

# Testing machine learning models for heuristic building damage assessment applied to the Italian Database of Observed Damage (DaDO)

Subash Ghimire<sup>1</sup>, Philippe Guéguen<sup>1</sup>, Adrien Pothon<sup>2</sup>, Danijel Schorlemmer<sup>3</sup>

<sup>1</sup>ISTerre, Université Grenoble Alpes/CNRS/IRD/Université Gustave Eiffel, Grenoble, CS40700 38058 Grenoble cedex 9, France.

<sup>2</sup>AXA Group Risk Management, GIE AXA - 21 Avenue Matignon - 75008 Paris, France.

<sup>3</sup>German Research Center for Geosciences, Telegrafenberg, 14473 Potsdam, Germany.

Correspondence to: Subash Ghimire (subash.ghimire@univ-grenoble-alpes.fr)

## Abstract

Assessing or forecasting seismic damage to buildings is an essential issue for earthquake disaster management. In this study, we explore the efficacy of several machine learning models for damage characterization, trained and tested on the database of damage observed after Italian earthquakes (DaDO). Six models were considered: regression- and classification-based machine learning models, each using random forest, gradient boosting and extreme gradient boosting. The structural features considered were divided into two groups: all structural features provided by DaDO or only those considered to be the most reliable and easiest to collect (age, number of storeys, floor area, building height). Macroseismic intensity was also included as an input feature. The seismic damage per building was determined according to the EMS-98 scale observed after seven significant earthquakes occurring in several Italian regions. The results showed that extreme gradient boosting classification is statistically the most efficient method, particularly when considering the basic structural features and grouping the damage according to the traffic-light based system used, for example, during the post-disaster period (green, yellow and red), 68% buildings were correctly classified. The results obtained by the machine learning-based heuristic model for damage assessment are of the same order of accuracy (error values were less than 17%) as those obtained by the traditional Risk-UE method. Finally, the machine learning analysis found that the importance of structural features with respect to damage was conditioned by the level of damage considered.

## Key Words

Earthquake building-damage, DaDO building damage database, Machine learning, RISK-UE, Seismic vulnerability of buildings, Italy.

## 38 1. Introduction

39 Population growth worldwide increases exposure to natural hazards, increasing consequences in terms  
40 of global economic and human losses. For example, between 1985 and 2014, the world's population  
41 increased by 50% and average annual losses due to natural disasters increased from US\$14 billion to  
42 over US\$140 billion (Silva et al., 2019). Among other natural hazards, earthquakes represent one-fifth  
43 of total annual economic losses and cause more than 20 thousand deaths per year (Daniell et al., 2017;  
44 Silva et al., 2019). To develop effective seismic risk reduction policies, decision-makers and  
45 stakeholders rely on a representation of consequences when earthquakes affect the built environment.  
46 Two main risk metrics generally considered at the global scale are associated with building damage:  
47 direct economic losses due to costs of repair/replacement and loss of life of inhabitants due to building  
48 damage. The damage is estimated by combining the seismic hazard, exposure models and  
49 vulnerability/fragility functions (Silva et al., 2019).

50 For scenario-based risk assessment, damage and related consequences are computed for a single  
51 earthquake defined in terms of magnitude, location, and other seismological features. Many methods  
52 have been developed to characterize the urban environment for exposure models. In particular, damage  
53 assessment requires vulnerability/fragility functions for all types of existing buildings, defined  
54 according to their design characteristics (shape, position, materials, height, etc.) and grouped in a  
55 building taxonomy (e.g. among other conventional methods FEMA, 2003; Grünthal, 1998; Guéguen  
56 et al., 2007; Lagomarsino & Giovinazzi, 2006; Mouroux & Le Brun, 2006; Silva et al., 2022). At the  
57 regional/country scale, damage assessment is therefore confronted with the difficulty of accurately  
58 characterizing exposure according to the required criteria and assigning appropriate  
59 vulnerability/fragility functions to building features. Unfortunately, the necessary information is often  
60 sparse and incomplete, and the exposure model development suffers from economic and time  
61 constraints.

62 Over the past decade, there has been growing interest in artificial intelligence methods for seismic risk  
63 assessment, due to their superiority computational efficiency, easy handling of complex problems, and  
64 the incorporation of uncertainties (e.g., Riedel et al., 2014, 2015; Azimi et al., 2020; Ghimire et al.,  
65 2022; Hegde and Rokseth, 2020; Kim et al., 2020; Mangalathu & Jeon, 2020; Morfidis & Kostinakis,  
66 2018; Salehi & Burgueño, 2018; Seo et al., 2012; Sun et al., 2021; Wang et al., 2021; Xie et al., 2020;  
67 Y. Xu et al., 2020; Z. Xu et al., 2020). In particular, several studies have tested the effectiveness of  
68 machine learning methods in associating damage degrees with basic building features and spatially-  
69 distributed seismic demand with acceptable accuracy compared with conventional methods or tested  
70 with post-earthquake observations (e.g., Riedel et al., 2014, 2015; Guettiche et al., 2017; Harirchian et  
71 al., 2021; Mangalathu et al., 2020; Roeslin et al., 2020; Stojadinović et al., 2021; Ghimire et al., 2022).  
72 In parallel, significant efforts have been made to collect post-earthquake building damage observations  
73 after damaging earthquakes (Dolce et al., 2019; MINVU, 2010; MTPTC, 2010; NPC, 2015). With more

74 than 10,000 samples compiled, the Database of Observed Damage (DaDO) in Italy, a platform of the  
 75 Civil Protection Department, developed by the Eucentre Foundation (Dolce et al., 2019), allows  
 76 exploration of the value of heuristic vulnerability functions calibrated on observations (Lagomarsino et  
 77 al., 2021), as well as the training of heuristic functions using machine learning models (Ghimire et al.,  
 78 2022) and considering sparse and incomplete building features.

79 The main objective of this study is to investigate the effectiveness of several machine learning models  
 80 trained and tested on information from the DaDO to develop a heuristic model for damage assessment.  
 81 The model may be classified as heuristic because it applies a problem-solving approach in which a  
 82 calculated guess based on previous experience is considered for damage assessment (as opposed to  
 83 applying algorithms that effectively eliminate the approximation). The damage is thus estimated in a  
 84 non-rigorous way defined during the training phase and the results must be validated and then tested  
 85 against observed damage. By analogy with psychology, this procedure can reduce the cognitive load  
 86 associated with uncertainties when making decisions based on damage assessment, by explicitly  
 87 considering the uncertainties in the assessment, being aware about the incompleteness of the  
 88 information and the accuracy level to make a decision. The dataset and methods are described in the  
 89 data and method sections, respectively. The fourth section presents the results of damage prediction  
 90 produced by machine learning models compared with conventional methods, followed by a discussion  
 91 and a conclusion section.

92

## 93 2. Data

94 The Database of Observed Damage (DaDO, Dolce et al., 2019) is accessible through a web-GIS  
 95 platform and is designed to collect and share information about building features, seismic ground  
 96 motions and observed damage following major earthquakes in Italy from 1976 to 2019 (with the  
 97 exclusion of the 2016-2017 Central-Italy earthquake for which data processing is ongoing). A  
 98 framework was adopted to homogenize the different forms of information collected and to translate the  
 99 damage information into the EMS-98 scale (Grunthal et al., 1998) using the method proposed by Dolce  
 100 et al. (2019). For this study, we selected building damage data from seven earthquakes summarized in  
 101 Table 1 and presented in Fig.1.

102

103 **Table 1.** Building-damage data from the DaDO for the seven earthquakes considered in this study. ‘Ref’  
 104 is the reference to the earthquake used in the manuscript. ‘DL’ is the number of the damage grade  
 105 available in DaDO. ‘NB’ is the number of buildings considered in this study. AeDES is the post-  
 106 earthquake damage survey form, first introduced in 1997 and which now have become the official  
 107 operational tool recognized by the Italian Civil Protection in 2002.

Ref	Earthquake	Event date	Mag.	Epicentre		Damage survey form	DL	NB
				Lat.	Long.			
E1	Irpinia-1980	23/11/1980	6.9	40.91	15.37	Irpinia-980	8	37,828

E2	Pollino-1998	09/09/1998	5.6	40.04	15.98	AeDES-1998	4	9,485
E3	Molise-Puglia-2002	31/10/2002	5.9	41.79	14.87	AeDES-2000	4	6,396
E4	Emilia-Romagna-2003	14/09/2003	5.3	44.33	11.45	AeDES-2000	4	239
E5	L'Aquila-2009	06/04/2009	6.3	42.34	13.34	AeDES-2008	4	37,999
E6	Emilia-Romagna-2012	20/05/2012	6.1	44.89	11.23	AeDES-2008	4	10,581
E7	Garfagnana-Lunigiana-2013	21/06/2013	5.3	44.15	10.14	AeDES-2008	4	1,474

108

109 The converted EMS-98 damage grade (DG) ranges from damage grade DG0 (no damage) to DG5 (total  
110 collapse). The building features are available for each individual building and relate to the shape and  
111 design of the building and the built-up environment (Tab. 2, Fig. 2), as follows:

112 **Building location** - the location of each building is defined by its latitude and longitude, assigned using  
113 either the exact address of the building if available or the address of the local administrative centre  
114 (Dolce et al., 2019).

115 **Numbers of storeys** - total number of floors above the surface of the ground.

116 **Age of building** - time difference between the date of the earthquake and the date of building  
117 construction/renovation.

118 **Height of building** - total height of the building above the surface of the ground, in m.

119 **Floor area** – average of the storey surface area, in m<sup>2</sup>.

120 **Ground slope condition** - four types of ground slope conditions are defined (flat, mild slope, steep  
121 slope, and ridge).

122 **Roof type** – four types of roofs are defined (thrusting heavy roof, non-thrusting heavy roof, thrusting  
123 light roof, and non-thrusting light roof).

124 **Position of building** - indication of the building's position in the block: isolated, extreme, corner, and  
125 intermediate.

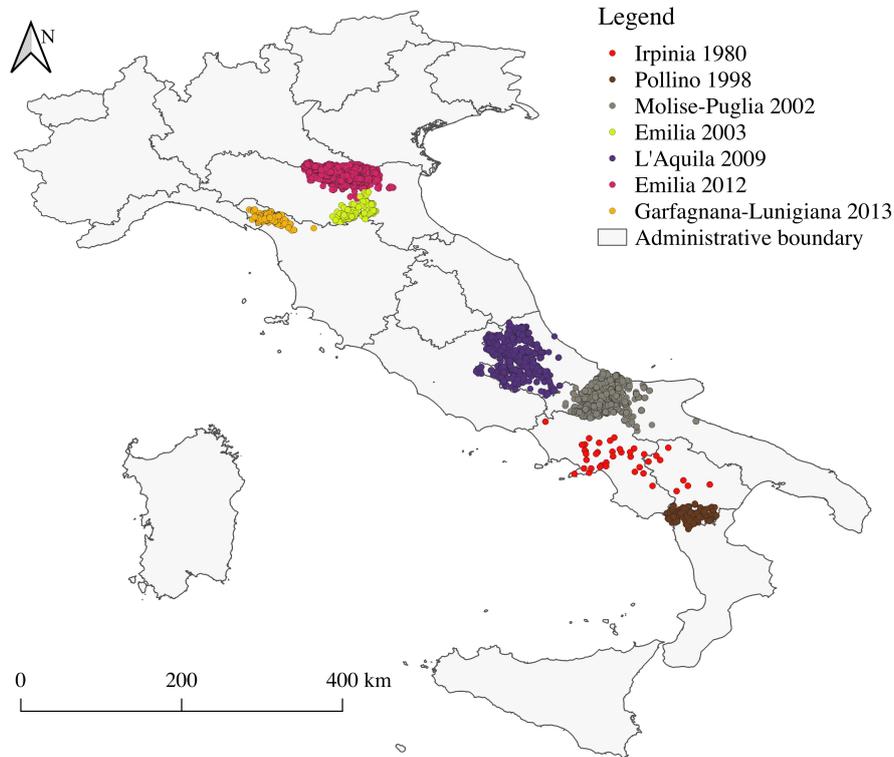
126 **Regularity**: building regularity in terms of plan and elevation, classified as either irregular or regular.

127 **Construction material**: vertical elements: good and poor-quality masonry, good and poor quality  
128 mixed frame masonry, reinforced concrete frame and wall, steel frame, and other.

129 For features defined as value ranges (e.g., date of construction/renovation, floor area, and building  
130 height), the average value was used. Furthermore, the Irpinia-1980 building damage portfolio (E1) was  
131 constructed using the specific Irpinia-1980 damage survey form, while the AeDES damage survey form  
132 was used for the others. The Irpinia-1980 dataset will therefore be analysed separately.

133 Building damage data from earthquake surveys other than the Irpinia-1980 earthquake damage survey  
134 primarily include damaged buildings. This is because the data was collected based on requests for  
135 damage assessments after the earthquake event (Dolce et al. 2019). The damage information in the  
136 DaDO database is still relevant for testing the machine learning models for heuristic damage  
137 assessment. Mixing these datasets to train machine learning models can lead to biased outcomes.  
138 Therefore, the machine learning models were developed on the other earthquake dataset excluding the  
139 Irpinia dataset, and the Irpinia earthquake dataset was used only in the testing phase.

140 The distribution of the samples is very imbalanced (Fig. 2): for example, there is a small proportion of  
 141 buildings in DG4+DG5 (7.59%), and a large majority of masonry (65.47%) compared to reinforced  
 142 concrete frame (21.31%) buildings. This imbalance should be taken into account when defining the  
 143 machine learning models.  
 144



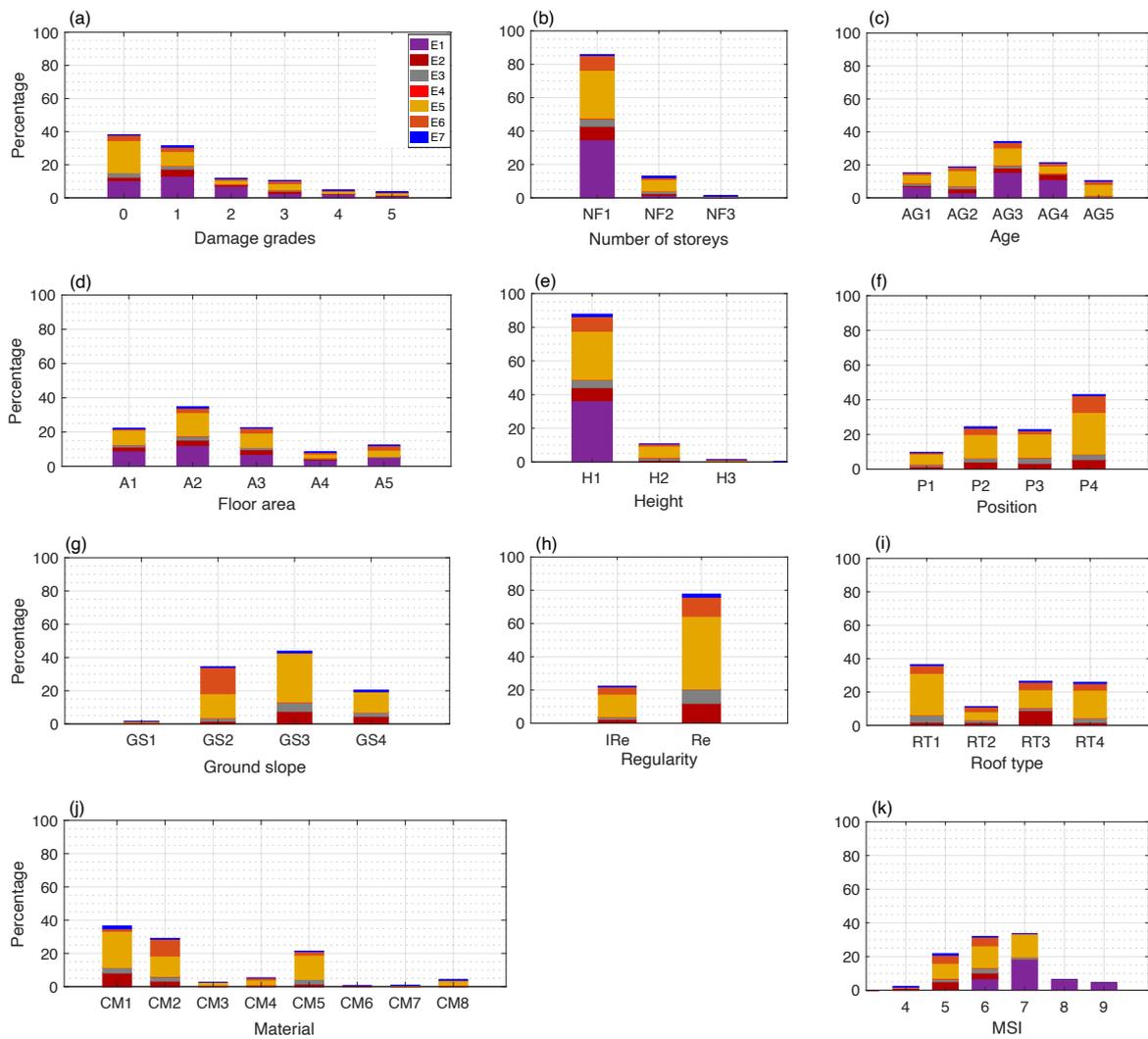
145  
 146 **Figure 1.** Geographic location of the buildings considered in this study.  
 147

148 To consider spatially-distributed ground motion, the original DaDO data are supplemented with the  
 149 main event macroseismic intensities (MSI) provided by the United States Geological Survey (USGS)  
 150 ShakeMap tool (Wald et al., 2005). Macroseismic intensities (MSI) given in terms of modified Mercalli  
 151 intensities are considered and assigned to buildings based on their location. The distribution of MSI  
 152 values in the database is shown in Fig. 2k.

153  
 154 **Table 2.** Distribution of the different features used in this study.

No.	Parameters	Data type	Distribution (%)	Remarks
1	No damage	DG0	43.63	Fig. 2a
	Slight damage	DG1	28.90	
	Moderate damage	DG2	7.41	
	Substantial damage	DG3	12.48	
	Very heavy damage	DG4	3.94	
	Total collapse	DG5	3.65	

2	Number of storeys	0-3	NF1	Numerical	85.81	Fig. 2b
		3-5	NF2		13.01	
		> 5	NF3		1.19	
3	Age (years)	0-20	AG1	Numerical	15.22	Fig. 2c
		21-40	AG2		18.81	
		41-60	AG3		34.15	
		61-80	AG4		21.34	
		>80	AG5		10.49	
4	Floor area (square metres)	0-50	A1	Numerical	22.16	Fig. 2d
		50-100	A2		34.73	
		100-150	A3		22.53	
		150-200	A4		8.32	
		> 200	A5		12.26	
5	Height (metres)	0-10	H1	Numerical	87.78	Fig. 2e
		10-15	H2		10.69	
		>15	H3		1.50	
6	Position	Corner	P1	Categorical	9.71	Fig. 2f
		Extreme	P2		24.47	
		Internal	P3		22.80	
		Isolated	P4		43.02	
7	Ground slope	Ridge	GS1	Categorical	2.62	Fig. 2g
		Plain	GS2		34.25	
		Moderate slope	GS3		43.74	
		Steep Slope	GS4		20.39	
8	Regularity	Irregular in plan and elevation	IR	Categorical	22.28	Fig. 2h
		Regular in plan and elevation	Re		77.72	
9	Roof type	Heavy no thrust	R1	Categorical	36.43	Fig. 2i
		Heavy thrust	R2		11.25	
		Light thrust	R3		26.48	
		Light no thrust	R4		25.83	
10	Material	Masonry poor quality	CM1	Categorical	36.51	Fig. 2j
		Masonry good quality	CM2		28.96	
		Mixed frame masonry poor quality	CM3		2.64	
		Mixed frame masonry good quality	CM4		5.21	
		Reinforced concrete frame	CM5		21.31	
		Reinforced concrete wall	CM6		0.42	
		Steel frame	CM7		0.09	
		Other	CM8		4.10	



156

157 **Figure 2.** Distribution of the different features in the database. E1, E2, E3, E4, E5, E6, and E7, representing  
 158 Irpinia-1980, Pollino-1998, Molise-Puglia-2002, Emilia-Romagna-2003, L'Aquila-2009, Emilia-Romagna-2012,  
 159 and Garfagnana-Lunigiana-2013 building damage portfolios, respectively. The y-axis is the percentage  
 160 distribution and the x-axis is (a) Damage grade, (b) Number of storeys (NF1: 0-3, NF2: 3-5, NF3: >5), (c) Building  
 161 age (AG1: 0-20, AG2: 21-40, AG3: 41-60, AG4: 61-80, AG5: >80), (d) Floor area (A1: 0-50, A2: 51-100, A3:  
 162 101-150, A4: 151-200, A5: >200), (e) Height (H1: 0-10, H2: 10-15, H3: >15), (f) Building position (P1: corner,  
 163 P2: extreme, P3: internal, P4: isolated), (g) Ground slope condition (GS1: ridge, GS2: plain, GS3: moderate slope,  
 164 GS4: steep slope), (h) Regularity in plan and elevation (IRe: irregular, Re: Regular), (i) Roof type (RT1: heavy  
 165 no thrust, RT2: heavy thrust, RT3: light no thrust, RT4: light thrust), (j) Construction material (CM1: poor-quality  
 166 masonry, CM2: good-quality masonry, CM3: poor-quality mixed frame masonry, CM4: good-quality mixed  
 167 frame masonry, CM5: reinforced concrete frame, CM6: reinforced concrete wall, CM7: steel frames, CM8: other),  
 168 and (k) macro-seismic intensity.

169

### 170 **3. Method**

#### 171 **3.1. Machine learning models**

172 Ghimire et al. (2022) applied classification- and regression-based machine learning models to the  
173 damage observed after the 2015 Gorkha Nepal earthquake (NPC, 2015). The main concepts for method  
174 selection, the definition of the dataset for training and testing, and the representation of model  
175 performance are presented here.

176 To develop the heuristic damage assessment model, the damage grades are considered as the target  
177 feature. The damage grades are discrete labels, from DG0 to DG5. Three most advanced classification  
178 and regression machine learning algorithms were selected: random forest (RFC) and regression (RFR)  
179 (Breiman, 2001), gradient boosting classification (GBC) and regression (GBR) (Friedman, 1999), and  
180 extreme gradient boosting classification (XGBC) and regression (XGBR) (Chen and Guestrin, 2016).  
181 A label (or class) was thus assigned to the categorical response variables (DG) for the classification-  
182 based machine learning models. For the regression-based machine learning models, DG is converted  
183 into a continuous variable to minimize misclassifications (Ghimire et al., 2022). For the regression-  
184 based machine learning models, DG is converted into a continuous variable as tested by  
185 Ghimire et al. (2022): first, the damage grades were ordered and considered as a continuous  
186 variable ranging between 0 (DG0) and 5 (DG5). Because the regression model outputs a real  
187 value between 1 and 5 and not an integer, we rounded the output (real number) to the nearest  
188 integer to plot the confusion matrix. However, the error matrices were computed without  
189 rounding the model outputs to the nearest integer.

190  
191 Building features and macroseismic intensities were considered as input features. A one-hot encoding  
192 technique was used to convert the categorical features (i.e., ground slope condition, building position,  
193 roof type, construction material) into binary values (1 or 0), resulting in 28 input variables (Tab. 2). No  
194 input features were removed from the dataset: some building features (e.g., number of storeys and  
195 height) may be correlated but we assumed that the presence of correlated features does not impact the  
196 overall performance of these machine learning methods (Ghimire et al., 2022). No specific data cleaning  
197 methods were applied to the DaDO database.

198 The machine learning algorithms from the Scikit-learn package developed in Python (Pedregosa et al.,  
199 2011) were applied. The machine learning models were trained and tested on the randomly selected  
200 training (60% of the dataset) and testing (40% of the dataset) subsets of data, considering a single  
201 earthquake dataset or the whole DaDO dataset. The testing subset was kept hidden from the model  
202 during the training phase.

203

#### 204 **3.2. Machine learning model efficacy**

205 The efficacy of the heuristic damage assessment model (i.e., its ability to predict damage to a  
206 satisfactory or expected degree) was analysed in three stages: comparison of the efficacy of the machine  
207 learning models using metrics; analysis of specific issues related to machine learning using the selected  
208 models; and application of the heuristic model to the whole DaDO dataset.

209

### 210 **3.2.1 First stage: model selection**

211 In the first stage, only the L'Aquila-2009 portfolio was considered for the training and testing phases.  
212 This is the largest dataset in terms of the number of buildings and was obtained using the AeDES survey  
213 format (Baggio et al., 2007; Dolce et al., 2019). Model efficacy was provided by the confusion matrix,  
214 which represents model prediction compared with the so-called "ground truth" value. Accuracy was  
215 then represented on the confusion matrix by the ratio of the number of correctly predicted DGs to the  
216 total number of observed values per DG ( $A_{DG}$ ).

217 Total accuracy ( $A_T$ ) was computed as the ratio of the number of correctly predicted DGs to the total  
218 number of observed values.  $A_T$  and  $A_{DG}$  values close to 1 indicate high efficacy. Moreover, the  
219 quantitative statistical error was also calculated as the mean of the absolute value of errors (MAE) and  
220 the mean squared error (MSE) (MAE and MSE values close to 0 indicate high efficacy). For  
221 classification-based machine learning models, the ordinal value of the DG was used to calculate the  
222 MAE and MSE scores directly. For the regression-based machine learning models, the output DG  
223 values were rounded to the nearest integer for the accuracy scores plotted for the confusion matrix, but  
224 not for the MAE and MSE value calculations.

225

### 226 **3.2.2 Second stage: machine learning related issues**

227 In the second stage, the best heuristic model for damage assessment was selected based on the highest  
228 efficacy, and used to analyse and test specific issues related to machine learning: (1) the imbalance  
229 distribution of DGs in the DaDO, (2) the performance of the selected model when only some basic, but  
230 accurately assessed, building features are considered (i.e., number of storeys, location, age, floor area),  
231 and (3) the simplification of the heuristic model, in the sense that DGs are grouped into a traffic-light-  
232 based classification (i.e., green, yellow and red, corresponding to  $DG_0+DG_1$ ,  $DG_2+DG_3$  and  
233  $DG_4+DG_5$ , respectively). In the second stage, the issues related to machine learning were first analysed  
234 using the L'Aquila-2009 portfolio. The whole DaDO dataset was then used.

235

### 236 **3.2.2 Third stage: application to the whole DaDO portfolio and comparison with Risk-UE**

237 In the third stage, several learning and testing sequences were considered, with the idea of moving to  
238 an operational configuration in which past information is used to predict damage from future  
239 earthquakes: either learning based on a portfolio of damage caused by one earthquake and tested on  
240 another portfolio, or learning based on a series of damage portfolios and tested on the portfolio of  
241 damage caused by an earthquake placed in the chronological continuity of the earthquake sequence

242 considered. In this stage, the efficacy of the heuristic damage assessment model was analysed by  
 243 comparing the prediction values with the so-called “ground truth” values through the error distribution,  
 244 as follows:

$$245 \quad \varepsilon_d(\%) = \left(\frac{n_e}{N}\right) * 100 \quad (1)$$

246 where  $n_e$  is the total number of buildings at a given error level (difference between observed and  
 247 predicted DGs),  $N$  is the total number of buildings in the damage portfolio.

248 In this stage, the efficacy of the heuristic damage assessment model was compared with the  
 249 conventional damage prediction framework proposed by the RISK-UE method (Milutinovic and  
 250 Trendafiloski, 2003). The RISK-UE method assigns a vulnerability index (IV) to a building, based on  
 251 its construction material and structural properties (e.g., height, building age, position, regularities,  
 252 geographic location, etc.). For a given level of seismic demand (MSI), the mean damage ( $\mu_d$ ) and the  
 253 probability,  $p_k$ , of observing a given damage level  $k$  ( $k = 0$  to  $5$ ) are given by:

$$254 \quad \mu_d = 2.5 \left[ 1 + \tanh \left( \frac{MSI + 6.25IV - 13.1}{2.3} \right) \right] \quad (2)$$

$$255 \quad p_k = \frac{5!}{k!(5-k)!} \left(\frac{\mu_d}{5}\right)^k \left(1 - \frac{\mu_d}{5}\right)^{5-k} \quad (3)$$

256  
 257  
 258 Herein, comparing the heuristic model and the RISK-UE method amounts to considering the following  
 259 steps, based on the equations given by RISK-UE:

260 **Step 1** - The buildings in the training and testing datasets are grouped into different classes according  
 261 to construction material.

262 **Step 2** - For a given building class in the training dataset, computation of

263 **Step 2.1** - mean damage ( $\mu_d$ ) using the observed damage distribution at a given MSI value by:

$$264 \quad \mu_d = \sum_{k=0}^5 p_k k \quad (4)$$

265  
 266 **Step 2.2** - vulnerability index (IV) with the  $\mu_d$  obtained in step 2.1 by:

$$267 \quad IV = \frac{1}{6.25} \left[ 13.1 - MSI + 2.3 \left( \tanh^{-1} \left( \frac{\mu_d}{2.5} - 1 \right) \right) \right] \quad (5)$$

268  
 269  
 270 **Step 3** - For the same building class in the test dataset, calculation of

271 **Step 3.1** - mean damage ( $\mu_d$ ) Eq. 2 for a given MSI value with the value of IV obtained in step  
 272 2.2;

273 **Step 3.2** - damage probability ( $p_k$ ) Eq. 3 with the value of  $\mu_d$  obtained in step 3.1;

275 **Step 3.3** - distribution of buildings in each damage grade within a range of MSI values observed  
276 in the test dataset as follows:

277

$$278 \quad N_{pred,k} = \sum_{MSI} p_k n_{obs,MSI} \quad (6)$$

279

280 where  $n_{obs,MSI}$  is the total number of buildings observed in the test set for a given MSI  
281 value;

282 **Step 3.4** –absolute error ( $\varepsilon_k$ ) in each damage level k, given by:

$$283 \quad \varepsilon_k = \left| \frac{N_{obs,k} - N_{pred,k}}{N} \right| \quad (7)$$

284

285 where,  $N_{obs,k}$  is the total number of buildings observed in the given damage grade k.

286

287 Similarly, the heuristic damage assessment model was also compared with the mean damage  
288 relationship (Eq. 4) applied to the test set. Thus, for each building class in the test set, the error value  
289 (Eq. 7) for each DG was computed from the  $\mu_d$  on the observed damage using Eq. (4), the probability  
290  $p_k$  of obtaining a given DG k (k= 0 to 5) using Eq. (3), and the distribution of buildings in each DG  
291  $N_{pred,k}$  for a given MSI value using Eq. (6).

292

## 293 **4. Result**

### 294 **4.1 First stage: model selection**

295 The efficacy of the regression (RFR, GBR, XGBR) and classification (RFC, GBC, XGBC) machine  
296 learning models trained and tested on the randomly selected 60% (training set) and 40% (test set) of the  
297 2009 -L’Aquila earthquake building damage portfolio is summarized in Table 3. The hyperparameters  
298 indicated in Tab. 3 were chosen after tests performed by Ghimire et al. (2022). The regression-based  
299 machine learning models RFR, GBR and XGBR yielded similar MSE scores (1.22, 1.22 and 1.21) and  
300 accuracy scores ( $A_T = 0.49, 0.50$  and  $0.50$ ), considering the five DGs of the EMS-98 scale. In the  
301 confusion matrix (Fig. 3a: RFR, Fig. 3b: GBR, and Fig. 3c: XGBR), the accuracy  $A_{DG}$  values show that  
302 the efficacy of these models is higher for the lower DGs (around 60% for DG0 and 55% for DG1) and  
303 lower for the higher DGs (6% and 1% of the buildings are correctly classified in DG4 and DG5,  
304 respectively).

305 For the classification-based machine learning models, the XGBC model ( $[MSE, A_T] = [1.78, 0.59]$ ) was  
306 more effective than the RFC ( $[MSE, A_T] = [1.86, 0.57]$ ) and GBC ( $[MSE, A_T] = [1.80, 0.58]$ ) models,  
307 considering the EMS-98 scale. In the confusion matrix (Fig. 3d: RFC, Fig. 3e: GBC, and Fig. 3f:  
308 XGBC), the accuracy  $A_{DG}$  values also show higher model efficacy for the lower DGs (86% for DG0  
309 and 39% for DG1) and lower efficacy for the higher DGs (5%, 23%, 12% and 17% buildings correctly  
310 classified in DG2, DG3, DG4 and DG5, respectively).

311

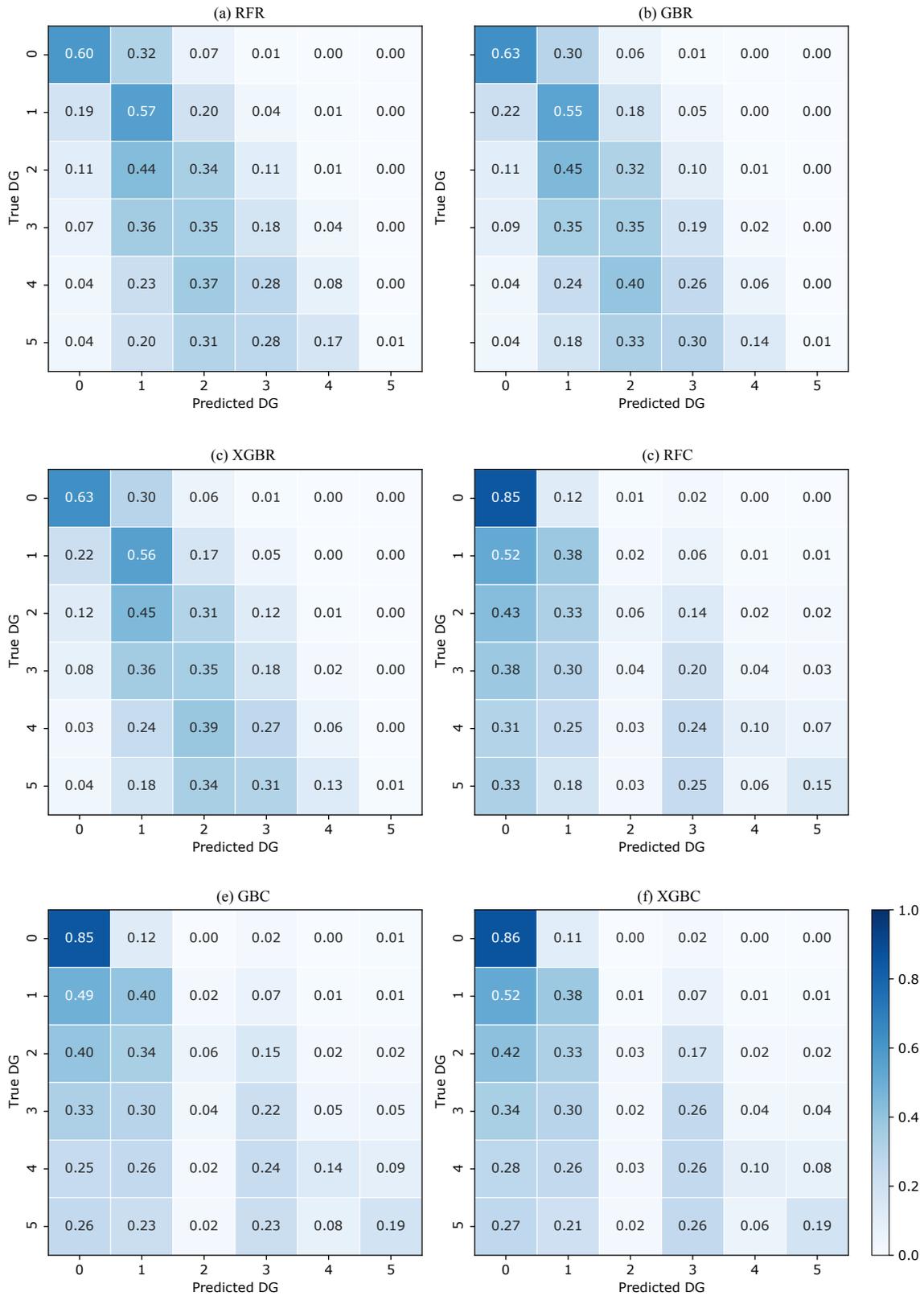
312 **Table 3.** Summary of optimized hyperparameters parameters, accuracy  $A_T$  and quantitative statistical  
 313 error values for the regression-based and classification-based machine learning methods in the test set.  
 314 The parameters are the hyperparameters chosen for the machine learning models (the other model  
 315 parameters not mentioned here are the default parameters in the Scikit-learn documentation (Pedregosa  
 316 et al., 2011)). The best accuracy and error values are indicated in bold. The optimum hyperparameters  
 317 were selected thanks to k-fold cross-validation (with 10-fold), by randomly select a % for training  
 318 and % for testing, for different combination of hyperparameters and the optimum evaluated in  
 319 terms of performance metrics on testing is finally selected.

Method	Parameters	Accuracy $A_T$	MSE	MAE
RFR	n_estimators = 1000 max_depth = 25	0.49	1.22	0.77
GBR	n_estimators = 1000 max_depth = 10 learning_rate = 0.01	0.50	1.22	0.77
XGBR	n_estimators = 1000 max_depth = 10 learning_rate = 0.01	0.50	<b>1.21</b>	0.76
RFC	no_estimators = 1000 max_depth = 25	0.57	1.86	0.77
GBC	no_estimators = 1000 max_depth = 10 learning_rate = 0.01	0.58	1.80	0.77
XGBC	n_estimators = 1000 max_depth = 10 learning_rate = 0.01	<b>0.59</b>	1.78	<b>0.74</b>

320

321 The classification-based machine learning models thus yielded slightly better predictive efficacy, but  
 322 still lower than recent studies applied to other datasets (Ghimire et al., 2022; Harirchian et al., 2021;  
 323 Mangalathu et al., 2020; Roeslin et al., 2020; Stojadinović et al., 2021). The high classification error in  
 324 the higher DGs could be related to the characteristics of the building portfolio and the imbalance of DG  
 325 distribution. Among the classification methods, the XGBC model showed slightly higher classification  
 326 efficacy; the XGBC model was therefore selected for the next stages 2 and 3.

327



328

329 **Figure 3.** Normalized confusion matrix between predicted and observed DGs. The values given in each main  
 330 diagonal cell are the accuracy scores  $A_{DG}$ . All values are also represented by the colour scale.

331

332 **4.2 Second stage: issues related to machine learning**

333 **4.2.1 Imbalance distribution of the DGs in the DaDO**

334 The efficacy of the heuristic damage assessment model depends on the distribution of target features in  
 335 the training dataset. This can lead to low prediction efficacy, especially for minority classes (Estabrooks  
 336 & Japkowicz 2001; Japkowicz & Stephen 2002; Branco et al. 2017; Ghimire et al., 2022). The previous  
 337 section reports significant misclassification associated with the highest DGs for all classification- and  
 338 regression-based models (Fig. 3), i.e., for the DGs with the lowest number of buildings (Fig. 2a). The  
 339 efficacy of the XGBC model is analysed below, addressing the class-imbalance issue with data  
 340 resampling techniques applied to the training phase and considering the L’Aquila-2009 portfolio.

341

342 Four strategies to solve the class imbalance issue were tested:

343 (a) random undersampling: randomly selecting the number of data entries in each class equal to the  
 344 number of data entries in the minority class (DG4 in our case);

345 (b) random oversampling: randomly replacing the number of data entries in each class equal to the  
 346 number of data entries in the majority class (DG0 in our case);

347 (c) Synthetic Minority Oversampling Technique (SMOTE): creating an equal number of data entries in  
 348 each class by generating synthetic samples by interpolating the neighbouring data in the minority class;

349 (d) a combination of oversampling and undersampling methods: oversampling of the minority class  
 350 using the SMOTE method, followed by the Edited Nearest Neighbours (ENN) undersampling method  
 351 to eliminate data that is misclassified by its three nearest neighbours (SMOTE-ENN).

352

353 Fig. 4 shows the confusion matrices of the four strategies considered for the class imbalance issue.  
 354 Compared with Fig. 3f (i.e., XGBC), the effects of addressing the issue of imbalance were as follows:

355 (a) undersampling (Fig. 4a):  $A_{DG}$  value increased by 20/22/26% for DG2/DG4/DG5 and decreased by  
 356 29% for DG0.

357 (b) oversampling (Fig. 4b):  $A_{DG}$  value increased by 11/16/18% for DG2/DG4/DG5 and decreased by  
 358 13% for DG0

359 (c) SMOTE (Fig. 4c):  $A_{DG}$  value increased by 4/1/4% for DG2/DG4/DG5 and decreased by 3% for  
 360 DG0

361 (d) SMOTE-ENN (Fig. 4d):  $A_{DG}$  value increased by 13/9/8% for DG2/DG4/DG5 and decreased by 25%  
 362 for DG0.

363 The  $A_T$ , MAE and MSE scores are given in Table 4 with the associated effects.

364

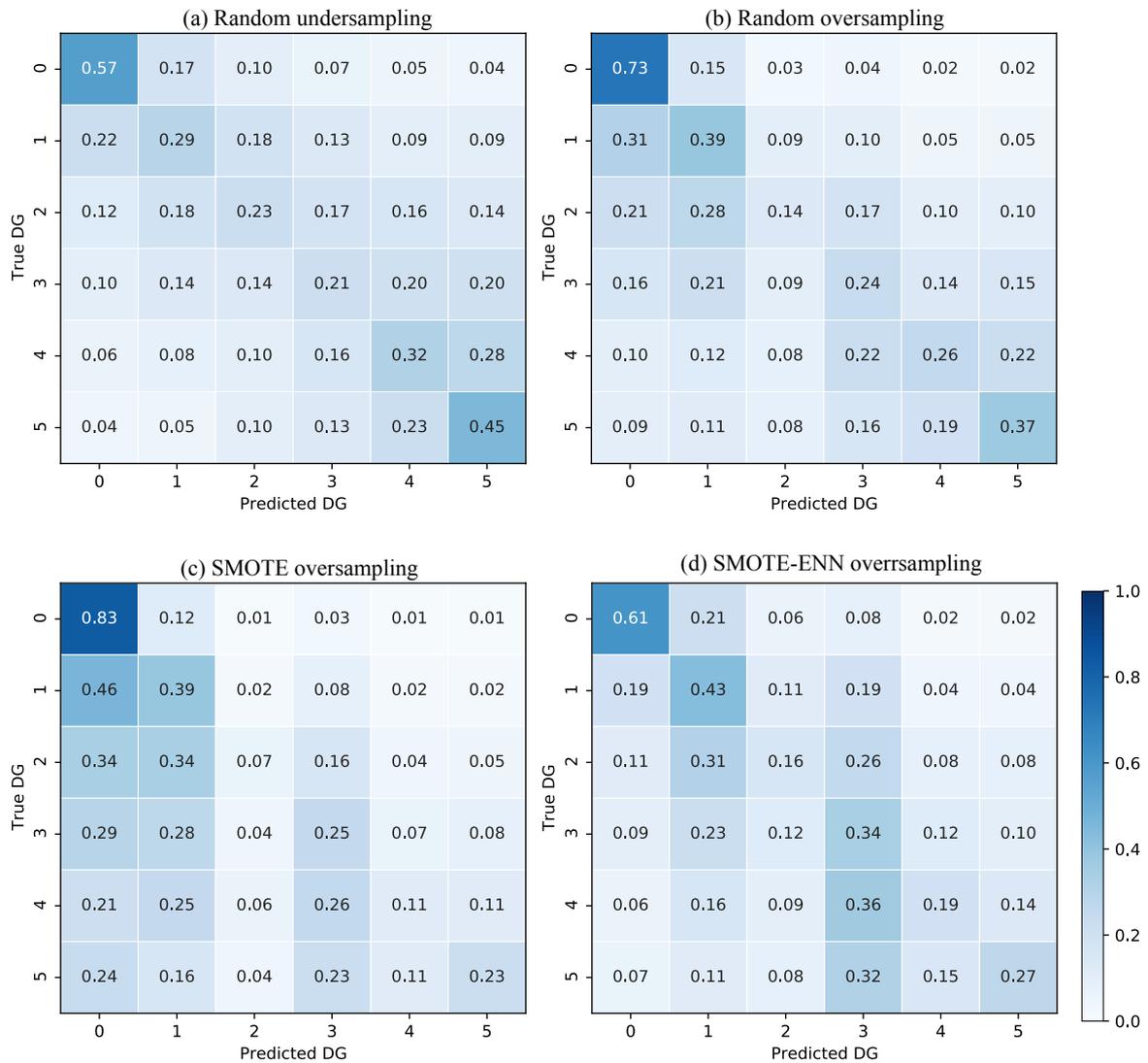
365 Table 4 – Scores of the accuracy  $A_T$ , MSE and MAE metrics in the test set considering the imbalance  
 366 issue and their variation  $\Delta$  compared with values without consideration of the imbalance.

Method	Accuracy $A_T$		MSE		MAE	
	Scores	$\Delta$	Score	$\Delta$	Score	$\Delta$

Undersampling	0.26	-0.33	1.24	-0.34	1.20	0.46
Oversampling	0.53	-0.06	2.13	0.35	0.86	0.12
SMOTE	0.57	-0.02	1.87	0.09	0.77	0.03
SMOTE-ENN	0.49	-0.10	2.28	0.50	0.93	0.19

367

368 In conclusion, the random oversampling method improves prediction in the minority class without  
 369 significantly decreasing prediction in the majority class. The random oversampling method was  
 370 therefore applied in this study.



371

372 **Figure 4.** Confusion matrices for the four methods to solve the DG imbalance issue in the DaDO. The values  
 373 given in each main diagonal cell are the accuracy scores  $A_{DG}$ . All values are also represented by the colour scale.  
 374

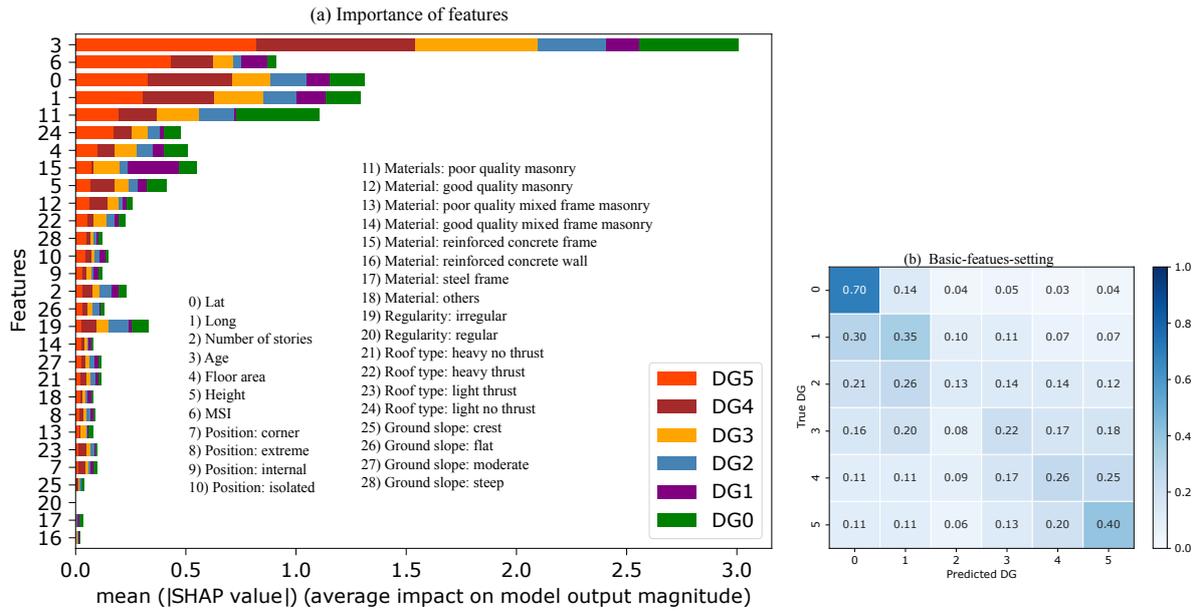
374

375 **4.2.2 Testing the XBGC model with basic features**

376 This section begins by exploring the importance of each feature in the heuristic damage assessment  
377 model applied to the L'Aquila-2009 portfolio. We used the Shapely Additive Explanations (SHAP)  
378 method developed by Lundberg and Lee (2017). The SHAP method compares the efficacy of the model  
379 with and without considering each input feature to measure its average impact, provided in terms of  
380 mean absolute SHAP values.

381 Figure 5a shows the average SHAP value associated with each feature considered in this study as a  
382 function of DG. The most weighted features are building age, location (latitude and longitude), material  
383 (poor quality masonry, RC frame), MSI, roof type, floor area, and height. Interestingly, the mean SHAP  
384 values are dependent on the DG, i.e., the weight of the feature is not linear depending on the DG  
385 considered; this is never taken into account in vulnerability methods. For example, Scala et al. (2022)  
386 and Del Gaudio et al. (2021) observed a decrease in the vulnerability of structures as construction year  
387 increases, without distinguishing the DG considered, which is not the case herein. Note also that the  
388 importance score associated with the location feature can indirectly capture variations in local  
389 geological properties and the spatially distributed vulnerability associated with the built-up area of the  
390 L'Aquila-2009 portfolio (e.g., the distinction between the historic town and more modern urban areas).  
391 Furthermore, the average SHAP value obtained for poor quality masonry buildings for DG3/DG4/DG5  
392 confirms the same high vulnerability of this typology as in the EMS-98 scale (Grünthal, 1998),  
393 regardless of DG.

394 Some basic features of the building (e.g., location, age, floor area, number of storeys, height) are  
395 observed with a high mean SHAP value (Fig. 5a). Compared with others, these five basic features can  
396 be easily collected from the field or provided by national census databases, for example. Fig. 5b shows  
397 the efficacy of the heuristic damage assessment model using XGBC trained with a set of easily  
398 accessible building features (i.e., basic-features-setting: geographic location, floor area, number of  
399 stories, height, age, MSI), after addressing the class-imbalance issue using the random oversampling  
400 method. Compared with Fig. 4b (considering all features and named as the full-features-setting), the  
401 XGBC model with the basic-features-setting (Fig. 5b) gives almost the same efficacy with only a 6%  
402 average reduction in the accuracy scores.



403

404

**Figure 5.** (a) Graphic representation of the importance scores associated with the different input features considered for the XGBC model. The features (the same as in Fig. 2) considered in this study are on the y-axis, and the x-axis is the mean SHAP score according to DG. (b) Confusion matrices considering the basic-features-setting. The values given in each main diagonal cell are the accuracy scores  $A_{DG}$ . All values are also represented by the colour scale.

408

409

410

#### 4.2.3 Testing the XBGC model with the traffic-light system for damage grades

411

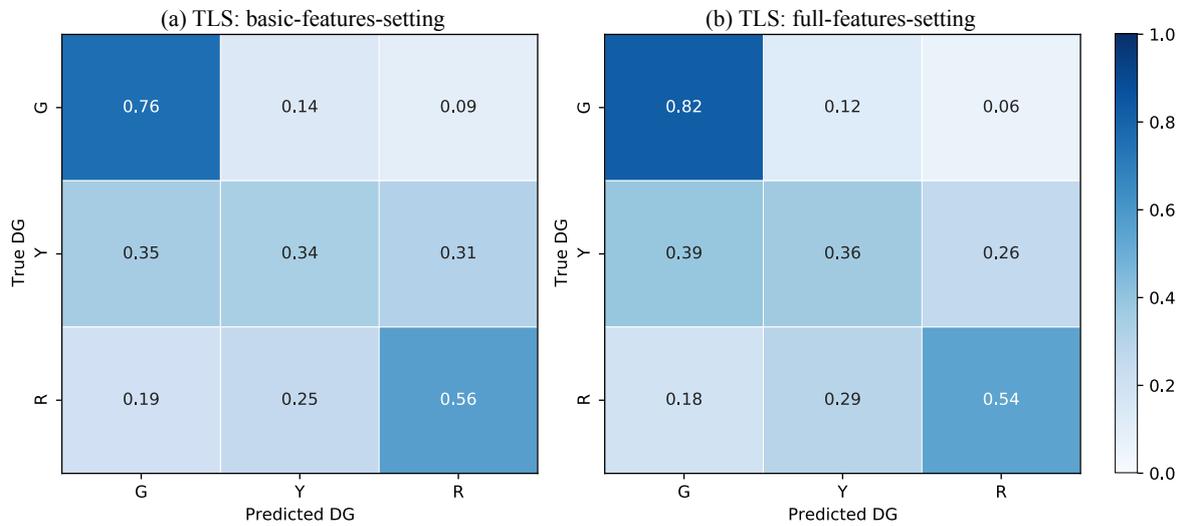
In this section, a simplified version of the DG scale was used, in the sense that the DGs are classified according to a traffic-light system (TLS) (i.e., green G, yellow Y and red R classes, corresponding to DG0+DG1, DG2+DG3 and DG4+DG5, respectively), as monitored during post-earthquake emergency situations (Mangalathu et al., 2020; Riedel et al., 2015; ATC, 2005; Bazzurro et al., 2004). For the TLS-based damage classification, the XGBC model (after oversampling to compensate of the imbalance issue) with the basic-features-setting applied to the L'Aquila-2009 portfolio (Fig. 6a) gives almost the same efficacy compared to the full-features-setting (Fig. 6b). For example, accuracy values  $A_{DG}$  using the basic-features-setting and the full-features-setting were 0.76/0.34/0.56 and 0.82/0.36/0.54 for G/Y/R classes, with the accuracy score  $A_T$  of 0.68 and 0.72, respectively. Mangalathu et al. (2020), Roslin et al., (2020), and Harirchian et al., (2021) reported similar damage grade classification accuracy values of 0.66, 0.67, and 0.65 respectively.

422

The efficacy of the heuristic damage assessment model using TLS-based damage classification indicates that classifying damage into three classes is much easier for the machine learning model compared with the six-class classification system (EMS-98 damage classification). This is also observed during damage surveys in the field, which sometimes find it hard to distinguish the intermediate damage grades, such as DG2 and DG3, or DG3 and DG4. Similar observations have been

426

427 reported in previous studies by Guettiche et al., (2017); Harirchian et al., (2021); Riedel et al., (2015);  
 428 Roeslin et al., (2020) and Stojadinović et al., (2021).  
 429



430  
 431 **Figure 6.** Confusion matrices for (a) the basic-features-setting and (b) the full-features-setting using the traffic-  
 432 light (TLS)-based classification, grouping the EMS-98 damage grades (DG) into three classes (green for no or  
 433 slight damage; yellow for moderate damage; and red for heavy damage). The values given in each main diagonal  
 434 cell are the accuracy scores  $A_{DG}$ . All values are also represented by the colour scale.

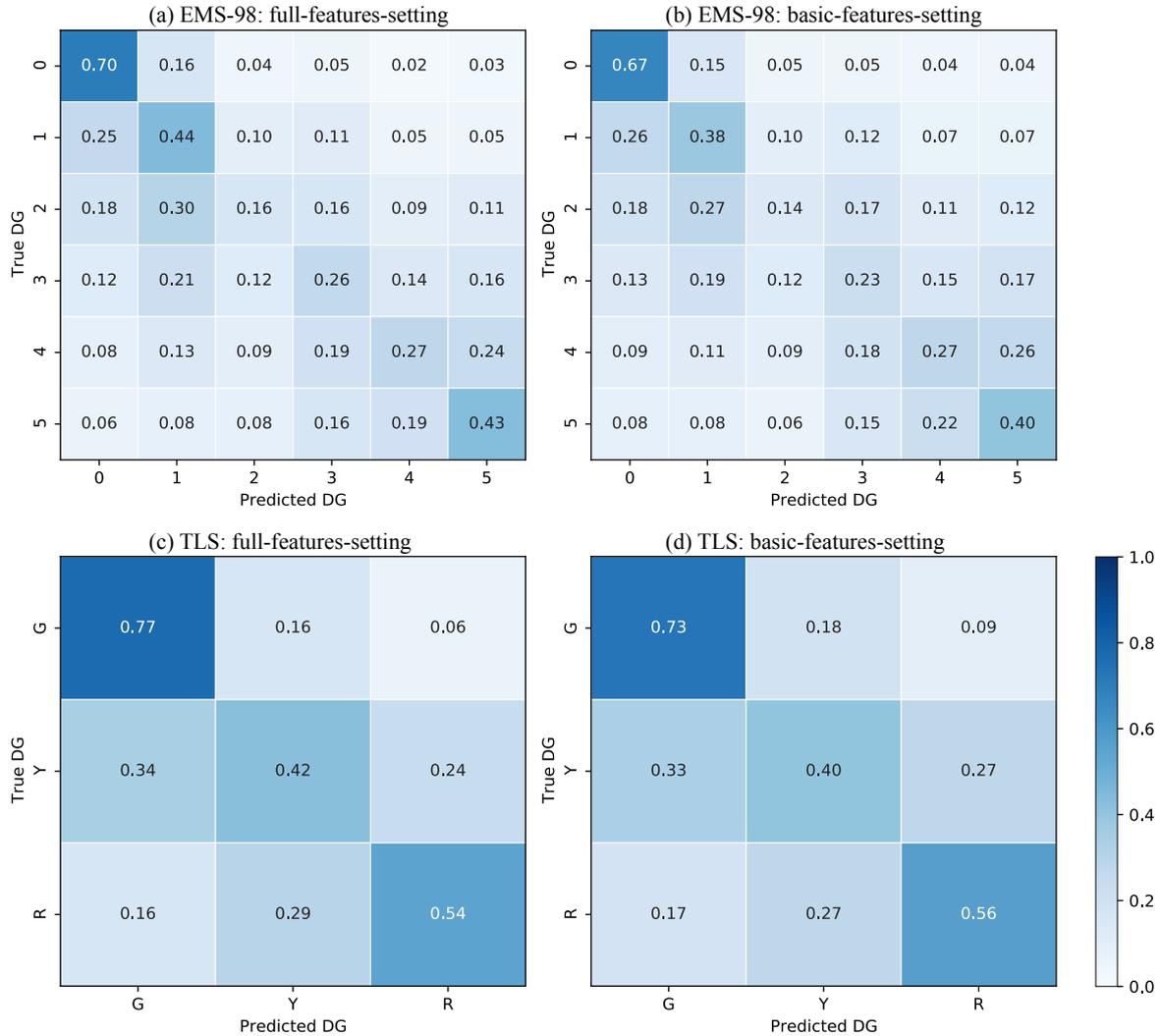
435  
 436 **4.2.4 Testing the XGBC model with the whole dataset**

437 The efficacy of the XGBC model was tested using a dataset with six building damage portfolios,  
 438 excluding the 1980-Irpinia building damage portfolio. The XGBC model was trained and tested on the  
 439 randomly selected 60% (training set) and 40% (test set) of the dataset for EMS-98/TLS damage  
 440 classification, with two sets of features (full-features-setting and basic-features-setting), applying the  
 441 random oversampling method to compensate for class-imbalance issues. Fig.7 shows the associated  
 442 confusion matrix.

443 The basic-features-setting resulted in a similar level of damage prediction compared with the full-  
 444 features setting for both EMS-98 and TLS-based damage classification systems. For EMS-98 damage  
 445 classification (Fig. 7a, b), the accuracy  $A_{DG}$  scores indicated in the confusion matrices are almost the  
 446 same for the basic-features-setting and the full-features-setting. Furthermore, the accuracy  $A_T$  and MAE  
 447 scores are also almost the same (0.45 and 1.08 for the basic-features-setting and 0.48 and 0.95 for the  
 448 full-features-setting).

449 Likewise, for TLS-based damage classification (Fig. 7c, d), the accuracy values  $A_{DG}$  for the basic-  
 450 features-setting and the full-features-setting are almost the same, with similar accuracy  $A_T$  and MAE  
 451 scores (0.63/0.45 and 0.67/0.39, respectively).

452  
 453



454  
455  
456  
457  
458  
459

**Figure 7.** Confusion matrices for EMS-98 (a, b) and TLS (c, d) damage classification systems using the basic- and full-features-settings (green for no or slight damage; yellow for moderate damage; red for heavy damage) with (c) the full-features-setting and (d) the basic-features-setting. The values given in each main diagonal cell are the accuracy scores  $A_{DG}$ . All values are also represented by the colour scale.

### 460 4.3 Third stage: application to the whole DaDO portfolio and comparison with Risk-UE

461 In this section, the efficacy of the heuristic damage assessment model was considered for building  
462 damage predictions, without respecting the time frame of the earthquakes. Two scenarios were  
463 considered: (1) a single building damage portfolio was used for training and the model was then tested  
464 on the others (named single-single), in situations using a single portfolio to predict future damage; and  
465 (2) some building damage portfolios were used for training but testing was performed on a single  
466 portfolio (named aggregate-single), i.e. a larger number of damage portfolios were used as a training  
467 set to predict the damage caused by the next earthquake. The model XGBC was applied with the basic-  
468 features-setting (number of storeys, building age, floor area, height, MSI for EMS-98) and EMS-98-  
469 and TLS-based damage classification.

470

#### 471 **4.3.1 Single-single scenario**

472 First, a series of building damage portfolios, concerning earthquakes occurring in northern or southern  
473 Italy and of different magnitudes, was used for training and testing:

474 (i) Training set: E3 – test set: E1, E5, E7.

475 (ii) Training set: E5 – test set: E1, E3, E7.

476 (iii) Training set: E7 – test set: E1, E3, E5.

477

478 Figure 8 shows the distribution of correct DG classification (i.e.,  $1 - \varepsilon_d$  in % given by Eq. 1) observed  
479 for each building for the EMS-98 damage grade (8a) and the TLS (8b) systems. The x-axis represents  
480 the incremental error in the damage grade (e.g., 1 corresponds to the delta of damage grade between  
481 observation and prediction, regardless of the DG considered).

482 For the EMS-98 damage scale, correct classification (x-value centred on 0) in the range of 31% to 48%  
483 was found, depending on the training/test data sets. The error distribution is quite wide with incorrect  
484 predictions of +/-1 DG in the range of +/- 13-35%. Remarkably, when considering the E1 portfolio  
485 (Irpinia-1980), for which the post-earthquake inventory was based on another form, as the test set, the  
486 error is larger. The predictions at +/-1 DG (i.e., the sum of the x-values Fig. 8a between -1 and +1) were  
487 70.5%, 69.9% and 72.8% with portfolios E3, E5 and E7 as the test set, respectively, for an average of  
488 71%. For the other portfolios, the average of the predictions at +/- 1 DG was 77%, 78% and 77%,  
489 respectively, for portfolios E5, E3 and E7 as the test set. This tendency was also observed for the TLS  
490 damage system (Fig. 8b). In this case, the classification of the E1 portfolio was correct on average  
491 (average of x-values centred on 0) at 63% and equal to 72%, 73%, and 70.5% for the test on portfolios  
492 E5, E3, and E7. For both damage scales, the distributions were skewed, with a larger number of  
493 predictions being underestimated (positive x-values), as certainly a consequence of the choice of  
494 machine learning models, their implementation (including imbalance issues), the distribution of input  
495 and target features considered, or all. The interest of machine learning model is also to have a relevant  
496 representation of the errors and limits of these methods.

497

#### 498 **4.3.2. Aggregate-single scenario**

499 Secondly, several aggregated building damage portfolio scenarios were considered to predict a single  
500 earthquake, thus testing whether the prediction was improved by increasing the number of post-  
501 earthquake damage observations. Three scenarios were tested. They are represented in Fig. 9 applying  
502 the EMS-98 damage grade (9a) and the TLS (9b):

503

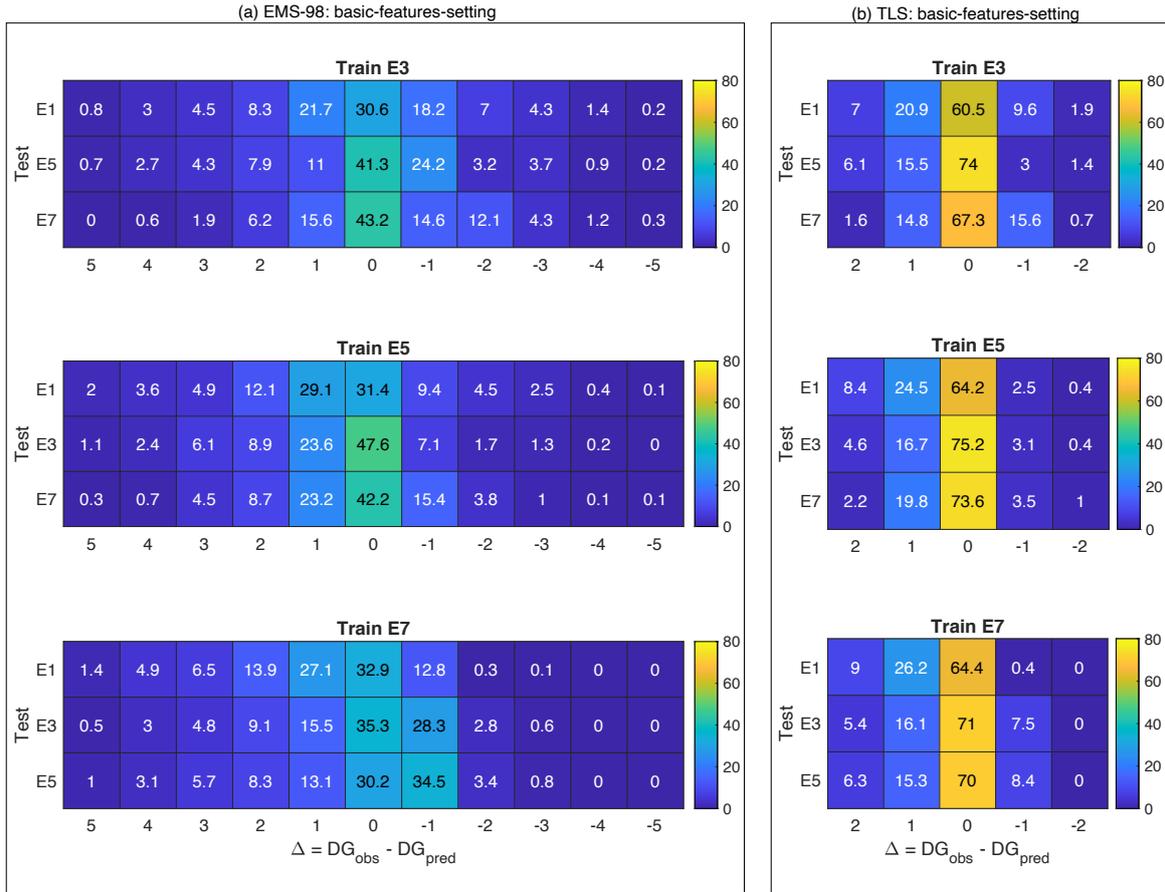
504 (i) Training set: E2+E3+E4+E6 (shown as E2346) – test set: E1, E5 and E7.

505 (ii) Training set: E2+E4+E5+E6 (shown as E2456) – test set: E1, E3 and E7.

506 (iii) Training set: E2+E4+E6+E7 (shown as E2467) – test set: E1, E3 and E5.

507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527

For the EMS-98 damage scale, correct classification ( $x$ -value centred on 0) in the range of 27% to 49% was found, depending on the training/test datasets. As in Fig. 8, using the E1 (Irpinia-1980) earthquake for testing scored lower regardless of the portfolio used for training (28.7%, 27.2% and 27.4% prediction accuracy). With E1 as the test set, the predictions at  $\pm 1$  DG (i.e., the sum of the  $x$ -values on Fig. 9a between -1 and +1) were 65.7%, 63.8% and 62.4% considering the E2346, E2456 and E2467 portfolios as the training set, respectively, for an average of 64% (compared with the 70% score for the single portfolio scenario, Fig. 8a). Other scenarios were also tested by aggregating the building damage portfolios differently (not presented herein), leading to the two main conclusions: (1) the quality and homogeneity of the input data (i.e., building features) affect the efficacy of the heuristic model and (2) this efficacy is limited and not improved by increasing the number of building damage observations, with a score (excluding E1) between 40% and 49% ( $x$ -value centred on 0), and up to 78% (average of the two scenarios, Fig. 8a and Fig. 9a) at  $\pm 1$  DG. Considering the TLS damage scale (Fig. 9b), a damage prediction efficacy of about 72% was obtained (compared with 72% in Fig. 8b), i.e., but no significant improvement was observed when the number of damaged buildings in the training portfolio was increased. For EMS-98 and TLS, the distributions were skewed, with a larger number of predictions being underestimated (positive  $x$ -values). Finally, in conclusion, the heuristic damage assessment model based on the XGBC model gives a better score for TLS damage assessment than for the EMS-98 damage scale. The TLS system also allows for quick assessment of damage on a large scale such as a city or region from an operational point of view.



528

529

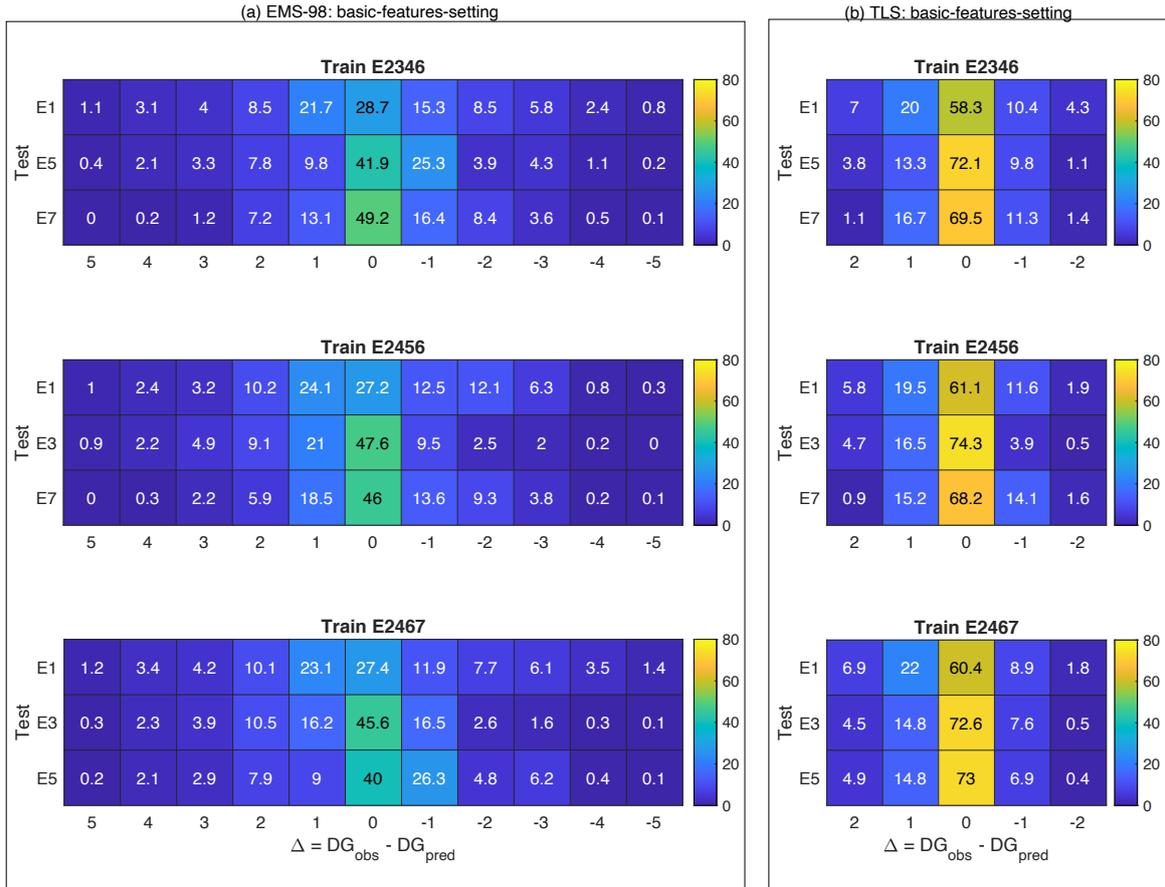
**Figure 8.** Distribution of the classification value ( $1 - \varepsilon_d$  in % given by Eq. 1) for (a) EMS-98- and (b) TLS-based damage classification using XGBC machine learning models and considering a single damage portfolio to predict a single portfolio (single-single scenario). The colour bar indicates the associated value in each cell. The x-values are the difference between the DG observed and the DG predicted, regardless of the DG considered.

530

531

532

533



534

535 **Figure 9.** Distribution of the classification value ( $1 - \varepsilon_d$  in % given by Eq. 1) for (a) EMS-98- and (b) TLS-based  
 536 damage classification using XGBC machine learning models and considering an aggregate damage portfolio to  
 537 predict a single portfolio (aggregate-single scenario). The colour bar indicates the associated value in each cell.  
 538 The x-values are the difference between the DG observed and the DG predicted, regardless of the DG considered.

539

#### 540 4.3.3 Comparing efficacy with the Risk-UE model

541 The efficacy of the heuristic damage assessment model was then compared with conventional damage  
 542 prediction methods, i.e., RISK-UE and mean damage relationship (Eq. 2 to 7), considering the basic-  
 543 features-settings. For RISK-UE, mean damage  $\mu_d$  (Eq. 4) was computed using the training set and the  
 544 vulnerability index  $IV$  for each building (Eq. 5). A vulnerability index was then attributed to all the  
 545 buildings in each class defined according to building features. The vulnerability indexes were then  
 546 attributed to every building in the test set, mean damage ( $\mu_d$ ) was computed with Eq. 2 and then DG  
 547 distribution with Eq. 3, before being compared with the damage portfolio used for testing. Finally, the  
 548 distribution of the mean damage observed (Eq. 4) was compared with the distribution of damage directly  
 549 on the test set, using Eq. 3.

550 Fig. 10 shows the distribution of absolute errors associated with the RISK-UE, mean damage  
 551 relationship, and XGBC methods (with and without compensation for the class-imbalance issue) trained

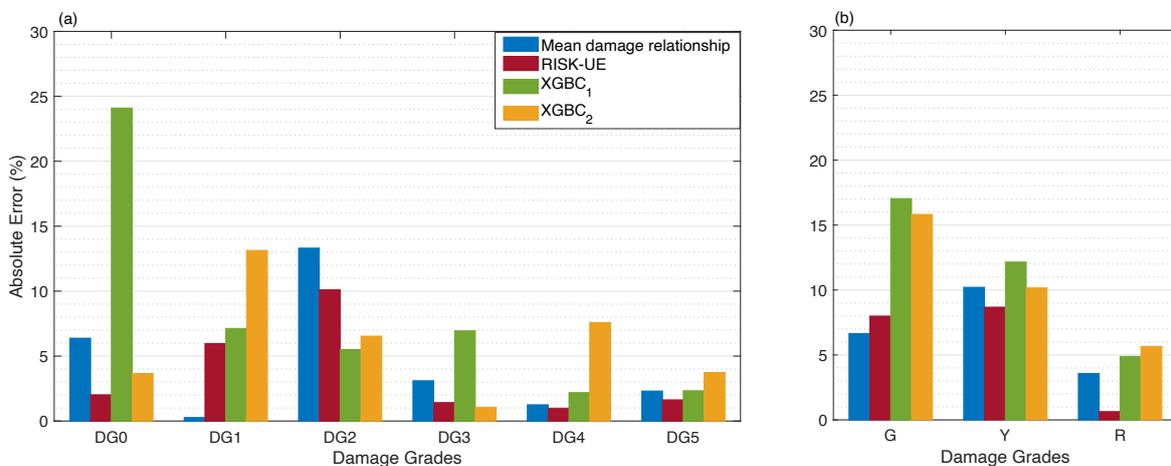
552 on earthquake building damage portfolio E5 and tested on E3. For EMS-98 damage classification (Fig.  
 553 10a), the XGBC model (without compensation for class-imbalance issues) resulted in a level of absolute  
 554 errors similar to that of the RISK-UE and/or mean damage relationship, except for DG0 (24%). Random  
 555 oversampling to compensate for the class-imbalance issues improved the distribution of errors for the  
 556 XGBC model (errors less than 8%, except for DG1: 13%).

557 For TLS-based damage classification, the XGBC model also resulted in a similar level of errors  
 558 compared with the mean damage relationship and/or RISK-UE methods (Fig. 10b), except for the green  
 559 class (no or slight damage, 17.04%). Compensation for class-imbalance issues slightly improved the  
 560 distribution of errors for the XGBC model with a 2% drop in errors for green (no/slight damage) and  
 561 yellow (moderate damage) classes.

562 Figure 11 shows the distribution of absolute errors trained using the E2456 portfolio and tested on the  
 563 E3 portfolio. For EMS-98 damage classification (Fig. 11a), the XGBC model (without compensation  
 564 for class-imbalance issues) resulted in a level of errors similar to that of the RISK-UE and/or mean  
 565 damage relationship; errors were highest for DG0 with 15.15%. With compensation for the class-  
 566 imbalance issues, the XGBC model achieved a slightly lower error distribution for DG0 (5%) and DG3  
 567 (4%); however, for other damage grades, the error value increased significantly (DG1: 11%, DG2: 12%  
 568 DG4: 7%, DG5: 2%). For TLS-based damage classification, the distribution of absolute errors was  
 569 similar for both the XGBC model and the mean damage relationship and/or RISK-UE methods (Fig.  
 570 11b). The highest absolute error value was associated with the green (no or slight damage) class of  
 571 buildings (16.40%). Compensation for the class-imbalance issues slightly increased the error  
 572 distribution for the XGBC model with nearly 5% for buildings in the green (no or slight) and red (heavy)  
 573 classes.

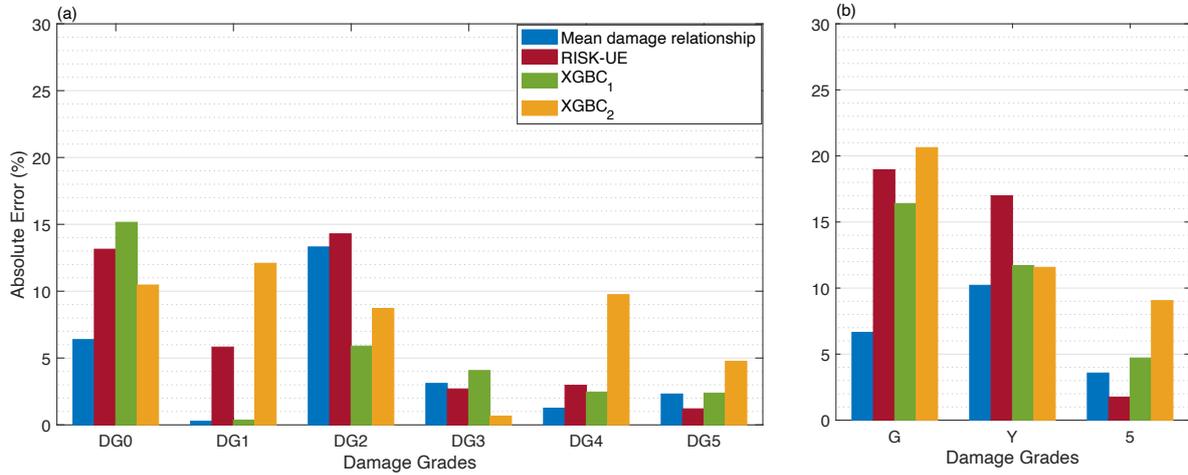
574 These results show that the heuristic building damage model based on the XGBC model, trained using  
 575 building damage portfolios with the basic-features-setting, provides a reasonable estimation of potential  
 576 damage, particularly with TLS-based damage classification.

577



578

579 **Figure 10.** Comparison of the efficacy of the heuristic model with the conventional model considering the DaDO  
 580 portfolio (training set: E5; test set: E3) for (a) EMS-98- and (b) TLS-based damage classification. The x-axis is  
 581 the damage grade and the y-axis is the percentage of absolute error ( $\epsilon_k$  in % given by Eq. 7). The blue bar  
 582 corresponds to the mean damage relationship, the red bar corresponds to the RISK-UE method, the green and  
 583 orange bars correspond to the heuristic model without (XGBC<sub>1</sub>) and with (XGBC<sub>2</sub>) compensation for the class-  
 584 imbalance issues, respectively.



585  
 586 **Figure 11.** Comparison of the efficacy of the heuristic model with the conventional model considering the DaDO  
 587 portfolio (training set: E2456; test set: E3) for (a) EMS-98- and (b) TLS-based damage classification. The x-axis  
 588 is the damage grade and the y-axis is the percentage of absolute error ( $\epsilon_k$  in % given by Eq. 7). The blue bar  
 589 corresponds to the mean damage relationship, the red bar corresponds to the RISK-UE method, the green and  
 590 orange bars correspond to the heuristic model without (XGBC<sub>1</sub>) and with (XGBC<sub>2</sub>) compensation for the class-  
 591 imbalance issues, respectively.

592  
 593 **5. Discussion**

594 Previous studies have aimed to test a machine learning framework for seismic building damage  
 595 assessment (e.g., Mangalathu et al., 2020; Roeslin et al., 2020; Harirchan et al., 2021; Ghimire et al.,  
 596 2022). They evaluated various machine learning and data balancing methods to classify earthquake  
 597 damage to buildings. However, these studies (Mangalathu et al., 2020, Roeslin et al., 2020, Harirchan  
 598 et al., 2021) had limitations such as limited data samples, damage classes, and building characteristics  
 599 limited to a spatial coverage and range of seismic demand values. Ghimire et al. (2022) also used a  
 600 larger building damage database, but did not investigate the importance of input features as a function  
 601 of damage levels and did not compare machine learning with conventional damage assessment methods.  
 602 This study aims to go beyond previous studies by testing advanced machine learning methods and data  
 603 resampling techniques using the unique DaDO dataset collected from several major earthquakes in Italy.  
 604 This database covers a wide range of seismic damage and seismic demands of a specific region,  
 605 including undamaged buildings. Most importantly, this study highlights the importance of input features  
 606 according to the degrees of damage and finally compares the machine learning models with a classical  
 607 damage prediction model (Risk-UE). The machine learning models achieved comparable accuracy to

608 the Risk-UE method. In addition, TLS-based damage classification, using red for heavily damaged,  
609 yellow for moderate damage, and green for no to slight damage, could be appropriate when the  
610 information for undamaged buildings is unavailable during model training.

611 Indeed, it is worth noting that the importance of the input features used in the learning process changes  
612 with the degree of damage: this indicates that each feature may have a contribution to the damage that  
613 changes with the damage level. Thus, the weight of each feature does not depend linearly on the degree  
614 of damage, which is not considered in conventional vulnerability methods.

615 The prediction of seismic damage by machine learning remains until now tested on geographically  
616 limited data. The damage distribution is strongly influenced by region-specific factors such as  
617 construction quality and regional typologies, implementation of seismic regulations and hazard level.  
618 Therefore, machine learning-based models can only work well in regions with comparable  
619 characteristics and a host-to-target transfer of these models should be studied. In addition, the  
620 distribution of damage is often imbalanced, impacting the performance of machine learning models by  
621 assigning higher weights to the features of the majority class. However, data balancing methods like  
622 random oversampling can reduce bias caused by imbalanced data during the training phase, but they  
623 may also introduce overfitting issues depending on the distribution of input and target features. Thus,  
624 integrating data from a wider range of input features and earthquake damage from different regions,  
625 relying on a host-to-target strategy, could help achieve a more natural balance of data sets and lead to  
626 less biased results. Moreover, the machine learning methods only train on the data available in the  
627 learning phase, that reflects the building portfolio in the study area. The importance of the features  
628 contributing to the damage could thus be modulated, and would require a host-to-target adjustment for  
629 the application of the model to another urban zone/seismic region.”

630 However, the machine learning models trained and tested on the DaDO dataset resulted in similar  
631 damage prediction accuracy values reported in existing literature using different models and datasets  
632 with different combinations of input features. This might suggest that the uncertainty related to building  
633 vulnerability in damage classification may be smaller than the primary source of uncertainty related to  
634 the hazard component (such as ground motion, fault rupture, slip duration, etc.).

635

636 In recent years, there has been a proliferation of open building data, such as the OpenStreetMap-based  
637 dynamic global exposure model (Schorlemmer et al., 2020) and building damage dataset after an  
638 earthquake (such as DaDO). We must therefore continue this paradigm shift initiated by Riedel et al.  
639 (2014, 2015) which consisted in identifying the exposure data available and as certain as possible, and  
640 in finding the most effective relationships for estimating the damage, unlike conventional approaches  
641 which proposed established and robust methods but relying on data not available and therefore difficult  
642 to collect. The global dynamic exposure model will make it possible to meet the challenge of modelling  
643 exposure on a larger scale on available data, using a tool capable of integrating this large volume of  
644 data. Machine learning methods are one such rapidly growing tool that can aid in exposure classification

645 and damage prediction by leveraging readily available information. It is therefore necessary to continue  
646 in this direction in order to evaluate the performance of the methods and their pros and cons for  
647 maximum efficacy of the prediction of damage.

648 Future works will therefore have to address several key issues that have been discussed here but that  
649 need to be further investigated. For example, the weight of the input features varies according to the  
650 level of damage, but one can question the systematization of this observation whatever the dataset and  
651 the feature considered. The efficiency of the selected models and the management of imbalance data  
652 remain to be explored, in particular by verifying regional independence. Taking advantage of the  
653 increasing abundance of exposure data and post-seismic observations, the imbalanced feature  
654 distribution and observed damage levels could be solved by aggregating datasets independent of the  
655 exposure and hazard contexts of the regions, once the host-to-target transfer of the models has been  
656 resolved. Finally, key input features (still not yet identified) describing hazard or vulnerability may be  
657 unexplored, and incorporating them into the models may improve the accuracy of damage classification.

658

## 659 **6. Conclusion**

660 In this study, we explored the efficacy of machine learning models trained using DaDO post-earthquake  
661 building damage portfolios. We compared six machine learning models: RFC, GBC, XGBC, RFR,  
662 GBR, and XGBR. These models were trained on numbers of building features (location, number of  
663 storeys, age, floor area, height, position, construction material, regularity, roof type, ground slope  
664 condition) and ground motion intensity defined in terms of macro-seismic intensity. The classification  
665 models performed slightly better than the regression methods and the XGBC model was ultimately  
666 found to be the most efficient model for this dataset. To solve the imbalance issue concerning observed  
667 damage, the random oversampling method was applied to the training dataset to improve the efficacy  
668 of the heuristic damage assessment model by rectifying the skewed distribution of the target features  
669 (DGs).

670 Surprisingly, we found that the weight of the most important building feature evolves according to DG,  
671 i.e., the weight of the feature for damage prediction changes depending on the DG considered, which is  
672 not taken into account in conventional methods.

673 The basic-features-setting (i.e., considering number of storeys, age, floor area, height and macroseismic  
674 intensity, which are accurately evaluated for the existing building portfolio) gave the same accuracy  
675 (0.68) as the full-features-settings (0.72) with the TLS-based damage classification method. For training  
676 and testing, the homogeneity of the information in the portfolios is a key issue for the definition of a  
677 highly effective machine learning model, as shown by the data from the E1 earthquake (Irpinia-1990).  
678 However, the efficacy of the model reaches a limit which is not improved by increasing the number of  
679 damaged buildings in the portfolio used as the training set, for example. For damage prediction, this  
680 type of heuristic model results in approximately 75% correct classification. Other authors (e.g., Riedel

681 et al., 2014, 2015; Ghimire et al. 2022) have already reached this same conclusion by increasing the  
682 percentage of the training set compared with the test set.

683 Despite this limit threshold, the level of accuracy achieved remains similar to that attained by  
684 conventional methods, such as Risk-UE and the mean damage relationship, for the basic-features-  
685 settings and TLS-based damage classification (error values less than 17 %). Machine learning models  
686 trained on post-earthquake building damage portfolios could provide a reasonable estimation of damage  
687 for a different region with similar building portfolios, after host-to-target adjustment.

688 Some variability may have been introduced into the damage prediction model due to the framework  
689 defined to translate the original damage scale to the EMS-98 damage scale and because in the DaDO  
690 database, the year of construction and the floor area of each building are provided as interval values,  
691 and missing locations of buildings were replaced with the location of local administrative centres. The  
692 latter can lead to a smoothing of the macro-seismic intensities to be considered for each structure and  
693 also affect the distance to the earthquake. Similarly, the building damage surveys were carried out after  
694 the seismic sequence, which includes aftershocks as well as the mainshock, whereas the MSI input  
695 corresponds to the mainshock from the USGS ShakeMap. All these issues may reduce the efficacy of  
696 the heuristic model and its limit threshold. Addressing these issues could improve the damage prediction  
697 performance of machine learning models.

698

#### 699 **Code availability**

700 The machine learning models were developed using Scikit-learn documentation and the value of  
701 hyperparameters used are provided in table 3.

#### 702 **Data availability**

703 The data used in this study is available in the Database of Observed Damage (DaDO) web-GIS platform  
704 of the Civil Protection Department, developed by the Eucentre Foundation.

705 [https://egeos.eucentre.it/danno\\_osservato/web/danno\\_osservato?lang=EN](https://egeos.eucentre.it/danno_osservato/web/danno_osservato?lang=EN).

706

#### 707 **Author contribution**

708 Subash Ghimire: Conceptualization, methodology, data preparation, investigation, visualization, draft  
709 preparation. Philippe Guéguen: Conceptualization, investigation, visualization, supervision, review and  
710 editing. Adrien Pothon: Conceptualization, supervision, review and editing draft. Danijel Schorlemmer:  
711 Conceptualization, supervision, review and editing draft.

712

#### 713 **Competing interests**

714 The authors declare that they have no conflict of interest.

715

#### 716 **Acknowledgment**

717 The author(s) disclosed receipt of the following financial support for the research, authorship, and/ or  
718 publication of this article: This study was funded by the URBASIS-EU project (H2020-MSCA- ITN-  
719 2018, Grant No. 813137). A.P. and P.G. thank the AXA Research Fund supporting the project New  
720 Probabilistic Assessment of Seismic Hazard, Losses and Risks in Strong Seismic Prone Regions. P.G.  
721 thanks LabEx OSUG@2020 (Investissements d’avenir- ANR10LABX56)

## References

- 722 ATC: ATC-20-1, Field Manual: Postearthquake Safety Evaluation of Buildings Second Edition,,  
723 Applied Technology Council, Redwood City, California., 2005.
- 724 Azimi, M., Eslamlou, A. D., and Pekcan, G.: Data-driven structural health monitoring and damage  
725 detection through deep learning: State-ofthe- art review, <https://doi.org/10.3390/s20102778>, 2020.
- 726 Baggio, C., Bernardini, A., Colozza, R., Pinto, A. V, and Taucer, F.: Field Manual for post-earthquake  
727 damage and safety assessment and short term countermeasures (AeDES) Translation from Italian:  
728 Maria ROTA and Agostino GORETTI, 2007.
- 729 Bazzurro, P., Cornell, C. A., Menun, C., and Motahari, M.: GUIDELINES FOR SEISMIC  
730 ASSESSMENT OF DAMAGED BUILDINGS, in: 13th World Conference on Earthquake  
731 Engineering, Vancouver, B.C. m Canada, 74–76, <https://doi.org/10.5459/bnzsee.38.1.41-49>, 2004.
- 732 Branco, P., Ribeiro, R. P., Torgo, L., Krawczyk, B., and Moniz, N.: SMOGN: a Pre-processing  
733 Approach for Imbalanced Regression, *Proc. Mach. Learn. Res.*, 74, 36–50, 2017.
- 734 Breiman, L.: Random Forests, *Mach. Learn.*, 5–32, 2001.
- 735 Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System, in: 22nd acm sigkdd  
736 international conference on knowledge discovery and data mining, 785–794,  
737 <https://doi.org/10.1145/2939672.2939785>, 2016.
- 738 Daniell, J. E., Schaefer, A. M., Wenzel, F., and Tsang, H. H.: The global role of earthquake fatalities  
739 in decision-making: earthquakes versus other causes of fatalities, *Proc. Sixth. world Conf. Earthq. Eng.*  
740 *Santiago, Chile*, 9–13, 2017.
- 741 Dolce, M., Speranza, E., Giordano, F., Borzi, B., Bocchi, F., Conte, C., Meo, A. Di, Faravelli, M., and  
742 Pascale, V.: Observed damage database of past italian earthquakes: The da.D.O. WebGIS, *Boll. di*  
743 *Geofis. Teor. ed Appl.*, 60, 141–164, <https://doi.org/10.4430/bgta0254>, 2019.
- 744 Estabrooks, A. and Japkowicz, N.: A mixture-of-experts framework for learning from imbalanced  
745 data sets, *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes*  
746 *Bioinformatics)*, 2189, 34–43, [https://doi.org/10.1007/3-540-44816-0\\_4](https://doi.org/10.1007/3-540-44816-0_4), 2001.
- 747 FEMA: Hazus –MH 2.1 Multi-hazard Loss Estimation Methodology Earthquake, 2003.
- 748 Friedman, J. H.: Greedy Function Approximation:A Gradient Boosting Machine, 1999.
- 749 Del Gaudio, C., Scala, S. A., Ricci, P., and Verderame, G. M.: Evolution of the seismic vulnerability  
750 of masonry buildings based on the damage data from L’Aquila 2009 event, *Bull. Earthq. Eng.*, 19,  
751 4435–4470, 2021.
- 752 Ghimire, S., Gueguen, P. ;, and Schorlemmer, D.: Earthquake Damage Prediction of Buildings in  
753 Nepal using Machine Learning Tools, in: VIKAS, A Journal of Development, Nepal’s Post

- 754 Earthquake Recovery and Reconstruction, special issue, Volume 1, P. 124-131, 2021.
- 755 Ghimire, S., Guéguen, P., Giffard-Roisin, S., and Schorlemmer, D.: Testing machine learning models  
756 for seismic damage prediction at a regional scale using building-damage dataset compiled after the  
757 2015 Gorkha Nepal earthquake, *Earthq. Spectra*, <https://doi.org/10.1177/87552930221106495>, 2022.
- 758 Grünthal, G.: Escala Macro Sísmica Europea EMS - 98, 101 pp., 1998.
- 759 Guéguen, P., Michel, C., and Lecorre, L.: A simplified approach for vulnerability assessment in  
760 moderate-to-low seismic hazard regions: Application to Grenoble (France), *Bull. Earthq. Eng.*, 5,  
761 467–490, <https://doi.org/10.1007/s10518-007-9036-3>, 2007.
- 762 Guettiche, A., Guéguen, P., and Mimoune, M.: Seismic vulnerability assessment using association  
763 rule learning: application to the city of Constantine, Algeria, *Nat. Hazards*, 86, 1223–1245,  
764 <https://doi.org/10.1007/s11069-016-2739-5>, 2017.
- 765 Harirchian, E., Kumari, V., Jadhav, K., Rasulzade, S., Lahmer, T., and Das, R. R.: A synthesized  
766 study based on machine learning approaches for rapid classifying earthquake damage grades to rc  
767 buildings, *Appl. Sci.*, 11, <https://doi.org/10.3390/app11167540>, 2021.
- 768 Hegde, J. and Rokseth, B.: Applications of machine learning methods for engineering risk assessment  
769 – A review, *Saf. Sci.*, 122, 104492, <https://doi.org/10.1016/j.ssci.2019.09.015>, 2020.
- 770 Japkowicz, N. and Stephen, S.: The class imbalance problem A systematic study fulltext.pdf, 6, 429–  
771 449, 2002.
- 772 Kim, T., Song, J., and Kwon, O. S.: Pre- and post-earthquake regional loss assessment using deep  
773 learning, *Earthq. Eng. Struct. Dyn.*, 49, 657–678, <https://doi.org/10.1002/eqe.3258>, 2020.
- 774 Lagomarsino, S. and Giovinazzi, S.: Macro seismic and mechanical models for the vulnerability and  
775 damage assessment of current buildings, *Bull. Earthq. Eng.*, 4, 415–443,  
776 <https://doi.org/10.1007/s10518-006-9024-z>, 2006.
- 777 Lagomarsino, S., Cattari, S., and Ottonelli, D.: The heuristic vulnerability model: fragility curves for  
778 masonry buildings, Springer Netherlands, 3129–3163 pp., [https://doi.org/10.1007/s10518-021-01063-](https://doi.org/10.1007/s10518-021-01063-7)  
779 7, 2021.
- 780 Lundberg, S. M. and Lee, S.-I.: A Unified Approach to Interpreting Model Predictions, in: 31st  
781 Conference on Neural Information Processing Systems, 2017.
- 782 Mangalathu, S. and Jeon, J.-S.: Regional Seismic Risk Assessment of Infrastructure Systems through  
783 Machine Learning: Active Learning Approach, *J. Struct. Eng.*, 146, 04020269,  
784 [https://doi.org/10.1061/\(asce\)st.1943-541x.0002831](https://doi.org/10.1061/(asce)st.1943-541x.0002831), 2020.
- 785 Mangalathu, S., Sun, H., Nweke, C. C., Yi, Z., and Burton, H. V.: Classifying earthquake damage to  
786 buildings using machine learning, *Earthq. Spectra*, 36, 183–208,  
787 <https://doi.org/10.1177/8755293019878137>, 2020.
- 788 Milutinovic, Z. and Trendafiloski, G.: Risk-UE An advanced approach to earthquake risk scenarios  
789 with applications to different european towns, Rep. to WP4 vulnerability Curr. Build., 1–83, 2003.
- 790 Ministry of Housing and Urbanism of Chile, Terremoto y Tsunami 27F 2010:
- 791 Morfidis, K. and Kostinakis, K.: Approaches to the rapid seismic damage prediction of r/c buildings  
792 using artificial neural networks, *Eng. Struct.*, 165, 120–141,

- 793 <https://doi.org/10.1016/j.engstruct.2018.03.028>, 2018.
- 794 Mouroux, P. and Le Brun, B.: Presentation of RISK-UE Project, *Bull. Earthq. Eng.* 2006 44, 4, 323–  
795 339, <https://doi.org/10.1007/S10518-006-9020-3>, 2006.
- 796 Ministère des Travaux Publics, Transports et Communications: Evaluation des Bâtiments:  
797 [https://www.mtpqc.gouv.ht/accueil/recherche/article\\_7.html](https://www.mtpqc.gouv.ht/accueil/recherche/article_7.html).
- 798 NPA: Police Countermeasures and Damage Situation associated with 2011Tohoku district - off the  
799 Pacific Ocean Earthquake Total burn down Inundated below floor level Partially damaged Property  
800 damages Damaged roads Partial burn down March 10, 2021,  
801 [https://doi.org/https://www.npa.go.jp/news/other/earthquake2011/pdf/higaijokyo\\_e.pdf](https://doi.org/https://www.npa.go.jp/news/other/earthquake2011/pdf/higaijokyo_e.pdf) (last access:  
802 22 March 2021), 2021.
- 803 2015 Nepal Earthquake: Open Data Portal: /eq2015.npc.gov.np/#/, last access: 22 March 2021.
- 804 Pedregosa, F., Varoquaux, G., Buitinck, L., Louppe, G., Grisel, O., and Mueller, A.: Scikit-learn,  
805 *GetMobile Mob. Comput. Commun.*, 19, 29–33, <https://doi.org/10.1145/2786984.2786995>, 2011.
- 806 Riedel, I. and Guéguen, P.: Modeling of damage-related earthquake losses in a moderate seismic-  
807 prone country and cost–benefit evaluation of retrofit investments: application to France, *Nat. Hazards*,  
808 90, 639–662, <https://doi.org/10.1007/s11069-017-3061-6>, 2018.
- 809 Riedel, I., Guéguen, P., Dunand, F., and Cottaz, S.: Macroscale vulnerability assessment of cities  
810 using association rule learning, *Seismol. Res. Lett.*, 85, 295–305, <https://doi.org/10.1785/0220130148>,  
811 2014.
- 812 Riedel, I., Guéguen, P., Dalla Mura, M., Pathier, E., Leduc, T., and Chanussot, J.: Seismic  
813 vulnerability assessment of urban environments in moderate-to-low seismic hazard regions using  
814 association rule learning and support vector machine methods, *Nat. Hazards*, 76, 1111–1141,  
815 <https://doi.org/10.1007/s11069-014-1538-0>, 2015.
- 816 Roeslin, S., Ma, Q., Juárez-García, H., Gómez-Bernal, A., Wicker, J., and Wotherspoon, L.: A  
817 machine learning damage prediction model for the 2017 Puebla-Morelos, Mexico, earthquake, *Earthq.*  
818 *Spectra*, 36, 314–339, <https://doi.org/10.1177/8755293020936714>, 2020.
- 819 Salehi, H. and Burgueño, R.: Emerging artificial intelligence methods in structural engineering, *Eng.*  
820 *Struct.*, 171, 170–189, <https://doi.org/10.1016/j.engstruct.2018.05.084>, 2018.
- 821 Scala, S. A., Del Gaudio, C., and Verderame, G. M.: Influence of construction age on seismic  
822 vulnerability of masonry buildings damaged after 2009 L’Aquila earthquake, *Soil Dyn. Earthq. Eng.*,  
823 157, 107199, <https://doi.org/10.1016/J.SOILDYN.2022.107199>, 2022.
- 824 Seo, J., Dueñas-Osorio, L., Craig, J. I., and Goodno, B. J.: Metamodel-based regional vulnerability  
825 estimate of irregular steel moment-frame structures subjected to earthquake events, *Eng. Struct.*, 45,  
826 585–597, <https://doi.org/10.1016/j.engstruct.2012.07.003>, 2012.
- 827 Silva, V., Crowley, H., Pagani, M., Monelli, D., and Pinho, R.: Development of the OpenQuake  
828 engine, the Global Earthquake Model’s open-source software for seismic risk assessment, *Nat.*  
829 *Hazards*, 72, 1409–1427, <https://doi.org/10.1007/s11069-013-0618-x>, 2014.
- 830 Silva, V., Pagani, M., Schneider, J., and Henshaw, P.: Assessing Seismic Hazard and Risk Globally  
831 for an Earthquake Resilient World, *Contrib. Pap. to GAR 2019*, 24 p., 2019.
- 832 Silva, V., Brzev, S., Scawthorn, C., Yepes, C., Dabbeek, J., and Crowley, H.: A Building

- 833 Classification System for Multi-hazard Risk Assessment, *Int. J. Disaster Risk Sci.*, 13, 161–177,  
834 <https://doi.org/10.1007/s13753-022-00400-x>, 2022.
- 835 Stojadinović, Z., Kovačević, M., Marinković, D., and Stojadinović, B.: Rapid earthquake loss  
836 assessment based on machine learning and representative sampling, *Earthq. Spectra*, 38, 152–177,  
837 <https://doi.org/10.1177/87552930211042393>, 2021.
- 838 Sun, H., Burton, H. V., and Huang, H.: Machine learning applications for building structural design  
839 and performance assessment: State-of-the-art review, *J. Build. Eng.*, 33, 101816,  
840 <https://doi.org/10.1016/j.jobe.2020.101816>, 2021.
- 841 Wald, D. J., Worden, B. C., Quitoriano, V., and Pankow, K. L.: ShakeMap manual: technical manual,  
842 user's guide, and software guide, *Techniques and Methods*, <https://doi.org/10.3133/tm12A1>, 2005.
- 843 Wang, C., Yu, Q., Law, K. H., McKenna, F., Yu, S. X., Taciroglu, E., Zsarnóczy, A., Elhaddad, W.,  
844 and Cetiner, B.: Machine learning-based regional scale intelligent modeling of building information  
845 for natural hazard risk management, *Autom. Constr.*, 122,  
846 <https://doi.org/10.1016/j.autcon.2020.103474>, 2021.
- 847 Xie, Y., Ebad Sichani, M., Padgett, J. E., and DesRoches, R.: The promise of implementing machine  
848 learning in earthquake engineering: A state-of-the-art review, *Earthq. Spectra*, 36, 1769–1801,  
849 <https://doi.org/10.1177/8755293020919419>, 2020.
- 850 Xu, Y., Lu, X., Cetiner, B., and Taciroglu, E.: Real-time regional seismic damage assessment  
851 framework based on long short-term memory neural network, *Comput. Civ. Infrastruct. Eng.*, 1–18,  
852 <https://doi.org/10.1111/mice.12628>, 2020a.
- 853 Xu, Z., Wu, Y., Qi, M. zhu, Zheng, M., Xiong, C., and Lu, X.: Prediction of structural type for city-  
854 scale seismic damage simulation based on machine learning, *Appl. Sci.*, 10,  
855 <https://doi.org/10.3390/app10051795>, 2020b.
- 856