1   # PMT: new analytical framework for automated evaluation of

2   # geo-environmental modelling approaches

3   Omid Rahmati [1,2*], Aiding Kornejady [3], Mahmood Samadi [4], Ravinesh C. Deo [5], Christian

4   Conoscenti [6], Luigi Lombardo [7], Kavina Dayal [8], Ruhollah Taghizadeh-Mehrjardi [9, 10], Hamid

5   Reza Pourghasemi [11], Sandeep Kumar [12], Dieu Tien Bui [13,14*]

6   [1] Geographic Information Science Research Group, Ton Duc Thang University, Ho Chi Minh City,

7   Vietnam

8   [2] Faculty of Environment and Labour Safety, Ton Duc Thang University, Ho Chi Minh City, Vietnam

9   [3] Young Researchers and Elite Club, Gorgan Branch, Islamic Azad University, Gorgan, Iran

10   [4] Faculty of Natural Resources, University of Tehran, Karaj, Iran

11   [5] School of Agricultural, Computational and Environmental Sciences, Centre for Sustainable Agricultural

12   Systems & Centre for Applied Climate Sciences, University of Southern Queensland, Springfield, QLD

13   4300, Australia

14   [6] Department of Earth and Marine Sciences (DISTEM), University of Palermo, Via Archirafi 22, 90123

15   Palermo, Italy

16   [7] Faculty of Geo-Information Science and Earth Observation (ITC), University of Twente, Enschede,

17   Netherlands

18   [8] CSIRO Agriculture and Food, 15 College Road, Sandy Bay, TAS 7005, Australia

19   [9] Department of Geosciences, Soil Science and Geomorphology, University of Tübingen, Tübingen,

20   Germany

21   [10] Faculty of Agriculture and Natural Resources, Ardakan University, Ardakan, Iran

22  [11] Department of Natural Resources and Environmental Engineering, College of Agriculture, Shiraz

23  University, Shiraz, Iran

24  [12] Department of Agronomy, Horticulture and Plant Science, South Dakota State University, Brookings,

25  USA

26  [13] Institute of Research and Development, Duy Tan University, Da Nang 550000, Vietnam

27  [14] Geographic Information System Group, Department of Business and IT, University of South-Eastern

28  Norway, N-3800 Bø i Telemark, Norway

29  * Corresponding authors' Email addresses: Orahmati68@gmail.com; Dieu.T.Bui@usn.no

30

31

# PMT: new analytical framework for automated evaluation of geo-environmental modelling approaches

32

33

## Abstract

34

35 Geospatial computation, data transformation to a relevant statistical software, and step-wise

36 quantitative performance assessment can be cumbersome, especially when considering that the

37 entire modelling procedure is repeatedly interrupted by several input/output steps, and the self-

38 consistency and self-adaptive response to the modelled data and the features therein are lost

39 while handling the data from different kinds of working environments. To date, an automated

40 and a comprehensive validation system, which includes both the cutoff-dependent and –

41 independent evaluation criteria for spatial modelling approaches, has not yet been developed for

42 GIS based methodologies. This study, for the first time, aims to fill this gap by designing and

43 evaluating a user-friendly model validation approach, denoted as Performance Measure Tool

44 (PMT), and developed using freely available Python programming platform. The considered

45 cutoff-dependent criteria include receiver operating characteristic (ROC) curve, success-rate

46 curve (SRC) and prediction-rate curve (PRC), whereas cutoff-independent consist of twenty-one

47 performance metrics such as efficiency, misclassification rate, false omission rate, F-score, threat

48 score, odds ratio, etc. To test the robustness of the developed tool, we applied it to a wide variety

49 of geo-environmental modelling approaches, especially in different countries, data, and spatial

50 contexts around the world including, the USA (soil digital modelling), Australia (drought risk

51 evaluation), Vietnam (landslide studies), Iran (flood studies), and Italy (gully erosion studies).

52 The newly proposed PMT is demonstrated to be capable of analyzing a wide range of

53    environmental modelling results, and provides inclusive performance evaluation metrics in a

54    relatively short time and user-convenient framework whilst each of the metrics is used to address

55    a particular aspect of the predictive model. Drawing on the inferences, a scenario-based protocol

56    for model performance evaluation is suggested.

57    *Keywords*: PMT, Spatial modelling, Goodness-of-fit, Validation, Performance analysis;

58    predictive model evaluation framework

59

## 60    Software and data availability

61    Name of tool:                    PMT (Performance Measure Tool)

62    Developers:                      Samadi M., Kornejady A., and Rahmati O.

63    Hardware required:               General-purpose computer (2 Gb RAM)

64    Software required:               ArcGIS 10.2

65    Programming languages:           Python$^©$ 2.7

66    Program size:                    120 KB

67    Availability and cost:           Freely available in GitHub

68    Web link:                        https://github.com/mahmoodsamadi/PMT

69    Year first available:            2018

70

## 71    1. Introduction

72    Spatially-applicable predictive models must include a mandatory step where different aspects

73    of the model performance can be quantitatively benchmarked. Without considering the

74    performance of such geo-environmental models, the users would not be confident about the

75    veracity of the modelling results, and is unlikely to utilize them for practical decision making

4

76 (Pullar and Springer, 2000; Glade, 2005; Beguería, 2006). The accuracy of predictive models,

77 which is a pertinent factor demonstrating the usefulness of the relevant models, can significantly

78 result in the misclassification costs of the approach depending on the error magnitudes and types

79 (Frattini et al., 2010). For example, in the modelling of natural hazards, the Error Type I (*i.e.*, false

80 positive) is likely to indicate that a stable part of a spatial region is classified as being unstable,

81 and therefore, it can lead to unnecessary control and risk mitigation measures that are

82 implemented. The Error Type II (*i.e.*, false negative) can imply that a given terrain unit is

83 susceptible to the hazard, and it can be incorrectly classified as being stable, and consequently,

84 this terrain region can be allowed to be occupied by people or infrastructure without a responsible

85 and actionable risk mitigation activity. These errors, if not assessed properly, can consequently

86 incur social and economic costs, depending on the vulnerability and economic value of the

87 elements at risk (*e.g.*, infrastructures, lives, etc.). In light of this need, a robust investigation of

88 such predictive errors in spatially-applicable models is highly warranted, to make the modelling

89 approaches and model results more viable for real-life usage, risk mitigation and implementation.

90     Over the past couple of decades, a number of susceptibility assessment models have been built,

91 each striving to portray the current and future spatial patterns of a specific phenomenon. Many

92 studies have included a "model comparison" or a "performance assessment" step that was aimed

93 to evaluate the spatial modelling result, and to select the most optimal spatially-relevant model.

94 These sorts of models, largely promulgated as an operational tool, have largely been reported in

95 different fields and applications, such as landslide susceptibility studies (*e.g.,* Kornejady et al.,

96 2017; Kavzoglu et al., 2019; Yan et al., 2019), flood susceptibility studies (*e.g.,* Rahmati and

97 Pourghasemi, 2017; Siahkamari et al., 2018; Choubin et al., 2019), forest fire modelling purposes

98 (*e.g.,* Arpaci et al., 2014;Tien Bui et al., 2017), groundwater potential modelling studies (*e.g.,*

99    Naghibi et al., 2017; Miraki et al., 2019), species distribution modelling tasks (*e.g.,* Bucklin et al.,

100   2015; Shabani et al., 2016; Quillfeldt et al., 2017), land subsidence modelling (*e.g.,* Abdollahi et

101   al., 2018; Ghorbanzadeh et al., 2018), soil digital mapping (*e.g.,* Minasny and McBratney, 2007;

102   Wiesmeier et al., 2011; Malone et al., 2017), gully-erosion susceptibility (*e.g.,* Akgün and Türk,

103   2011; Conoscenti et al., 2014; Garosi et al., 2018). The evaluation of predictive models with

104   different statistical metrics and their implemented approaches, especially in such a diverse range

105   of studies, clearly warrant automated and coherent scientific strategies where performance

106   evaluations are implemented by means of a universally acceptable and statistically robust tool.

107       A review of published literature in this respect reveals significant advancements in predictive

108   model performance evaluations where the context of application and the respective model type

109   were seen to play a pivotal role in how these evaluation tools were implemented. Recently, the

110   study of Rahmati and Pourghasemi (2018) compared the performance of ten different advanced

111   machine learning models for the modelling of landslide susceptibility, while the study of Fukuda

112   et al. (2013) applied and compared seven different data-driven models for developing species

113   distribution maps. These authors considered the receiver operating characteristic (ROC) curves

114   and a number of cutoff-dependent methods for judging the capability of their model, and

115   consequently, in preparing and transporting the results to their statistical software, although this

116   was a relatively time-consuming task. Particularly, one must note that when susceptibility maps

117   are supposed to be directly incorporated into land-use planning, the best performing model are

118   likely to be highly favored for practical decision-making tasks (Youssef et al., 2016; Siahkamari et

119   al., 2018). This is primarily because the model performance assessments provide immensely

120   useful insights into the optimal structure of such models, and the possibility of their practical

121   implementation for perceived risk mitigation (Van Westen, 2006).

122    Most performance evaluation metrics that are designed to evaluate the overall learning skill of

123    the predictive model, and the validity of the generated results from them are based on comparing

124    the predicted patterns in spatial models with the actual observation datasets (Chung and Fabbri,

125    2008). In a somewhat different approach to the traditional model evaluation approaches (*e.g.,*

126    graphical check of the model's susceptibility maps in respect to the ground-truth datasets), the new

127    generation of model performance metrics is mainly applicable for quantifying the traditional terms

128    and the models' functionality. According to a general consensus, the performance indices in a

129    predictive model can be classified into two different categories: cutoff-dependent metrics (*e.g.*,

130    Cohen's Kappa, sensitivity, and specificity) and the -independent metrics (*e.g.* receiver operating

131    characteristic, ROC method) (Frattini et al., 2010). These approaches have been used in a number

132    of spatial modelling sub-fields.

133    Meanwhile, there is little doubt that the ArcGIS software, by virtue of its wide flexibility,

134    portability and the relevance in spatial modelling approaches (e.g. geostatistics, mapping tools,

135    variogram, kriging, and local/global scale metrics), has been unceasingly used by many

136    researchers to implement the most basic as well as the more complex spatial functions and

137    statistical criterion that are available. In spite of this widespread usage of ArcGIS software as a

138    spatial modeling platform , the absence of a dedicated GIS-based tool and its non-availability to

139    aspiring researchers and practitioners who are outside of the major subscribed users and

140    institutions, is still very challenging (Scott and Janikas, 2010). Furthermore, the GIS users need to

141    employ cumbersome step-by-step procedures in order to calculate each of their performance

142    indices, and occasionally, they need to reach out for additional commercial and/or freely available

143    software platforms (*e.g.*, Microsoft Excel, SPSS, and R packages). These types of external model

144    evaluation frameworks and largely the expensive software that need to be used to analyze these

145    data outside of GIS platforms, represent a challenging task when aiming at optimizing any

146    modeling workflow.

147        In respect to these arguments for more robust evaluation of spatially-relevant predictive

148    models, some of the freely available software, such as the R package in the form of "cvTools"

149    (Alfons, 2012) or "CrossValidate" (Coombes, 2018), and the relevant modelling platforms in the

150    R software have partially satisfied the need to compute these metrics. However, these add-in tools

151    also seem to be relatively deficient in terms of their inclusiveness in the respective modelling

152    approach, and also sometimes, they may require additional external coding skills, which in some

153    cases may not available to the users. Furthermore, each of these add-in software are likely to

154    include only some of the cutoff-dependent and/or –independent evaluation criteria, and not include

155    the others (as necessary) within a universally desirable manner, and therefore, the external

156    software may be less flexible and attractive to the novice modeler and other non-scientific

157    stakeholders, practitioner and decision-makers.

158        To address inherent limitations posed by existing approaches adopted in the evaluation of

159    spatial models, this research study aims to propose and construct a new, robust and comprehensive

160    GIS-based package, denoted as the Performance Measure Tool (PMT), to scrutinize in a

161    statistically sound manner the performance of spatially-relevant predictive models. The merits of

162    the proposed PMT, augmented by its extensive validation in diverse regions, contextual

163    applications and global studies, are likely to enable modelers and risk mitigation practitioners to

164    calculate practically useful performance metrics (both cutoff-dependent and the –independent

165    category). The PMT is designed in such way that it has the ability to provide information in a

166    tabular and graphical format with a relatively simple platform and self-explanatory user interface.

167    This proposed tool is likely to be useful for any spatially-relevant model, various types of end-

168 users—from the beginner who are not familiar with advanced coding, to those who are

169 comfortable with a 'click-based procedure' and also practitioners in any scientific sub-field who

170 need to implement decisions about the model's versatility. To further ensure credibility and

171 generalizability of the software, the proposed PMT has been benchmarked rigorously to evaluate

172 its relative performance in different geo-environmental modelling contexts and in different parts

173 of the world including studies in Australia, Asia, Europe, and America.

## 2.      Basic Design Framework of the Performance Measure Tool

175    Implementing the notion that performance evaluation of a spatially-relevant predictive model

176 must be an important cornerstone of any spatial modelling attempt, in this study different cutoff-

177 dependent and –independent evaluation criteria, elaborated later in greater detail, have been

178 proposed. A brief review of recent literature shows that most of these analyses are underpinned by

179 a matrix-wise calculation, termed as the confusion matrix (and also, sometimes known as the table

180 of confusion, error matrix, or the matching matrix) and the contingency table (also known as the

181 cross tabulation or crosstab). Some researchers have interchangeably used these two names in

182 their studies and considered the confusion matrix largely as a special derivative of the contingency

183 matrix. Other researchers, however, pointed out a delicate, and logical difference, in that the

184 former is more suitable for evaluating the performance of different classifiers (*i.e.*, more common

185 in data mining models), while the latter is used to evaluate the rules of association and

186 interrelations between any two variables (Powers, 2011). However, the name "matching matrix" is

187 well-adapted in unsupervised machine learning algorithms, whereas the confusion matrix is used

188 in supervised learning (*i.e.*, input data fed by the training instances).

189    In this research, the confusion matrix has been considered as a way to describe the primary

190    basis for constructing the proposed spatially-relevant model evaluation tool. Consequently, a $2 \times 2$

191    confusion matrix is created where the rows are the instances in an actual class (*i.e.*, the

192    observations) and the columns are the instances in the predicted class, as illustrated in Table 1. As

193    the name "confusion" implies, the matrix is able to examine the degree of mislabeling one state (as

194    another) by means of directly comparing the predictions and the observations. The statistics

195    derived from the matrix are therefore all presented as either the row-wise (*e.g.*, positive and

196    negative predictive values) or the column-wise (*e.g.*, sensitivity and specificity) in the

197    implemented PMT tool.


198    **Table 1** HERE


199    It should be emphasized that the process and various stages of model performance assessments

200    can be rather a time-consuming and a complex task for the performance measures in a traditional

201    approach must be calculated separately using the geo-statistical techniques. This is particularly the

202    case for novice end-users (*e.g.*, risk mitigation practitioners who may be unfamiliar with various

203    mathematical and statistical knowledge). More importantly, to the best of the authors' knowledge,

204    there is hardly any reliable, comprehensive and end-user-friendly tool currently available that can

205    be used to consider the most relevant performance metrics, particularly in the widely adopted

206    ArcGIS environment. Considering this deficit, this paper aims to develop an efficient and

207    automated approach that operates in a quick, reliable and organized manner, and also presents a

208    relatively effective framework providing a user-friendly interface. The PMT has deliberately been

209    written in the freeware, the Python programming environment using a portability feature that

210    enables it to be installed easily within a geo-processing framework found in the ArcToolbox of the

211    ArcGIS 10.2 software.

212    Figure S1 (refer to supplementary information) illustrates the graphical user interface and the

213    execution process of the proposed PMT.

214                                          **Fig. S1** HERE

215    To illustrate the operational mechanism of the proposed PMT, one part of the Python code used

216    for calculating the evaluation criteria is displayed in Figure S2. The required inputs used to

217    execute the tool and the relevant outputs files are given in Tables 2 and 3, respectively. It is

218    important to note that the PMT extension allows the end-users to evaluate the accuracy of the

219    predictive model in both steps, composed of training/calibration and the validation phase. End-

220    users can also adopt both parts of the training and validation process to check the accuracy of their

221    predictive models, although investigating the accuracy of the model in the training step can also be

222    left unchecked in this particular tool. This option is added because most of the interest is usually

223    focused on the validation component, as it guarantees the viability of the model to be used for the

224    prediction and decision-making process. Conversely, calibration is a component uniquely voted to

225    build the reference model, and to evaluate the covariate effects, although these can be subjected to

226    some degree of overfitting (Lombardo et al., 2018).   These stages make the model easy-to-use

227    with no special skills required to run the proposed tool.

228                                          **Fig. S2** HERE

229                                          **Table 2** HERE

230                                          **Table 3** HERE

231    **3. Statistical background of the performance metrics**

232    **3.1. Confusion matrix**

233    In what follows next, the authors outline the kinds of information these metrics are able to

234    convey regarding the model performance. In order to construct a confusion matrix from a spatial

235    model, the users should define a cutoff (in percentile units) to split the spatial map into two

236    distinct classes in which the PMT can calculate the cutoff-dependent performance metrics. This is

237    the analogous operation to splitting a probability distribution into two distinct classes, although in

238    our case, this is performed directly within ArcGIS into map form. In this process, the first class

239    (*i.e.*, the lower percentage of susceptibility/ suitability map) is considered as the absence areas

240    (*e.g.,* the landslide-free areas) and the upper part as the presence locations (*e.g.*, the landslide

241    affected areas). For instance, let us assume a 50% cutoff for a landslide susceptibility map of

242    particular interest with 20 landslides located within the lower 50% (*i.e.*, low to moderate

243    susceptible areas). In this case, those 20 samples will be considered as error sources (denoted as

244    the 'false negative error', that has been discussed later) by the proposed tool and consequently, it

245    can reduce the performance of the predictive model since the landslides that have already occurred

246    are supposed to be located within the areas with the highest susceptibility values. The 50% cutoff

247    value is also quite common in existing literature, especially for the equally balanced

248    presence/absence datasets (*e.g.,* Lombardo and Mai, 2018). However, the prevalence can be

249    considered as the best alternative since it is able to represent the inherent predominance of a

250    phenomenon and it is not controlled by the experimenter. Additionally, quantifying the prevalence

251    of a natural phenomenon is somewhat problematic (discussed in Section *5.3*). Most of the data

252    mining models can circumvent this issue by calculating the prevalence by means of estimating the

253    best possible distribution of an event using generalized algorithms which is common in the

254    presence-only models (*e.g.* Maximum entropy model).

255    **3.2. Cutoff-dependent Approach**

256    Cutoff-dependent metrics include True Positive Rate (TPR), True Negative Rate (TNR), False

257    Positive Rate (FPR), False Negative Rate (FNR), Misclassification rate, Accuracy, Positive

258    Predictive Value (PPV), False Discovery Rate (FDR), Negative Predictive Value (NPV), False

259    Omission Rate (FOR), F-score, Matthews Correlation Coefficient (MCC), Informedness

260    (Bookmaker informedness; BM), Markedness (MK), Threat Score, Equitable threat score, True

261    skill statistic, Heidke's skill score, Odds ratio, Odd ratio skill score, and Cohen's kappa. Table 4

262    details the equations for all of the cutoff-dependent metrics.

263                                        **Table 4** HERE

264    The TPR, also termed as the sensitivity, recall, or hit rate, represents the probability of correctly

265    predicting the positives as observed in reality (given as True positives (TP) / total number of

266    positives (P)). The TNR, termed as the specificity, aims to quantify the probability of correctly

267    predicting the negatives as observed in reality (given as true negatives (TN)/ total number of

268    negatives (N)). The FPR, also known as the "1–specificity" or fall-out, aims to indicate the

269    probability of incorrectly predicting a non-event location as an event (given as false positives

270    (FP)/ total number of negatives (N)). Furthermore, the FNR, also denoted as the miss rate,

271    indicates the probability of incorrectly predicting an event location as a non-event (given as false

272    negatives (FN)/ total number of positives (P)). This quantity is used to express how often the

273    model wrongly predicts absences. Misclassification rate undertakes both the false negative and

274    false positive values and therefore reflects an overall error rate ((FP+FN)/total). The accuracy (or

275    the model efficiency) is the opposite metric compared to the misclassification rate, since it is able

276    to highlights the overall success of the predictive model ((*i.e.*, TP+TN)/total). Overall, this metric

277    shows how often the predictive model is correct. The PPV, also denoted as the confidence or the

278    precision in data mining approaches, or as Powers (2011) analogously calls it as the accuracy of

279    predicted positives, is used to measure the proportion of predicted presences that correctly

280    represent the real presence. As a complement component of the PPV, a false discovery rate is

281    applied to conceptualize the Type I errors (*i.e.*, rejection of a true null hypothesis) (Benjamini

282    and Hochberg, 1995). In accordance with the PPV, the NPV is used to measure the precision of

283    the predictive model in predicting the absence (or non-event) locations. However, this metric

284    largely ignores how well the model is able to handle the presence locations and that the FOR

285    simply is the complement of the NPV. The F-score is also called the harmonic mean of the

286    precision and the recall (i.e., sensitivity) where it reaches its best values at 1 (*i.e.*, best precision

287    and recall) and the worst at 0. In essence, MCC is a correlation coefficient metric computed

288    between the observed and the predicted binary classifications, and it is able to undertake a true

289    and a false positive and negative value. The terms *informedness* and *markedness*, implemented in

290    the PMT, were introduced initially by Powers (2011). Informedness, however, is likely to be the

291    only unbiased indicator in the confusion matrix and it measures the probability that an informed

292    decision that is being made rather than guessing, either the correct or the incorrect decision (due

293    to overtraining, atypical data, or even deliberately) (Powers, 2011). Markedness, also referred to

294    as *deltaP* in psychology, is the complementary pair of informedness indicating the probability

295    that an outcome is marked by the predictor (marker). Threat Score also penalizes the rare events

296    since some success of correct predictions of a less frequent event might be resulted out of

297    random chance. Although Threat Score uses different statistics in conjunction, the actual sources

298    of misclassification error are not discernible. Equitable Threat Score also known as the Gilbert's

299    skill score (Gilbert, 1884; Schaefer, 1990), the equitable threat score functions as per above

300    based on critical success score, but it is also used to eliminate the hit rates (*i.e.*, true positive

301    rates) originated by random chance. True skill statistic (TSS) (also called the Hanssen and

302    Kuipers discriminant or Pierces skill score), is applied to measure the ability of a predicted value

303    to discriminate between the events and the non-events, using all of the elements in the confusion

304    matrix. The Heidke's Skill Score operates according to the accuracy level but it is also used to

305    improve its meaning by eliminating the true positive rates that would be expected to occur by

306    chance (Heidke, 1926). Odd Ratio is used to measure the odds that an event (or an outcome) will

307    occur given a particular exposure, compared to the odds of the event occurring in the absence of

308    that exposure (Pepe et al., 2004). Odd Ratio Skill Score (also known as the Yule's Q) rescales

309    the values of the odds ratio into the -1 and the +1 range. In addition, Kappa is essentially a

310    measure of how well the model has performed as compared to how well it would have performed

311    purely by chance, and this would enable the modeler to better understand the true outcome of the

312    model in respect to the random occurrence of that value.**3.3. Cutoff-independent approach**

313    This approach, included in the PMT, includes two different methods that can be categorized as:

314    (1) receiver operating characteristic (ROC) curve, and (2) success-rate curve (SRC) and

315    prediction-rate curve (PRC).

316    **3.3.1. ROC curve**

317    The ROC curve, used typically in risk assessment through predictive model results, simply

318    plots the sensitivity (*i.e.*, true positive rates) on the *Y*-axis against the 1–specificity (*i.e.*, false

319    positive rate) on the *X*-axis (Gorsevski et al., 2006). The area under the ROC curve (denoted as

320    AUROC, bounded by [0, 1]), is the actual measure of the model evaluation since it generates a

321    quantitative value of the performance (Pontius and Schneider, 2001; Mas et al., 2013; Swets,

322    2014). The closer the AUROC is to unity, the better is the performance. The ROC curve can be

323    interpreted differently depending on the dataset; it can address the learning capability (or the so-

324 called goodness-of-fit) of the model if the training set is used for plotting; it can also infer the

325 predictive skill of the model if the validation set is used (Fawcett, 2006; Lombardo and Mai 2018).

326     In this regard, the proportion between training and validation samples is highly relevant. A

327 70:30% split is quite common among the researchers (Pradhan and Lee, 2010). Although different

328 partitions have also been used, such as 80:20% (*e.g.* Lipovetsky, 2009), 70:30% (*e.g.* Choubin et

329 al., 2019) or even 50:50% (*e.g.* Deo et al., 2016; Deo et al., 2017), there is no empirical consensus

330 on the best partition since this is more of an expert-user based decision. Irrespective of this, having

331 a large amount of inventory data (*i.e.*, number of events), one can assign a greater percentage of

332 such data to train the predictive model and a lesser percentage for validation. Opting for a suitable

333 approach to partition the training and validation sets is yet another crucial matter that has been the

334 subject of many studies, e.g. Kornejady et al. (2017). In this regard, the random sampling, self-

335 organizing maps for input selection, Mahalanobis distance, excerpting separate training/validation

336 areas, and temporal partitioning are all some of the common sample partitioning approaches. For

337 more details, readers can refer to the references therein.

### 338   3.3.2. Success-Rate Curve (SRC) and Prediction-Rate Curve (PRC)

339     The SRC is a measure of the learning capability of the model, while the PRC is able to examine

340 the predictive power. Although the SRC and the PRC may share some common features with the

341 ROC, the ROC in particular uses almost all the elements of the confusion matrix. This includes

342 positive (TPR and TNR) and negative (FPR and FNR) aspects of the model, while the SRC and

343 the PRC are calculated independently from the confusion matrix. In fact, the SRC represents the

344 cumulative areal percentage of the susceptibility classes (*i.e*, from the highest values to the lowest)

345 on the *X*-axis against the areal cumulative percentage of the training set located within those

16

346 susceptibility classes on the Y-axis (Chung, 2006; Blahut et al., 2010). In terms of its physical

347 interpretation, a steeper SRC curve is used to indicate that more events fall within the highly

348 susceptible classes; *i.e.*, a good learning skill. The PRC curve, however, follows the same plotting

349 process as the SRC, but the training data are replaced by the validation set.

**4. Testing the Efficacy of PMT: Selected Case Studies**

351 In this section, the proposed PMT is applied to 5 distinct, real geo-environmental modelling

352 tasks and case studies in order to robustly investigate its credibility and generalizability, and also

353 to demonstrate the potential benefits in considering different evaluation criteria promulgated by

354 the PMT. It is imperative to note that the selected case studies exhibited various noticeable

355 characteristics in terms of the issue under investigation, the modelling strategies, the overall

356 frameworks and the predictive model type, spatial or temporal scales considered and the

357 geographical and climatic conditions that influence the results and implementation of the model.

358 To provide a robust evaluation of the proposed PMT, the most relevant and a relatively diverse

359 range of data sets were obtained from most recently conducted research studies and also some

360 newly implemented models based on: (1) gully erosion prediction mapping in two small

361 catchments of central-western Sicily, Italy (Conoscenti et al., 2018) (2) flood hazard modelling in

362 the Galikesh region, Iran (Rahmati and Pourghasemi, 2017) (3) drought risk modelling in south-

363 east Queensland, Australia (Dayal, 2018; Dayal et al. 2018) (4) landslide susceptibility modelling

364 in the Kon Tum province, Vietnam (5) soil digital modelling in South Dakota, USA (Fig. 1). Each

365 of these studies employed a range of geo-spatial models where the PMT is used to provide a

366 consolidated assessment of its efficacy in providing greater insights into the practicality of the

367 modelling various frameworks.

368    An overall description of the study areas and the applied models are provided as follows whereas

369    further details of the modelling approaches are provided in the references therein.

370    A detailed flowchart of the various studies is shown in Fig. 2.

371                                    **Fig. 1** HERE

372                                    **Fig. 2** HERE

373    **4.1. Gully Erosion Modelling (Italy)**

374    Intense farming activities in two small catchments of central-western Sicily, Italy, have

375    expedited many erosion processes. In particular, the gully erosion has led to the landscape

376    dissection and massive soil loss (Conoscenti et al., 2018). The gullies in the study area have

377    developed as a result of the interrelation of several geo-environmental factors and human activities

378    such as access roads, parcel borders, wheel tracks, and plow furrows. In addition to the

379    multivariate adaptive regression splines (MARS) model already utilized by Conoscenti et al.

380    (2018) for gully erosion prediction mapping, in this paper we used the generalized linear model

381    (GLM) to conduct a fair comparison of their approach (Fig. 3).

382                                    **Fig. 3** HERE

383    **4.2. Flood Hazard Modelling (Iran)**

384    Over the last few decades, the Galikesh region, located in the Golestan province, in the north-

385    east of Iran, has witnessed severe flood events due to the particular climatic and topo-hydrological

386    conditions that resulted in many economic losses and causalities attributable to environmental

387    mismanagement (*e.g.*, deforestation, overgrazing, and over-exploitation). Since flood-inundation

388    has been one of the major issues of the urban areas in Golestan province for decades, Rahmati and

18

389 Pourghasemi (2017) used evidential belief function (EBF) to investigate the flood-prone hotspots

390 (Fig. 4). In this paper, we have implemented the proposed PMT as a statistical and decision-

391 support tool to provide an inclusive performance evaluation of their model.

392                                        **Fig. 4** HERE

393 **4.3. Drought Risk Modelling (Australia)**

394     An area located in the south-east of Queensland, Australia, encompasses intensive agricultural

395 activities, such as grazing, horticulture, and animal production, other than the densely populated

396 localities, which require a reliable water supply. As the study area is affected by severe and

397 frequent drought events, Dayal (2018) and Dayal et al. (2018) attempted to develop a spatial

398 drought risk map by employing the Bayes' theorem (*i.e.*, classifying spatial indicators), fuzzy

399 logic (*i.e.*, standardizing spatial indicators), and fuzzy GAMMA overlay (*i.e.*, aggregating drought

400 vulnerability, exposure, and hazard indices) technique (Fig. 5). Employing the findings of that

401 study, in this paper we utilized their final drought risk map as a potential input to the proposed

402 PMT, enabling us to examine the different aspects of its performance over the geospatial scale. In

403 order to investigate the influence of the cutoff values on the performance analysis, three different

404 cutoffs, *i.e.*, 50%, 70%, and 90% were taken into account and the results were compared, as

405 illustrated in Fig. 6.

406                                        **Fig. 5** HERE

407                                        **Fig. 6** HERE

408 **4.4. Landslide Susceptibility Modelling (Vietnam)**

19

409       Landslides are the dominant geo-hazardous elements in the Kon Tum province of Vietnam.

410 Hence, this study has used two novel data mining models including maximum entropy (MaxEnt)

411 and a recently developed model named as BayGmmKda (Bayesian-based ensemble of Gaussian

412 mixture model and radial-basis-function Fisher discriminant analysis) (Tien Bui and Hoang, 2017)

413 (Fig. 7). This study also uses the proposed PMT to highlight the potential asymmetries among the

414 performance metrics.

415                                                 **Fig. 7** HERE

416    **4.5. Soil Digital Modelling (USA)**

417       Soil digital modelling has received significant attention amongst scientists in recent years,

418 where computer-assisted pedometric-predictive mapping of soil properties has led to the creation

419 of an inclusive geographically-referenced soil database. To this end, an attempt is carried out to

420 map the soil bulk density (BD) predictive distribution in South Dakota, USA, by obtaining soil

421 bulk density samples of the study area and using two data mining models, namely the artificial

422 neural network (ANN) and decision tree (DAT) (Fig. 8). We have delineated the need for

423 rendering quantitative suitability maps into probability values to be able to use the proposed PMT

424 for further assessing the models' performance. In general, there is a few differences between

425 models' requirements. For example, DAT model does not require a separate dataset to optimize

426 parameters and just uses the training dataset for model building (i.e., learning and predicting),

427 whereas ANN model uses both the training and validation datasets for model building, validation,

428 and reevaluation and tuning parameters. Therefore, in ANN model, soil inventory dataset was

429 divided into three subsets: training (50% of input data) and 25% each for validation and testing.

430 For comparison sake, the same 25% testing dataset was kept in a vault and used for assessing the

431 generalization power of both the ANN and DAT models.

432

## 5. Results and Discussion

435    The following results and the subsequent discussions are based on Table 5, containing all the

436    previously-described performance metrics that have been calculated by means of the newly

437    proposed GIS-based PMT extension system. After a preliminary diagnosis of the models in each

438    of the aforementioned case studies, a detailed comparison of the performance metrics is provided.[1]

439    

### 5.1. Gully Erosion Modelling, Italy

441    According to the AUROC values, both the GLM and the MARS model show excellent

442    performance where the differences in the AUROC values were almost negligible. According to

443    Conoscenti et al. (2018), the excellent performance of these two models is indebted to a well-

444    investigation of the gullies in the study area and opting the main controlling factors that best

445    defined the occurrence mechanism. This process has been carried out by building a base model

446    comprised of the slope gradient and the contribution area and is then fed by nine other geo-

447    environmental factors one at a time. Moreover, the exemplary features of the chosen model have

448    also led to a significantly good performance, defined by measures such as the handling of all types

449    of factors (*i.e.*, both categorical and continuous) and well detecting the interactions among the

450    factors and also between the factors and the response event. Notably, Gómez-Gutiérrez et al.

---

[1] Note: the discussion provided here follows a particular way as the inferences derived from each case study is modified or reemphasized perpetually on the basis of the collective information obtained from different case studies and modelling scenarios. It is tried to be err on the side of caution to avoid raising any misleading points and engaging in dogmatic defense of one approach to the detriment of another.

451 (2015) also applied the MARS model to predict the gully occurrence in a relatively close (ca. 85

452 km) catchment with similar characteristics; however, the AUROC values stood at the range of

453 about 0.75-0.85, which was lower than that of Conoscenti et al. (2018). This highlights the

454 importance of making a well-structured input data and the calibration/ validation techniques. To

455 this point, both models seem to have rather similar performances.

456 However, a greater discrimination between models become apparent, as present in the results,

457 after breaking down these overall precision metrics into smaller components (*i.e.*, considering

458 simpler indices) that explain the efficacy of the approach more elaborately. Considering the

459 misclassification rate of both models, it is evident that the GLM approach has most likely

460 misclassified the presence and the absence more than the MARS model. Also, accuracy, as

461 understood to be the opposite concept of misclassification rate, attested the same pattern, where

462 the MARS model exhibited a higher accuracy in the classification of the presence and the absence

463 localities generated by the spatially-relevant model.

464 Further exploring the confusion matrix, it becomes evident that the higher value of the

465 misclassification rate in the GLM approach is directly rooted in the false negative rate. That is, the

466 GLM approach appears to have misclassified a number of 'presence locations' as the 'absence

467 locations' (in fact, this happened almost 13 folds greater than the MARS model). This indicates

468 that the GLM approach has somewhat failed to locate the gullies in notable study areas, and

469 therefore, may require further careful consideration prior to its application for real-life decisions.

470 In fact, the present analysis shows that this error appears to have also spread out to the other

471 metrics such as the sensitivity, F-score, NPV, and the FOR. The reason for the high AUROC value

472 for the GLM approach is plausibly due to that the latter is a cutoff independent metric, while the

473 confusion matrix elements have been calculated based on a 50% cutoff value. However, this does

474    not justify the GLM's underperformance at misclassifying the absence locations, since both

475    predictive models are compared under the same situation.

476    As explained in the *Theory* section, in such situations, the MCC may be the best representative

477    of the model's performance regarding the agreement between the observations and predictions.

478    One reason for this is because, as opposed to AUROC, AUSRC, and AUPRC, the MCC values the

479    cost of error and attempts to avoid to circumvent or truncate any error sources. Expectedly, the

480    MCC has well differentiated the performance of both MARS and GLM approaches, where the

481    MARS model with a value close to 1 almost represents a perfect model, while the GLM approach

482    with a value below 0.5 has shown a lesser degree of agreement between the observations and

483    predictions. This notion raises the possibility of some randomness (*i.e.*, being closer to zero). The

484    underperformance of the GLM approach highlights the disadvantages of using a predictive model

485    that is built on linear functions. Such a model is largely incapable of considering the nonlinear

486    interactions between the causal factors and the response event, may be sensitive to the number of

487    predictors, and more importantly, it could be sensitive to the outliers which are robustly handled

488    by non-linear basis functions in the MARS model. Given that the asymmetries of the cutoff-

489    dependent and –independent metrics are now more evident, a greater degree of scrutinization is

490    perhaps required, as provided by a more extensive discussion in the following real-life case

491    studies.

492    **5.2. Flood hazard modelling, Iran**

493    Recently, Evidential Belief Function (EBF), as a bivariate statistical model underpinned by the

494    Dempster-Shafer theory (Shafer 1976), has been adopted for flood inundation and susceptibility

495    mapping in Iran (Rahmati and Pourghasemi, 2017). Starting with the AUROC values, the overall

496  performance is acceptable, with respectively, 0.86 as the learning capability (obtained from the

497  training set) and 0.78 as a predictive skill (obtained from the validation set). Higher learning skill

498  compared to the predictive capability is common, and generally expected since the model's

499  parameters have been calibrated on a much larger data sample compared to the validation set.

500  However, this might question the possibility of overfitting, where a statistical model begins to

501  describe the random error in the data rather than the relationships between variables; that is, the

502  model becomes accustomed to the pre-used set of data. In this regard, simple statistical

503  assumptions have been identified as one of the main sources of overfitting issues, especially in

504  bivariate statistical models. This can negatively influence the generalization power and the

505  transferability of the model's results to the validation set/ areas/ time periods.

506  Considering the results presented here, all of the favorable qualities of the model (*i.e.*, all the

507  performance metrics highlighting the success of the model) have deteriorated to some extent in the

508  model validation stage. Although according to the AUROC classifications provided by Hosmer

509  and Lemeshow (2000), the values greater than 0.7 and 0.8, respectively, indicate an acceptable and

510  excellent performances, which in turn somewhat addresses the possibility of overfitting. This is

511  also evident in the AUSRC and AUPRC values, indicating that the predictive model is

512  respectively well-performing in terms of both the learning capability and the predictive skill. As

513  for the AUSRC and AUPRC values, the differences are discernable when compared to the

514  training- and validation-derived AUROC values. These differences are conceivable, given that the

515  AUSRC and AUPRC are calculated merely based on the presence localities. Therefore, by using

516  the AUSRC and AUPRC, the potential error sources (*i.e.*, polluting the presence population to

517  some absences which are incorrectly classified as positives) are left unclear and some degree of

518  success (*i.e.*, correctly detecting the absence locations) are also not acknowledged and not

519   included in the final calculation. This makes using the AUSRC and AUPRC less favorable to use

520   due to their erroneous behavior (Frattini et al., 2010).

521     A closer scrutinization appears to shed more light on the randomly-driven performances and

522   consequently, the weakness of the model or the input data. Considering the MCC—so far

523   suggested as an all-inclusive metric in this study—the values greater than zero (*i.e.*, random

524   agreement) reveals a promising level of precision; however, the values may not be high enough

525   (*i.e.,* far from a perfect precision to be certain of a non-random performance. In particular, the

526   level of disagreement between the observed and the predicted values appears to increase in the

527   validation stage. Other comprehensive measures, such as the true skill statistic, informedness, and

528   markedness are also in concurrence with the MCC value.

529     The Heidke's Skill Score, well-known for providing a robust accuracy value by diminishing the

530   TPR values generated by random chance, shows how the preliminary accuracy values (i.e.,

531   efficiency) is likely to decay. Similarly, the Cohen's Kappa aims to address the random aspect of

532   the model performance and provides new values in agreement with the latter. However, as stated

533   in our recent discussion, one should be cautious when using the cutoff-dependent metrics.

534   Drawing relevance from a report given by Frattini et al. (2010), the score-based metrics, despite

535   providing valuable insights, highly relies on certain cutoff values. That is, different cutoff values

536   might result in different performance values. However, this assumption still does not contradict

537   using the score-based metrics for a comparison purpose, since, as stated above, all the predictive

538   models were supposed to be compared under the same cutoff value(s) (*e.g.*, the Italian case study).

539   To test this concern, we have applied three different cutoffs for assessing the performance of a

540   drought risk map developed in the south-east region of Queensland, Australia.

541    To elaborate further, we provide two assumptions regarding the reduction in the accuracy of the

542    EBF metric. The first assumption pertains to the model's structure. Bivariate statistical models

543    have long been criticized for ignoring the interactions among the predictors, which can have direct

544    (and largely negative) influence on both the learning and the predictive skills. Moreover, as stated

545    by Ruspini et al. (1992), and more recently Reineking (2014), a need for categorizing factors with

546    continuous nature and also presenting a generalized probabilistic reasoning limit the application of

547    the EBF metric only to some specific problems (*e.g.*, detecting the uncertainty sources) rather than

548    a general use. However, a review of the previous work of Rahmati and Pourghasemi (2017)

549    reveals that the two other well-known data mining models (*i.e.*, boosted regression trees and the

550    random forest) have been used in addition to the EBF and surprisingly, we noted that the EBF

551    outperformed both of the data mining models, although the differences were negligible (*i.e.*,

552    AUROCs= 0.73-0.78), which leaves us with the second hypothesis.

553    Regarding the latter, the input data can be responsible for such limited performances of all three

554    models. Reviewing the model input data shows that only 63 flooding points were used as an input

555    for the modelling process in the period of 2001-2009, let alone that they were categorized into two

556    sets of 47 (training) and 16 (validation) locations which seems to be rather small to build a proper

557    predictive model. Complementing the inventory map by collecting more data from a broader time

558    period would provide a larger information matrix for the models to rely on. This highlights a note

559    given by Ruspini et al. (1992); "*the alleged lack of decision-support and counterintuitive nature of*

560    *evidential belief models, in fact, indicates the lack of basic informational shortcomings*".

561    **5.3. Drought Risk Spatial Attribution and Modelling, Australia**

562    For a drought risk map produced in the south-east of Queensland, Australia, the following

563    inferences can be derived from the validation stage only in order to focus on the alteration of the

564    performance metric values. The question mentioned above regarding the liability of the cutoff-

565    dependent metrics is answered by means of producing three cutoffs thresholds, *i.e.*, 50%, 70%,

566    and 90%, and then comparing these results.

567    It was evident that the AUROC and AUSRC expectedly yielded intact performance values

568    through all of the three cutoffs (Table 5). Based on this, the predictive skill of the fuzzy model

569    appears to be well performing. However, the values of all the cutoff-dependent metrics drastically

570    change at each cutoff. It is evident that by a transition from 50% to 90% cutoff, the area of danger

571    zone appears to shrink (as illustrated in Fig. 8). Moreover, at each cutoff threshold, a different

572    population of the negatives and the positives appears to fall within the safe and danger zones.

573    The direct impact of these transitions on the results is transparent in Table 5.  As appears,

574    moving from 50% to 70% cutoff, the *FN* error decreases to a certain level and adds to the *TN*,

575    serving as an advantage point for the model, while the false positives and true positives have

576    remained intact. Moreover, a vivid increase is also discernible in the values of the cutoff-

577    dependent metrics. However, another step towards the 90% cutoff backfires, where—similar to the

578    previous transition—although the *FN* value decreases and adds to *TN*, most of the *TP* population

579    migrates to *FP* category. This expectedly decreases the values of some cutoff-dependent metrics

580    such as F-score and PPV. Although 70% cutoff performed better than 50% and 90% cutoffs. Such

581    a choice would not be advisable for the other study areas and certainly not for the other predictive

582    models, because it is only in favor of this particular predictive model and the specific distribution

583    of the positive/negative points throughout the study area.

584    As previously mentioned in the *Theory* section, the only suitable substitute for the cutoff value

585    is the prevalence of the phenomenon, which again is difficult to ascertain, unless one constructs an

586    inclusive archive of the 'presence-absence locations' by visiting numerous sites. This type of data

587    compilation is more common in species distribution assessment, whereas, in natural hazard-related

588    studies, extracting absence locations are executed as an additional stage after inventory mapping,

589    based on random selection or other analytical strategies. Drawing on these inferences, it is

590    reasonable to ascertain that using cutoff-dependent performance metrics may not be practical for

591    individual model assessment, unless it is accompanied by mentioning the cutoff value from which

592    the metrics' values are extracted (*i.e.*, 50% for Iran, Italy, and all the following case studies), or it

593    is carried out by setting the prevalence as the cutoff value.

594    As with the case of Iran, the AUROC yielded the most accurate performance value that a

595    spatial modeler can rely on. Thus, based on current arguments, we confirm the second assumption

596    in which the incapability of the models (*i.e.*, EBF, BRT, and RF) to progress is due to the

597    unsatisfactory input data (*i.e.*, either scarce inventory, inadequate spatial indicators or spatial

598    resolution) rather than the models' structure. Analogously, the AUROC and AUPRC values are

599    more representative for the fuzzy model's performance in Queensland, Australia. Also, they are

600    comparatively in accordance with the validation method of Dayal (2018) and Dayal et al. (2018),

601    based on which the correlation of the drought risk map and the soil moisture/ rainfall departure

602    maps confirmed plausible predictive skills.

603    Comparing the different predictive models (*i.e.*, choosing the premier model among the many

604    alternatives) or different scenarios of a specific model (*i.e.*, opting the best scenario from different

605    sample partitioning techniques, different spatial resolution, and so forth), is still feasible by using

606    the cutoff-dependent metrics as they do provide valuable information that can lead to a more

607 transparent distinction between the choices. In particular, the cutoff-dependent indices can assist

608 us with distinguishing the features of the GLM and the MARS models for the case study in Italy.

609 Hence, in the following case studies, the cutoff-metrics are used only for a comparison and

610 selection of the better-performing predictive model.

611 **5.4. Landslide Susceptibility Modelling, Vietnam**

612    In accordance with the analytical evidence from the results of previous case studies, this study

613 avers that the use of the cutoff-dependent metrics can be informative for a predictive model

614 comparison. The inferences of this case study are interesting in several ways, showing that how

615 one should interpret the latter with some degree of caution. According to the AUROC and

616 AUPRC values of MaxEnt and BayGmmKda models tested in Vietnam (Table 5), the MaxEnt

617 appears to slightly excel in predictive skill, although both models show an excellent performance

618 (AUROC> 0.8). On the other hand, asymmetries are evident in the values of the cutoff-dependent

619 metrics, as we have categorized them as the ROC-accordant and -discordant metrics (see Table 6).

620                                          **Table 6** HERE

621    According to Table 6 and the relevant equations provided in Table 4, both categories support

622 high *TP* and *TN* values, while there is a subtle difference that makes them oppose. In fact, a

623 model's success in *FP* stage is highly favored in the ROC-accordant metrics, while the discordant

624 group leans towards penalizing a model's downfall in the *FN* stage. This is evident in the

625 confusion matrix of the MaxEnt and BayGmmKda, in which the MaxEnt shows an outstanding

626 performance with a zero *FP* value, while the *FN* population is drastically increased in such a way

627 that it even surpasses the *FN+FP* population in BayGmmKda model. In this case, the

628 BayGmmKda has well balanced the FP and FN population that accords to Table 7. As previously

629 mentioned in the *Theory* section, although a zero FP (Type I error) in MaxEnt results cause no

630 infrastructural and study costs, a drastic increase in FN (Type II error) values can cause massive

631 casualties via misrepresenting an area as a safe location.

632                                  **Table 7** HERE

633     Considering the structure of these predictive models, as opposed to the presence-absence nature

634 of the BayGmmKda, MaxEnt is considered as a presence-only model where some randomly

635 chosen pseudo-absence locations (*i.e.*, background samples) help the model differentiate the

636 presence locations and eventually predict an occurrence pattern. Therefore, presence-absence-

637 based validation metrics (*i.e.*, all the metrics provided in this study) may not be a good fit for the

638 performance assessment of MaxEnt. This being the case, AUPRC might be the best fit for MaxEnt

639 and in fact, it has clearly distinguished the performance of both models. However, according to

640 Phillips et al. (2006), at least, those background locations should be considered as 'pure absences'

641 to be able to graph a ROC curve, and also to calculate the metrics derived from confusion matrix.

642 This is an inevitable process for the MaxEnt. Another critical inference of this case study

643 underlines that although cutoff-dependent metrics are valuable metrics for comparing different

644 models, they are not necessarily supposed to be in line with cutoff-independent metrics. This is the

645 reason why MaxEnt and BayGmmKda both excel, but in different areas. Therefore, relying on

646 what we have conceived so far, each cutoff-dependent or -independent metric has a unique

647 indication of a model's performance.

648     There is a consensus that selecting the best predictive model can be a matter of the user

649 preference and study area's goals, which has been previously well-delineated in Goetz et al.

650 (2015). This can be carried out by relying on a pros and cons list for all the metrics and assessing

651 whether they work in agreement with the objective(s) of the project. Taking aside the

652 disadvantages of cutoff-dependent metrics, some critics have also been moved towards AUROC

653 (Lobo et al., 2008). The main complains pertain to ignoring the PPV (addressed earlier in *Theory*

654 section) and equally weighting omission (not recording some instances) and commission (miss-

655 recording some instances) errors. However, this directly stems from predefining a series of

656 thresholds and the presence-absence fabric of AUROC which is not only specific to AUROC but

657 rather all the performance metrics. Furthermore, these limitations do not question the metric itself,

658 but rather the application of them. For instance, ROC curves were first employed in the study of

659 *"discriminator systems for the detection of radio signals in the presence of noise in the 1940s"*,

660 following the attack on Pearl Harbor, USA (Garrett et al., 2008). Even the use of AUROC in

661 clinical biochemistry is carried out under a presence-absence condition (Obuchowski et al., 2004).

662 Therefore, in order to employ AUROC and other cutoff/prevalence- independent metrics in a

663 probabilistic environmental modelling context, their limitation should be accepted in favor of their

664 valuable outcomes regarding the performance evaluation.

665     Under these premises, we aver that the project study goal can assist the decision maker with

666 opting the well-performing model. For instance, if the number of opposing metrics matters the

667 most, the BayGmmKda would be the well-performing one. In particular, many municipal

668 authorities may decide in favor of public safety, which in turn can end in an immediate rejection

669 of the MaxEnt due to having considerable Type II error that can also cause notable fatalities.

670 Comparatively, if the uncertain nature of the cutoff value is in question, one can choose the

671 decisive judgment of the AUROC.

672 **5.5. Soil Digital Modelling, USA**

673    As previously mentioned, this case study represents a unique application of the proposed PMT

674    for performance assessment of the Bulk Density (BD) lateral distribution in South Dakota, USA.

675    In contrast to the previous applications of data mining methods that deal with predicting the

676    probability of an occurrence, in this study we employed the ANN and DT approaches for

677    predicting an actual quantity of BD whose actual amounts can be measured in the field. Measuring

678    the BD samples from different location of the study area, root mean square error (RMSE) can be a

679    good metric to test the accuracy of the results (*i.e.*, an approximated standard deviation of data) if

680    the data are Gaussian (*i.e.*, rich data) and devoid of any outliers (Chai and Draxler, 2014).

681    However, RMSE or accuracy, in general, can be biased and may not reflect the total precision of a

682    predictive model, warranting the need for a consolidated list of model evaluation metrics that

683    provide more extensive insights into the predictive performance.

684    In respect to the above discussion, the proposed PMT approach can be a good alternative, but

685    the nature of the prediction map should be rendered into its probability terms or at least as an

686    indication of the probability. That is, the higher values of the prediction map can indicate the

687    greater probability of having higher BD values, and vice versa. By doing so, the cutoff-dependent

688    and -independent metrics have been calculated based on which, almost all the indices congruently

689    introduce ANN as a better-performing model compared to DAT; the rest of opposing metrics (e.g.

690    specificity and PPV) show negligible differences. This is in agreement with those reported by

691    Taghizadeh-Mehrjardi et al. (2017) where the ANN was seen to outperform the support vector

692    machine (SVM) model in the mapping of soil organic matter distribution.

693    **6. Synthesis and Conclusion**

32

694    This paper provides a novel scientific contribution towards the design and implementation of an

695    adaptive, largely automated and user-friendly GIS-based spatial model assessment system,

696    denoted as the Performance Measure Tool (PMT). PMT can be used to address existing challenges

697    in pragmatic evaluation of predictive models in diverse contexts, and generally, for any scientific

698    branch where information has a spatial connotation. The PMT encloses the relevant mathematical

699    formulations to make it an easy-to-use software; it has the added capability to evaluate the

700    accuracy of the spatial modelling approach based on the different cutoff-dependent and -

701    independent evaluation criteria. The PMT is considerably flexible, and hence, it can be widely

702    applicable in multiple scientific and engineering applications where spatially-relevant predictive

703    models are tested. The approach has the potential to be applied in diverse contexts, as verified in

704    this research study, to extend its usage from geo-environmental spatial models to fields such as

705    medical geography and epidemiology where data-driven approaches are adopted to generate

706    predictive models and such models require robust comparison with several benchmark models and

707    real-life (observed) datasets.

708    In context of proposing an additional GIS-based predictive model assessment tool, the

709    consolidated metrics that are generated and evaluated by the proposed PMT, certainly provides a

710    new practical pathway for real-life decision-makers who are seeking a better performing

711    predictive model (relative to any other comparative model). Based on contested reasons, and

712    evaluations of PMT with several studies collated in this research paper, real-life decision-makers

713    can deduce the grounds on which their predictive models performs better than the others prior to

714    implementing them for practical use. By accommodating multiple types of real-time geo-

715    environmental modelling instances in this study, the take-home messages are as follows. The use

716    of a merely row-wise or a column-wise calculated index from the confusion matrix is not a robust

717  approach for model selection as this can ignore the more practical concepts considered by their

718  counterpart tools.

719     In contrast, some of the model evaluation indices (*i.e.* cutoff-dependent and −independent ones)

720  generally use a collective information of the matrix in such a way that a set of multiple statistics

721  are used in conjunction with each other. Notwithstanding this, some cutoff-dependent metrics may

722  infer the same connotation which they can be used interchangeably (*e.g.*, threat score and

723  equitable threat score, or the odds ratio and the odds ratio skill score). Moreover, the choice of

724  using the cutoff-dependent metrics over each other without a prior knowledge can also constitute

725  an unjust approach since each metric is able to tackle a different aspect of the model performance.

726  However, all metrics can be highly sensitive to the cutoff values so, they should be suggested only

727  for the model comparison.

728     As demonstrated in the theory of PMT and relevant case studies, it becomes unambiguous that

729  the measurement of the prevalence of the studied phenomenon is highly advisable in order to

730  ascertain reliable cutoff-dependent values. Doing so, they are likely to be applicable even for the

731  performance assessment of an individual model, and also, they could be comparable with cutoff-

732  independent metrics.

733     On the other hand, the cutoff-independent metrics (*i.e.*, AUROC, AUSRC, and AUPRC) can

734  decisively screen the premier model regardless of the changes in their cutoff values. However, the

735  AUROC is also underpinned by some specific assumptions so that using it would require

736  accepting its mathematical fabric. Furthermore, AUSRC and AUPRC only support presence

737  locations, they show an erroneous behavior and in particular may result in an underestimation of

738  performance compared to AUROC. Moreover, all cutoff-dependent and -independent metrics can

739  occasionally mislead by providing different results and consequently different model ranks. In

740  such case, selecting the reference model is strictly tied to the aim of the research and specific

741  aspect(s) of interest. We also concluded that compartmentalizing models in different performance

742  categories is not feasible since the matter of performance itself is quite relative.

743      We also propose the following scenario-based decision-making inferences:

744  I.    Italy and USA case studies: having more than one model→ if AUROC values converge

745       and the changes are negligible→ using other cutoff-dependent metrics to derive the

746       better-performing model.

747  II.   Iran and Australia case studies: having one model→ no access to prevalence value

748       change→ cutoff-dependent metrics change drastically by altering cutoff values→ use

749       AUROC as the decisive metric.

750  III.  Vietnam case study: more than one model→ metrics are opposing and taking different

751       parts (i.e. each selecting a different model) → decision should be made based on the

752       project goal by making pros and cons list for all the metrics.

753     As our final upshot, ROC and AUC are metrics that tend to lump together the prediction as a

754  whole; however, studying confusion matrices, accuracy and precision of a model ensure a better

755  insight on a model hit and misses. This is something that can be rarely found in the literature,

756  despite its great importance. The PMT quickly provides a full suite of performance metrics

757  allowing the users to better evaluate their spatial model and supporting a more critical judgment,

758  which in turn can promote better decision-making procedures.

759  **Acknowledgments**

769

## References

771    Abdollahi, S., Pourghasemi, H. R., Ghanbarian, G. A., Safaeian, R. 2018. Prioritization of
772        effective factors in the occurrence of land subsidence and its susceptibility mapping using
773        an SVM model and their different kernel functions. B. Eng. Geol. Environ.. doi:
774        10.1007/s10064-018-1403-6.

775    Akgün, A., Türk, N., 2011. Mapping erosion susceptibility by a multivariate statistical method: a
776        case study from the Ayvalık region, NW Turkey. Comput. Geosci. 37(9), 1515-1524.

777    Alfons, A. 2012. Package "cvTools": Cross-validation tools for regression models. https://cran.r-
778        project.org/web/packages/cvTools/index.html

779    Allouche, O., Tsoar, A., Kadmon, R., 2006. Assessing the accuracy of species distribution
780        models: prevalence, kappa and the true skill statistic (TSS). J. Appl. Ecol. 43(6), 1223-
781        1232.

782    Arpaci, A., Malowerschnig, B., Sass, O., Vacik, H., 2014. Using multi variate data mining
783        techniques for estimating fire susceptibility of Tyrolean forests. Appl. Geogr. 53, 258-
784        270.

785    Beguería, S., 2006. Validation and evaluation of predictive models in hazard assessment and risk
786        management. Nat. Hazards 37 (3), 315–329

787    Benjamini, Y., Hochberg, Y. 1995. Controlling the false discovery rate: a practical and powerful
788            approach to multiple testing. Journal of the royal statistical society. Series B
789            (Methodological), 289-300.

790    Blahut, J., van Westen, C. J., Sterlacchini, S. 2010. Analysis of landslide inventories for accurate
791            prediction of debris-flow source areas. Geomorphology 119(1-2), 36-51.

792    Boughorbel, S., Jarray, F., El-Anbari, M. 2017. Optimal classifier for imbalanced data using
793            Matthews Correlation Coefficient metric. PloS One 12(6), e0177678.

794    Bucklin, D.N., Basille, M., Benscoter, A.M., Brandt, L.A., Mazzotti, F.J., Romanach, S.S.,
795            Speroterra, C., Watling, J.I., 2015. Comparing species distribution models constructed
796            with different subsets of environmental predictors. Divers. Distrib. 21(1), 23-35.

797    Chai, T., Draxler, R. R. 2014. Root mean square error (RMSE) or mean absolute error (MAE)?–
798            Arguments against avoiding RMSE in the literature. Geosci. Model Dev. 7(3), 1247-
799            1250.

800    Chapi, K., Singh, V.P., Shirzadi, A., Shahabi, H., Bui, D.T., Pham, B.T., Khosravi, K., 2017. A
801            novel hybrid artificial intelligence approach for flood susceptibility assessment. Environ.
802            Modell. Softw. 95, 229-245.

803    Choubin, B., Moradi, E., Golshan, M., Adamowski, J., Sajedi-Hosseini, F., Mosavi, A., 2019. An
804            Ensemble prediction of flood susceptibility using multivariate discriminant analysis,
805            classification and regression trees, and support vector machines. Sci. Total Environ. 651,
806            2087-2096.

807    Chung, C. J. 2006. Using likelihood ratio functions for modelling the conditional probability of
808            occurrence of future landslides for risk assessment. Comput. Geosci. 32(8), 1052-1068.

809    Chung, C. J., Fabbri, A. G. 2008. Predicting landslides for risk analysis—spatial models tested
810            by a cross-validation technique. Geomorphology 94(3-4), 438-452.

811    Cohen, J. 1960. A coefficient of agreement for nominal scales. Educ. Psychol. Meas. 20(1), 37–
812            46.

813    Conoscenti, C., Agnesi, V., Cama, M., Caraballo-Arias, N. A., Rotigliano, E. 2018. Assessment
814            of gully erosion susceptibility using multivariate adaptive regression splines and
815            accounting for terrain connectivity. Land Degrad. Dev. 29(3), 724-736.

816    Conoscenti, C., Angileri, S., Cappadonia, C., Rotigliano, E., Agnesi, V., Märker, M., 2014. Gully
817            erosion susceptibility assessment by means of GIS-based logistic regression: a case of
818            Sicily (Italy). Geomorphology 204, 399-411.

819    Coombes, K.R., 2018. Package 'CrossValidate': Classes and Methods for Cross Validation of
820        "Class Prediction" Algorithms. https://cran.r-
821        project.org/web/packages/CrossValidate/index.html

822    Dayal, K. S., Deo, R. C., Apan, A. A. 2018. Spatio-temporal drought risk mapping approach and
823        its application in the drought-prone region of south-east Queensland, Australia. Nat.
824        Hazards, 1-25.

825    Dayal, K.S., 2018. Development of Statistical and Geospatial-Based Framework For Drought-
826        Risk Assessment. PhD Thesis. University of Southern Queensland, Australia. 248pp.

827    Deo, R. C., Samui, P., Kim, D. 2016. Estimation of monthly evaporative loss using relevance
828        vector machine, extreme learning machine and multivariate adaptive regression spline
829        models. Stoch. Env. Res. Risk. A. 30(6), 1769-1784.

830    Deo, R. C., Tiwari, M. K., Adamowski, J. F., Quilty, J. M. 2017. Forecasting effective drought
831        index using a wavelet extreme learning machine (W-ELM) model. Stoch. Env. Res. Risk.
832        A. 31(5), 1211-1240.

833    Fawcett, T., 2006. An introduction to ROC analysis. Pattern Recogn. Lett. 27(8), 861-874.

834    Frattini, P., Crosta, G., Carrara, A. 2010. Techniques for evaluating the performance of landslide
835        susceptibility models. Eng. Geol. 111(1-4), 62-72.

836    Garosi, Y., Sheklabadi, M., Pourghasemi, H.R., Besalatpour, A.A., Conoscenti, C., Van Oost, K.,
837        2018. Comparison of differences in resolution and sources of controlling factors for gully
838        erosion susceptibility mapping. Geoderma 330, 65-78.

839    Garrett, P. E., Lasky, F. D., Meier, K. L. 2008. User protocol for evaluation of qualitative test
840        performance; approved guideline. CLSI.

841    Ghorbanzadeh, O., Blaschke, T., Aryal, J., Gholaminia, K., 2018. A new GIS-based technique
842        using an adaptive neuro-fuzzy inference system for land subsidence susceptibility
843        mapping. J. Spat. Sci. 1-17.

844    Gilbert, G. K. 1884. Finley's tornado predictions. American Meteorological Journal. A Monthly
845        Review of Meteorology and Allied Branches of Study (1884-1896), 1(5), 166.

846    Glade, T. (Ed.), 2005. Landslide Hazard and Risk. Wiley, Chichester, pp. 41–74.

847    Goetz, J. N., Brenning, A., Petschko, H., Leopold, P. 2015. Evaluating machine learning and
848        statistical prediction techniques for landslide susceptibility modelling. Comput.
849        Geosci. 81, 1-11.

850 Gómez-Gutiérrez, Á., Conoscenti, C., Angileri, S. E., Rotigliano, E., Schnabel, S. 2015. Using
851     topographical attributes to evaluate gully erosion proneness (susceptibility) in two
852     mediterranean basins: advantages and limitations. Nat. Hazards 79(1), 291-314.

853 Gorsevski, P.V., Gessler, P.E., Foltz, R.B., Elliot, W.J., 2006. Spatial prediction of landslide
854     hazard using logistic regression and ROC analysis. T. GIS 10(3), 395-415.

855 Hanssen, A. W., Kuipers, W. J. A. 1965. On the Relationship Between the Frequency of Rain
856     and Various Meteorological Parameters:(with Reference to the Problem Ob Objective
857     Forecasting). Staatsdrukkerij-en Uitgeverijbedrijf.

858 Heidke, P., 1926: Berechnung des Erfolges und der Gute der Windstarkevorhersagen im
859     Sturmwarnungsdienst (Measures of success and goodness of wind force forecasts by the
860     gale warning service) Geogr. Ann. 8, 301–349

861 Hosmer, D. W., Lemeshow, S. 2000. Applied Logistic Regression., 2nd edn.(Wiley: New
862     York.). NY, USA.

863 Kavzoglu, T., Colkesen, I. and Sahin, E.K., 2019. Machine Learning Techniques in Landslide
864     Susceptibility Mapping: A Survey and a Case Study. In Landslides: Theory, Practice and
865     Modelling (pp. 283-301). Springer, Cham.

866 Kornejady, A., Ownegh, M., Bahremand, A. 2017. Landslide susceptibility assessment using
867     maximum entropy model with two different data sampling methods. Catena 152, 144-
868     162.

869 Lipovetsky, S., 2009. Pareto 80/20 law: derivation via random partitioning. Int. J. Math. Educ.
870     Sci. Technol. 40(2), 271-277.

871 Lobo, J. M., Jiménez-Valverde, A., Real, R. 2008. AUC: a misleading measure of the
872     performance of predictive distribution models. Global Ecol. Biogeogr. 17(2), 145-151.

873 Loiselle, B.A., Howell, C.A., Graham, C.H., Goerck, J.M., Brooks, T., Smith, K.G., Williams,
874     P.H., 2003. Avoiding pitfalls of using species distribution models in conservation
875     planning. Conserv. Biol. 17, 1591–1600.

876 Lombardo, L., Mai, P. M. 2018. Presenting logistic regression-based landslide susceptibility
877     results. Eng. Geol. 244, 14-24.

878 Lombardo, L., Opitz, T., Huser, R. 2018. Point process-based modelling of multiple debris flow
879     landslides using INLA: an application to the 2009 Messina disaster. Stoch. Env. Res.
880     Risk. A. 32(7), 2179-2198.

881   Malone, B.P., Minasny, B., McBratney, A.B., 2017. Combining Continuous and Categorical
882       Modelling: Digital Soil Mapping of Soil Horizons and Their Depths. In Using R for
883       Digital Soil Mapping (pp. 231-244). Springer, Cham.

884   Mas, J.F., Soares Filho, B., Pontius, R.G., Farfán Gutiérrez, M., Rodrigues, H., 2013. A suite of
885       tools for ROC analysis of spatial models. ISPRS Int. Geo-Inf. 2(3), 869-887.

886   Minasny, B., McBratney, A.B., 2007. Incorporating taxonomic distance into spatial prediction
887       and digital mapping of soil classes. Geoderma 142(3-4), 285-293.

888   Miraki, S., Zanganeh, S.H., Chapi, K., Singh, V.P., Shirzadi, A., Shahabi, H., Pham, B.T., 2019.
889       Mapping Groundwater Potential Using a Novel Hybrid Intelligence Approach. Water
890       Resour. Manage. 33(1), 281-302.

891   Naghibi, S.A., Ahmadi, K., Daneshi, A., 2017. Application of support vector machine, random
892       forest, and genetic algorithm optimized random forest models in groundwater potential
893       mapping. Water Resour. Manage. 31(9), 2761-2775.

894   Obuchowski, N. A., Lieber, M. L., Wians, F. H. 2004. ROC curves in clinical chemistry: uses,
895       misuses, and possible solutions. Clin. Chem. 50(7), 1118-1125.

896   Peirce, C. S. 1884. The numerical measure of the success of predictions. Science 4(93), 453-454.

897   Pepe, M.S., Janes, H., Longton, G., Leisenring, W., Newcomb, P., 2004. Limitations of the odds
898       ratio in gauging the performance of a diagnostic, prognostic, or screening marker. Am. J.
899       Epidemiol. 159(9), 882-890.

900   Perruchet, P., Peereman, R. 2004. The exploitation of distributional information in syllable
901       processing. J. Neurolinguist. 17(2-3), 97-119.

902   Phillips, S. J., Anderson, R. P., Schapire, R. E. 2006. Maximum entropy modelling of species
903       geographic distributions. Ecol. Modell. 190(3-4), 231-259.

904   Pontius Jr, R.G., Schneider, L.C., 2001. Land-cover change model validation by an ROC method
905       for the Ipswich watershed, Massachusetts, USA. Agr. Ecosyst. Environ. 85(1-3), 239-
906       248.

907   Pourghasemi, H.R., Rahmati, O., 2018. Prediction of the landslide susceptibility: Which
908       algorithm, which precision?. Catena 162, 177-192.

909   Powers, D. M. 2011. Evaluation: from precision, recall and F-measure to ROC, informedness,
910       markedness and correlation. http://hdl.handle.net/2328/27165.

911   Pradhan, B., Lee, S., 2010. Landslide susceptibility assessment and factor effect analysis:
912       backpropagation artificial neural networks and their comparison with frequency ratio and

bivariate logistic regression modelling. Environmental Modelling & Software, 25(6), 747-759.

Pullar, D., Springer, D., 2000. Towards integrating GIS and catchment models. Environ. Modell. Softw. 15(5), pp.451-459.

Quillfeldt, P., Engler, J.O., Silk, J.R., Phillips, R.A., 2017. Influence of device accuracy and choice of algorithm for species distribution modelling of seabirds: a case study using black-browed albatrosses. J. Avian Biol. 48(12), 1549-1555.

Rahmati, O., Pourghasemi, H. R. 2017. Identification of critical flood prone areas in data-scarce and ungauged regions: A comparison of three data mining models. Water Resour. Manage. 31(5), 1473-1487.

Reineking, T. 2014. Belief functions: Theory and algorithms. https://pdfs.semanticscholar.org/eb77/3cd7c84617bfd9e3abbb7695e113e94c9524.pdf

Ruspini, E. H., Lowrance, J. D., Strat, T. M. 1992. Understanding evidential reasoning. Int. J. Approx. Reason. 6(3), 401-424.

Sahana, M., Ganaie, T.A., 2017. GIS-based landscape vulnerability assessment to forest fire susceptibility of Rudraprayag district, Uttarakhand, India. Environ. Earth Sci. 76(20), p.676.

Schaefer, J. T. 1990. The critical success index as an indicator of warning skill. Weather Forecast. 5(4), 570-575.

Scott, L. M., Janikas, M. V. 2010. Spatial statistics in ArcGIS. In Handbook of applied spatial analysis (pp. 27-41). Springer, Berlin, Heidelberg.

Shabani, F., Kumar, L., Ahmadi, M., 2016. A comparison of absolute performance of different correlative and mechanistic species distribution models in an independent area. Ecol. Evol. 6(16), 5973-5986.

Shafer, G. 1976. A mathematical theory of evidence (Vol. 42). Princeton university press.

Siahkamari, S., Haghizadeh, A., Zeinivand, H., Tahmasebipour, N., Rahmati, O., 2018. Spatial prediction of flood-susceptible areas using frequency ratio and maximum entropy models. Geocarto Int. 33(9), 927-941.

Swets, J. A. 2014. Signal detection theory and ROC analysis in psychology and diagnostics: Collected papers. Psychology Press.

Tien Bui, D., Bui, Q.T., Nguyen, Q.P., Pradhan, B., Nampak, H., Trinh, P.T., 2017. A hybrid artificial intelligence approach using GIS-based neural-fuzzy inference system and

945       particle swarm optimization for forest fire susceptibility modelling at a tropical area. Agr.
946       Forest Meteorol. 233, 32-44.

947 Tien Bui, D., Hoang, N. D. 2017. A Bayesian framework based on a Gaussian mixture model
948       and radial-basis-function Fisher discriminant analysis (BayGmmKda V1. 1) for spatial
949       prediction of floods. Geosci. Model Dev. 10 (9), 3391-3409

950 Van Westen, C. J., Van Asch, T. W., Soeters, R. 2006. Landslide hazard and risk zonation—why
951       is it still so difficult?. B. Eng. Geol. Environ. 65(2), 167-184.

952 Wiesmeier, M., Barthold, F., Blank, B., Kögel-Knabner, I., 2011. Digital mapping of soil organic
953       matter stocks using Random Forest modelling in a semi-arid steppe ecosystem. PlantSoil
954       340(1-2), 7-24.

955 Yan, F., Zhang, Q., Ye, S., Ren, B., 2019. A novel hybrid approach for landslide susceptibility
956       mapping integrating analytical hierarchy process and normalized frequency ratio methods
957       with the cloud model. Geomorphology 327, 170-187.

958 Yule, G. U. 1900. VII. On the association of attributes in statistics: with illustrations from the
959       material of the childhood society, &c. Phil. Trans. R. Soc. Lond. A, 194(252-261), 257-
960       319.

961

962

963

964

965

966

967

968

969

970

**Figure Captions**

972     **Fig. 1** Study sites on the world map

973     **Fig. 2** Methodological flowchart adopted in this study

974     **Fig. 3** Gully erosion prediction maps of the central-western Sicily (Italy) generated by using the

975     GLM (a) and MARS (b) models

976     **Fig. 4** Flood-inundation susceptibility map of the Galikesh region (Iran) obtained from the EBF

977     model

978     **Fig. 5** Drought risk map of the south-east of Queensland (Australia) produced by using fuzzy

979     GAMMA overlay technique

980     **Fig. 6** Effects of 50% (a), 70% (b) and 90% (c) cutoff values on the extent of safe/danger zones

981     and classification of presence/absence samples in south-east of Queensland

982     **Fig. 7** Landslide susceptibility maps of the Kon Tum province (Vietnam) obtained from

983     BayGmmKda (a) and MaxEnt (b) models

984     **Fig. 8** Bulk density predictive distribution maps of South Dakota (USA) generated from ANN

985     (a) and DT (b) models

**Table 1** Confusion matrix elements.

| Observed | Predicted | |
|---|---|---|
| | Class stable (−) | Class unstable (+) |
| Class stable (−)* | (−|−) True negative (TN) | (+|−) False positive (FP; Error Type I) |
| Class unstable (+)** | (−|+) False negative (FN; Error Type II) | (+|+) True positive (TP) |

* Absence areas    ** Presence areas

**Table 2** The PMT input files

| ID | Setting | Description | ID | Setting | Description |
|----|---------|-------------|----|---------|-------------|
| 1 | Input raster layers | The raster maps generated by any spatial model representing the susceptibility or suitability of a phenomenon over an area (you can add different maps for the same area as many as desired). | 5 | Validation positives | Import the shapefile of all the validation samples of the phenomenon of interest (discarded dataset in the training stage). |
| 2 | Cutoff | An a-priori cutoff percentage to split the input raster into two segments (50% is set as default). | 6 | Validation negatives | Import the shapefile of the non-event validation locations. |
| 3 | Training positives | Import the shapefile of all the training samples of the phenomenon of interest. | 7 | Output workspace | The pass to contain the outputs (a folder address). |
| 4 | Training negatives | Import the shapefile of the absence training locations (should be prepared beforehand by different methods mentioned in the text) | 8 | Number of classes (for SRC and PRC curves) | The number into which the spatial raster is to be reclassified (100 classes are set as default). The reclassification method is based on an equal interval. A higher number of classes will result in smoother SRC and PRC curves with more precise AUSRC and AUPRC values. |

**Table 3** The PMT output files

| ID | Setting | Description |
|----|---------|-------------|
| 1 | Html file | It explains the main results of the performance analyses including confusion matrix, cutoff-dependent metrics, and cutoff–independent metrics. ROC, SRC, and PRC curves are other parts of this html file. In addition, all results were classified into two groups of cutoff-dependent and cutoff-independent approaches with some useful explanations regarding these approaches. |
| 2 | Microsoft excel file | This file summarize all of quantitative results (without explanations) |

**Table 4** Equations of cutoff-dependent performance metrics

| Performance metric | Equation | Performance metric | Equation |
|---|---|---|---|
| True positive rate (TPR; sensitivity) | $\dfrac{TP}{P} = \dfrac{TP}{TP+FN}$ | Matthews correlation coefficient (MCC) | $\dfrac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$ |
| False positive rate (FPR; fall-out; 1–specificity) | $\dfrac{FP}{N} = \dfrac{FP}{FP+TN} = 1 - \dfrac{TN}{TN+FP}$ | Informedness (Bookmaker informedness; BM) | $TPR + TNR - 1$ |
| True negative rate (TNR; specificity) | $\dfrac{TN}{N} = \dfrac{TN}{TN+FP}$ | Markedness (MK) | $PPV + NPV - 1$ |
| False negative rate (miss rate) | $\dfrac{FN}{P} = \dfrac{FN}{FN+TP} = 1 - TPR$ | Threat score | $\dfrac{TP}{TP+FN+FP}$ |
| Efficiency (accuracy) | $\dfrac{TP+TN}{T}$ | Equitable threat score | $\dfrac{TP - TP_{random}}{TP+FN+FP - TP_{random}}$<br><br>$where, \quad TP_{random} = \dfrac{(TP+FN)(TP+FP)}{T}$ |
| Misclassification rate | $\dfrac{FP+FN}{T}$ | True skill statistic (Pierce's skill score) | $\dfrac{TP}{TP+FN} - \dfrac{FP}{FP+TN} = \text{Sensitivity} + \text{Specificity} - 1$ |
| Positive predictive value (PPV; precision) | $\dfrac{TP}{TP+FP}$ | Heidke's skill score | $\dfrac{TP+TN-E}{T-E}$<br><br>$where\ E = \dfrac{1}{T}[(TP+FN)(TP+FP) + (TN+FN)(TN+FP)]$ |
| False discovery rate (FDR) | $1 - PPV = \dfrac{FP}{FP+TP}$ | Odds ratio | $\dfrac{TP \times TN}{FN \times FP}$ |
| Negative predictive value (NPV) | $\dfrac{TN}{TN+FN}$ | Odd ratio skill score (Yule's Q) | $\dfrac{(TP \times TN) - (FP \times FN)}{(TP \times TN) + (FP \times FN)}$ |
| False omission rate (FOR) | $1 - NPV = \dfrac{FN}{FN+TN}$ | Cohen's kappa | $\dfrac{(TP+TN) - [\{TP+FN)(TP+FP) + (FN+TN)(FP+TN)]/T}{T - [\{(TP+FN)(TP+FP) + (FN+TN)(FP+TN)\}/T\}}$ |
| F-score | $2\dfrac{PPV.TPR}{PPV+TPR} = \dfrac{2TP}{2TP+FP+FN}$ | - | - |

**Table 5** Performance metrics calculated for each case study

| Country | Subject | Model | Modelling step | Efficiency (accuracy) | True positive rate (TPR) | False positive rate (FPR) | Threat score | Equitable threat score | Hedke skill score | Odds ratio | Odd ratio skill score |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Australia | Drought risk mapping | Fuzzy function: 50% cutoff | Validation | 0.625 | 0.580 | 0.222 | 0.545 | 0.142 | 0.25 | 4.8462 | 0.657 |
| | | Fuzzy function: 70% cutoff | | 0.85 | 0.818 | 0.111 | 0.75 | 0.538 | 0.7 | 36 | 0.945 |
| | | Fuzzy function: 90% cutoff | | 0.625 | 1 | 0.428 | 0.25 | 0.142 | 0.25 | 0 | 1 |
| Iran | Flood inundation mapping | EBF | Training | 0.808 | 0.891 | 0.245 | 0.647 | 0.446 | 0.617 | 25.33 | 0.924 |
| | | | Validation | 0.718 | 0.769 | 0.315 | 0.526 | 0.28 | 0.437 | 7.22 | 0.756 |
| USA | Distribution of soil organic matters | DAT | Validation | 0.442 | 0.431 | 0 | 0.431 | 0.014 | 0.028 | 0 | 1 |
| | | ANN | Validation | 0.730 | 0.625 | 0.1 | 0.588 | 0.315 | 0.48 | 15 | 0.875 |
| Italy | Gully susceptibility mapping | MARS | Training | 0.970 | 0.963 | 0.022 | 0.942 | 0.888 | 0.940 | 1151 | 0.998 |
| | | | Validation | 0.976 | 0.970 | 0.016 | 0.954 | 0.910 | 0.953 | 1885 | 0.998 |
| | | GLM | Training | 0.656 | 0.592 | 0 | 0.592 | 0.185 | 0.312 | 0 | 1 |
| | | | Validation | 0.674 | 0.605 | 0 | 0.605 | 0.211 | 0.348 | 0 | 1 |
| Vietnam | Landslide susceptibility mapping | MaxEnt | Validation | 0.601 | 0.556 | 0 | 0.556 | 0.112 | 0.202 | 0 | 1 |
| | | BayGmmKda | | 0.739 | 0.731 | 0.2521 | 0.592 | 0.314 | 0.478 | 8.08 | 0.779 |

**Table 5 (**continued)

| Country | Subject | Model | Modelling step | True skill statistic | Cohen's kappa | True negative rate (TNR) | False negative rate (miss rate) | Misclassification rate | Positive predictive value (PPV) | False discovery rate (FDR) | Negative predictive value (NPV) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Australia | Drought risk mapping | Fuzzy function: 50% cutoff | Validation | 0.358 | 0.25 | 0.778 | 0.419 | 0.375 | 0.900 | 0.100 | 0.350 |
| | | Fuzzy function: 70% cutoff | | 0.707 | 0.7 | 0.889 | 0.182 | 0.150 | 0.900 | 0.100 | 0.800 |
| | | Fuzzy function: 90% cutoff | | 0.571 | 0.25 | 0.571 | 0.000 | 0.375 | 0.250 | 0.750 | 1.000 |
| Iran | Flood inundation mapping | EBF | Training | 0.646 | 0.617 | 0.754 | 0.108 | 0.192 | 0.702 | 0.298 | 0.915 |
| | | | Validation | 0.453 | 0.437 | 0.684 | 0.231 | 0.281 | 0.625 | 0.375 | 0.813 |
| USA | Predictive distribution of soil bulk density | DAT | Validation | 0.431 | 0.028 | 1.00 | 0.569 | 0.558 | 1.000 | 0.000 | 0.033 |
| | | ANN | Validation | 0.525 | 0.48 | 0.90 | 0.375 | 0.269 | 0.909 | 0.091 | 0.600 |
| Italy | Gully susceptibility mapping | MARS | Training | 0.941 | 0.940 | 0.978 | 0.037 | 0.030 | 0.978 | 0.022 | 0.963 |
| | | | Validation | 0.953 | 0.953 | 0.983 | 0.030 | 0.024 | 0.983 | 0.017 | 0.970 |
| | | GLM | Training | 0.592 | 0.312 | 1.000 | 0.407 | 0.344 | 1.000 | 0.000 | 0.313 |
| | | | Validation | 0.605 | 0.348 | 1.000 | 0.394 | 0.326 | 1.000 | 0.000 | 0.349 |
| Vietnam | Landslide susceptibility mapping | MaxEnt | Validation | 0.556 | 0.202 | 1.000 | 0.444 | 0.399 | 1.000 | 0.000 | 0.203 |
| | | BayGmmKda | | 0.479 | 0.478 | 0.748 | 0.269 | 0.261 | 0.757 | 0.243 | 0.722 |

**Table 5 (**continued)

| Country | Subject | Model | Modelling step | False omission rate (FOR) | F-score | Matthews correlation coefficient (MCC) | Informedness (Bookmaker informedness; BM) | Markedness (MK) | AUROC | AUSRC | AUPRC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Australia | Drought risk mapping | Fuzzy function: 50% cutoff | Validation | 0.650 | 0.706 | 0.299 | 0.358 | 0.250 | 0.873 | - | 74.400 |
| | | Fuzzy function: 70% cutoff | | 0.200 | 0.857 | 0.704 | 0.707 | 0.700 | 0.873 | | 74.400 |
| | | Fuzzy function: 90% cutoff | | 0.000 | 0.400 | 0.378 | 0.571 | 0.250 | 0.873 | | 74.400 |
| Iran | Flood inundation mapping | EBF | Training | 0.085 | 0.786 | 0.632 | 0.646 | 0.617 | 0.866 | 79.710 | - |
| | | | Validation | 0.188 | 0.690 | 0.445 | 0.453 | 0.438 | 0.787 | - | 75.209 |
| USA | Predictive distribution of soil bulk density | DAT | Validation | 0.967 | 0.603 | 0.120 | 0.431 | 0.033 | 0.839 | - | 77.620 |
| | | ANN | Validation | 0.400 | 0.741 | 0.517 | 0.525 | 0.509 | 0.879 | - | 79.630 |
| Italy | Gully susceptibility mapping | MARS | Training | 0.037 | 0.971 | 0.941 | 0.941 | 0.941 | 0.992 | 99.141 | - |
| | | | Validation | 0.030 | 0.977 | 0.953 | 0.953 | 0.953 | 0.995 | - | 99.285 |
| | | GLM | Training | 0.687 | 0.744 | 0.430 | 0.593 | 0.313 | 0.987 | 97.134 | - |
| | | | Validation | 0.651 | 0.754 | 0.460 | 0.606 | 0.349 | 0.992 | - | 97.542 |
| Vietnam | Landslide susceptibility mapping | MaxEnt | Validation | 0.797 | 0.715 | 0.336 | 0.556 | 0.203 | 0.889 | - | 0.855 |
| | | BayGmmKda | | 0.278 | 0.744 | 0.479 | 0.479 | 0.479 | 0.819 | - | 69.460 |

**Table 6** Opposing performance metrics for Vietnam's case study

| ROC-accordant | ROC-discordant |
| --- | --- |
| Informedness | Markedness |
| PPV | MCC |
| TNR | NPV |
| TSS | Misclassification rate |
| 1-Specificity | FNR |
| FDR | Cohen's Kappa |
| | F-score |
| | Hedke skill score |
| | Equitable threat score |
| | Threat score |
| | Sensitivity |
| | Accuracy |
| | FOR |

**Table 7** Comparing confusion matrix variants of MaxEnt and BayGmmKda models as implemented in Vietnam

| Observed | Models | |
| --- | --- | --- |
| | MaxEnt | BayGmmKda |
| TN | 330 | 1175 |
| TP | 1627 | 1231 |
| FN | 1297 | 452 |
| FP | 0 | 396 |

**Central-Western Sicily (Italy)**

**Kon Tum (Vietnam)**

**BSR Watershed South Dekota (USA)**

**Galikesh (Iran)**

**South-East Queensland (Australia)**

**Fig. 1**

**Fig. 2**

Fig. 3

**Fig. 4**

**Fig. 5**

**Fig. 6**

**Fig. 7**

**Fig. 8**