

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/320559210>

# Earthquake Prediction in California Using Regression Algorithms and Cloud-based Big Data Infrastructure

Article in Computers & Geosciences · November 2017

DOI: 10.1016/j.cageo.2017.10.011

---

CITATIONS

105

READS

5,659

4 authors:



Gualberto Asencio Cortés

University of Pablo de Olavide

78 PUBLICATIONS 1,306 CITATIONS

[SEE PROFILE](#)



Antonio Morales-Estebaran

University of Seville

99 PUBLICATIONS 1,832 CITATIONS

[SEE PROFILE](#)



Xueyi Shang

35 PUBLICATIONS 802 CITATIONS

[SEE PROFILE](#)



Francisco Martínez-Álvarez

Pablo de Olavide University

209 PUBLICATIONS 6,353 CITATIONS

[SEE PROFILE](#)

1      Earthquake Prediction in California Using Regression Algorithms and  
2      Cloud-based Big Data Infrastructure

3      G. Asencio-Cortés<sup>a</sup>, A. Morales-Esteban<sup>b</sup>, X. Shang<sup>c</sup>, F. Martínez-Álvarez<sup>a</sup>

4      <sup>a</sup>*Division of Computer Science, Pablo de Olavide University of Seville, Spain – {guaaescor,fmaralv}@upo.es*

5      <sup>b</sup>*Department of Building Structures and Geotechnical Engineering, University of Seville, Spain – ame@us.es*

6      <sup>c</sup>*School of Resources and Safety Engineering, Central South University, China – shangxueyi@csu.edu.cn*

---

7      **Abstract**

Earthquake magnitude prediction is a challenging problem that has been widely studied during the last decades. Statistical, geophysical and machine learning approaches can be found in literature, with no particularly satisfactory results. In recent years, powerful computational techniques to analyze big data have emerged, making possible the analysis of massive datasets. These new methods make use of physical resources like cloud based architectures. California is known for being one of the regions with highest seismic activity in the world and many data are available. In this work, the use of several regression algorithms combined with ensemble learning is explored in the context of big data (1 GB catalog is used), in order to predict earthquakes magnitude within the next seven days. Apache Spark framework, *H<sub>2</sub>O* library in R language and Amazon cloud infrastructure were been used, reporting very promising results.

8      **Key words:** Earthquake prediction, big data analytics, cluster computing, regression, ensemble learning

---

9      **1. Introduction**

10     Modern societies are threatened by natural risks and demand a proper preparation to reduce their  
11     impact. During the last years many initiatives have merged from modern societies in order to minimize the  
12     economical and human impact of natural disasters.

13     Natural risk is a concept embedded in the collective consciousness of modern societies. Against expected,  
14     an objective and universal definition of risk is yet to be established [1]. Nevertheless, it can be defined as a  
15     measure of the combined likelihood of occurrence of a threatening event and its potential consequences.

16     Natural disasters occur when a probable hazard turns into a real event. Then, potential consequences  
17     become real human and economic losses. Among natural disasters such as earthquakes, tsunamis, volcanic  
18     eruptions, hurricanes, tornadoes, floods and others, earthquakes stand out due to their devastating effects  
19     [2]. Earthquakes arrive suddenly and can destroy a whole city or region within seconds causing lost of  
20     lives or injures, property damage, social and economic breaks or environmental damage [3]. Moreover, many  
21     populated areas stand on seismic zones. Besides, earthquakes can produce correlated effects such as tsunamis  
22     [4], landslides [5] and liquefaction [6].

23     Seismic risk is a combination of seismic hazard and seismic vulnerability [7]. On the one hand, seismic  
24     hazard represents a potentially damaging seismic event that can cause damage. On the other hand, the  
25     potential consequences are the existing vulnerabilities that show the susceptibility to the damaging effect of  
26     the hazard.

27     Big data analytics has emerged as a very powerful technique. It is typically used to examine huge  
28     datasets in order to extract useful information and discover patterns [8]. When such big datasets must  
29     be analyzed, computational resources increase and traditional machine learning algorithms require new  
30     parallelized implementations that must be launched in clusters [9].

31     For all the aforementioned, there is a worldwide trend to enhance our understanding of earthquakes  
32     in order to increase our ability to manage them [10]. In this paper, earthquake prediction in one of the  
33     most seismic and populated areas of the world -California- is explored. So far, standard machine learning

algorithms were been applied to earthquake prediction. However, studied datasets' sizes were typically no bigger than several MB [11]. For this purpose, a 1 GB catalog was been generated, retrieving events from 1970 to 2017. Four regressors (linear models, gradient boosting machines, deep learning and random forests) and ensembles of them were been applied for predicting the maximum magnitude in the coming seven days. Due to the high computational resources required, the use of big data technologies and infrastructures was necessary. Spark distributed computing framework and  $H_2O$  library for cluster computing in R language were been used in this approach. Finally, Amazon cloud infrastructures were been also used.

The rest of the paper is structured as follows. Section 2 reports and discusses the relevant works related to big data and earthquake prediction. Section 3 details the proposed methodology to apply regression algorithms and using cloud-based big data infrastructure to predict earthquakes in California. Reported results along with illustrative comments can be found in Section 4. Finally, the conclusions drawn from this work are summarized in Section 5.

## 2. Related works

Earthquake generation is an extraordinarily stochastic process. Determining the time, the location and the magnitude of the next coming earthquake is an extremely difficult task. Moreover, considering that no direct measure of the accumulated stress and the strength of the material can be currently made. This search has made researchers to develop many different models [12].

These models can be classified in two groups [13]: the probabilistic methods (based on analysing the seismicity distribution) and the methods based on artificial intelligence (based on learning from the time series data).

Among all the initiatives done, the Regional Earthquake Likelihood Model (RELM) generated up to 18 different models. Based on that the probability of the occurrence of an earthquake follows a Poisson distribution, Petersen et al. [14] created a time-independent model. A 24-hours forecast method was proposed by [15]. This probability model uses foreshock/aftershock statistics. An intermediate to long time probabilistic forecast model was developed by [16]. Simple methods for determining the long-term average seismicity were created by [17]. The most fruitful author was Ward [18] who proposed five methods. The first one is based on smoothed seismicity. The second one considers GPS derived strain and the Kostrov's formula. The third one uses geological fault slip-data. The next one is an average of the previous methods and the last one uses an earthquake simulator. Two different models for a 24-hours forecast were proposed by [19]. These considered that an earthquake can be triggered by earlier shakes or can trigger later events. A five-year smoothed-seismicity model for  $M \geq 5.0$  for southern California was developed by [20]. Similarly, a smoothed-seismicity model was presented by [21] and it is based on small earthquakes for mapping large earthquakes. Ebel et al. [22] generated a 5-year forecast for  $M \geq 5.0$ . Moreover, they also presented two one-day forecast methods for earthquakes larger than or equal to 4.0. Finally, Rhoades [23] introduced a method for long-range forecasting, based on preceding minor earthquakes for forecasting large events.

Despite some successful predictions such as the Xiuyen prediction [24] and that at the Sanriku area in Japan on November 2001 [25], failures are dominant. The most significant failure is probably the one at Parkfield [26], due to the great effort and financial resources employed.

Recently, some promising models based on data mining have been proposed. In [27] the authors used clustering techniques to model seismic temporal data. This research was based on the previous work by [28]. Later, the M5' algorithm was used for relating the b-value and the occurrence of large earthquakes [29]. In [30] artificial neural networks were used for predicting earthquakes in Chile (a survey on artificial neural networks application to earthquake prediction can be found in [31]). A model for earthquake prediction on the seismogenic areas of the Iberian Peninsula was presented in [32]. The authors in [33] determined the best set of seismicity indicators to predict earthquakes. This work analyzed four zones of Chile (the most seismic country in the world) and it is based on the previous works of [34] and [30]. In [35] the authors produced a method to test the validity of the seismicity indicators used, closing the existing gap between geologists and data mining experts.

**Random forests algorithm (RF) has been extensively used in literature. In [36], RF was used to predict the remaining time before the next failure derived from earthquakes. In that**

84 work, RF identifies two classes of signals and uses them to predict failure: shear stress and  
85 dynamic strain encompassing two failure events, and a zoom of dynamic strain when failure is  
86 in the distant future. The work proposed in [37] addressed the prediction of large earthquakes  
87 (higher than or equal to magnitude 5.5) in the Hindukush region of Pakistan. Authors used  
88 artificial neural networks (ANN), recurrent neural networks, RF and a linear combination  
89 of tree classifiers named LPBoost, which maximizes a margin between training instances of  
90 different classes (binary classes).

91 A large set of classifiers were put in comparison in [38] for earthquake prediction in Italy.  
92 Specifically, Logit Boost, Bagging, Naive Bayes (NB), Bayes Net, Logistic regression, SV-Cm  
93 (a deep learning algorithm based on a supervised contractive map), MLP-Bp, C4.5, RF, KNN  
94 and Linear Regression. Despite the vast number of algorithms, the classification accuracy was  
95 between 30% and 40% for earthquakes with magnitudes larger than 3.0. In [39], five classifiers  
96 (SVM, ANN, KNN, C4.5 and NB) were used to analyze the predictability of earthquake  
97 datasets previously grouped by clustering. Such work proposed two studies: the first one  
98 analyzes the ability of the different groups to train general prediction models, the second  
99 analyzes the diversity and representativeness covered by the samples of the clusters of each  
100 group.

101 Classification of large earthquakes via ensemble learning was addressed in [40] up to  
102 magnitudes larger than 7.0 for Chile datasets. Different unbalanced techniques were analyzed,  
103 including undersampling and oversampling as preprocessing, along with ensembles as boosting  
104 and bagging. Finally, in [41], generalized linear models (GLM) were used in combination with a  
105 model fitting based on the AIC measure for predicting the probability of near-fault earthquake  
106 ground motion pulses.

107 The study of the state of the art reveals that machine learning algorithms have been recently used.  
108 However, explored datasets included limited information due to space and computational limitations that  
109 since machines exhibit. Therefore, the exploration of big data analytics in this context is justified and it is  
110 expected to serve as seed for future research works.

### 111 3. Methodology

112 The methodology used to carry out the proposed regression study is described in this section. It is a  
113 methodology through which a large amount of earthquake events are retrieved, divided and used to train  
114 and test a set of machine learning algorithms in a comparative way. The whole procedure is sustained by a  
115 big data infrastructure placed in a public cloud.

116 In Section 3.1 the entire procedure is summarized. The following five sections (3.2 to 3.6) explain in  
117 detail every phase of the methodology. Finally, Section 3.7 describes the underlying IT infrastructure used,  
118 which is based on public cloud and big data technologies.

#### 119 3.1. Overall methodology

120 The methodology carried out in this work is shown in a schematic way in Figure 1. First, a large catalog  
121 of earthquake events is retrieved from a public place. The purpose of selecting a large catalog was to prove  
122 the ability of the Big Data infrastructure to address large amount of events and process them.

123 The catalog acquired corresponds to the state of California from 1970 to 2017. A grid of latitudes and  
124 longitudes was established covering the state of California. Such grid produces a cell matrix where each cell  
125 has a size of  $0.5 \times 0.5$  (latitude  $\times$  longitude).

126 Every cell of the grid has a number of events ( $L_{i,j}$ ) and a magnitude range ( $\min_{i,j} - \max_{i,j}$ ). In order to  
127 select a subset of cells to perform a comparative regression study, two thresholds,  $\mu_1$  and  $\mu_2$ , were applied  
128 to filter cells. Specifically, cells with  $L_{i,j} \geq \mu_1$  and  $\max_{i,j} \geq \mu_2$  were selected producing 27 cells used to  
129 feed the study.

130 A set of 16 seismic features are generated from each selected cell resulting on a set of 27 regression  
131 datasets. The target prediction is the maximum magnitude in the next seven days. Every regression

132 dataset, which is sorted ascending by time, is divided in two parts: the first 75% for training models and  
133 the rest (25%) for testing purposes. Thus, 27 pairs training/test were produced.

134 A set of four machine learning-based regressors were used to carry out the regression study: generalized  
135 linear models (GLM), gradient boosting machines (GBM), deep learning (DL) and random forests (RF). The  
136 best performance was achieved by RF and, for such reason, four ensembles based on the stacking technique  
137 were built using RF as base learner: RF-GLM, RF-GBM, RF-DL and ALL (RF-GLM-GBM-DL).

138 Every regressor is trained using each training dataset producing a regression model. Models were then  
139 applied to each testing dataset resulting on earthquake predictions. All predictions are compared with  
140 their corresponding actual values producing a set of deviations. Such deviations are evaluated resulting on  
141 absolute and relative errors. Those errors are shown and discussed in Section 4.

### 142 *3.2. Data acquisition and preparation*

143 Earthquake data was acquired from the FTP site of the ANSS Composite Earthquake Catalog  
144 (<ftp://www.ncedc.org/pub/catalogs/anss>), through the Northern California Earthquake Data Center  
145 (NCEDC) [42] (last accessed on Apr 15, 2017).

146 Table 1 shows the characteristics of the catalog of events downloaded from the FTP site. The catalog  
147 has a size of 917.7 MB of decompressed text files. The catalog contains a file for each month in CNSS  
148 format. A time period from Jan, 1970 to Apr, 2017 was considered, resulting on a set of 568 files containing  
149 earthquake events for each month in the considered period.

150 The variables used from the catalog entries were the latitude, the longitude and the magnitude of the  
151 events. The catalog was filtered according to a minimum magnitude  $M_0 = 2.5$ . Thus, only events with at  
152 least such magnitude will be considered from this point to the rest of the work. This filtering resulted in  
153 63,960 events with magnitude greater than or equal to  $M_0$ .

154 Figure 2 shows the location of the considered events that occur more frequently. Events with the highest  
155 frequency are colored in red. It can be noticed that there are up to seven zones of high seismic activity.  
156 Three of them are particularly recurrent: the Joseph D. Grant County Park (near to the city of San Jose),  
157 the Sierra National Forest (near to the city of Bishop) and the San Bernardino National Forest (near to the  
158 city of San Bernardino).

159 According to the magnitude of the events, Figure 3 shows a map with the location of the earthquakes  
160 and a color scale from 2.5 (yellow) to 7.3 (red) degrees in the Richter scale. More significative areas in terms  
161 of high magnitude are the same highlighted in Figure 2 plus the area of the Humboldt Redwoods State Park,  
162 at the north of the state of California.

### 163 *3.3. Catalog grid and filtering*

164 From the catalog of earthquakes described in the previous subsection, a grid of cells defined by latitudes  
165 and longitudes was built. Table 2 shows the grid configuration used in this work. Specifically, cells of the  
166 grid are squared and they have a fixed size of 0.5 degrees (grid granularity).

167 The grid starts in the coordinates computed by the floor of the minimum latitude and longitude from  
168 the entire catalog of events. Specifically, the minimum latitude and longitude of catalog events are 32.55  
169 and -124.37, respectively. Therefore, the grid starts at coordinates (32, -125), or in GPS coordinates (32 N,  
170 125 W). From this point, cells are counted each 0.5 degrees in both latitude and longitude directions until  
171 the maximum point of the catalog (41.99, -114.56), forming the grid of study.

172 Note that the events were previously filtered by  $M_0$  and the geographic distribution of events is not  
173 uniform. For such reason, void cells (without events) can appear and, therefore, they were removed from  
174 the grid resulting on a set of 177 cells.

175 In order to analyze the most significant earthquakes in the presented regression study, only cells which  
176 contain at least 500 events and one or more events with more than 5 degrees of magnitude were considered  
177 (these filtering specifications are summarized in Table 2). Finally, as result of the indicated filters, 27 cells  
178 are considered for the regression study. Those cells contain from 538 to 5575 events inside.

179 A map of locations for the 27 cells of study is provided in Figure 4. These locations are placed in the  
180 previously described highest seismicity areas in the state of California from 1970 to 2017. Specifically, these  
181 locations are indicated in detail in Table 3.

182 The dataset names of the selected cells are assigned according to their corresponding cell coordinates.  
 183 For example, dataset named 2-19 contains the events in the cell (2,19) of the grid. Table 3 shows for each  
 184 dataset both the latitude and the longitude ranges (minimum and maximum) and its centroid of points  
 185 (specified in columns Lat.Cen and Lon.Cen).

186 The magnitude distribution of events in the considered datasets is summarized in Table 4. Columns size,  
 187 Q1, median, mean, Q3 and max show the number of events, the first quartile of magnitudes, median, mean,  
 188 third quartile and the maximum magnitude of each dataset, respectively. As it can be noticed, datasets 6-18  
 189 and 17-2 have the highest values of earthquake magnitudes.

#### 190 3.4. Feature generation

191 For each selected cell from the grid a propositional dataset was built containing a set of features to be  
 192 served as input to further regression models. Along with these features an outcome variable (continuous class)  
 193 is included in datasets indicating the maximum magnitude of events in the next week. In this subsection,  
 194 the seismic input features are described.

195 In Table 5 the set of seismic features are enumerated. The definitions of all used seismic indicators were  
 196 been taken from two previous works [34, 30]. Specifically, the features  $x_1, x_2, x_3, x_4, x_5, x_6$  and  $x_7$  were  
 197 firstly introduced in [30] and  $b, a, \eta, \Delta M, T, \mu, c, dE^{1/2}$  and  $M_{mean}$  were proposed in [34]. To assess the  
 198 features of a given event, the previous  $n$  events are calculated. In this work  $n$  was been set to 50 events, as  
 199 suggested in [43] and successfully used in [33, 32, 30]. All the attributes were been normalized between 0  
 200 and 1.

201 The seismic indicator  $b$  corresponds to the Gutenberg-Richter law's b-value [44]. The authors in [34]  
 202 used the least squares method for calculating the b-value. Due to the lack of robustness of this method  
 203 when large infrequent earthquakes happen, [30] used the maximum likelihood method which is described in  
 204 Equation 1:

$$b = \frac{\log e}{(1/n) \sum_{j=0}^{n-1} M_{i-j} - M_0} \quad (1)$$

205 The parameters involved in the Equation 1 are the number of events considered prior to  $e_i$ ,  $n$ ; the  
 206 magnitude of the event  $e_{i-j}$ ,  $M_{i-j}$ ; and the cutoff magnitude of the seismic zone,  $M_0$ . The rest of the  
 207 seismic features ( $x_1, x_2, x_3, x_4, x_5, x_6, x_7, a, \eta, \Delta M, T, \mu, c, dE^{1/2}$  and  $M_{mean}$ ) are defined as they  
 208 were proposed in [34, 30].

#### 209 3.5. Regression algorithms

210 In this section, the machine learning-based algorithms used for the regression in the presented study are  
 211 described. Five different approaches were considered: generalized linear models, gradient boosting machine,  
 212 deep learning, random forests and stacking ensembles.

213 Generalized Linear Models (GLM) provides a flexible generalization of the ordinary multiple linear  
 214 regression with error distribution models other than a Gaussian distribution. GLM unifies various other  
 215 statistical models, including Poisson, linear, logistic, and others when using L1 and L2 regularization. Due  
 216 to the problem nature of earthquakes magnitude prediction, response variable is continuous and, therefore,  
 217 a Gaussian distribution was used.

218 GLM belongs to the most commonly-used models for many types of data analysis use cases. Some  
 219 problems, specially linear ones, can be addressed successfully using GLM, but others may not be as accurate  
 220 if the variables are more complex. Namely, when the response variable has a non-linear distribution or  
 221 the effect of the input variables is not linear, GLM can produce results less accurate than other non-linear  
 222 models.

223 Gradient Boosting Machines (GBM) is a decision tree-based algorithm which is an ensemble method.  
 224 GBM makes iteratively more than one decision tree combining their outputs. Boosting is the ensemble  
 225 technique used in GBM to produce and select the different decision trees. Boosting-based techniques give  
 226 more importance to the harder-to-learn training data, it tends to reduce bias in its predictions. GBM  
 227 focuses its attention on the difficult instances in the training data, the ones that are hard to learn. That is

228 favorable, but it can also be risky. If there is one outlier that each tree keeps getting wrong it is going to  
229 get boosted and boosted achieving the maximum importance. If such outlier is a real unusual event then  
230 learning procedure is adequate, but if it is a measuring error it can distort the GBM accuracy.

231 GBM produces a prediction model in the form of an ensemble of weak prediction models, fitting  
232 consecutive trees where each solves for the net error of the prior trees. GBM builds the model in a stage-  
233 wise fashion and is generalized by allowing an arbitrary differentiable loss function. Results of new trees are  
234 applied partially to the entire solution. GBM often produces the best possible model but it is necessary to  
235 set proper stopping points to avoid overfitting, because GBM is sensitive to noise and extreme values.

236 Whereas underlying concept of linear models-based GLM is mathematics and decision trees-based GBM  
237 is logic, Deep Learning (DL) is a black box inspired by the human brain. DL models create high-level  
238 abstractions in data by using non-linear transformations in a layered-based iterative procedure. DL can  
239 address both supervised and unsupervised learning, which can use unlabeled data and it is widely used for  
240 pattern recognition in images or speech.

241 DL is based on a Deep Neural Network (DNN) which, in turn, is an Artificial Neural Network (ANN) with  
242 multiple hidden layers between input and output layers. DNNs can model complex non-linear relationships  
243 and generate compositional models where patterns in data are expressed as a layered composition of  
244 primitives. The extra layers enable composition of features from previous layers giving the potential of  
245 modeling complex data. DNNs are typically designed as feedforward networks, but other architectures like  
246 recurrent or convolutional neural networks were also applied.

247 DL implementation used in this work is based on a multi-layer feed-forward ANN that is trained with  
248 stochastic gradient descent using back-propagation. The network contains a large number of hidden layers  
249 consisting of neurons with tanh activation function. Despite its high computational cost, it scales well  
250 with big data, because each compute node in the cluster trains asynchronously a copy of the global model  
251 parameters on its local data with multi-threading and contributes periodically to the global model across  
252 the network.

253 Random Forests (RF) is an ensemble of decision trees based on the bagging technique using bootstrap  
254 aggregation. The idea is avoiding overfitting in complex data sets considering wide sets of trees and using of  
255 them to perform predictions on new data. RF combine multiple decision trees, each fit to a random sample  
256 of the original data. For classification the most frequent response is returned, for regression the mean of  
257 each tree response is used.

258 RF do not train every tree with all training data, instead of this, different random samples of rows and  
259 columns are given to each tree. RF is non-linear and it is robust for noisy data. It is able to reduce variance  
260 when predicting non-seen data with minimal increase in bias. Moreover, RF has the advantage of having  
261 few parametrization, mainly the number of trees used in the ensemble.

262 Apart of ensemble-based methods GBM and RF, the stacking technique of ensemble learning was  
263 considered in the present study. Both boosting and bagging of GBM and RF, respectively, are ensembles  
264 that take a collection of weak learners and forms a single strong learner. Stacking technique involves the  
265 training of a second-level metalearner to ensemble a group of base learners. The metalearner algorithm  
266 learns the optimal combination of the base learner fits. Unlike boosting and bagging, the goal in stacking is  
267 to ensemble strong and diverse sets of learners together.

268 Stacking builds the ensemble by training each of a set of  $B$  base algorithms on the training set. It then  
269 performs a k-fold cross-validation on each of these base algorithms and collect the cross-validated predicted  
270 values from each of the  $B$  algorithms. The  $M$  cross-validated predicted values from each of the  $B$  algorithms  
271 are combined to form a dataset called *level-one* with  $M$  instances and  $B$  predictors plus the original outcome  
272 variable. Then the metalearner is trained with such dataset. The ensemble model consists of the  $B$  base  
273 learning models and the metalearner model, which can then be used to produce predictions on test sets.

### 274 3.6. Validation and evaluation

275 In order to analyze and compare the performance of the regression algorithms on the different datasets,  
276 a hold-out scheme of validation was used. Specifically, each dataset was split in two parts: the first one  
277 includes the 75% of data and it is used to train the models, the second one includes the remaining 25% and  
278 it is used to test the models.

279 Since 27 datasets were splitted, there were 27 training sets and 27 test sets, each one of  
 280 them for each dataset. Regressors were trained with each training set separately producing 27  
 281 models, one for each training set. Each model was tested also separately with its corresponding  
 282 test set, producing its predictions.

283 Evaluation metrics computed for predictions made with different models are described as follows. Due  
 284 to the regression nature of the problem of earthquakes magnitude prediction, metrics computed are absolute  
 285 (MAE) and Relative Errors (RE). Equations 2 and 3 show the formulas of such types of error.

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (2)$$

$$RE = \frac{1}{n \times \max(y_i)} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (3)$$

286 In equations 2 and 3,  $n$  is the number of predicted test instances,  $y_i$  is the actual outcome value (the  
 287 magnitude of the maximum event in the next week) and  $\hat{y}$  is the predicted value for the outcome variable.  
 288 Due to non-significant events were removed ( $y_i \geq M_0, \forall i = 1..n$ ), some instances do not have any event in  
 289 the next week. In such cases, the outcome variable is zero. For this reason, standard relative errors like  
 290 RAE or MAPE cannot be used because they divide deviations between actual values. Therefore, the RE  
 291 is divided by the maximum magnitude in the dataset, which is enough significant information about the  
 292 proportion of the magnitude range that errors represents.

### 293 3.7. Big Data infrastructure

294 To perform the regression study, all phases of the methodology were implemented on a cloud-based Big  
 295 Data infrastructure. The use of Big Data technologies was necessary due to the high number of earthquake  
 296 events considered in the study. The infrastructure implemented in this work is drawn in Figure 5.

297 Amazon Web Services (AWS) were selected to provide a platform as a service in which all software  
 298 components needed are deployed in the proposed methodology. The first step is the data ingestion, that  
 299 consists in loading the whole catalog of earthquake events into the AWS cloud. Such procedure was carried  
 300 out importing ANSS Composite catalogs from FTP to Amazon S3 service directly (without passing by any  
 301 local filesystems) using the import tool of the S3 on-line control panel.

302 Catalog files were parsed using an eventual instance of Amazon EC2 to feed structured data to an Amazon  
 303 Redshift database. Fields parsed were the timestamp, the longitude, the latitude and the magnitude of the  
 304 earthquake events. An Amazon EMR 5.5.0 cluster was launched with Hadoop 2.7.3 and Apache Spark 2.1.0.  
 305 Once the cluster is running, R 3.4.0 and the  $H_2O$  3.10.4.6 library were installed using the bundle Sparkling  
 306 Water distribution for Apache Spark.

307 An R script was developed to build the grid of cells from the event catalog tables in Amazon Redshift.  
 308 Such grid was filtered according to the procedure described in Section 3.3 producing the set of 27 selected  
 309 cells of study. Those cells are stored in Amazon Redshift in a new table. Next, seismic features are generated  
 310 for each selected cell and resulting datasets were split in training and test. Later, they were stored in a next  
 311 table of Amazon Redshift. Such process was run by other R script inside the Amazon EMR cluster.

312 Once training and test datasets were built, the machine learning-based regressors (GLM, GBM, DL,  
 313 RF and Stacking ensembles) implemented in the  $H_2O$  library were executed to train models from training  
 314 datasets and to predict the earthquakes of the test splits. Finally, absolute and relative errors were computed  
 315 from an R script inside the Amazon EMR cluster. Error tables were stored both in Amazon S3 using LATEX  
 316 format and in Amazon RDS. Different error plots were produced and stored in a specific Amazon S3 output  
 317 bucket for plots.

## 318 4. Results and discussion

319 Prediction results of the proposed regression study are shown and discussed in this section. In order  
 320 to measure and compare the effectiveness achieved by the different machine learning-based algorithms on

321 the 27 datasets described in the methodology section, the evaluation metrics previously introduced in such  
322 section were computed.

323 Tables 6 and 7 show the absolute and the relative errors, respectively, produced by the different methods  
324 for each dataset of study. Lowest errors are shown in bold text for each dataset. RF achieved the best  
325 performance on average with a mean absolute error of 0.74 degrees in the Richter scale. Moreover, RF was  
326 the most effective for each dataset in the 81% of cases.

327 GLM and DL had the worst performance, producing up to 2.50 and 4.04 absolute deviations, respectively.  
328 To see their behaviors in detail, Figure 6 shows the dispersion of errors using a boxplot representation for  
329 each regressor. As it can be noticed, GLM and DL had the highest error dispersion.

330 GBM achieved competitive performance when predicting earthquakes magnitudes in regression, obtaining  
331 the highest accuracy in three of the datasets (3-18, 5-16 and 9-9). Moreover, its error dispersion is low and  
332 similar to that achieved by RF, as it can be seen in Figure 6.

333 All stacking ensemble combinations were carried out, but only those which includes RF were shown.  
334 All other ensembles (GLM-GBM, GLM-DL, GBM-DL and GLM-GBM-DL) performed worse and, to be  
335 concise, they were omitted from the study. All presented ensembles performed similarly, they achieved the  
336 same absolute error average (0.76) and showed very similar error dispersion.

337 In order to analyze in detail the sign of deviations produced by regressors in their predictions, Figure 7  
338 shows histograms of errors for each regressor. All RF-based algorithms (RF, GLM-RF, GBM-RF, DL-RF,  
339 ALL) show the same behavior in the sign of their errors. They produced higher quantity of positive errors  
340 than negative ones (positive bias), as it can be noticed in the histograms of Figure 7. Precisely, stacking  
341 versions of RF produced higher positive errors than RF and lower negative errors.

342 RF-based stacking ensembles did not overcome to the base RF algorithm. This could be due to a lack of  
343 error complementation among RF and other methods. Specifically, instances in the cross-validated procedure  
344 of stacking formation predicted inaccurately by RF were not better predicted by any other algorithm in the  
345 ensemble. For such reason, stacking ensembles performed very similar to the base RF in the test sets.

346 GBM shown less bias in its error signs. They are more equilibrated around zero error and it approaches  
347 to a Gaussian distribution of its errors, which is desirable for a regressor. This behavior was expected due  
348 to nature of boosting (GBM) versus bagging (RF) ensembles. Boosting methods tend to decrease bias while  
349 bagging techniques tend to reduce variance increasing its bias, as it is mentioned in Section 3.5.

350 GLM and, specially, DL produced eventually very high errors (error outliers), up to -7.86 in GLM and  
351 -16.21 in DL. GLM could be too simple for this problem, because future earthquake magnitudes depend on  
352 very complex non-linear relationships of input seismic features.

353 The case of DL is different, it is the model with the highest complexity of the comparative. Its low  
354 accuracy could be due to its high parametrization. It has the highest number of parameters among the  
355 analyzed regressors. It could be improved performing a previous parameter tuning using training sets. In  
356 this work, the default configuration of deep learning implementation in the *H2O* library was used.

357 A detailed analysis of predictions was performed according to the different magnitude values  
358 in order to see the performance of the different regressors on larger magnitudes, which are  
359 more complex to be accurately predicted. To carry out this study, all predictions performed  
360 for 25%-test splits of the 27 studied datasets were considered. Then, predictions were divided  
361 by regressor and classified by a set of intervals of the actual magnitude, ranging from the  
362 interval [0,3] to [7,8] (six intervals of size 1).

363 The complete set of predictions analyzed contains 90,920 samples and the maximum actual  
364 magnitude was 7.3. Tables 8 and 9 show mean absolute and mean squared errors, respectively;  
365 all averaged for each regressor and magnitude interval (mean and standard deviations were  
366 computed). The last two rows of these tables indicate the number (#) and percentage (%) of  
367 samples predicted, respectively, for each magnitude interval. Note that magnitudes lower than  
368 5 suppose the 92.74% of the samples, leaving remaining events with larger magnitudes under-  
369 represented and, therefore, specially difficult to be predicted by machine learning techniques.

370 As it can be seen in Tables 8 and 9, the range of magnitudes [3,7] was reasonably well  
371 predicted with mean absolute errors lower than 0.6 with RF. Specially, range [4,7] was the  
372 most accurate in which errors were lower than or equal to 0.26. Extreme intervals [0,3)

373 and [7,8] were predicted with higher errors (up to 2.03 of MAE with RF in the case of  
374 the last interval). Mean squared errors (Table 9) reveal that ensemble methods were the best  
375 performance regressors for extreme intervals [0,4) and [7,8]. Such result could suggest possible  
376 improvement for large-magnitude predictions (magnitudes larger than 7) using more complex  
377 ensembles.

378 Figure 8 shows the scatter plots of actual-vs-predicted points for each regressor. A blue  
379 line was included in the scatter plots indicating the mean predicted value for each actual  
380 magnitude. As it can be seen, both GLM and DL perform worst due to their predicted values  
381 were significantly lower than the actual ones for large magnitudes. By contrast, GBM, RF  
382 and ensembles were sensible to higher magnitudes showing high correlation between actual  
383 and predicted values (up to  $R^2 = 0.80$  with RF).

384 Table 10 shows execution times consumed for each process carried out in the proposed  
385 methodology. All regressors are parallelized in  $H_2O$  and they were executed in batch mode as  
386 Amazon EMR tasks. As expected, training process was the most consuming task (6 hours).  
387 DL and DL-RF were the slowest regressors in their training phase, consuming 2.4 hours  
388 together. The faster regressor was GBM consuming only 10 minutes. RF, the most effective  
389 regressor previously analyzed, was also one of the faster methods, taking only 18 minutes to  
390 train the models for all datasets.

## 391 5. Conclusions

392 Data from California earthquakes from 1970 to 2017 were been analyzed in this work. A total of 1 GB  
393 of information, divided into 27 datasets that identify cells of  $0.5^\circ \times 0.5^\circ$ , were been processed by means of  
394 cloud based infrastructure. In particular, distributed implementations from  $H_2O$  package of four popular  
395 regression methods were been used to predict earthquakes magnitude within the next seven days. As  
396 next step, stacking-based ensemble learning was been applied, reporting relative errors verging on 10% and  
397 absolute errors verging on 0.5. Methods based on trees performed better and these methods reached, in  
398 general terms, lower regression errors. In conclusion, the use of big data analytics in the field of earthquakes  
399 magnitude prediction opens a very promising research line that may help in simultaneously processing  
400 massive data with huge number of variables.

## 401 Acknowledgements

402 The authors would like to thank the Spanish Ministry of Economy and Competitiveness, Junta de  
403 Andalucía for the support under projects TIN2014-55894-C2-R and P12-TIC-1728, respectively.

## 404 References

- 405 [1] T. Aven. On how to define, understand and describe risk. *Reliability Engineering and System Safety*, 95(6):623–631, 2010.
- 406 [2] E. Florido, F. Martínez-Álvarez, A. Morales-Esteban, J. Reyes, and J. L. Aznarte. Detecting precursory patterns to  
enhance earthquake prediction in Chile. *Computers and Geosciences*, 76:112–120, 2015.
- 407 [3] A. Spicák and J. Vanek. Earthquake swarms reveal submarine magma unrest induced by distant mega-earthquakes:  
Andaman Sea region. *Journal of Asian Earth Sciences*, 116:155–163, 2016.
- 408 [4] C. Cecioni, G. Bellotti, A. Romano, A. Abdolali, P. Sammarco, and L. Franco. Tsunami Early Warning System based on  
Real-time Measurements of Hydro-acoustic Waves. *Procedia Engineering*, 70:311–320, 2014.
- 409 [5] D. K. Keefer. Landslides caused by earthquakes. *Bulletin of the Seismological Society of America*, 95(4):406–421, 1984.
- 410 [6] R. Verdugo and J. González. Liquefaction-induced ground damages during the 2010 Chile earthquake. *Soil Dynamics and  
Earthquake Engineering*, 79(B):280–295, 2015.
- 411 [7] L. F. Sá, A. Morales-Esteban, and P. Durand. A Seismic Risk Simulator for Iberia. *Bulletin of the Seismological Society  
of America*, 106(3):1198–1209, 2016.
- 412 [8] C. W. Tsai, C. F. Lai, H. C. Chao, and A. V. Vasilakos. Big data analytics: a survey. *Journal of Big Data*, 2:21:1–32,  
2015.
- 413 [9] J. C. Jackson, V. Vijayakumar, A. Quadir, and C. Bharathi. Survey on Programming Models and Environments for  
Cluster, Cloud, and Grid Computing that Defends Big Data. *Procedia Computer Science*, 50:517–523, 2015.

- 421 [10] X. Romão, E. Paupério, and N. Pereira. A framework for the simplified risk analysis of cultural heritage assets. *Journal*  
 422 *of Cultural Heritage*, 20:696–708, 2014.
- 423 [11] Q. Wang, D. D. Jackson, and Y. Y. Kagan. California earthquakes, 1800-2007: A unified catalog with moment magnitudes,  
 424 uncertainties, and focal mechanisms. *Seismological Research Letters*, 80(3):446–457, 2009.
- 425 [12] K. F. Tiampo and R. Shcherbakov. Seismicity-based earthquake forecasting techniques: Ten years of progress.  
 426 *Tectonophysics*, 522-523:89–121, 2012.
- 427 [13] G. Asencio-Cortés, F. Martínez-Álvarez, A. Morales-Esteban, and A. Troncoso. Medium-large earthquake magnitude  
 428 prediction in Tokyo with artificial neural networks. *Neural Computing and Applications*, 28(5):1043–1055, 2017.
- 429 [14] M. D. Petersen, T. Cao, K. W. Campbell, and A. D. Frankel. Time-independent and time-dependent seismic hazard  
 430 assessment for the state of California: Uniform California earthquake rupture forecast model 1.0. *Seismological Research*  
 431 *Letters*, 78(1):99–109, 2007.
- 432 [15] M. C. Gerstenberger, L. M. Jones, and S. Wiemer. Short-term aftershock probabilities: Case studies in California.  
 433 *Seismological Research Letters*, 78(1):66–77, 2007.
- 434 [16] Z. Z. Shen, D. D. Jackson, and Y. Y. Kagan. Implications of geodetic strain rate for future earthquakes, with a five-year  
 435 forecast of M5 earthquakes in southern California. *Seismological Research Letters*, 78(1):116–120, 2007.
- 436 [17] P. Bird and Z. Liu. Seismic hazard inferred from tectonics: California. *Seismological Research Letters*, 78(1):37–48, 2007.
- 437 [18] S. N. Ward. Methods for evaluating earthquake potential and likelihood in and around California. *Seismological Research*  
 438 *Letters*, 78(1):121–133, 2007.
- 439 [19] R. Console, M. Murru, F. Catalli, and G. Falcone. Real time forecasts through an earthquake clustering model constrained  
 440 by the rate-and-state constitutive law: comparison with a purely stochastic ETAS model. *Seismological Research Letters*,  
 441 78(1):49–56, 2007.
- 442 [20] Y. Y. Kagan, D. D. Jackson, and Y. Rong. A testable five-year forecast of moderate and large earthquakes in southern  
 443 California based on smoothed seismicity. *Seismological Research Letters*, 78(1):94–98, 2007.
- 444 [21] A. Helmstetter, Y. Y. Kagan, and D. D. Jackson. High-resolution time-independent grid-based forecast for M=5  
 445 earthquakes in California. *Seismological Research Letters*, 78(1):78–86, 2007.
- 446 [22] J. E. Ebel, D. W. Chambers, A. L. Kafka, and J. A. Baglivo. Non-poissonian earthquake clustering and the hidden Markov  
 447 model as bases for earthquake forecasting in California. *Seismological Research Letters*, 78(1):57–65, 2007.
- 448 [23] D. A. Rhoades. Application of the EEPAS model to forecasting earthquakes of moderate magnitude in southern California.  
 449 *Seismological Research Letters*, 78(1):110–115, 2007.
- 450 [24] S. D. Zhang. *The 1999 Xiuyen-Haichung, Liaoning, M5.4 earthquake*. Beijing Seismological Press, 2004.
- 451 [25] T. Matsuzawa, T. Igarashi, and A. A. Hasegawa. Characteristic small-earthquake sequence off Sanriku, northeastern  
 452 Honshu, Japan. *Geophysical Research Letters*, 29(11):1543–147, 2002.
- 453 [26] W. H. Bakun, B. Aagaard, B. Dost, W. L. Ellsworth, J. L. Hardebeck, R. A. Harris, C. Ji, M. J. S. Johnston, J. Langbein,  
 454 J. J. Lienkaemper, A. J. Michael, J. R. Murray, R. M. Nadeau, P. A. Reasenberg, M. S. Reichle, E. A. Roeloffs, A. Shakal,  
 455 R. W. Simpson, and F. Waldhauser. Implications for prediction and hazard assessment from the 2004 Parkfield earthquake.  
*Nature*, 437:969–974, 2005.
- 456 [27] A. Morales-Esteban, F. Martínez-Álvarez, A. Troncoso, J. L. de Justo, and C. Rubio-Escudero. Pattern recognition to  
 457 forecast seismic time series. *Expert Systems with Applications*, 37(12):8333–8342, 2010.
- 458 [28] P. Nuannin, O. Kulhanek, and L. Persson. Spatial and temporal b value anomalies preceding the devastating off coast of  
 459 nw sumatra earthquake of december 26, 2004. *Geophysical Research Letters*, 32, 2005.
- 460 [29] F. Martínez-Álvarez, A. Troncoso, A. Morales-Esteban, and J. C. Riquelme. Computational intelligence techniques for  
 461 predicting earthquakes. *Lecture Notes in Artificial Intelligence*, 6679(2):287–294, 2011.
- 462 [30] J. Reyes, A. Morales-Esteban, and F. Martínez-Álvarez. Neural networks to predict earthquakes in Chile. *Applied Soft  
 463 Computing*, 13(2):1314–1328, 2013.
- 464 [31] E. Florido, J. L. Aznarte, A. Morales-Esteban, and F. Martínez-Álvarez. Earthquake magnitude prediction based on  
 465 artificial neural networks: A survey. *Croatian Operational Research Review*, 7(2):687–700, 2016.
- 466 [32] A. Morales-Esteban, F. Martínez-Álvarez, and J. Reyes. Earthquake prediction in seismogenic areas of the Iberian  
 467 Peninsula based on computational intelligence. *Tectonophysics*, 593:121–134, 2013.
- 468 [33] F. Martínez-Álvarez, J. Reyes, A. Morales-Esteban, and C. Rubio-Escudero. Determining the best set of seismicity  
 469 indicators to predict earthquakes. Two case studies: Chile and the Iberian Peninsula. *Knowledge-Based Systems*, 50:198–  
 470 210, 2013.
- 471 [34] A. Panakkat and H. Adeli. Neural network models for earthquake magnitude prediction using multiple seismicity indicators.  
 472 *International Journal of Neural Systems*, 17(1):13–33, 2007.
- 473 [35] G. Asencio-Cortés, F. Martínez-Álvarez, A. Morales-Esteban, and J. Reyes. A sensitivity study of seismicity indicators in  
 474 supervised learning to improve earthquake prediction. *Knowledge-Based Systems*, 101:15–30, 2016.
- 475 [36] Bertrand Rouet-Leduc, Claudia Hulbert, Nicholas Lubbers, Kipton Barros, Colin Humphreys, and Paul A Johnson.  
 476 Machine learning predicts laboratory earthquakes. *arXiv preprint arXiv:1702.05774*, 2017.
- 477 [37] KM Asim, F Martínez-Álvarez, A Basit, and T Iqbal. Earthquake magnitude prediction in hindukush region using machine  
 478 learning techniques. *Natural Hazards*, 85(1):471–486, 2017.
- 479 [38] P. M. Buscema, G. Massini, and G. Maurelli. Artificial Adaptive Systems to predict the magnitude of earthquakes.  
 480 *Bollettino di Geofisica Teorica ed Applicata*, 56(2):227–256, 2015.
- 481 [39] Gualberto Asencio-Cortés, Sanja Scitovski, Rudolf Scitovski, and Francisco Martínez-Álvarez. Temporal analysis of  
 482 croatian seismogenic zones to improve earthquake magnitude prediction. *Earth Science Informatics*, pages 1–18, 2017.
- 483 [40] Manuel Jesús Fernández-Gómez, Gualberto Asencio-Cortés, Alicia Troncoso, and Francisco Martínez-Álvarez. Large  
 484 earthquake magnitude prediction in chile with imbalanced classifiers and ensemble learning. *Applied Sciences*, 7(6):625,  
 485

- 486            2017.
- 487 [41] Shrey K Shahi and Jack W Baker. Regression models for predicting the probability of near-fault earthquake ground  
488 motion pulses, and their period. *Applications of Statistics and Probability in Civil Engineering*, 30(4):459, 2011.
- 489 [42] UC Berkeley Seismological Laboratory, United States Geological Survey (USGS), Calpine and Unocal Corporations.  
490 Northern California earthquake data center, 2014.
- 491 [43] P. Nuannin. The potential of b-value variations as earthquake precursors for small and large events. Technical Report  
492 183, Uppsala University, Sweden, 2006.
- 493 [44] B. Gutenberg and C. F. Richter. Frequency of earthquakes in California. *Bulletin of the Seismological Society of America*,  
494 34:185–188, 1944.

**495    Tables**

<i>Zone</i>	State of California
<i>Period</i>	1970-2017
<i>Source</i>	ANSS Composite Earthquake Catalog
<i>Size</i>	917.7 MB
<i>Events</i>	1,421,691 events
<i>Max. magnitude</i>	7.3

Table 1: Specifications of the original catalog used in this work from whose the different datasets for the regression study were extracted.

<i>Lat/Long. granularity</i>	0.5 degrees
<i>Max. magn. threshold</i>	5
<i>Min. events/cell</i>	500 events
<i>Num. selected cells</i>	27 cells
<i>Resulting size range</i>	538 - 5,575 events

Table 2: Definition of the grid established over the original catalog of California. Each selected cell from the grid produces a dataset to perform earthquake predictions in the presented regression study.

Dataset	Lat.Min	Lat.Cen	Lat.Max	Lon.Min	Lon.Cen	Lon.Max
2-19	32.50	32.83	33.00	-116.00	-115.73	-115.50
3-18	33.00	33.29	33.50	-116.50	-116.29	-116.00
3-19	33.00	33.13	33.50	-116.00	-115.68	-115.50
4-17	33.50	33.81	34.00	-117.00	-116.73	-116.50
4-18	33.50	33.89	34.00	-116.50	-116.28	-116.00
5-13	34.00	34.32	34.50	-119.00	-118.62	-118.50
5-14	34.00	34.30	34.50	-118.50	-118.38	-118.00
5-16	34.00	34.13	34.50	-117.50	-117.25	-117.00
5-17	34.00	34.22	34.50	-117.00	-116.78	-116.50
5-18	34.00	34.22	34.50	-116.50	-116.38	-116.00
6-17	34.50	34.74	35.00	-117.00	-116.72	-116.50
6-18	34.50	34.68	35.00	-116.50	-116.32	-116.00
8-8	35.50	35.70	36.00	-121.50	-121.10	-121.00
8-9	35.50	35.68	36.00	-121.00	-120.80	-120.50
8-15	35.50	35.80	36.00	-118.00	-117.70	-117.50
9-9	36.00	36.23	36.50	-121.00	-120.76	-120.50
9-10	36.00	36.21	36.50	-120.50	-120.28	-120.00
9-15	36.00	36.12	36.50	-118.00	-117.79	-117.50
10-7	36.50	36.88	37.00	-122.00	-121.60	-121.50
10-8	36.50	36.67	37.00	-121.50	-121.24	-121.00
11-7	37.00	37.25	37.50	-122.00	-121.72	-121.50
11-13	37.00	37.45	37.50	-119.00	-118.73	-118.50
12-7	37.50	37.75	38.00	-122.00	-121.85	-121.50
12-13	37.50	37.61	38.00	-119.00	-118.85	-118.50
12-14	37.50	37.56	38.00	-118.50	-118.44	-118.00
14-5	38.50	38.80	39.00	-123.00	-122.79	-122.50
17-2	40.00	40.36	40.50	-124.50	-124.23	-124.00

Table 3: Location of the 27 datasets analyzed in the study. Each dataset was named using the number of the cell in a grid drawn over the state of California. The latitude and longitude ranges (minimum and maximum) are shown along with the mean of latitudes and longitudes (centroid) of the events inside each dataset (Lat.Cen, Lon.Cen).

Dataset	Size	Q1	Median	Mean	Q3	Max
2-19	2195	2.63	2.81	2.95	3.15	5.80
3-18	1051	2.62	2.79	2.91	3.09	5.43
3-19	1950	2.62	2.81	2.94	3.12	6.60
4-17	1065	2.59	2.73	2.84	2.97	6.00
4-18	1386	2.60	2.77	2.90	3.07	6.10
5-13	1022	2.62	2.80	2.97	3.15	6.70
5-14	1281	2.68	2.99	3.11	3.40	6.60
5-16	889	2.59	2.74	2.87	3.01	5.60
5-17	2326	2.61	2.78	2.92	3.08	6.30
5-18	3013	2.60	2.76	2.90	3.01	7.30
6-17	1606	2.60	2.75	2.85	2.98	5.26
6-18	1827	2.62	2.80	2.93	3.10	7.10
8-8	763	2.64	2.85	2.98	3.23	6.50
8-9	717	2.63	2.80	2.92	3.10	5.00
8-15	1281	2.61	2.77	2.89	3.02	5.75
9-9	1346	2.65	2.85	2.96	3.14	5.40
9-10	1840	2.65	2.87	3.00	3.21	6.70
9-15	1366	2.61	2.77	2.90	3.04	5.30
10-7	1940	2.69	2.92	3.02	3.22	5.40
10-8	5575	2.66	2.89	3.00	3.20	5.50
11-7	1615	2.63	2.82	2.95	3.11	6.90
11-13	1595	2.68	2.93	3.00	3.20	6.10
12-7	807	2.62	2.80	2.91	3.06	5.80
12-13	4002	2.68	2.95	3.04	3.26	6.20
12-14	724	2.65	2.89	3.04	3.28	6.40
14-5	3156	2.60	2.75	2.86	3.00	5.01
17-2	538	2.64	2.84	2.94	3.10	7.20

Table 4: Magnitude distributions and sizes of the 27 datasets analyzed in the study. The size is shown as the number of events included in each dataset. The magnitude distribution is expressed as its first quartile, median, mean, third quartile and the maximum magnitude for each dataset.

Feature	Description
$b$	Gutenberg-Richter law's b-value
$x_1$	Increment of $b$ between the events $i$ and $i - 4$
$x_2$	Increment of $b$ between the events $i - 4$ and $i - 8$
$x_3$	Increment of $b$ between the events $i - 8$ and $i - 12$
$x_4$	Increment of $b$ between the events $i - 12$ and $i - 16$
$x_5$	Increment of $b$ between the events $i - 16$ and $i - 20$
$x_6$	Maximum magnitude from the events recorded during the last week (OU's law)
$x_7$	Probability of recording an event with magnitude larger or equal to 6.0 using a probability density function
$a$	Gutenberg-Richter law's a-value
$\eta$	Mean square deviation
$\Delta M$	Magnitude deficit
$T$	Elapsed time
$\mu$	Mean time
$c$	Coefficient of variation
$dE^{1/2}$	Rate of square root of seismic energy
$M_{mean}$	Mean magnitude

Table 5: Set of features computed from catalogs to produce the datasets used to train and test the regression methods of the study.

Dataset	GLM	GBM	DL	RF	GLM-RF	GBM-RF	DL-RF	ALL
2-19	0.93	0.62	0.79	<b>0.56</b>	<b>0.56</b>	<b>0.56</b>	<b>0.56</b>	<b>0.56</b>
3-18	1.38	<b>1.33</b>	1.34	1.34	1.40	1.39	1.39	1.38
3-19	1.10	0.87	1.36	0.81	<b>0.80</b>	0.82	<b>0.80</b>	0.82
4-17	1.33	1.09	1.23	<b>1.05</b>	1.08	1.08	1.08	1.08
4-18	0.81	0.60	0.84	<b>0.59</b>	<b>0.59</b>	<b>0.59</b>	<b>0.59</b>	<b>0.59</b>
5-13	0.62	0.53	0.78	<b>0.46</b>	0.48	0.47	<b>0.46</b>	0.48
5-14	2.50	0.59	4.04	<b>0.57</b>	0.58	<b>0.57</b>	0.88	0.84
5-16	1.38	<b>1.30</b>	1.44	1.31	1.36	1.36	1.36	1.36
5-17	1.27	0.59	2.25	<b>0.55</b>	<b>0.55</b>	<b>0.55</b>	<b>0.55</b>	0.56
5-18	0.66	0.43	0.74	<b>0.37</b>	<b>0.37</b>	0.39	<b>0.37</b>	0.39
6-17	0.62	0.44	0.70	<b>0.41</b>	0.43	0.42	0.42	0.43
6-18	0.58	0.45	0.64	<b>0.40</b>	0.41	0.41	0.41	0.41
8-8	0.74	0.50	0.67	<b>0.46</b>	0.49	0.48	<b>0.46</b>	0.48
8-9	1.29	0.56	0.90	<b>0.53</b>	0.55	0.56	0.55	0.54
8-15	0.83	0.69	1.07	<b>0.65</b>	0.67	0.68	0.67	0.68
9-9	1.06	<b>0.87</b>	1.21	0.89	0.91	0.90	0.92	0.90
9-10	0.81	<b>0.59</b>	0.77	<b>0.59</b>	0.61	<b>0.59</b>	0.61	<b>0.59</b>
9-15	1.07	0.91	1.14	<b>0.89</b>	0.91	0.90	0.91	0.90
10-7	1.12	1.04	1.08	<b>1.02</b>	1.04	1.03	1.03	1.03
10-8	0.76	0.69	0.78	<b>0.64</b>	<b>0.64</b>	0.65	0.65	0.65
11-7	1.14	1.11	1.17	<b>1.08</b>	1.09	1.12	1.09	1.10
11-13	1.05	0.87	1.25	<b>0.84</b>	0.85	0.85	0.85	0.85
12-7	1.16	1.06	1.21	<b>1.05</b>	1.06	1.06	1.06	1.06
12-13	0.81	0.68	0.79	<b>0.61</b>	<b>0.61</b>	0.62	<b>0.61</b>	0.62
12-14	0.55	0.33	0.58	<b>0.30</b>	<b>0.30</b>	0.31	<b>0.30</b>	0.31
14-5	1.10	<b>1.09</b>	1.11	1.10	<b>1.09</b>	<b>1.09</b>	<b>1.09</b>	<b>1.09</b>
17-2	1.12	0.99	1.09	<b>0.95</b>	1.00	1.02	<b>0.95</b>	0.96
<b>Average</b>	1.03	0.77	1.15	<b>0.74</b>	0.76	0.76	0.76	0.76

Table 6: Mean absolute error of the different regressors analyzed when they predict the 27 earthquake datasets of the study. A hold-out validation scheme was applied splitting the first 75% of samples for training and the last 25% for testing. Regressors are generalized linear models (GLM), gradient boosting machines (GBM), deep learning (DL), random forests (RF) and 4 stacked ensembles including RF: GLM-RF, GBM-RF, DL-RF and GLM-GBM-DL-RF (ALL).

Dataset	GLM	GBM	DL	RF	GLM-RF	GBM-RF	DL-RF	ALL
2-19	0.16	0.11	0.14	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>
3-18	0.25	<b>0.24</b>	0.25	0.25	0.26	0.26	0.26	0.25
3-19	0.17	0.13	0.21	<b>0.12</b>	<b>0.12</b>	<b>0.12</b>	<b>0.12</b>	<b>0.12</b>
4-17	0.22	<b>0.18</b>	0.21	<b>0.18</b>	<b>0.18</b>	<b>0.18</b>	<b>0.18</b>	<b>0.18</b>
4-18	0.13	<b>0.10</b>	0.14	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>
5-13	0.09	0.08	0.12	<b>0.07</b>	<b>0.07</b>	<b>0.07</b>	<b>0.07</b>	<b>0.07</b>
5-14	0.38	<b>0.09</b>	0.61	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>	0.13	0.13
5-16	0.25	<b>0.23</b>	0.26	<b>0.23</b>	0.24	0.24	0.24	0.24
5-17	0.20	<b>0.09</b>	0.36	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>
5-18	0.09	0.06	0.10	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>
6-17	0.12	<b>0.08</b>	0.13	<b>0.08</b>	<b>0.08</b>	<b>0.08</b>	<b>0.08</b>	<b>0.08</b>
6-18	0.08	<b>0.06</b>	0.09	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>
8-8	0.11	0.08	0.10	<b>0.07</b>	<b>0.07</b>	<b>0.07</b>	<b>0.07</b>	<b>0.07</b>
8-9	0.26	<b>0.11</b>	0.18	<b>0.11</b>	<b>0.11</b>	<b>0.11</b>	<b>0.11</b>	<b>0.11</b>
8-15	0.14	0.12	0.19	<b>0.11</b>	0.12	0.12	0.12	0.12
9-9	0.20	<b>0.16</b>	0.22	<b>0.16</b>	0.17	0.17	0.17	0.17
9-10	0.12	<b>0.09</b>	0.12	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>	<b>0.09</b>
9-15	0.20	<b>0.17</b>	0.21	<b>0.17</b>	<b>0.17</b>	<b>0.17</b>	<b>0.17</b>	<b>0.17</b>
10-7	0.21	<b>0.19</b>	0.20	<b>0.19</b>	<b>0.19</b>	<b>0.19</b>	<b>0.19</b>	<b>0.19</b>
10-8	0.14	0.13	0.14	<b>0.12</b>	<b>0.12</b>	<b>0.12</b>	<b>0.12</b>	<b>0.12</b>
11-7	<b>0.16</b>	<b>0.16</b>	0.17	<b>0.16</b>	<b>0.16</b>	<b>0.16</b>	<b>0.16</b>	<b>0.16</b>
11-13	0.17	<b>0.14</b>	0.21	<b>0.14</b>	<b>0.14</b>	<b>0.14</b>	<b>0.14</b>	<b>0.14</b>
12-7	0.20	<b>0.18</b>	0.21	<b>0.18</b>	<b>0.18</b>	<b>0.18</b>	<b>0.18</b>	<b>0.18</b>
12-13	0.13	0.11	0.13	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>
12-14	0.09	<b>0.05</b>	0.09	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>
14-5	<b>0.22</b>							
17-2	0.16	0.14	0.15	<b>0.13</b>	0.14	0.14	<b>0.13</b>	<b>0.13</b>
<b>Average</b>	0.17	0.13	0.19	<b>0.12</b>	0.13	0.13	0.13	0.13

Table 7: Relative error of the different regressors analyzed when they predict the 27 earthquake datasets of the study. The relative error is computed as the mean absolute error divided by the maximum magnitude of each dataset. A hold-out validation scheme was applied splitting the first 75% of samples for training and the last 25% for testing. Regressors are generalized linear models (GLM), gradient boosting machines (GBM), deep learning (DL), random forests (RF) and 4 stacked ensembles including RF: GLM-RF, GBM-RF, DL-RF and GLM-GBM-DL-RF (ALL).

Regressor	[0,3)	[3,4)	[4,5)	[5,6)	[6,7)	[7,8]
GLM	$1.25 \pm 0.88$	$0.67 \pm 0.64$	$0.57 \pm 0.58$	$0.77 \pm 0.81$	$1.15 \pm 0.92$	$3.85 \pm 0.49$
GBM	$1.20 \pm 0.81$	$0.62 \pm 0.67$	$0.30 \pm 0.53$	$0.32 \pm 0.62$	$0.34 \pm 0.76$	$2.57 \pm 1.99$
DL	$1.24 \pm 0.87$	$0.79 \pm 0.72$	$0.60 \pm 0.60$	$0.80 \pm 0.87$	$1.36 \pm 1.00$	$3.49 \pm 0.39$
RF	<b><math>1.20 \pm 0.77</math></b>	<b><math>0.58 \pm 0.66</math></b>	<b><math>0.22 \pm 0.53</math></b>	<b><math>0.24 \pm 0.61</math></b>	<b><math>0.26 \pm 0.75</math></b>	<b><math>2.03 \pm 1.81</math></b>
GLM-RF	$1.20 \pm 0.78$	$0.58 \pm 0.68$	$0.25 \pm 0.53$	$0.26 \pm 0.60$	$0.29 \pm 0.74$	$2.07 \pm 1.85$
GBM-RF	$1.20 \pm 0.79$	$0.59 \pm 0.66$	$0.25 \pm 0.52$	$0.27 \pm 0.60$	$0.28 \pm 0.75$	$2.19 \pm 1.83$
DL-RF	$1.20 \pm 0.78$	$0.59 \pm 0.67$	$0.24 \pm 0.53$	$0.26 \pm 0.60$	$0.28 \pm 0.74$	$2.03 \pm 1.83$
ALL	$1.20 \pm 0.78$	$0.59 \pm 0.66$	$0.25 \pm 0.52$	$0.27 \pm 0.60$	$0.28 \pm 0.75$	$2.20 \pm 1.83$
# samples	35,608	27,168	21,544	5,936	616	48
% samples	39.16	29.88	23.7	6.53	0.68	0.05

Table 8: Mean absolute errors (mean  $\pm$  std. deviation) produced for each regressor according to the different magnitude intervals.

Regressor	[0,3)	[3,4)	[4,5)	[5,6)	[6,7)	[7,8]
GLM	$2.34 \pm 2.79$	$0.85 \pm 1.62$	$0.66 \pm 1.63$	$1.25 \pm 3.15$	$2.16 \pm 3.87$	$15.04 \pm 3.74$
GBM	$2.08 \pm 2.40$	$0.83 \pm 1.61$	$0.37 \pm 1.44$	$0.49 \pm 2.06$	$0.69 \pm 2.93$	$9.91 \pm 10.22$
DL	$2.28 \pm 2.71$	$1.11 \pm 1.89$	$0.72 \pm 1.64$	$1.40 \pm 3.33$	$2.84 \pm 4.09$	$12.28 \pm 2.67$
RF	$2.07 \pm 2.29$	$0.80 \pm 1.60$	<b><math>0.33 \pm 1.42</math></b>	<b><math>0.42 \pm 2.15</math></b>	<b><math>0.62 \pm 2.66</math></b>	$6.97 \pm 7.72$
GLM-RF	$2.05 \pm 2.25$	<b><math>0.78 \pm 1.53</math></b>	$0.34 \pm 1.45$	$0.43 \pm 2.11$	$0.62 \pm 2.68$	$7.15 \pm 7.85$
GBM-RF	$2.04 \pm 2.23$	<b><math>0.78 \pm 1.53</math></b>	$0.34 \pm 1.41$	$0.43 \pm 2.05$	$0.63 \pm 2.73$	$7.58 \pm 8.13$
DL-RF	$2.05 \pm 2.25$	$0.79 \pm 1.55$	$0.33 \pm 1.44$	$0.43 \pm 2.11$	$0.62 \pm 2.68$	<b><math>6.88 \pm 7.57</math></b>
ALL	<b><math>2.03 \pm 2.25</math></b>	$0.79 \pm 1.53$	$0.34 \pm 1.42$	$0.43 \pm 2.06$	$0.63 \pm 2.73$	$7.61 \pm 8.14$
# samples	35,608	27,168	21,544	5,936	616	48
% samples	39.16	29.88	23.7	6.53	0.68	0.05

Table 9: Mean squared errors (mean  $\pm$  std. deviation) produced for each regressor according to the different magnitude intervals.

Process	Execution time
Data acquisition	15 minutes
Cell building	4 minutes
Cell selection	3 minutes
Feature generation	6 minutes
Dataset splitting	10 seconds
Training of the regressors	6 hours
Prediction of the test subsets	2 minutes
Statistics computation	1 minute

Table 10: Execution times of each process of the proposed methodology. These values are obtained with the IT infrastructure described in Section 3.7.

496 **Figures**

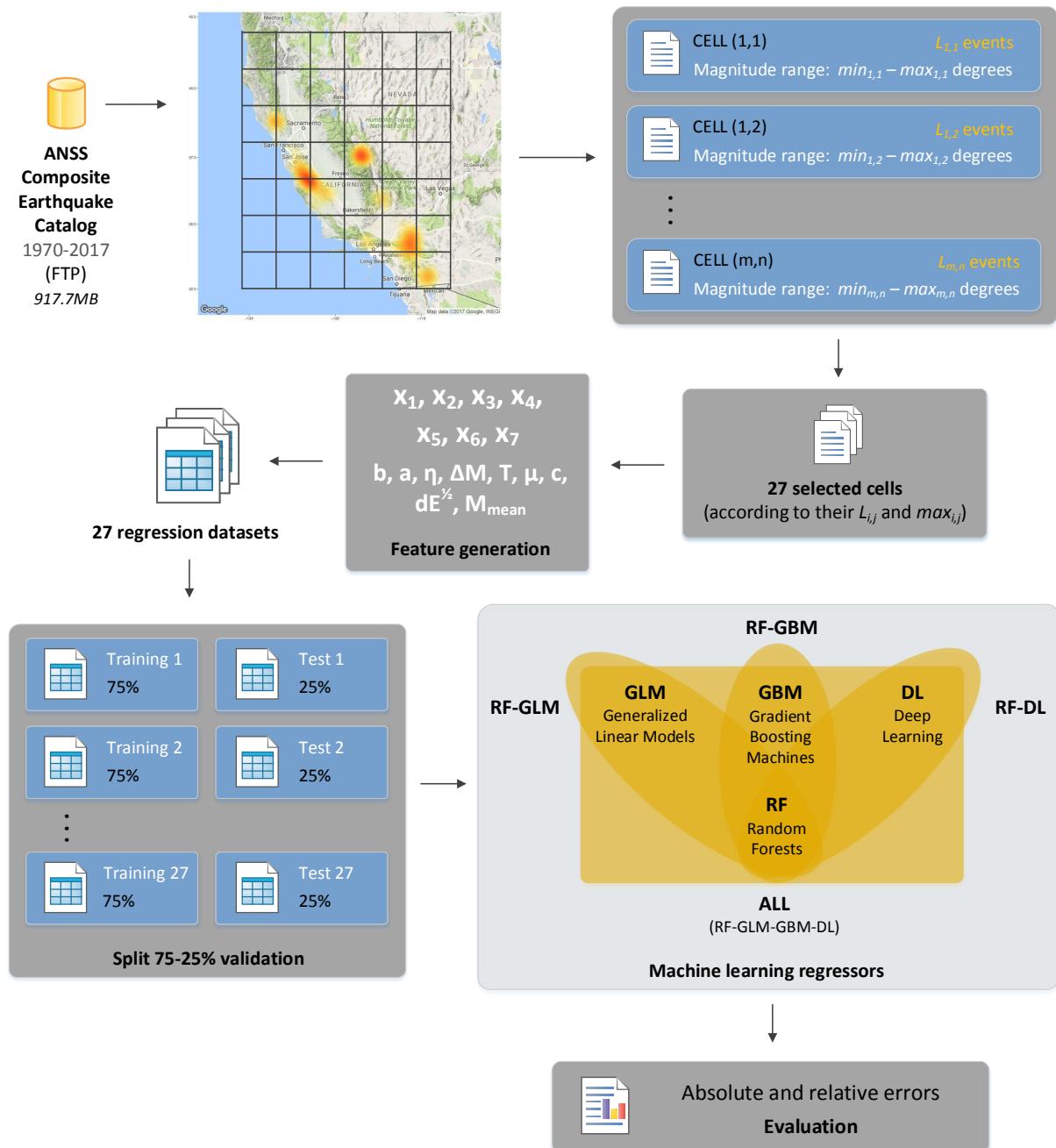


Figure 1: Schematic diagram of the methodology used to retrieve, divide, train and test the set of machine learning methods used to perform the proposed regression study.

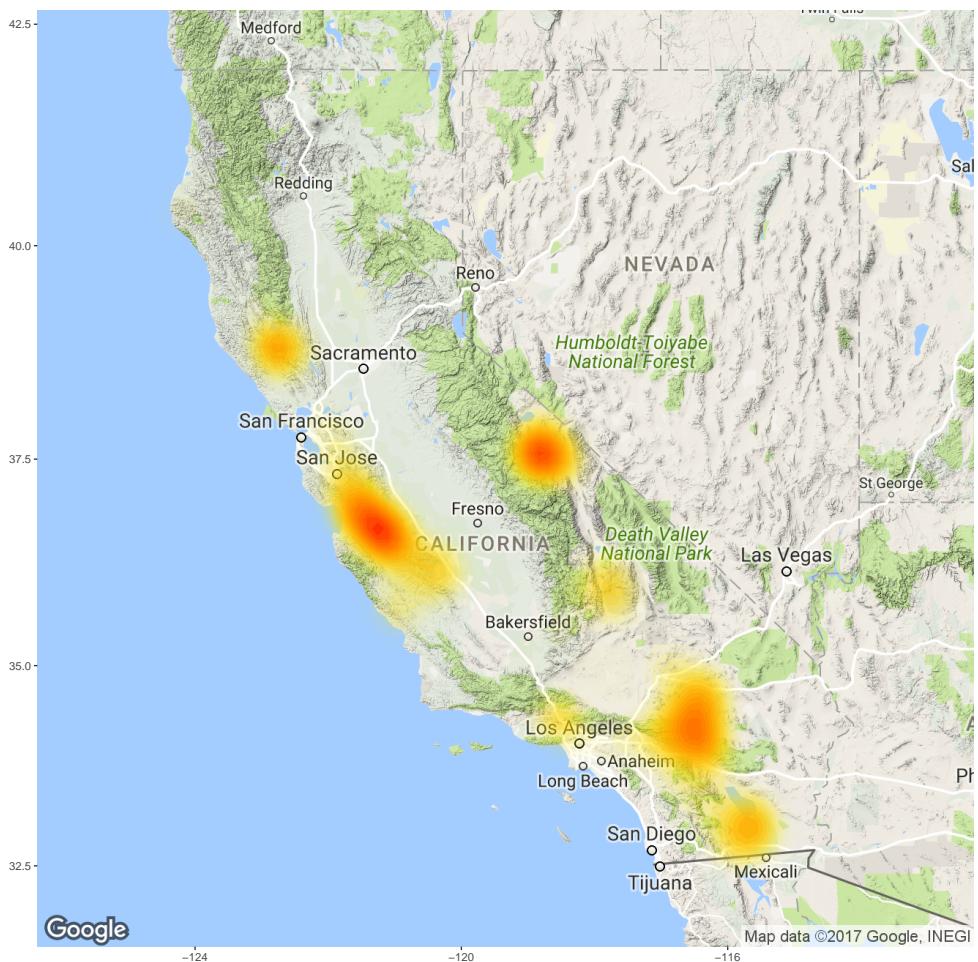


Figure 2: Earthquake events studied in California between 1970-2017 colored by the number of occurrences. Data were obtained from the ANSS Composite Earthquake Catalog.

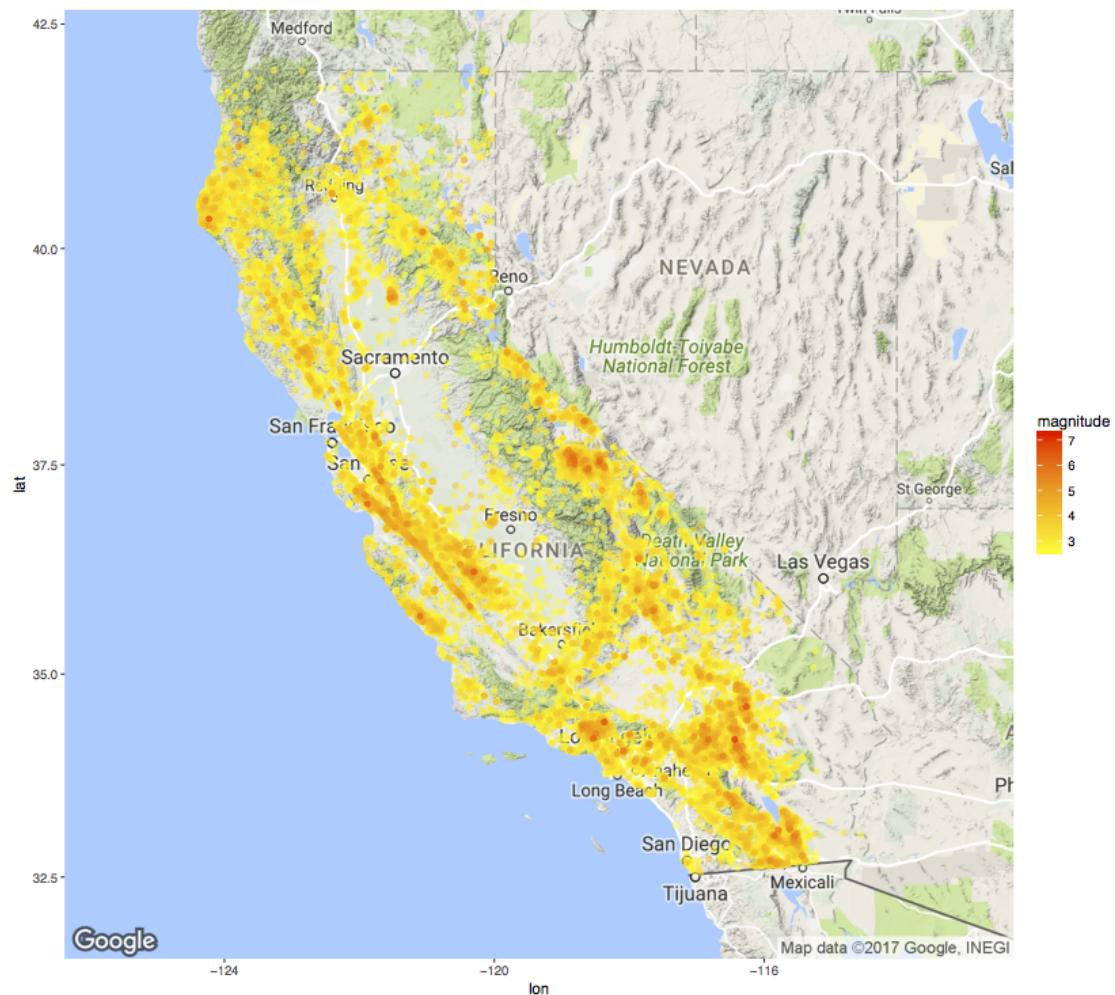


Figure 3: Earthquake events studied in California between 1970-2017 colored by their magnitude. Data were obtained from the ANSS Composite Earthquake Catalog.

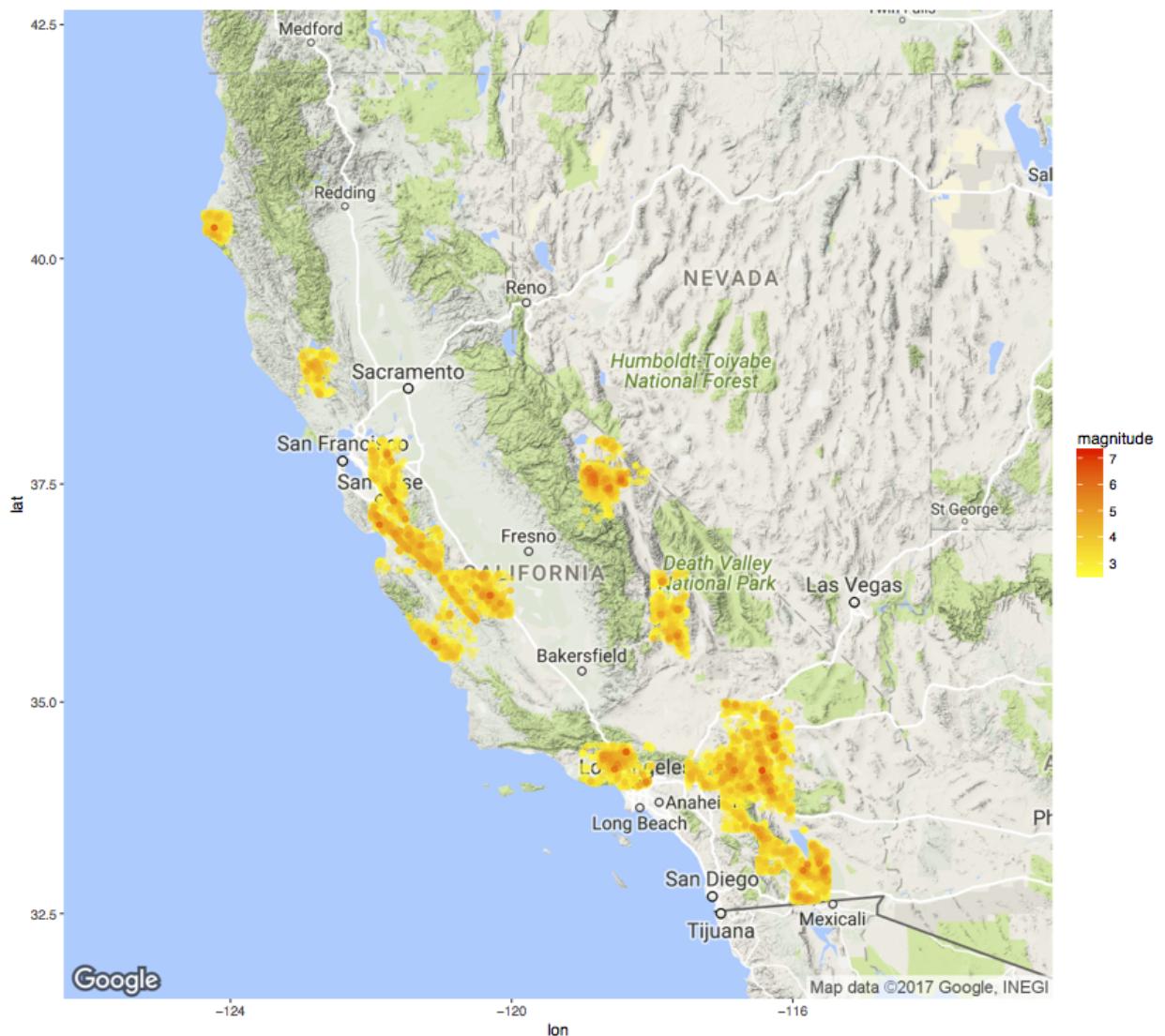


Figure 4: Earthquake events of the 27 filtered cells colored by their magnitude.

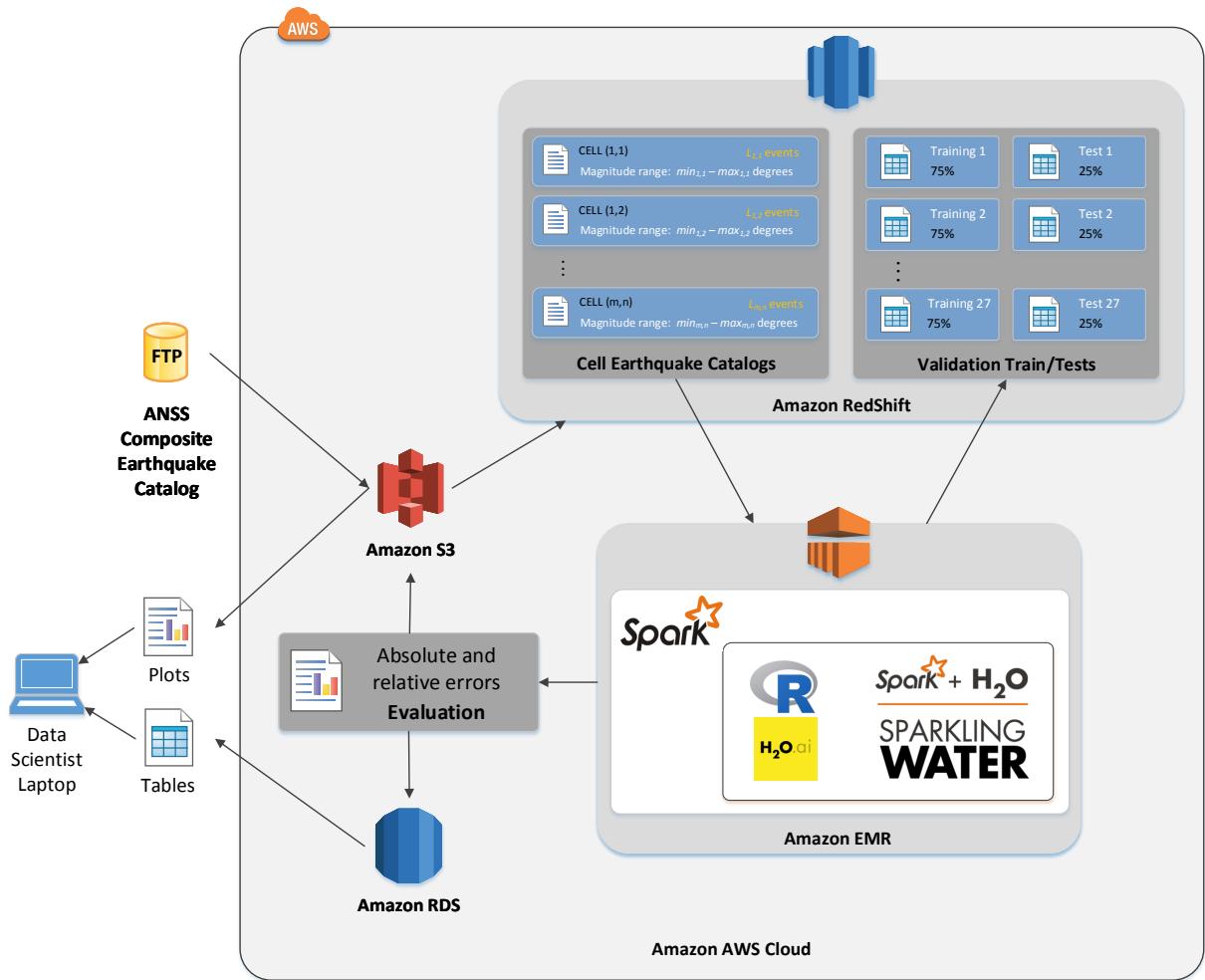


Figure 5: The cloud-based Big Data IT infrastructure implemented for the earthquake analysis and prediction study. Amazon Web Services were chosen to provide the platform as a service needed to deploy all software components of the proposed methodology.

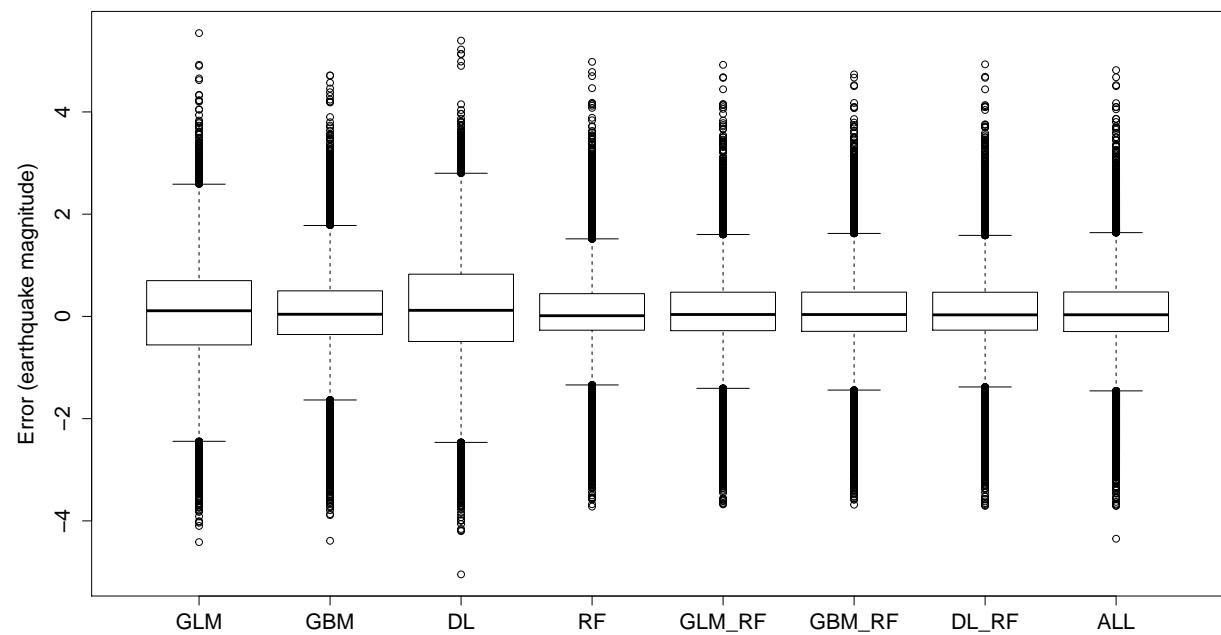


Figure 6: Boxplot of the errors produced by the regression algorithms when predicting the 27 datasets of the study.

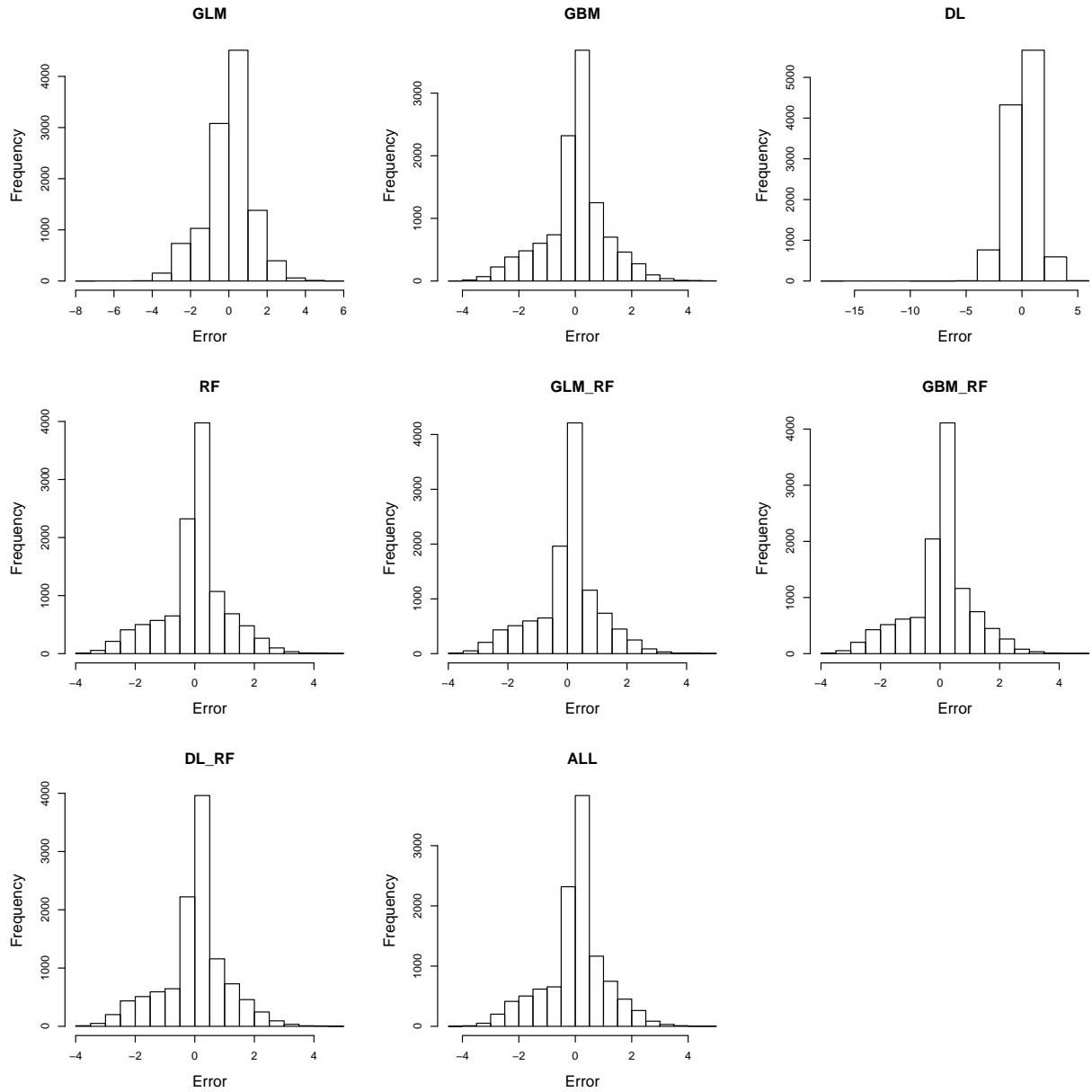


Figure 7: Histograms of the errors produced by the regression algorithms when predicting the 27 datasets of the study.

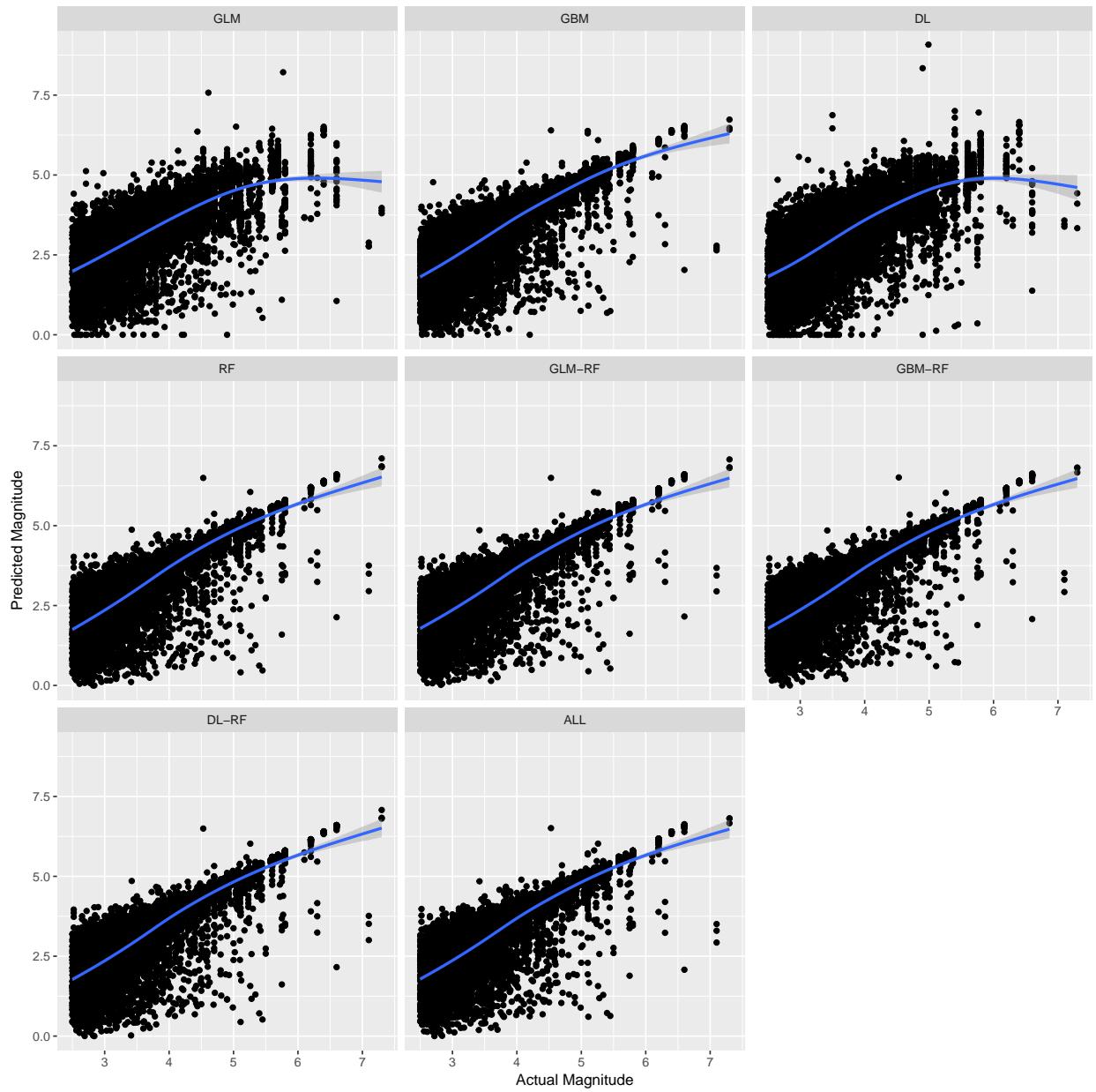


Figure 8: Actual versus predicted magnitude values for each regression algorithm when predicting the 27 datasets of the study.