith, J., Graham, C.H., Anderson, R.P., et al., 2006. Novel methods improve prediction of species' distributions from occurrence data. *Ecography* 29, 129–151.
- Ellison, A.M., 2004. Bayesian inference in ecology. *Ecology Letters* 7, 509–520.
- Ermini, L., Catani, F., Casagli, N., 2005. Artificial neural networks applied to landslide susceptibility assessment. *Geomorphology* 66, 327–343.
- Etzelmüller, B., Ødegård, R.S., Berthling, I., Sollid, J.L., 2001. Terrain parameters and remote sensing data in the analysis of permafrost distribution and periglacial processes: principles and examples from southern Norway. *Permafrost and Periglacial Processes* 12, 79–92.
- Etzelmüller, B., Heggem, E.S.F., Sharkhuu, N., et al., 2006. Mountain permafrost distribution modelling using a multi-criteria approach in the Hövsgöl area, northern Mongolia. *Permafrost and Periglacial Processes* 17, 91–104.
- Falaschi, F., Giacomelli, F., Federici, P.R., et al., 2009. Logistic regression versus artificial neural networks: landslide susceptibility evaluation in a sample area of the Serchio River valley, Italy. *Natural Hazards* 50, 551–569.
- Frattini, G.B., Crosta, A., Carrara, A., Agliardi, F., 2008. Assessment of rockfall susceptibility by integrating statistical and physically-based approaches. *Geomorphology* 94, 419–437.
- Freund, Y., Schapire, R.E., 1996. Experiments with a new boosting algorithm. In: Saitta, L. (Ed.), *Machine Learning: Proceedings of the Thirteenth International Conference*. Morgan Kaufman, San Francisco, CA, pp. 148–156.
- Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *Annals of Statistics* 29, 1189–1232.
- Friedman, J.H., 2002. Stochastic gradient boosting. *Computational Statistics and Data Analysis* 38, 367–378.
- Friedman, J.H., Meulman, J.J., 2003. Multiple additive regression trees with application in epidemiology. *Statistics in Medicine* 22, 1365–1381.
- Friedman, J., Hastie, T., Tibshirani, R., 2000. Additive logistic regression: a statistical view of boosting. *Annals of Statistics* 38, 337–374.
- Fronzek, S., Luoto, M., Carter, T.R., 2006. Potential effect of climate change on the distribution of palsas in subarctic Fennoscandia. *Climate Research* 32, 1–12.
- Fronzek, S., Carter, R., Räsänen, J., Ruokolainen, L., Luoto, M., 2010. Applying probabilistic projections of climate change with impact models: a case study for subarctic palsas in Fennoscandia. *Climatic Change* 99, 515–534.
- Gautam, M., Watanabe, K., Saegusa, H., 2000. Runoff analysis in humid forest catchment with an artificial neural network. *Journal of Hydrology* 235, 117–136.
- Goodchild, M.F., 1986. *Spatial Autocorrelation: Concepts and Techniques in Modern Geography*. Geo Books, Norwich, 49 pp.
- Goovaerts, P., 1999. Geostatistics in soil science: state-of-the-art and perspectives. *Geoderma* 89, 1–45.
- Goudie, A.S., 1995. *The Changing Earth. Rates of Geomorphological Processes*. Blackwell, Oxford, 302 pp.
- Guisan, A., Zimmermann, N.E., 2000. Predictive habitat distribution models in ecology. *Ecological Modelling* 135, 147–186.
- Guisan, A., Thuiller, W., 2005. Predicting species distribution: offering more than simple habitat models. *Ecology Letters* 8, 993–1009.

- Guzzetti, F., Carrara, A., Cardinali, M., Reichenbach, P., 1999. Landslide hazard evaluation: a review of current techniques and their application in a multi-scale study, Central Italy. *Geomorphology* 31, 181–216.
- Guzzetti, F., Reichenbach, P., Cardinali, M., Galli, M., Ardizzone, F., 2005. Probabilistic landslide hazard assessment at the basin scale. *Geomorphology* 72, 272–299.
- Harris, C., Arenson, L.U., Christiansen, H.H., et al., 2009. Permafrost and climate in Europe: monitoring and modelling thermal, geomorphological and geotechnical responses. *Earth-Science Reviews* 92, 117–171.
- Hastie, T., Tibshirani, R., 1986. Generalized additive models. *Statistical Science* 1, 297–318.
- Hastie, T.J., Tibshirani, R.J., 1990. Generalized Additive Models. Chapman and Hall, London, 335 pp.
- Hastie, T., Tibshirani, R., Friedman, J., 2001. The Elements of Statistical Learning: Data Mining, Inference and Prediction. Springer, New York, NY, 533 pp.
- Heikkilä, R.K., Luoto, M., Araújo, M.B., et al., 2006. Methods and uncertainties in bioclimatic envelope modelling under climate change. *Progress in Physical Geography* 30, 1–17.
- Hjort, J., 2006. Environmental Factors Affecting the Occurrence of Periglacial Landforms in Finnish Lapland: A Numerical Approach. Shaker Verlag, Aachen, 162 pp.
- Hjort, J., Luoto, M., 2006. Modelling patterned ground distribution in Finnish Lapland: an integration of topographical, ground and remote sensing information. *Geografiska Annaler* 88A, 19–29.
- Hjort, J., Luoto, M., 2009. Interaction of geomorphic and ecologic features across altitudinal zones in a subarctic landscape. *Geomorphology* 112, 324–333.
- Hjort, J., Luoto, M., 2011. Novel theoretical insights into geomorphic process–environment relationships using simulated response curves. *Earth Surface Processes and Landforms* 36, 363–371.
- Hjort, J., Marmion, M., 2008. Effects of sample size on the accuracy of geomorphological models. *Geomorphology* 102, 341–350.
- Hjort, J., Marmion, M., 2009. Periglacial distribution modelling with a boosting method. *Permafrost and Periglacial Processes* 20, 15–25.
- Hjort, J., Luoto, M., Seppälä, M., 2007. Landscape scale determinants of periglacial features in subarctic Finland: a grid-based modelling approach. *Permafrost and Periglacial Processes* 18, 115–127.
- Ibanez, J.M., Benítez, C., Gutierrez, L.A., et al., 2009. The classification of seismo-volcanic signals using Hidden Markov Models as applied to the Stromboli and Etna volcanoes. *Journal of Volcanology and Geothermal Research* 187, 218–226.
- Kohonen, T., 1984. Self-Organization and Associative Memory. Springer, Berlin, 312 pp.
- Lamelas, M.T., Marimon, O., Hoppe, A., Riva, J., 2008. Doline probability map using logistic regression and GIS technology in the central Ebro Basin (Spain). *Environmental Geology* 54, 963–977.
- Leathwick, J.R., Elith, J., Francis, M.P., Hastie, T., Taylor, P., 2006. Variation in demersal fish species richness in the oceans surrounding New Zealand: an analysis using boosted regression trees. *Marine Ecology Progress Series* 321, 267–281.
- Lee, S., 2007. Landslide susceptibility mapping using an artificial neural network in the Gangneung area, Korea. *International Journal of Remote Sensing* 28, 4763–4783.
- Lee, S., Ryu, J.-H., Lee, M.-J., Won, J.-S., 2003. Use of an artificial neural network for analysis of the susceptibility to landslides at Boun, Korea. *Environmental Geology* 44, 820–833.
- Lee, S., Ryu, J.H., Won, J.S., Park, H.J., 2004. Determination and application of the weights for landslide susceptibility mapping using an artificial neural network. *Engineering Geology* 71, 289–302.
- Lees, B.G. (Ed.), 1996. Neural Network Applications in the Geosciences. Computers and Geosciences 22, 955–1052.
- Lehmann, A., Overton, J.M., Leathwick, J.R., 2002. GRASP: generalized regression analysis and spatial prediction. *Ecological Modelling* 157, 189–207.
- Lek, S., Guégan, J.F., 1999. Artificial neuronal networks as a tool in ecological modelling, an introduction. *Ecological Modelling* 120, 65–73.
- Leverington, D., Duguay, C., 1997. A neural network method to determine the presence or absence of permafrost near Mayo, Yukon Territory, Canada. *Permafrost and Periglacial Processes* 8, 205–215.
- Lewkowicz, A.G., Ednie, M., 2004. Probability mapping of Mountain Permafrost using the BTS method, Wolf Creek, Yukon Territory, Canada. *Permafrost and Periglacial Processes* 15, 67–80.
- López-Moreno, J.I., Nogués-Bravo, D., 2005. A generalized additive model for the spatial distribution of snowpack in the Spanish Pyrenees. *Hydrological Processes* 19, 3167–3176.
- López-Moreno, J.I., Nogués-Bravo, D., Chueca-Cía, J., Julián-Andrés, A., 2006. Glacier development and topographic context. *Earth Surface Processes and Landforms* 31, 1585–1594.
- Luoto, M., 2007. New insights into factors controlling drainage density in subarctic landscapes. *Arctic, Antarctic, and Alpine Research* 39, 117–126.
- Luoto, M., Hjort, J., 2004. Generalized linear models in periglacial studies: terrain parameters and patterned ground. *Permafrost and Periglacial Processes* 15, 327–338.
- Luoto, M., Hjort, J., 2005. Evaluation of current statistical approaches for predictive geomorphic mapping. *Geomorphology* 67, 299–315.
- Luoto, M., Hjort, J., 2006. Scale matters – a multi-resolution study of the determinants of patterned ground activity in subarctic Finland. *Geomorphology* 80, 282–294.
- Luoto, M., Hjort, J., 2008. Downscaling of coarse-grained geomorphic data. *Earth Surface Processes and Landforms* 33, 75–89.
- Luoto, M., Seppälä, M., 2002. Modelling the distribution of palsas in Finnish Lapland with logistic regression and GIS. *Permafrost and Periglacial Processes* 13, 17–28.
- Luoto, M., Marmion, M., Hjort, J., 2010. Assessing the spatial uncertainty in predictive geomorphological mapping: a multi-modelling approach. *Computers and Geosciences* 36, 355–361.
- Manel, S., Dias, J.-M., Ormerod, S., 1999. Comparing discriminant analysis, neural networks and logistic regression for predicting species distribution: a case study with a Himalayan river bird. *Ecological Modelling* 120, 337–347.
- Marmion, M., Hjort, J., Thuiller, W., Luoto, M., 2008. A comparison of predictive methods in modelling the distribution of periglacial landforms in Finnish Lapland. *Earth Surface Processes and Landforms* 33, 2241–2254.
- Marmion, M., Hjort, J., Thuiller, W., Luoto, M., 2009. Statistical consensus methods for improving predictive geomorphology maps. *Computers and Geosciences* 35, 615–625.
- McCullagh, P., Nelder, J.A., 1989. Generalized Linear Models, Second ed. Chapman and Hall, London, 511 pp.
- McKillop, R.J., Clague, J.J., 2007. Statistical, remote sensing-based approach for estimating the probability of catastrophic drainage from moraine-dammed lakes in southwestern British Columbia. *Global and Planetary Change* 56, 153–171.
- Menke, S.B., Holway, D.A., Fisher, R.N., Jetz, W., 2009. Characterizing and predicting species distributions across environments and scales: argentine ant occurrences in the eye of the beholder. *Global Ecology and Biogeography* 18, 50–63.
- Melchiorre, C., Matteucci, M., Azzoni, A., Zanchi, A., 2008. Artificial neural networks and cluster analysis in landslide susceptibility zonation. *Geomorphology* 94, 379–400.
- Moore, I.D., Grayson, R.B., Ladson, A.R., 1991. Digital terrain modelling: a review of hydrological, geomorphological, and biological applications. *Hydrological Processes* 5, 3–30.
- Murray, A.B., Lazarusa, E., Ashton, A., et al., 2009. Geomorphology, complexity, and the emerging science of the Earth's surface. *Geomorphology* 103, 496–505.
- Nefeslioğlu, H.A., Duman, T.Y., Durmaz, S., 2008. Landslide susceptibility mapping for a part of tectonic Kelkit Valley (Eastern Black Sea region of Turkey). *Geomorphology* 94, 401–418.
- Nelder, J.A., Wedderburn, R.W.M., 1972. Generalized linear models. *Journal of the Royal Statistical Society* 135A, 370–384.
- Neuland, H., 1976. A prediction model of landslips. *Catena* 3, 215–230.
- Oh, H.J., Lee, S., Soedradjat, G.M., 2010. Quantitative landslide susceptibility mapping at Pemalang area, Indonesia. *Environmental Earth Sciences* 60, 1317–1328.
- Olden, J.D., Jackson, D.A., 2002. Illuminating the 'black box': a randomization approach for understanding variable contributions in artificial neural networks. *Ecological Modelling* 154, 135–150.
- Oreskes, N., Shrader-Frechette, K., Belitz, K., 1994. Verification, validation, and confirmation of numerical models in the Earth Sciences. *Science* 263, 641–646.
- Ott, R.L., Longnecker, M., 2010. An Introduction to Statistical Methods and Data Analysis, Sixth ed. Cengage, Belmont, 1284 pp.
- Park, N.W., Chi, K.H., 2008. Quantitative assessment of landslide susceptibility using high-resolution remote sensing data and a generalized additive model. *International Journal of Remote Sensing* 29, 247–264.
- Pearson, R.G., Dawson, T.P., Berry, P.M., Harrison, P.A., 2002. SPECIES: a spatial evaluation of climate impact on the envelope of species. *Ecological Modelling* 154, 289–300.
- Phillips, J.D., 2003. Sources of nonlinearity and complexity in geomorphic systems. *Progress in Physical Geography* 27, 1–23.
- Phillips, J.D., 2009. Changes, perturbations, and responses in geomorphic systems. *Progress in Physical Geography* 33, 1–14.

- Phillips, S.J., Dudik, M., Schapire, R.E., 2006. Maximum entropy modeling of species geographic distributions. *Ecological Modelling* 190, 231–259.
- Reimann, C., Filzmoser, P., Garrett, R.G., Dutter, R., 2008. Statistical Data Analysis Explained. Applied Environmental Statistics with R. Wiley, Chichester, 362 pp.
- Remondo, J., Oguchi, T., 2009. GIS and SDA applications in Geomorphology. *Geomorphology* 111, 1–3.
- Ridgeway, G., 1999. The state of boosting. *Computing Sciences and Statistics* 31, 172–181.
- Ripley, B., 1996. Pattern Recognition and Neural Networks. Cambridge University Press, New York, NY, 416 pp.
- Rossi, M., Guzzetti, F., Reichenbach, P., Mondini, A.C., Peruccacci, S., 2010. Optimal landslide susceptibility zonation based on multiple forecasts. *Geomorphology* 114, 129–142.
- Rowbotham, D.N., Dudycha, D., 1998. GIS modelling of slope stability in Phewa Tal watershed, Nepal. *Geomorphology* 26, 151–170.
- Sarangi, A., Bhattacharya, A.K., 2005. Comparison of artificial neural network and regression models for sediment loss prediction from Banha watershed in India. *Agricultural Water Management* 78, 195–208.
- Smith, M., 1993. Neural Networks for Statistical Modeling. Van Nostrand Reinhold, New York, NY, 235 pp.
- Snelder, T.H., Lamouroux, N., Leathwick, J.R., et al., 2009. Predictive mapping of the natural flow regimes of France. *Journal of Hydrology* 373, 57–67.
- Sokal, R.R., Rohlf, F., 1995. Biometry. Third ed. WH Freeman, New York, NY, 887 pp.
- Stockwell, D., Peters, D., 1999. The GARP modelling system: problems and solutions to automated spatial prediction. *International Journal of Geographical Information Science* 13, 143–158.
- Swets, J.A., 1988. Measuring the accuracy of diagnostic systems. *Science* 240, 1285–1293.
- Thuiller, W., 2003. BIOMOD – optimizing predictions of species distributions and projecting potential future shifts under global change. *Global Change Biology* 9, 1353–1362.
- Thuiller, W., Lafourcade, B., Engler, R., Aráujo, M.B., 2009. BIOMOD – a platform for ensemble forecasting of species distributions. *Ecography* 32, 1–5.
- Tu, J.V., 1996. Advantages and disadvantages of using artificial neural networks versus logistic regression for predicting medical outcomes. *Journal of Clinical Epidemiology* 49, 1225–1231.
- Van Den Eeckhaut, M., Vanwalleghem, T., Poesen, J., et al., 2006. Prediction of landslide susceptibility using rare events logistic regression: a case-study in the Flemish Ardennes (Belgium). *Geomorphology* 76, 392–410.
- Venables, W.N., Ripley, B.D., 2002. Modern Applied Statistics with S. Fourth ed. Springer, New York, NY, 495 pp.
- Whittingham, M.J., Stephens, P.A., Bradbury, R.B., Freckleton, R.P., 2006. Why do we still use stepwise modelling in ecology and behavior? *Journal of Animal Ecology* 75, 1182–1189.
- Wood, S.N., 2006. Generalized Additive Models. An Introduction with R. Chapman and Hall/CRC, New York, 422 pp.
- Wu, Q., Xu, H., Pang, W., 2008. GIS and ANN coupling model: an innovative approach to evaluate vulnerability of karst water inrush in coalmines of north China. *Environmental Geology* 54, 937–943.
- Yee, T.W., Mitchell, N.D., 1991. Generalized additive models in plant ecology. *Journal of Vegetation Science* 2, 587–602.

Biographical Sketch



Jan Hjort is a professor in physical geography at the Department of Geography, University of Oulu, Finland. He focuses on spatial and statistical analysis of geomorphic and related phenomena. His main interest fields in physical geography are geomorphology, geodiversity, and biogeomorphology of Arctic and sub-Arctic areas.



Miska Luoto is a professor in physical geography at the Department of Geosciences and Geography, University of Helsinki, Finland. His research falls principally within the fields of geomorphology and biogeography, with present emphasis on the development of robust spatial models for global change impact assessments.