

the literature has witnessed a radical change in the spectrum of possible answers to the same question. Specifically, the following decades welcomed contributions that explored more numerical-oriented approaches.

With the advent of geographical information system software, information on landscape characteristics associated with potential failures became available in digital form. This first era of innovation produced a wave of contributions centered around heuristic models (Luckman, 1987; Leoni et al., 2009). These were later replaced by bivariate statistical approaches. Among these, tools such as certainty factor (Wislocki and Bentley, 1991) or weight of evidence (WoE – Bonham-Carter et al., 1990; Agterberg et al., 1990; Luzi et al., 2000) offered good performance and easy to implement models that were used until recent times (e.g., Regmi et al., 2010). However, these methods suffered from a lack of quantitative outputs, which diminished model interpretability.

Statistical and machine learning susceptibility models provided a more quantitatively robust method of estimating landslide susceptibility, which led to its dominance in the field up to the present. The introduction of binomial statistical model by Atkinson and Massari (1998) increased susceptibility model performance and interpretation (Huang et al., 2020). Nevertheless, these models only allow for the relations between dependent and independent variables to be estimated in a linear fashion, an assumption that may not hold in many cases. For this reason, nonlinear extensions such as generalized additive models (GAMs – Brenning, 2008) have superseded them.

In the early 2000s, machine learning also made its appearance in the landslide susceptibility science, with a number of applications that extended widely to encompass (i) tree-based models (Yeon et al., 2010) and their three main derivatives: random forest (RF – Ho, 1995; Breiman, 2011; Hong et al., 2019), boosted regression trees and XGBoost (Zeng et al., 2023), ii) support vector machines (Huang and Zhao, 2018), (iii) artificial neural networks (ANNs – Amato et al., 2023; Bragagnolo et al., 2020). Historically, statistical and machine learning approaches occupied very distinct areas, with the former being sought after for its interpretability and uncertainty estimation components (Di Napoli et al., 2023), whereas the latter was used in performance-oriented applications (Marjanović et al., 2011). Only recently have these differences become more blurred, with statistical models incorporating spatial dependence information (Chalkias et al., 2020) and machine learning approaches offering tools to facilitate their interpretation (Dahal and Lombardo, 2023).

Notably, with the introduction of all these models, the landslide susceptibility community entered a somewhat dormant era between the years 2010 and 2020. During this time, a plethora of publications aimed at comparing certain models against others, each time taking a different set of tools under consideration and a different dataset to build such comparison. As a result, hundreds of articles appeared with no explicit research question other than comparing a set of models and a set of data in a particular context or setting. This practice does not allow for the systematic comparability required for general advancements in the landslide susceptibility field, and it is precisely with this idea in mind that the present work proposes a standard for a benchmark dataset (refer to Section 3).

In addition to exploring the effectiveness of different susceptibility model types, much work has investigated the effects of the ratio of landslide presence and absence data. The standard definition of landslide susceptibility, as given by Carrara et al. (1995) or Guzzetti et al. (1999), does not formally require retrieving an absolute probability as it is often the case in statistics. Conversely, the susceptibility definition corresponds to a relative probability between different mapping units. In other words, susceptibility assessment seeks to define which locations are more prone than others to experience a slope failure rather than assigning them an exact probability value (Akgun et al., 2008; Sterlacchini et al., 2011).

To strictly compute probabilities, a model should be fit with all the available presence/absence data. However, the common approach in the

literature is to keep all the presences while subsampling the absences (Huang et al., 2024). Many examples can be found where a balanced sampling strategy is pursued (e.g., Erener et al., 2017; Lucchese et al., 2021). In other contributions, the absences are still subsampled from the whole dataset but kept at a greater proportion with respect to the presences (Heckmann et al., 2014; Moreno et al., 2024; Bornaetxea et al., 2018). Importantly, sampling strategies vary between machine learning and statistical modeling.

If we consider statistical modeling, the implication of varying the proportion of the presence/absence label can essentially be seen in the global intercept. The first work to refer to this effect is by Petschko et al. (2014). There, the authors note the effect of sampling a subset of the absences on the global intercept and propose an equation to correct for this effect, thus effectively bringing the obtained relative probabilities to the standard strictly prescribed in statistics. The same line of research has been further explored (refer to the supplementary materials in Lombardo and Mai, 2018), where the effects on the global intercept have been demonstrated in a simulation exercise. In short, a dataset with much fewer presences than absences would estimate very negative global intercept values. This, in turn, applies a constant probability shift towards the left side of the susceptibility distribution. In other words, a balanced sampling choice returns probability values shaped according to a Gaussian or near-Gaussian distribution centered at around 0.5. Thus, as the absence proportion progressively increases, the distribution of susceptibility values becomes more and more positively skewed (Lombardo and Tanyas, 2022). A more positively skewed or even heavy-tailed susceptibility distribution matches the reality, with few locations being highly susceptible and most of the landscape is considered stable (Jia et al., 2021). The artificial transposition of this shape towards a normal distribution implies that any landscape is approximately split into two sides, 50 % to be considered stable and 50 % to be considered unstable. This is obviously not what happens in reality and it is also the reason why susceptibility values are almost always presented in a reclassified form. In such a way, grouping probability into low-to-high susceptibility classes removes the differences induced by values concentrated either in the bulk or tails of the distribution. This is, therefore, another area where many differences exist in the literature. In this sense, the Jenks natural break classification is quite common (Mărgărint et al., 2013; Elia et al., 2023), and alternatives can be found in an equal interval (Kavzoglu et al., 2014; Chen et al., 2016) or quantile descriptions of the susceptibility range (Steger et al., 2020; Wang et al., 2022).

The sampling strategies using machine/deep learning tools are more regulated. Machine learning largely prescribes that users select balanced sampling strategies (e.g., Batista et al., 2004), unless custom-made loss functions are used to account for data imbalance (Prakash et al., 2020; Dahal et al., 2024). For this reason, many fewer studies explore absence selection effects (e.g., Hong et al., 2019; Liang et al., 2021; Rabby et al., 2023).

3. Methods and data

The first action of this study was devising a tentative dataset, and publishing a call for expressions of interest in participating in a quantitative experiment comparing susceptibility methods on a proposed benchmark dataset. We proposed this experiment as a topical session at the annual European Geosciences Union General Assembly 2023.¹ Participants presented their calculations to obtain landslide susceptibility in the study area using the proposed dataset. The approaches of the 11 participating groups were different in many respects. In Sections 4.1–4.11, we report for each participant group information about (i) type of model, (ii) variable selection, (iii) calibration/validation approach, and (iv) performance assessment. Section 4.12 summarizes

¹ <https://meetingorganizer.copernicus.org/EGU23/session/47046>.

similarities and differences of the participants' results.

The second set of actions, collectively aimed at devising a final dataset, was based on the results of the previous step and is described in [Section 5](#). Feedback from the previous step suggested that the dataset should be updated to remove collinearity, which was an issue for most of the contributions. Moreover, it was clear that the workflow applied to obtain LSMs should be standardized both for a meaningful comparison of results from different methods and for benchmarking independent calculations against the results presented here. To this end, a final benchmark dataset was obtained adding new predictors and removing collinearity by reducing the number of variables, as described in [Section 5.1](#). The updated data were distributed to the contributors, with well-defined requirements for cross-validation (CV), so that the individual groups produced a more informative set of results.

3.1. Data

We selected a slope–unit (SU)–based dataset because SUs have a meaningful correspondence with topography ([Guzzetti et al., 2006](#)). We extracted a subset of the dataset used by [Loche et al. \(2022\)](#) for landslide susceptibility maps in Italy, who adopted an SU set previously optimized for Italy by [Alvioli et al. \(2020\)](#).

Out of the entire SU map of Italy, containing about 330,000 polygons, we selected a subset of 7360 units encompassing an area of 4,095 km² in Central Italy. [Fig. 1](#) shows the spatial location of the area of interest. The data had an attribute table containing several different morphometric and thematic variables. The morphometric variables were calculated using the European digital elevation model EU–DEM with 25-m resolution.² A few variables were obtained from the SoilGrids global dataset ([Hengl et al., 2017](#)). [Table 1](#) lists the full set of variables.

The SoilGrids dataset is an application of machine–learning models trained on over 230,000 soil profile observations from the world soil information WoSIS database ([ISRIC, 2024](#)). Lower and upper limits of a 90 % prediction interval quantify prediction uncertainty. Global datasets and models are increasingly being used to make use of data–hungry, high–performance, large–scale, machine–learning models. The accuracy of global datasets depends on the density and quality of data points used for building the models. Freely available global products are useful both in the context of this work, aiming at becoming a reference dataset for landslide susceptibility mapping, and for similar datasets, developed in different areas.

The original landslide location map in [Loche et al. \(2022\)](#) contained eight different presence/absence flags, corresponding to the point locations (highest point of landslide crown) of eight types of landslides from the Italian National landslide database assembled by the Italian Geological Survey (ISPRA; [Trigila et al. \(2010\)](#)). For this work, we selected only presence/absence of translational landslides.

To provide the contributors with two different landslide presence scenarios, we flagged landslide presence with two attribute fields, called p₁ and p₂, which is similar to flagging an SU as unstable if it contains a minimum landslide area ([Guzzetti et al., 2006; Schlögel et al., 2018](#)).

To define p₁, we selected SUs labeled as “without landslide” (p₁ flag: 0) where an SU contained no points at all, in 3766 cases (1,443.1 km²), and as “with landslides” (p₁ flag: 1) in the remaining 3594 cases (2,652.1 km²). For p₂, we selected SUs labeled as “without landslides” (p₂ flag: 0) where an SU contained up to one point, in 5089 cases (2,087.1 km²), and as “with landslides” (p₂ flag: 1) in the remaining 2271 cases (2,008.2 km²).

Note that, using p₁ as landslide presence, one would have an approximately balanced dataset with respect to the number of zeros/ones; using p₂, instead, one would have an approximately balanced dataset with respect to the total surface area covered by the SUs labeled

either with zero, or one. [Fig. 2](#) shows the spatial distribution of SUs labeled as positive/negative in the two cases. In such a varied methodological landscape pertaining to sampling ratios ([Section 2](#)), a detailed exploration of the selection of non–landslide data is beyond the scope of this work.

We invited participants to consider both landslide presence flags to produce two different LSMs for the study area. Moreover, we invited them to use their best strategy, or the strategy that best fits their model of choice, to produce a result for a landslide susceptibility index – a float number ranging from zero to unity – and an associated uncertainty, where possible.

4. Preliminary assessment of the benchmark dataset for landslide susceptibility

The following describe the methods applied in step one of the experiment, by each participating group. In each contribution, we have distinguished model selection, variable selection, calibration–validation approach, and model evaluation. In cases where no exclusion of variables is described, all variables were retained for the group's results.

4.1. Group 1. Application of the LAND–SUITE multi–model software

This contribution was presented as the (EGU) abstract by [Bornaetxea et al. \(2023b\)](#).

Model selection. Group 1 (G1) utilized the LAND–SUITE software ([Rossi et al., 2022](#)), a suite of R script modules ([R Core Team, 2021](#)) designed to support the landslide susceptibility inference process. LAND–SUITE contains several statistically driven approaches, including linear discriminant analysis (LDA), logistic regression (LR), and quadratic discriminant analysis (QDA). Additionally, the software provides an option to combine the outputs of the selected statistical methods into a single combination forecast model (CFM), where LR is used to determine the best fit among the original outputs ([Rossi and Reichenbach, 2016](#)). They tested all of the mentioned approaches (LDA, LR, QDA, and CFM), considering the two proposed landslide presence scenarios (p₁ and p₂), along with the explanatory variables, resulting in eight LSMs.

Variable selection. Each modeling process was preceded by an exploration phase to identify possible correlations among pairs of explanatory variables and, in such cases, to select only the most significant variable. Including highly correlated variables in the model training would likely inflate the model error and uncertainty estimate, thus negatively affecting the overall performance and the interpretation of variable effects ([Amato et al., 2019](#)). This in turn would increase computational time and, in some extreme collinearity cases, may even hinder the model convergence, especially for statistical models. To address this, G1 computed mutual correlation coefficients among the 26 explanatory variables (including SU area) and considered two variables as highly correlated if the Pearson correlation coefficient exceeded |0.7|. A leave-one-out (LOO) test assessed the individual significance of each variable with respect to the others ([Gong, 2006; Sin Yin et al., 2010](#)). With *n* preliminary runs of each model, excluding one variable at a time, G1 determined which variable's absence resulted in the largest loss in model performance, and they excluded the less significant variable from each highly correlated pair. This approach is referred to in different ways depending on the user background, including terms such as jackknife tests in ecology ([Shcheglovitova and Anderson, 2013](#)), or ablation studies in computer science ([Aguilera et al., 2022](#)). Additionally, for the LR outputs, G1 verified the *p*–value corresponding to each variable. Variables with *p*–values >> 0.05 were excluded from the analysis. This process was performed for each of the provided target variables (p₁ and p₂). The original values of the explanatory variables were scaled between their minimum/maximum values.

Calibration–Validation approach. Every experiment (LDA, QDA, LR, and CFM) used the same training and validation data partition, which

² <https://www.eea.europa.eu>.

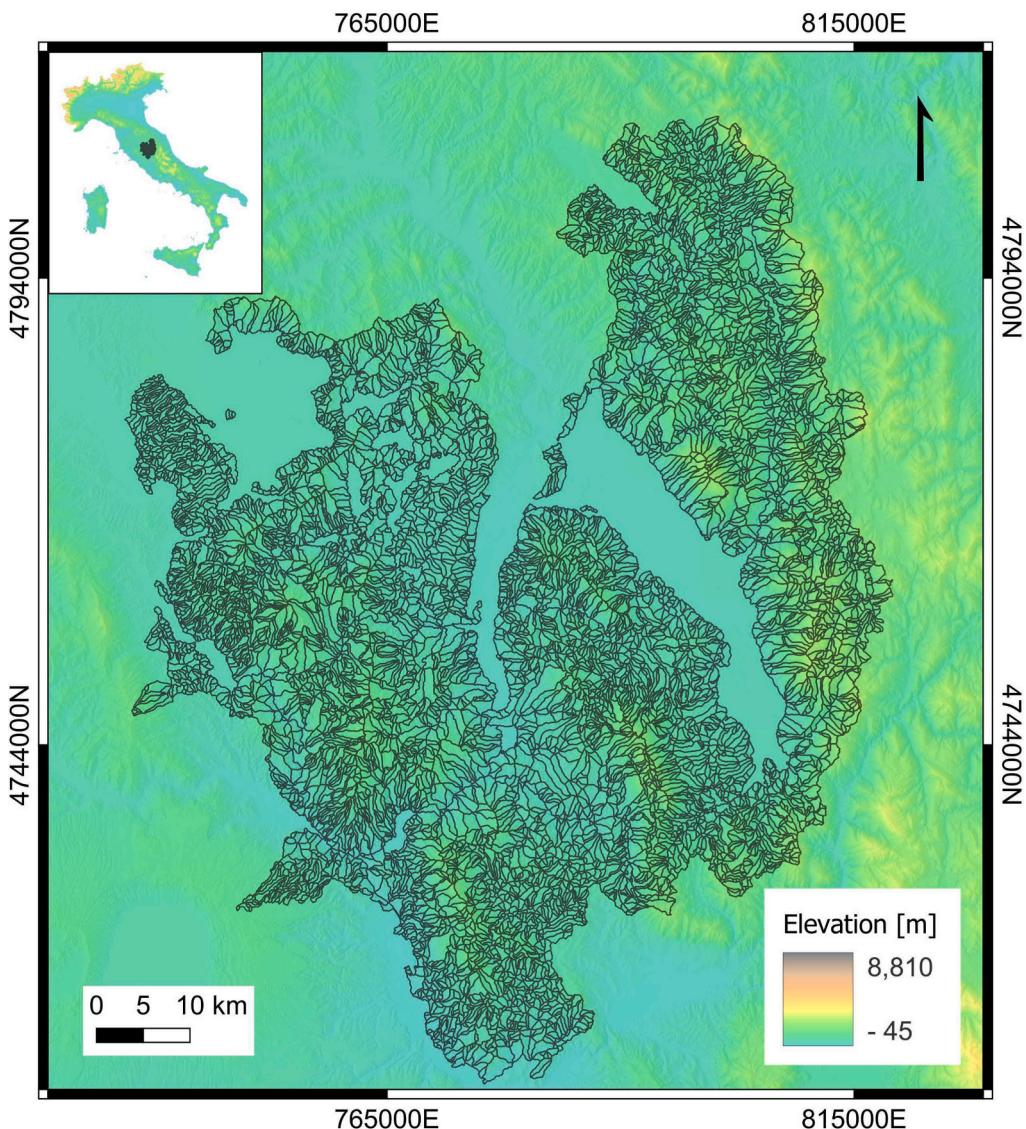


Fig. 1. Geographical location (inset) of the area covered by the slope unit set (main figure) selected in this work as a benchmark dataset for landslide susceptibility zonation. The dataset is a subset of the slope unit map obtained by Alvioli et al. (2020), and used by Loche et al. (2022) for a landslide susceptibility map of Italy. In the dataset proposed here, we selected point locations of translational landslides from the Italian national inventory known as 'IFFI' (Trigila et al., 2010). Map is in EPSG:32632 projected reference system.

involved a simple (one-fold) random CV approach. Group 1 allocated 75 % of the dataset for training and the remaining 25 % for validation. They made sure that a balanced number of landslide presence (1) and absence (0) were found in both the training and validation sets. To assess the internal variability of the results due to the randomly obtained training samples, G1 used the bootstrap resampling method (Davison and Hinkley, 1997). They conducted 100 resample iterations for each tested model and plotted the average and standard deviation of the results on variability plots (Rossi and Reichenbach, 2016).

Model evaluation. Validation of the results was performed using the area under the receiver operating characteristic (ROC) curves (AUC_{ROC} , Fawcett (2006)), calculated for both the training and validation samples. Group 1 obtained four-fold and histogram plots to visualize the overall agreement of the model compared to the observed results in the validation dataset. In all figures and tables, results corresponding to this paragraph are labeled as LDA, QDA, LR₁, and CFM (cf. Fig. 3).

4.2. Group 2. Generalized additive models with shrinkage option and geomorphological plausibility check

This contribution was presented as the EGU abstract by Camera and Bajni (2023).

Model selection. Group 2 (G2) applied GAMs, using the `mgcv` library in R (Wood, 2017). This class of models was selected because these models are easily interpretable and widely applied in recent literature with good results (e.g., Goetz et al. (2011); Bajni et al. (2023); Fang et al. (2024); Wang et al. (2024)).

Variable selection. An exploratory correlation analysis was carried out between the 27 independent variables (SU area included). Variable selection was done during the GAM fitting through shrinkage, which consists in removing the variables that explain a small part of model variance (usually variables highly correlated with others). This approach is quite intuitive in the linear case, with a penalization term used to shrink the regression coefficient values towards zero, checking at each time whether the shrinkage leads to loss in performance or not (Ranstam and Cook, 2018). In a nonlinear case, the penalization is executed in two dimensions, both for the regression coefficients as well

4.12. Overview and results of methods chosen by contributors in step one

The above sections outlined a plethora of approaches and choices intrinsic to an LSM exercise. Fig. 3 shows answers to specific questions, provided by the contributors along with results of their calculations for step one of the experiment. One can observe that the only features common to all of the groups were to produce a vector output and calculate AUC_{ROC} as a performance metric. Many groups also provided an uncertainty along with their results, most of them checked for collinearity and dropped a few variables, using only a subset of the initial variables for their classification, with different mixtures. Most of the contributors performed the frequently used training-validation data split, and a few implemented spatial cross validation.

Table 2 lists all the models considered at step one of this experiment, including model names, group who applied the models, whether they calculated results for p₁ and/or p₂, and a one-line description of the model.

Overall, we identified a few key points in which these approaches differ, and that can potentially affect the susceptibility values and model performance. Firstly, both statistical (e.g., LR, LDA, QDA, and WoE) and machine learning models (e.g., ANN, RF, and XGBoost) were used. Secondly, a few authors have included all variables, and others have used analytical or heuristic approaches prior to – or during the modeling exercise – to remove some variables. Thirdly, the approaches differed in the way data were split to calibrate and validate the model. Although all authors performed separate calibration and validation with variable fractions of input data, the use of CV and spatial CV was not systematic. In particular, only G4 and G11 applied spatial CV.

A few groups selected more than one method, two groups (G1 and G11) proposed ensemble modeling with different ways of combining a few results into a single one. A few groups selected the same method: LR and Bayesian LR were used by G1, G7, and G8; generalized additive models were adopted by G2, G3, and G9; a tree boosting system called XGBoost was used by G4, G8, and G11; RF was adopted by G5, G6, and G11.

Although many contributing teams chose to remove some variables, they differed in how they decided to do so, ranging from pre-modeling heuristics (e.g., based on a VIF-thresholding, G8 and G10), and penalization-based methods during modeling (e.g., G2). Likewise, although many groups performed a CV, they differed in the number of folds/repetitions.

As expected, all of the participants used AUC_{ROC} as a performance metric. In addition, G1 considered explicitly the graphical representation of four-fold plots (not shown here); G5 and G7 utilized success/prediction rate to evaluate performance on the basis of SU ranking of the results of classification; G8 suggested the use of Brier score, equivalent to the mean of squared differences between the (probabilistic) prediction and the target variable in each SU; G11 used several performance metrics, including accuracy, precision, F1, and Cohen's Kappa.

Due to the large number of methods applied in the first part of the experiment, and the heterogeneity of the application approaches, we do not show susceptibility maps in this section. We will show maps corresponding to the final benchmark dataset devised in this experiment, instead, described in the next section.

Fig. 4 shows a pairwise comparison of results of step one, calculated as follows:

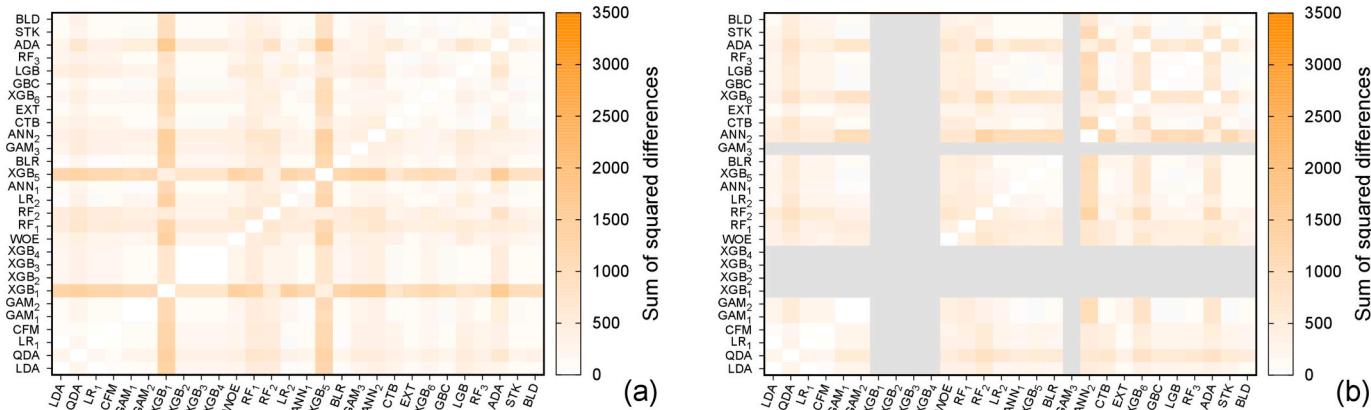


Fig. 4. Pairwise comparison of the results of different methods, in step one of the experiment, described in Section 4, calculated as in Eq. (2). Panels (a) and (b) correspond to the target variables p₁ and p₂, respectively. Names of the different methods are as in Fig. 3, Sections 4.1–4.11. Grey color denotes missing data.

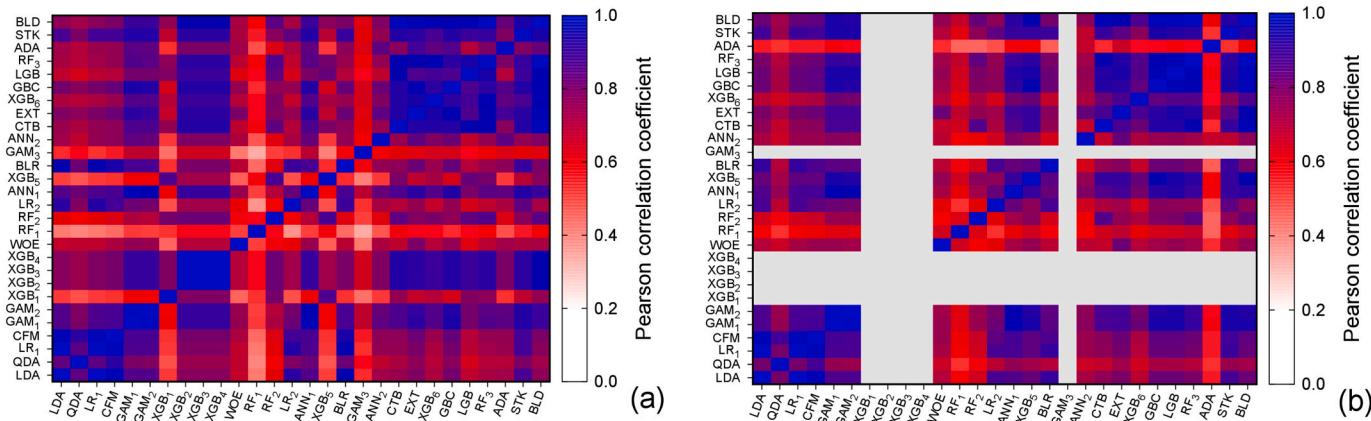


Fig. 5. As in Fig. 4, but for pairwise correlations between results for different methods. Grey color denotes missing data.

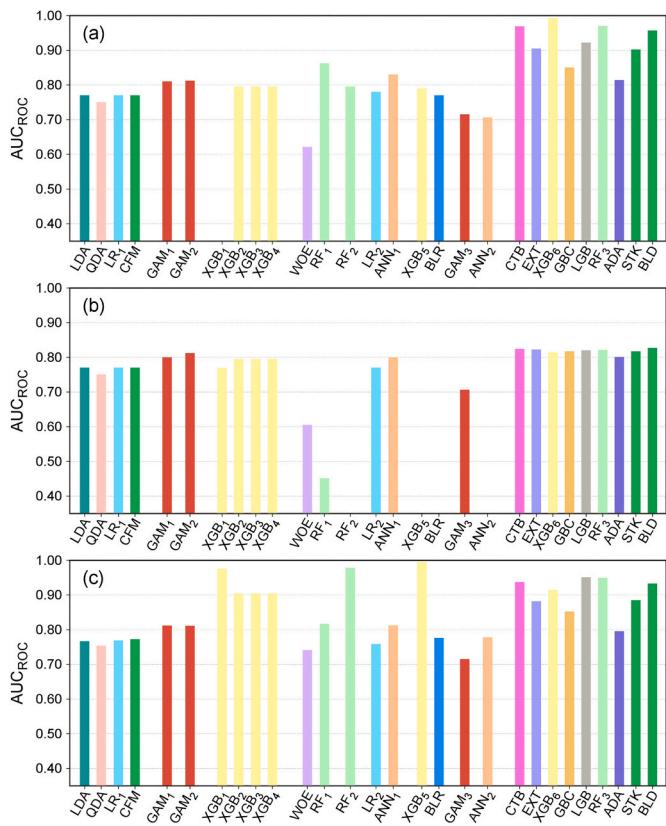


Fig. 6. Comparison of area under the receiver operating characteristic (AUC_{ROC}) values for the different methods considered in step one of the experiment, limited to the p1 landslide flag. Both in calibration (a) and in validation (b) the same method gives different results, due to the different workflow of application by different research groups. The plot in (c) was obtained by the organizers of the experiment, calculating AUC_{ROC} from the numerical results provided by contributors. Names of the different methods are as in Fig. 3, Sections 4.1–4.11.

$$D_{ij} = \sum_{k=1}^{N_{SU}} (S_i^k - S_j^k)^2, \quad (2)$$

where i and j run in the set of methods, k labels the SUs in the dataset, and $0 \leq S_i^k \leq 1$ is the susceptibility value of the i -th SU, calculated with the k -th method. One can observe from Fig. 4 that adoption of the same method did not necessarily lead to very similar results. For example, the largest differences, according to the criterion of Eq. (2), are between XGB₁, XGB₅ and all other methods; the difference between RF₁, RF₂, ADA and other methods stands out as well (Fig. 4(a), corresponding to p1). For the target label p2, Fig. 4(b), the classification from XGB_{1...4} gave the (exactly) same results in this case, and for p1 XGB₁ differs slightly from the XGB_{2..4}: they differ from the former only for the spatial CV strategy which, thus, seems to be marginally relevant, here. A few values are missing in the figures because no results were provided by the participants at step one. On the other hand, we note that results for the GAM method (p1) in the variants GAM₁ and GAM₂ are similar, although they somewhat differ from the variant GAM₃. The reason for that may reside in the different way of using predictors in the application of GAM₃, in which the role of SU size (*area*) was emphasized.

In both panels of Fig. 4 (for p1–p2), the results by G1 (LDA, QDA, LRM, and CFM with the software LAND–SUITE) and by G11 (CTB, EXT,

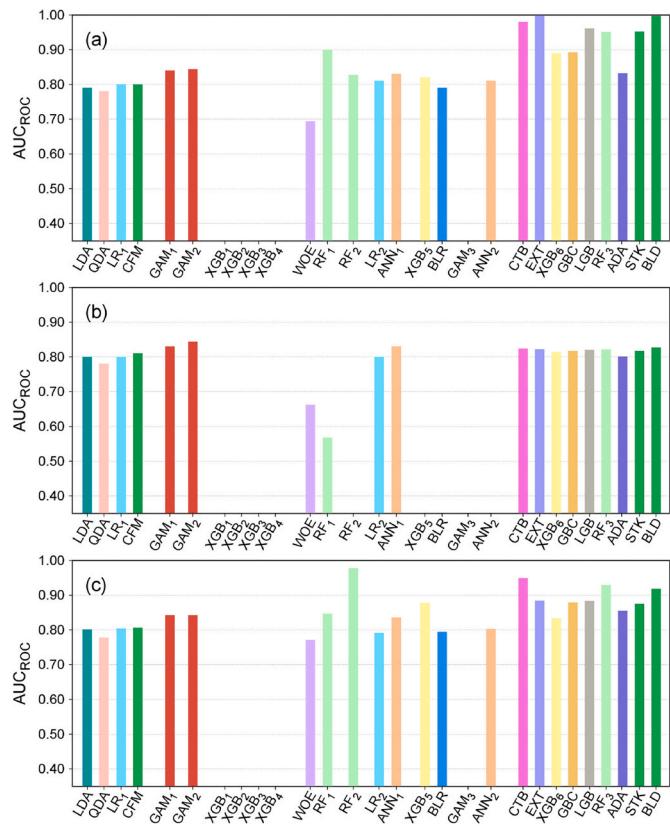


Fig. 7. As in Fig. 6, but for the p2 landslide presence flag.

XGB₆, GBC, LGB, RF₃, ADA, STK, and BLD with an ensemble machine learning) show similarities within each subset, indicating that different methods applied by the same author somehow produce more similar results. This may reside in pre-processing of data, variable selection, and training-validation strategy. We reached the same conclusion looking at Fig. 5, which shows pairwise correlations between results, calculated with the *corrplot* function of the *corrplot* package in R (Wei and Simko, 2021).

Fig. 5 also shows a few high-correlation blocks (e.g., the blueish blocks between methods GAM₁, GAM₂, and XGB_{2..4}; the methods within the ensemble machine learning by G11 for the p1 case; and similar cases for p2). Interpretation of this occurrence is not straightforward, because they correspond to multiple methods applied by different authors. Results for XGB_{1..4} and GAM₃ were not provided for p2 (grey bands in the figures).

Figs. 6 and 7 show AUC_{ROC} values for all the different methods, for training and validation reported by the authors. Moreover, we obtained AUC_{ROC} independently, from susceptibility values in the attribute tables provided by the users and the target values p1–p2, using the *roc()* function of the *pROC* R package (Robin et al., 2011). In the figures, a few missing entries are due to a few authors providing only results for p1 and/or only for the training step. Colors are consistent between the same method applied by different groups; combined models (CFM, STK, and BLD) also share the same color (dark green).

A general comment about AUC_{ROC} results is that the validation values are systematically lower than the training (fit) values, which is by design. We observe that the calculated performance is often different from the values provided by the users, both in excess or in deficiency.

This is also partially expected because most authors classified the final maps using a different subsample of the provided dataset. However, the XGB_{1...4}, WOE, RF₂, and XGB₅ (_{p1}) models resulted in higher AUC_{ROC} measured by the organizers than both training and validation values, denoting some effect on how the final maps were assembled, possibly other than a combination of training (fit) and validation (predicted) results.

Lower values for the validation case are more prominent for a few methods. They are prominent in RF₁, for both landslide presence scenarios (likely due to the stable SU sampling strategy), and in CTB, XGB₆, RF₃, and BLD, to a lesser degree. This indicates potential overfitting of the training data and diminished performance on unseen data. However, ensemble methods from G11 (including CTB, XGB₆, RF₃, and BLD) exhibit the highest validation scores among all methods (Figs. 6 and 7). Thus, the results obtained from the spatial CV scheme appear consistent and generalizable in these cases, despite the slight overfitting observed in the models.

The findings described in this section allowed us to draw the conclusion that a truly useful benchmark dataset for LSMs would be complemented by a minimal and well-defined set of prescriptions about application of the classification methods. As a result, the variability in the output LSMs include a substantial component due to such choices. In the next section, we describe both changes in the proposed dataset and a set of prescriptions for such choices, aimed at minimizing the effect of methodological choices to obtain a meaningful benchmark.

5. Final assessment of the benchmark dataset for landslide susceptibility

The selection of input variables – along with the type of model applied – has a large effect on the final susceptibility results. Thus, to improve the comparability of the models, we introduced a second step (step two) whereby we updated the input variables and asked the contributing teams to include the entire, updated dataset in their susceptibility model.

We updated the dataset firstly by including lithological information from the geo-mechanical lithology map of Italy by Bucci et al. (2022) because many contributors asked to include such information. We first calculated the percentage presence of lithological classes in the whole area, and selected the five classes with largest percentage cover, namely: alluvial deposits (Al, 12 %), unconsolidated sedimentary rocks (Ucr, 27 %), marlstone (M, 4 %), schistose metamorphic rocks (Ssr, 35 %), and carbonate rocks (Cr, 18 %). Total percentage was 96 %. Fig. 8 shows a simple description of the new variables. Lithological classes were provided as areal percentage in each SU polygon, which includes information about SUs containing different lithologies.

Fully comparable results also required each participant to adopt the same workflow concerning training/validation steps. To do that, the organizers had contributors apply a 10-fold CV procedure with mutual exclusion. That amounts to splitting the dataset into 10 numerically balanced parts. Training would be applied 10 times, on 90 % of the data, and validation would be performed on the remaining 10 %. Iterating the procedure 10 times on the 10 different splits for validation provides a fully validated LSM. The advantage of this procedure is that susceptibility values in each SU are calculated with a model trained with independent data.

5.1. Removing correlations

The updated dataset contains the following variables: slope, curvatures (morphometric) northerness/easterness, elevation, TWI (topographic wetness index), max distance and MD/ \sqrt{A} (max distance over the square root of SU area), bulk density, clay/sand/silt content (related to soil properties and texture), and percentage of lithology classes.

A few of these quantities are probably strongly correlated, and we

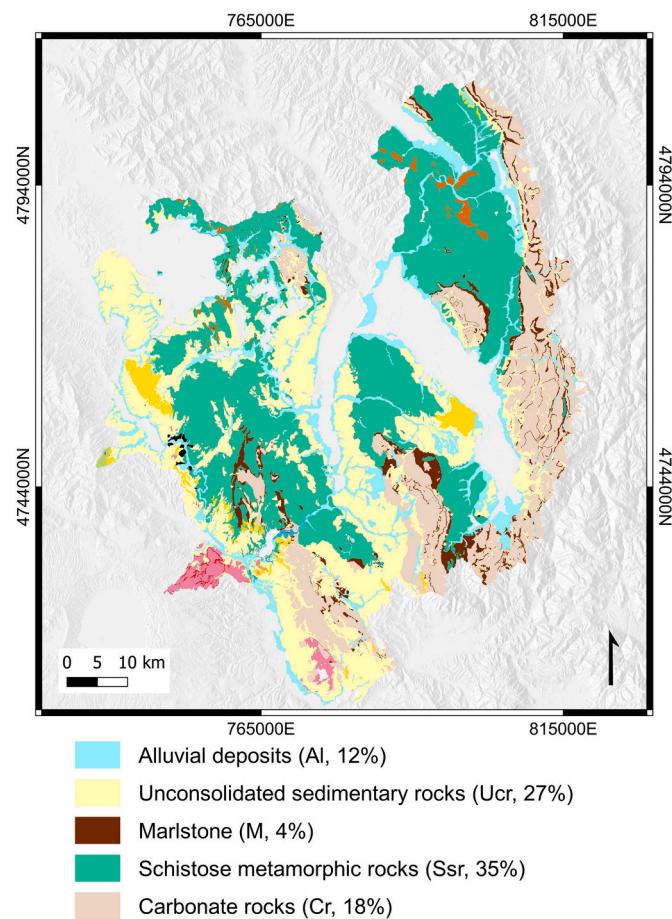


Fig. 8. Lithological classes included in the final version of the benchmark dataset. The figure shows a subset of the map prepared for the whole of Italy by Bucci et al. (2022). For this study, we considered only the five most representative classes as predictors, covering 96 % of the study area. Shaded relief map as in Fig. 2.

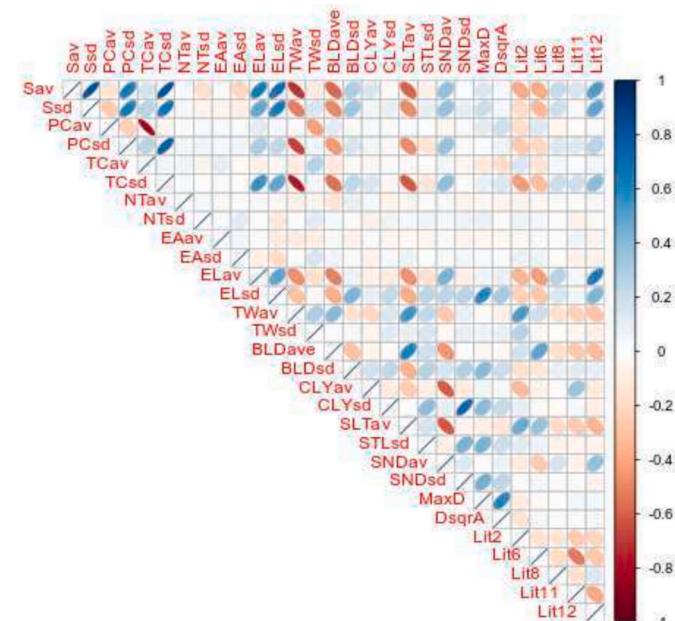


Fig. 9. Linear correlation plot of the 29 candidate variables for the final benchmark dataset. Short names of variables are as in Table 1. Symbols with larger ellipticity correspond to a higher degree of correlation between the two considered variables.

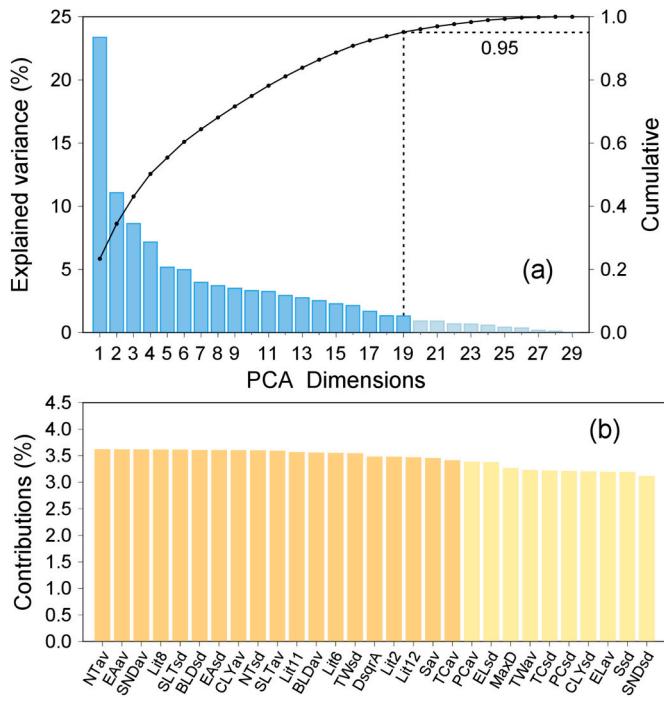


Fig. 10. (a)-(b) Result of principal component analysis (PCA) on the full dataset: (a) scree plot, i.e., contributions to the total variance explained by each principal component (PC); (b) contribution of each variable to the PC decomposition in (a).

would like to obtain a benchmark dataset free from major correlations. That is because the classification performance of a few methods (for example LR) would be penalized by correlations, while the performance of most pure machine learning models would not be affected. Moreover, we want to reduce the chance of overfitting the data and reduce the overall dimensionality of the problem.

The process of defining and removing correlations should not be linked to a specific classification method, so that we would end up with data that can be equally useful for a fair comparison of a range of different landslide susceptibility models. It should not be linked to the values of the target variable either because we have two target variables (p_1, p_2) and this would give performance advantages to some methods.

The final dataset was prepared to minimize correlations between variables and to mitigate any bias in the results of part two of this experiment. We first excluded depth to bedrock and SU area in the final dataset. The depth to bedrock variable was hardly changing across the study area, due to a low number of data points used to build the (global) SoilGrids model. Slope unit area had a peculiar behavior in at least two of the independent studies conducted in step one by G7 and G9.

After removing these variables, we checked for correlations between each pair of quantities in the dataset, using the standard `corrplot` R package (Wei and Simko, 2021). A graphical representation of the test is in Fig. 9, and the figure clearly shows a large degree of correlations. To visually explore correlations beyond linear, we used the `GGally` R package (Schloerke et al., 2022) to plot two predictors against each other, for all possible pairs, and check the general trend of data with respect to the binary variables p_1 and p_2 . As the full dataset is rather large, to visualize the results we split the data into five pieces and prepared figures for each pair combination, for a total of 10 figures for each presence variable. As this is only a visual inspection of data, we presented them in the supplementary material (Figs. S1–S10). It is clear from these figures that most variable pairs have a complex relationship with one another. Moreover, it is hard at this stage to establish which variable has a distinct behavior with respect to the absence or presence of landslides, despite a few differences described by the diagonal plots.

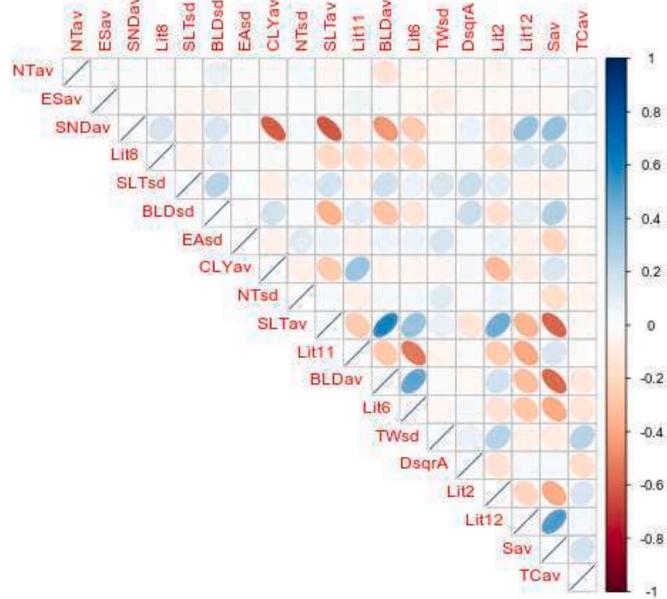


Fig. 11. Linear correlation plot of the selected 19 variables in the final benchmark dataset. Ordering is according to the decreasing contribution from principal component analysis (PCA), as in Fig. 10(b). Symbols with larger ellipticity correspond to a higher degree of correlation between the two considered variables. One can observe that, after variable reduction, most linear correlations were removed.

The lithology variables clearly stand out with respect to the others, which is expected, as they are different in nature from the other variables.

We reduced the correlation between variables using information gained from principal component analysis (PCA). We chose this method because it is not biased towards any modeling method and can measure overall correlations rather than just pairwise linear correlations. Principal component analysis is an orthogonal linear transformation of data to new coordinates that concentrate the largest variance in a smaller number of axes with respect to the original coordinates (Jolliffe, 2002).

Table 3

Attribute table of the final dataset. After principal component analysis (PCA) analysis, we retained the 19 variables contributing to 95 % of the total variance in the full dataset. In the table, SD stands for standard deviation and SU stands for slope unit.

Column name	Variable	Short name
id	Unique slope unit identifier	–
slope_aver	Mean Slope Steepness [deg]	Sav
tcurv_aver	Mean Profile Curvature	PCav
nthsns_aver	Mean Northernness	NTav
nthsns_stdd	SD of Northernness	NTsd
easns_aver	Mean Easternness	EAav
easns_stdd	SD of Easternness	EAsd
twi_stddev	SD of Topographic Wetness Index	TWsd
BLDFIE_ave	Mean Bulk density [kg/m ³]	BLDAv
BLDFIE_std	SD of Bulk density [kg/m ³]	BLDsd
CLYPPT_ave	Mean Weight % of clay particles	CLYav
SNDPPT_ave	Mean Weight % of sand particles	SNDav
SLTPPT_ave	Mean Weight % of silt particles	SLTav
SLTPPT_std	SD of Weight % of silt particles	SLTsd
D_sqrt_A	Maximum Distance/ \sqrt{SUArea}	MaxD
Litho2	Alluvial deposits (%)	Al / Lit2
Litho6	Unconsolidated sedimentary rocks	Usr / Lit6
Litho8	Marlstone	M / Lit8
Litho11	Schistose metamorphic rocks	Ssr / Lit11
Litho12	Carbonate rocks	Cr / Lit12
presence1	Binary landslide presence flag	p1
presence2	Binary landslide presence flag	p2

At variance with standard applications of PCA, we eventually did not use the transformed orthogonal variables singled out by PCA, but only the reduced set of original variables with maximal information content.

First, we performed PCA on the whole set of data, amounting to 29 variables (refer to Table 1). We used the R package factoextra (Kassambara and Mundt, 2020) to visualize results as in Fig. 10.

Data were normalized (R function `scale`) before feeding them to `pca()`. The percent contributions to principal components (PCs) calculated with the full dataset (Fig. 10(a)) is useful for determining how many PCs must be retained to explain a given fraction of the total variance in the data. We opted for a threshold of 95%; accumulating the contributions, we conclude that the first 19 PCs are enough for the purpose.

Second, we considered the contributions of all original variables to the first 19 PCs (Fig. 10(b)). Explaining the required 95% total variance would require at least 19 PCs, and we would need a minimum of 19 variables to do so with a linear combination, as in PCA. Thus, we decided to select the set of 19 variables that contributed to the 19 most relevant PCs, going from left to right in Fig. 10(b).

We validated the results of the PCA delimited dataset measuring the pairwise linear correlations among the reduced set of 19 variables. Fig. 11 shows that few correlations are left. Table 3 lists the final set of variables; note that all of the five lithological variables were retained by the PCA-based procedure.

6. Results with final dataset and constrained workflow

Results for the final dataset (Section 5, Table 3) present LSMs obtained with 16 different models. They correspond to one map for each of the 11 groups participating in the benchmark for each scenario p1 and p2, except for G1 (Section 4.1), who contributed with four results as in the first step (LDA, QDA, LRM, and CFM), G7 (Section 4.7) with two results (LR and ANN₁), and G8 (Section 4.8) with two results (XGB₅ and BLR). All the remaining groups either presented one result from the very beginning, or decided to show only their best result. G11 (Section 4.11) presented results for two different methods for p1 (STK) and p2 (RF₃).

To compare the 16 susceptibility values, we calculated the sum of squared differences as in Eq. (2) (Fig. 12), the pairwise Pearson correlation coefficient (Fig. 13), AUC_{ROC} (Figs. 14 and 15, for p1 and p2, respectively) and Brier score as in Eq. (1), (Fig. 16).

In general, metrics are much more similar to each other than in the first step; this is largely expected, given the prescription for training/CV and use of same data. Figs. 12 and 13 indicate similar considerations. A few methods stand out as more dissimilar from the others, for example QDA, RF₂ and STK (p1) and QDA, RF₁, and RF₂ (p2), in both figures. The correlation plots also show a peculiar pattern within the block of the first four results (all by G1), particularly for the combined model, CFM (which is expected).

The difference between the participant reported AUC_{ROC} values and values of AUC_{ROC} measured by the organizers are less variable compared to step one of this experiment (Figs. 14 and 15). All of the AUC_{ROC} values are between 0.7 and 0.75, are slightly higher for p2, and highest for the

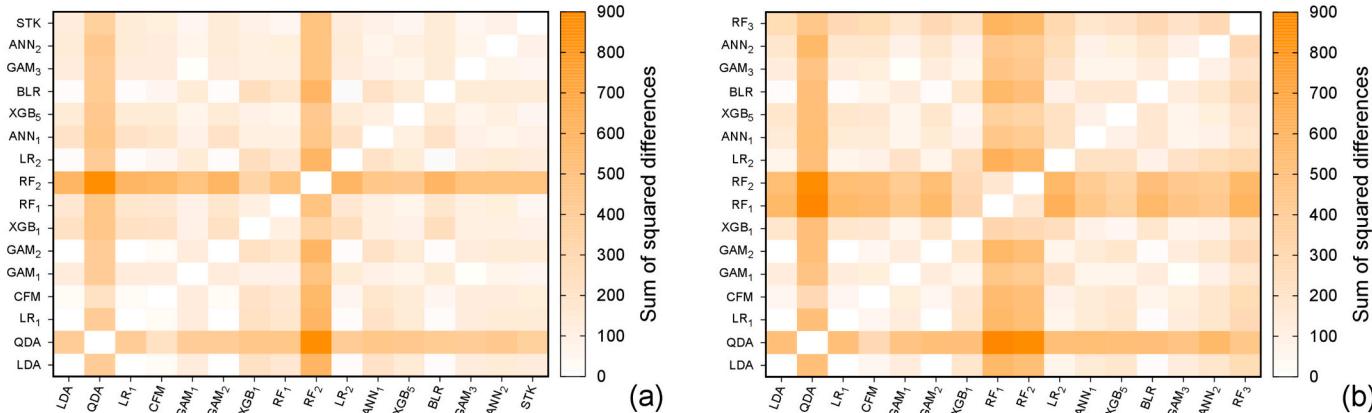


Fig. 12. Pairwise comparison of the results of different methods applied with same workflow, during step two of the experiment, described in Section 5. Numerical values calculated as in Eq. (2). Panels (a) and (b) correspond to the target variables p1 and p2, respectively. Names of the different methods are as in Fig. 3, Sections 4.1–4.11.

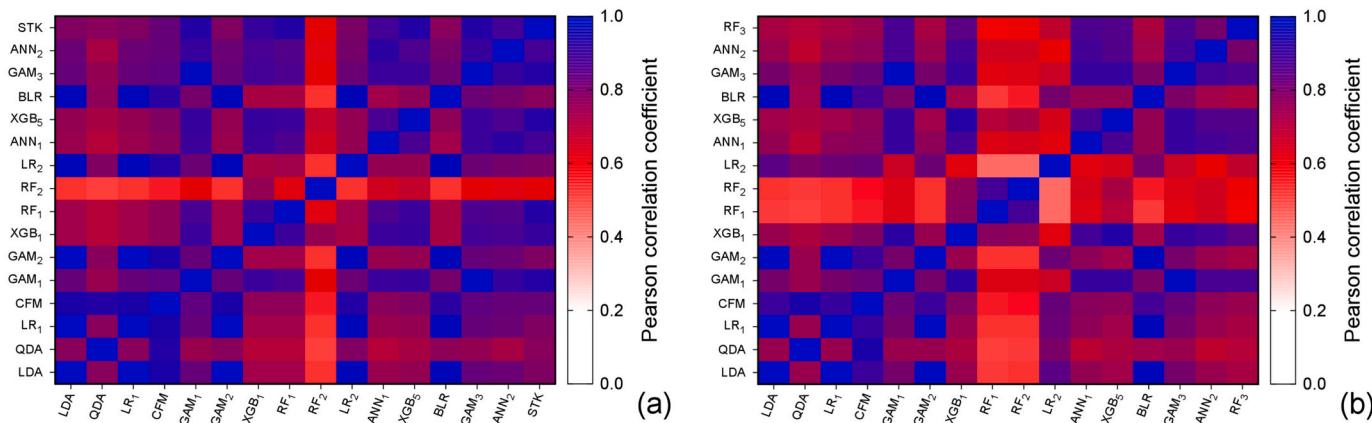


Fig. 13. As in Fig. 12, but for pairwise Pearson correlations.

values of class breaks give an alternative description of the susceptibility distributions, with respect to both maps (Figs. 17 and 18) and to box-plots (Fig. 19).

7. Discussion

We presented the first cross-examination of susceptibility modeling approaches with truly independent tests using a benchmark dataset. We surveyed interested contributors in the landslide science community, asking to prepare susceptibility maps with their methods of choice for a given dataset. The experiment was designed to be a rather specific one but, aside from quantitative calculations, we obtained a range of subjective responses on how to approach the problem, and the proposed experiment resulted in multiple outcomes.

Besides the use of different classification methods, the research groups participating in the experiment tried to answer alternative questions in addition to the simple, technical ones asked in the original survey. This triggered a second iteration in response to findings from the initial experiment, during which the dataset underwent modifications and a specific workflow was singled out.

We discuss separately the results corresponding to the preliminary dataset, presented in Sections 4.1–4.11, and the final version of the dataset, *i.e.*, our proposal for a benchmark dataset for landslide susceptibility assessment, presented in Section 5, with results in Section 6.

7.1. Discussion: preliminary dataset and results

The first part of the experiment, in which the contributors could play freely with methods and spatial variables, showed how the same algorithm could return quite different outputs, due to different model implementation techniques that mostly went unreported. This highlights the need to understand how these techniques influence model outcomes.

Variable selection was performed by 8 out of 11 participant groups, with several different methodologies, and the number of considered variables ranged from one (SU area) to all. We can distinguish different methods in two classes: (1) in a few cases groups would only use data to perform the selection, (2) in other cases, groups would also require performance assessment (LOO, VIF, other statistical tests), making the selection specific to the modeling approach. Selection based solely on data always consisted of removing collinear variables, although with different approaches for identifying them, mostly pairwise. These considerations helped in shaping the final dataset of the benchmark experiment, as described in Section 5.

One notable point was the relevance attributed to the SU area variable. In fact, in the original experiment, area was not even intended to be considered as a predictor. Nevertheless, G9 explicitly investigated the effect of using SU area alone, or in combination with other predictors, *versus* the case in which it was excluded. Slope unit area was a meaningful predictor for landslide occurrence. However, this correlation was assessed as a random effect rather than a causal relation. Group 7 discussed instead a related issue, namely the aggregation of variables over SU polygons (*i.e.*, zonal statistics: the process of calculating mean, standard deviation, and percentages) and the possible confounding effect of different SU area. Group 2 kept SU area in the analysis, treating it as any other covariate. The results of the performed k-fold CV showed a wide dispersion of the mean decrease in deviance explained (*i.e.*, variable importance) for both p₁ and p₂ (similar interquartile range from 2.7 % to 7.2 %), indicating a random effect of the variable on the model output.

At the data level, it is tricky to differentiate the SU area effect from potential causal contributions of other covariates because the SU area is inherently present in all covariates. Even after explicitly eliminating the area variable, SU area still controls the distribution of other covariates due to aggregation within SU polygons. This issue of aggregating covariates and target variables across non-uniform aerial units gives rise to the modifiable aerial unit problem, known as MAUP (Openshaw, 1984), mentioned in Alvioli et al. (2020).

We suggest that accounting for SU area effects could enhance the practical applicability and interpretation of the LSM. This is also relevant when assessing performance. We explore this point more fully later in this section. However, a model structure where landslide occurrence is represented as presence/absence does not seem ideal to visualize this relationship, and an investigation of the use of this variable in models targeting landslide counts *versus* landslide density is beyond the scope of this work and could be considered in a separate study.

Two participant groups considered the relevance of the geographical extent and location of the benchmark dataset. In fact, the dataset proposed here is an excerpt from a larger dataset used in a previous study at national level (Italy). The larger dataset, in turn, included landslide information from a national inventory. This triggered two different approaches from G2, which independently built a similar dataset for a different region in Northern Italy, and from G3, who compared the results for the benchmark subset with those at the national scale.

In the first case (G2), the average AUC_{ROC} values in the Italian Western Alps were higher than in Central Italy (up to 0.08 AUC_{ROC} points), but the average loss of performance between the training and test phases were lower in Central Italy than in the Alps (0.01 and 0.06 AUC_{ROC} points, respectively). Moreover, a different number of predictors was kept as significant in the two areas (26 in Central Italy, and only 9 in the Italian Western Alps). This indicates the need to develop more consistent models to account for spatial CV variability and to develop the model for the unique attributes within the study area.

The varying degree of completeness of the national inventory (Trigila et al., 2010) is discussed in great detail by Loche et al. (2022). This emphasizes that multiple benchmark datasets based on international open-source data are useful not only to investigate the comparative performance of mapping methods but also to explore and discuss geographical differences. For example, the landslides within this part of central Italy are dominated by translational slides. Consequently, the predictors included in this benchmark dataset were chosen to describe the attributes of the terrain that influence these landslide types within this region of Italy. Other variables can affect translational landslides, *e.g.*, terrain attitude or rock structure, which were not available for this study. Other geographic areas or other landslide types may require a different combination of predictors for a meaningful susceptibility model comparison to be conducted.

In the second case, G3 compared the national susceptibility maps (Loche et al., 2022) with the results of the benchmark dataset, and evaluated the response of fixed and random effects, as they were modeled in the same way in the two studies. They found that the non-linear variables (slope_aver, Max_Distan and D_sqrt_A) behave in the very same way in the two cases. Moreover, they calculated the Pearson correlation coefficient between both presence scenarios and all of Italy, which resulted in values of 0.81 and 0.83, for p₁ and p₂, respectively. This confirms the relationship between the susceptibility zonations applied across different scales, from our small study area for the benchmark dataset up to the entire nation of Italy.

G4 and 11 explored the effects of different CV methods. Group 4 implemented several strategies, including non-spatial CV, spatial block

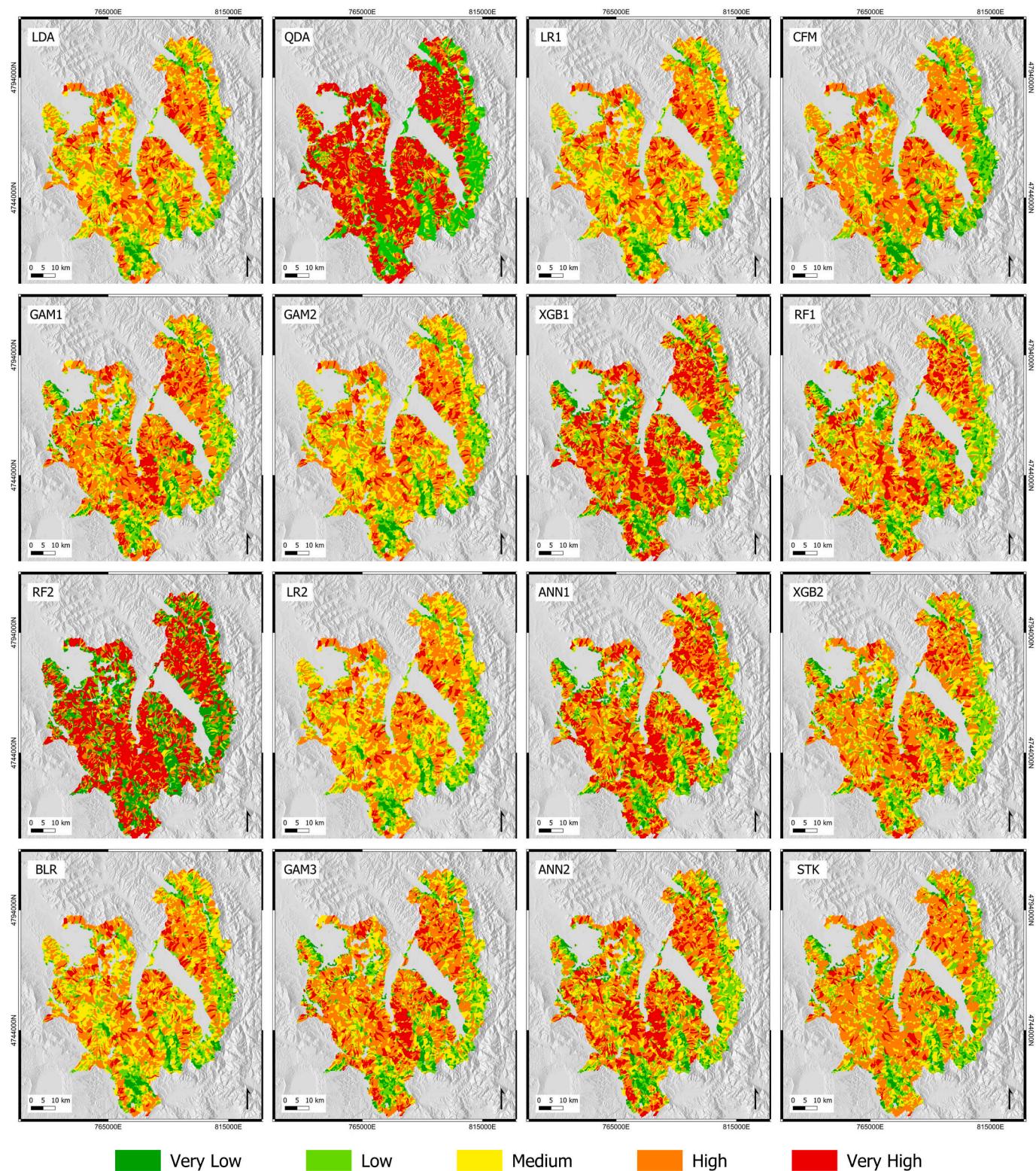


Fig. 17. Susceptibility maps corresponding to the case p1; model names as in Fig. 3 and Table 2. Shaded relief map as in Fig. 2.

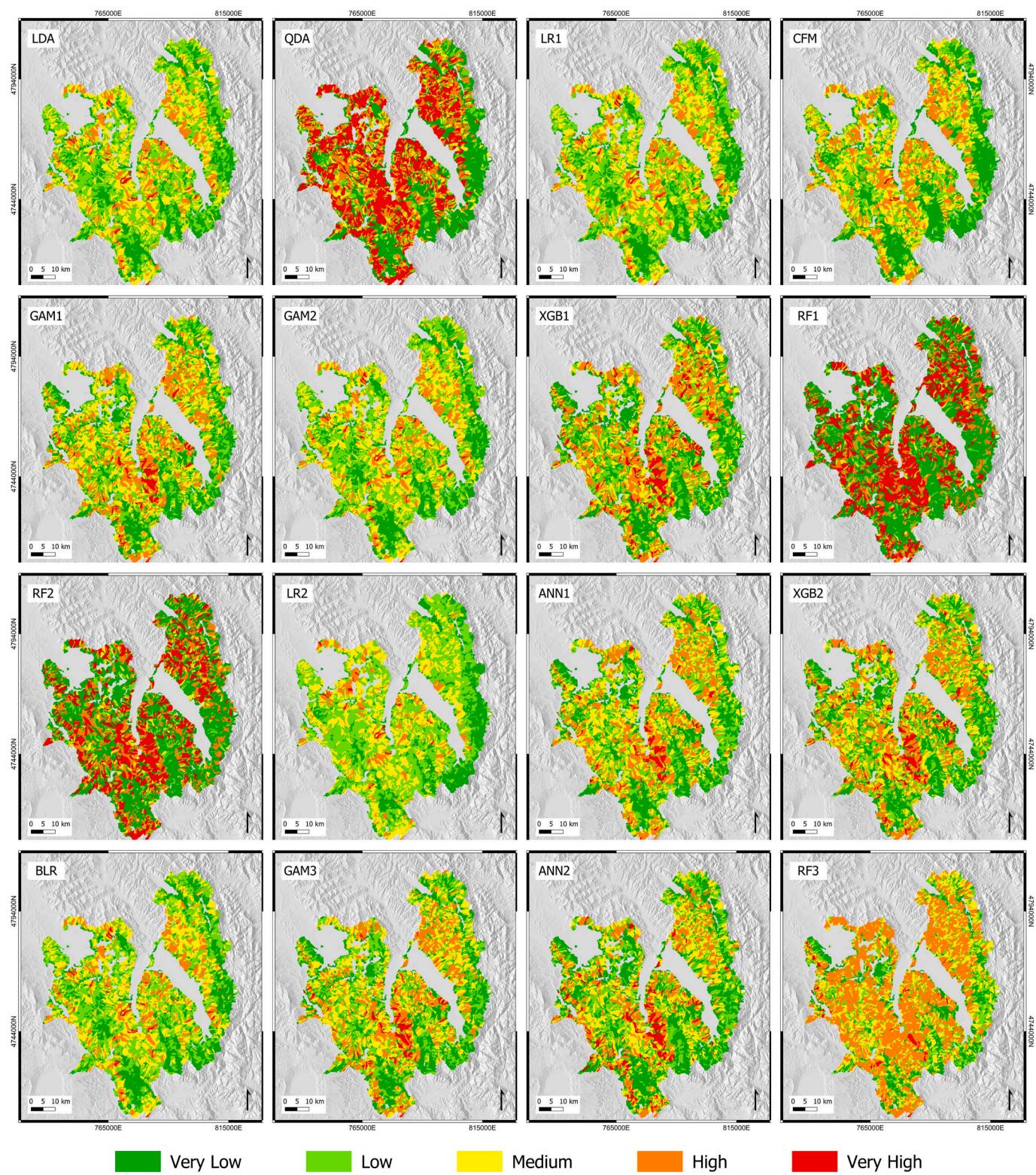


Fig. 18. Susceptibility maps corresponding to the case p2; model names as in Fig. 3 and Table 2. Shaded relief map as in Fig. 2.

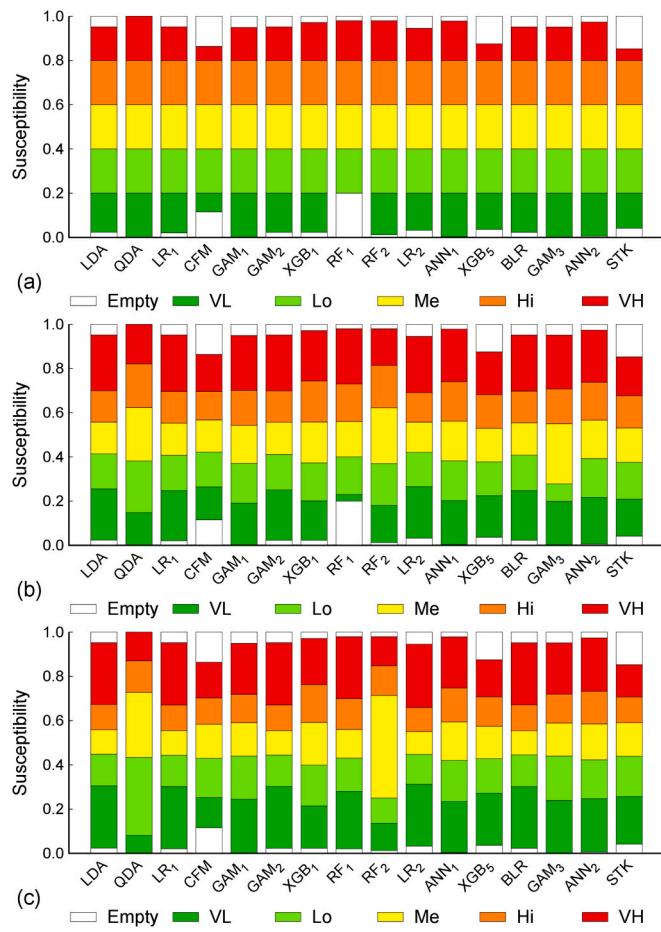


Fig. 20. Stacked bars representing different class breaks, for the p1 scenario, using (a) equal intervals (same breaks, different number of values in each class; used in Fig. 17), (b) Jenks natural breaks (different breaks and different number of values in each class), and (c) percentiles of the distributions (different breaks, but each class contains the same number of values, *i.e.*, of slope units). In all cases, we show the actual lower and upper limits in the distribution of susceptibility values with a white band.

point to not change the presence indicator in p2 may also increase the variance in the model results. This is because the attributes of SUs containing one or many landslides may be similar, hampering a model's ability to differentiate between landslide and non-landslide SUs. The p1 indicator avoids this issue by not allowing any landslide points within the SUs labeled as a non-landslide slope unit. The reduced area indicating landslide presence (Fig. 2) also explains the reduction in probabilities using the p2 flag. In summary, these results do not favor one method of categorizing SUs as containing a landslide. However, the chosen method may have notable effects on the model results.

Finally, we stress that we did not intend to select the “best” method here, as our only aim is to establish a benchmark dataset and workflow, that could be useful as a standard reference for calculations by other scholars. Different methods may be more useful in diverse settings, and different predictors may be available in other areas, with respect to those considered here. We note that the standardized workflow leveled out model performances, as Figs. 14–16 show, with respect to Figs. 6 and

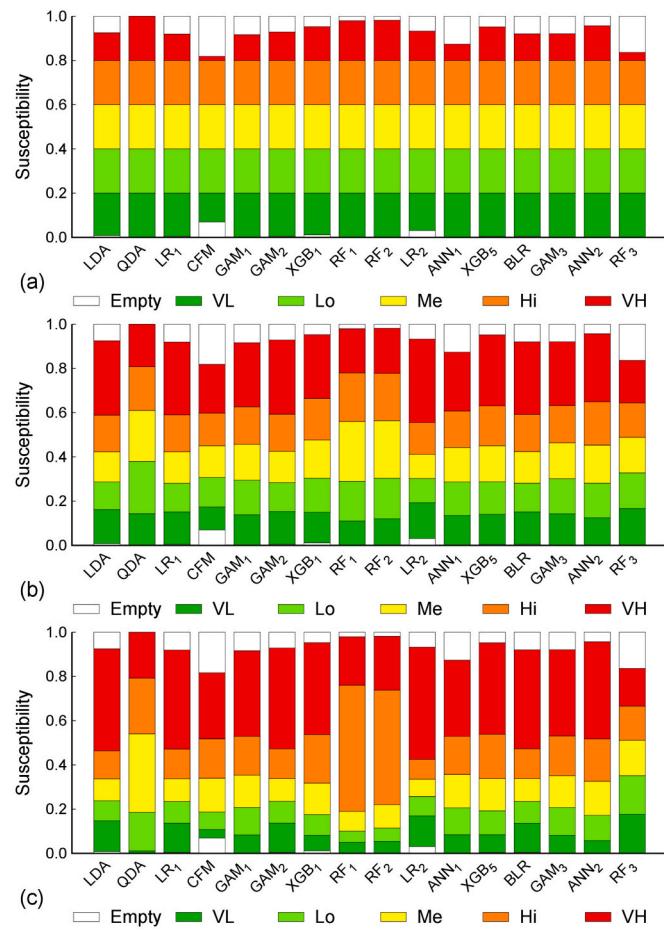


Fig. 21. As in Fig. 20, but for the p2 scenario; breaks in (a) correspond to the classification of Fig. 18.

7. Nevertheless, the spatial patterns in Figs. 17 and 18 show clear differences that one could not grasp from the numerical values of the performance metrics. The residual differences can be ascribed to different predictive abilities of the models in the study area presented here as a benchmark. Fig. 22 shows the variance of results across the 16 different methods, in each SU, colorized in five classes. We observe that variability is larger for the second scenario, and has a different spatial pattern, which highlights the relevance of the method adopted to quantify landslide presence, starting from point landslide locations.

The pronounced leveling with the final dataset combined with the observation that results are most similar within each contributing group highlights that user-caused variability in model performance and quality is significant. Different software packages, coding styles, user error, and unclear workflows can all lead to significant model differences that are difficult to parse out, as evidenced by step one of this experiment. This further highlights the benefit of a benchmark dataset for directly comparing future susceptibility modeling approaches and emphasizes that direct comparison between models produced by different researchers without a standardized workflow should be done with caution.

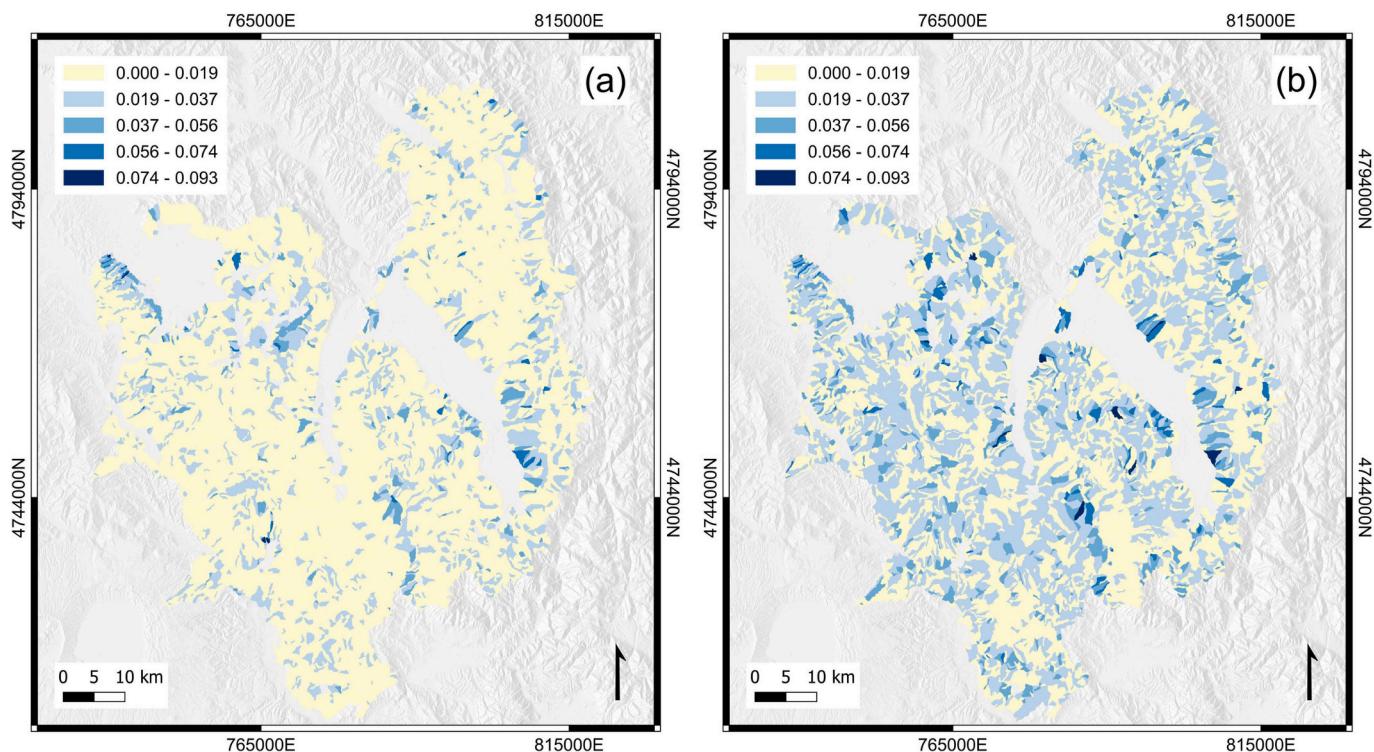


Fig. 22. Variance of the susceptibility values in each slope unit, obtained with the benchmark dataset proposed here for 16 different classification methods. (a) Variance of results shown in Fig. 17, for the landslide presence scenario p_1 , and (b) scenario p_2 , shown in Fig. 18. Shaded relief map as in Fig. 2.

8. Conclusions

We presented the first participatory experiment to systematically investigate current methods and practices in landslide susceptibility modeling. Our results clearly highlight the benefit for landslide scientists and practitioners to benchmark their results against known data, methods, and workflows.

We demonstrated that two fundamental steps for obtaining robust and reproducible LSMs are (1) a critical selection of model input and detailed rules for application of a classification method (referred to as “workflow” throughout this work), and (2) a critical evaluation of the output. Our findings also highlight the benefit not only of benchmark datasets, but also of very specific, unambiguous, and shared guidelines on how to build an LSM. These warrant being included in guidance ranging from data collection, data selection and pre-processing, to selection and application of methods and, eventually, to the assessment of results with proper metrics.

We designed the dataset to serve exactly this purpose, and the results of the experiment demonstrate the success of our objectives. Moreover, similar datasets and results for different regions of the world, specific to landslide science, would improve our understanding of landslide susceptibility predictions, and would enable their real-world applicability.

Data and code availability

The final benchmark dataset corresponds to the same SUs set as in the preliminary dataset, with a modified attribute table. In addition, we provide in separate vector layers the results for all of the methods used in step two of the experiment, used to prepare Figs. 11–21. We also share code to (1) calculate susceptibility value within the GAM model of Section 4.3, and (2) calculate AUC_{ROC} values and Brier scores from the attribute table of the vector layer with results.

The benchmark dataset, results, and code are publicly available for download at the main CNR IRPI SU project page, at <https://geomorphology.irpi.cnr.it/tools/slope-units>,

under the section Data → Benchmark Dataset. We provide the datasets in OGC GeoPackage vector format (GeoPackage is an open, standards-based, platform-independent, portable, self-describing, compact format for transferring geospatial information) and in ESRI Shapefile format. The maps are identical, both of them have the same attributes, and are in EPSG:32632 - WGS 84 / UTM zone 32 N projected reference system. Code is a combination of bash scripting, GRASS GIS, and R.

Authors contributions

M. Alvioli: designed the experiment, prepared the dataset (DD), collected results from participants, wrote the first draft (WD), revised the manuscript (RM). M. Loche: DD, contributed as G3, and RM. L. Jacobs: WD and RM. C. H. Grohmann: RM. M. T. Abraham: G10 and RM. K. Gupta and N. Satyam: G10. G. Scaringi: G3. T. Bornetxtea and M. Rossi: G1. I. Marchesini: G3 and RM. L. Lombardo: G3. M. Moreno: G9 and RM. S. Steger: G9. C. Camera: G2 and RM. G. Bajni: G2. G. Samodra, E. E. Wahyudi and N. Susyanto: G4 and RM. M. Sinčić: G5 and RM. S. Bernat Gazibara: G5. F. Sirbu: G6 and RM. J. Torizin and N. Schüßler: G7 and RM. B. Mirus and J. Woodard: G8, RM, proof-read English and took care of internal USGS revision. Héctor Aguilera Alonso: G11 and RM. J. S. Rivera-Rivera: G11.

Disclaimer

Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the U.S. Government.

Declaration of competing interest

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

- landslide susceptibility maps. *Geomorphology* 262, 8–23. <https://doi.org/10.1016/j.geomorph.2016.03.015>.
- Steger, S., Schmaltz, E., Glade, T., 2020. The (f)utility to account for pre-failure topography in data-driven landslide susceptibility modelling. *Geomorphology* 354, 107041. <https://doi.org/10.1016/j.geomorph.2020.107041>.
- Steger, S., Mair, V., Kofler, C., Pittore, M., Zebisch, M., Schneiderbauer, S., 2021. Correlation does not imply geomorphic causation in data-driven landslide susceptibility modelling — benefits of exploring landslide data collection effects. *Sci. Total Environ.* 776, 145935 <https://doi.org/10.1016/j.scitotenv.2021.145935>.
- Sterlacchini, S., Ballabio, C., Blahut, J., Masetti, M., Sorichetta, A., 2011. Spatial agreement of predicted patterns in landslide susceptibility maps. *Geomorphology* 125, 51–61. <https://doi.org/10.1016/j.geomorph.2010.09.004>.
- Süzen, M.L., Doyuran, V., 2004. Data driven bivariate landslide susceptibility assessment using geographical information systems: a method and application to Asarsuyu catchment, Turkey. *Eng. Geol.* 71, 303–321. [https://doi.org/10.1016/S0013-7952\(03\)00143-1](https://doi.org/10.1016/S0013-7952(03)00143-1).
- Thai Pham, B., Prakash, I., Dou, J., Singh, S., Phan Trong, T., Trung Tran, H., Minh Le, T., Van Phong, T., Dang Kim, K., Shirzadi, A., Tien Bui, D., 2020. A novel hybrid approach of landslide susceptibility modelling using rotation forest ensemble and different base classifiers. *Geocarto Int.* 35, 1267–1292. <https://doi.org/10.1080/10106049.2018.1559885>.
- Tien Bui, D., Thai Pham, B., Quoc Nguyen, P., Hoang, N.D., 2016. Spatial prediction of rainfall-induced shallow landslides using hybrid integration approach of Least-Squares Support Vector Machines and differential evolution optimization: a case study in Central Vietnam. *Int. J. Digit. Earth* 9, 1077–1097. <https://doi.org/10.1080/17538947.2016.1169561>.
- Torizin, J., Schüßler, N., 2023. Exploring the benchmark dataset for tasks related to landslide susceptibility assessment. In: EGU General Assembly 2023. <https://doi.org/10.5194/egusphere-egu23-13362>. Vienna, Austria, 24–28 April.
- Trigila, A., Iadanza, C., Spizzichino, D., 2010. Quality assessment of the Italian landslide inventory using GIS processing. *Landslides* 7, 455–470. <https://doi.org/10.1007/s10346-010-0213-0>.
- UCI, 2024. Machine learning repository. <https://archive.ics.uci.edu/>. Accessed: July 19, 2024.
- Wang, N., Cheng, W., Marconcini, M., Bachofer, F., Liu, C., Xiong, J., Lombardo, L., 2022. Space-time susceptibility modeling of hydro-morphological processes at the Chinese national scale. *Eng. Geol.* 301, 106586 <https://doi.org/10.1016/j.enggeo.2022.106586>.
- Wang, T., Dahal, A., Fang, Z., van Westen, C., Yin, K., Lombardo, L., 2024. From spatio-temporal landslide susceptibility to landslide risk forecast. *Geosci. Front.* 15, 101765.
- Wei, T., Simko, V., 2021. R Package ‘Corrplot’: Visualization of a Correlation Matrix. URL: <https://github.com/taiyun/corrplot>.
- Wislocki, A., Bentley, S., 1991. An expert system for landslide hazard and risk assessment. *Comput. Struct.* 40, 169–172. [https://doi.org/10.1016/0045-7949\(91\)90469-3](https://doi.org/10.1016/0045-7949(91)90469-3).
- Wood, S., 2017. Generalized Additive Models – An Introduction with R, 2nd edition. In: Mathematics & Statistics. Chapman and Hall/CRC. <https://doi.org/10.1201/9781315370279>.
- Wood, S.N., Augustin, N.H., 2002. GAMs with integrated model selection using penalized regression splines and applications to environmental modelling. *Ecol. Model.* 157, 157–177. [https://doi.org/10.1016/S0304-3800\(02\)00193-X](https://doi.org/10.1016/S0304-3800(02)00193-X).
- Yeon, Y.K., Han, J.G., Ryu, K.H., 2010. Landslide susceptibility mapping in Injae, Korea, using a decision tree. *Eng. Geol.* 116, 274–283. <https://doi.org/10.1016/j.enggeo.2010.09.009>.
- Yesilnacar, E., Topal, T., 2005. Landslide susceptibility mapping: a comparison of logistic regression and neural networks methods in a medium scale study, Hendek region (Turkey). *Eng. Geol.* 79, 251–266. <https://doi.org/10.1016/j.enggeo.2005.02.002>.
- Yong, C., Jinlong, D., Guo, F., Bin, T., Tao, Z., Hao, F., Li, W., Qinghua, Z., 2022. Review of landslide susceptibility assessment based on knowledge mapping. *Stoch. Env. Res. Risk A.* 36, 2399–2417. <https://doi.org/10.1007/s00477-021-02165-z>.
- Zeng, T., Wu, L., Peduto, D., Glade, T., Hayakawa, Y.S., Yin, K., 2023. Ensemble learning framework for landslide susceptibility mapping: different basic classifier and ensemble strategy. *Geosci. Front.* 14, 101645 <https://doi.org/10.1016/j.gsf.2023.101645>.
- Zuur, A.F., Ieno, E.N., Walker, N., Saveliev, A.A., Smith, G.M., 2009. Mixed Effects Models and Extensions in Ecology with R. Springer. <https://doi.org/10.1007/978-0-387-87458-6>.