# Transfer Learning to Build a Scalable Model for the Declustering of Earthquake Catalogs

Florent Aden-Antoniow<sup>1</sup>, William Benjamin Frank<sup>1</sup>, and Leonard Seydoux<sup>2</sup>

<sup>1</sup>Massachusetts Institute of Technology <sup>2</sup>Université Grenoble Alpes

November 21, 2022

#### Abstract

The rate of background seismicity, or the earthquakes not directly triggered by another earthquake, in active seismic regions is indicative of the stressing rate of fault systems. However, aftershock sequences often dominate the seismicity rate, masking this background seismicity. The identification of aftershocks in earthquake catalogs, also known as declustering, is thus an important problem in seismology. Most solutions involve spatio-temporal distances between successive events, such as the Nearest-Neighbor-Distance algorithm widely used in various contexts. This algorithm assumes that the space-time metric follows a bi-modal distribution with one peak related to the background seismicity and another peak representing the aftershocks. Constraining these two distributions is key to accurately identify the aftershocks from the background events. Recent work often uses a linear-splitting based on nearest-neighbor distance threshold, ignoring the overlap between the two populations and resulting in a mis-identification of background earthquakes and aftershock sequences. We revisit this problem here with both machine-learning classification and clustering algorithms. After testing several popular algorithms, we show that a random forest trained with various synthetic catalogs generated by an Epidemic Type Aftershock Sequence model outperforms approaches such as K-means, Gaussian-mixture models, and Support Vector Classifications. We evaluate different data features and discuss their importance in classifying aftershocks.

We then apply our model to two different actual earthquake catalogs, the relocated Southern California Earthquake Center catalog and the GeoNet catalog of New Zealand. Our model capably adapts to these two different tectonic contexts, highlighting the differences in aftershock productivity between crustal and intermediate depth seismicity.

## Transfer Learning to Build a Scalable Model for the Declustering of Earthquake Catalogs

1

2

3

## F. Aden-Antoniów<sup>1</sup>, W. B. Frank<sup>1</sup>, L. Seydoux<sup>2</sup>

 <sup>1</sup>Department of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA
 <sup>2</sup>Équipe Ondes et Structures, ISTerre, Université Grenoble-Alpes, UMR 5375, Gières, France

Corresponding author: Florent Aden-Antoniow, faden@mit.edu

#### $_{7}$ Abstract

The rate of background seismicity, or the earthquakes not directly triggered by another 8 earthquake, in active seismic regions is indicative of the stressing rate of fault systems. q However, aftershock sequences often dominate the seismicity rate, masking this background 10 seismicity. The identification of aftershocks in earthquake catalogs, also known as declus-11 tering, is thus an important problem in seismology. Most solutions involve spatio-temporal 12 distances between successive events, such as the Nearest-Neighbor-Distance algorithm widely 13 used in various contexts. This algorithm assumes that the space-time metric follows a bi-14 modal distribution with one peak related to the background seismicity and another peak rep-15 resenting the aftershocks. Constraining these two distributions is key to accurately identify 16 the aftershocks from the background events. Recent work often uses a linear-splitting based 17 on nearest-neighbor distance threshold, ignoring the overlap between the two populations 18 and resulting in a mis-identification of background earthquakes and aftershock sequences. 19 We revisit this problem here with both machine-learning classification and clustering al-20 gorithms. After testing several popular algorithms, we show that a random forest trained 21 with various synthetic catalogs generated by an Epidemic Type Aftershock Sequence model 22 outperforms approaches such as K-means, Gaussian-mixture models, and Support Vector 23 Classifications. We evaluate different data features and discuss their importance in classify-24 ing aftershocks. We then apply our model to two different actual earthquake catalogs, the 25 26 relocated Southern California Earthquake Center catalog and the GeoNet catalog of New Zealand. Our model capably adapts to these two different tectonic contexts, highlighting 27 the differences in aftershock productivity between crustal and intermediate depth seismicity. 28

#### <sup>29</sup> Plain Language Summary

The seismic catalog of earthquakes that have occurred in a given fault zone is a window 30 into the tectonic processes that occur at depth. These earthquakes rupture faults when the 31 fault can no longer support the stress built-up by tectonic motion. When an earthquake oc-32 curs spontaneously from tectonic stresses, this mainshock will trigger aftershocks, which can 33 themselves trigger even more aftershocks. Aftershock sequences often dominate catalogs due 34 to the sudden increase in the number of earthquakes. Distinguishing between mainshocks 35 and aftershocks requires an understanding of the connection between earthquakes. An ac-36 curate classification of mainshocks and aftershocks in a catalog allow notably a more precise 37 investigation of the evolution of accumulated stress on a fault and is often a necessary tool 38 to estimate the seismic hazard of a region. In this work, we develop a machine-learning 39 algorithm to achieve such classification. We tested our model on large earthquake catalogs 40 of Southern California and New Zealand to demonstrate the effectiveness of our approach. 41 We show that our approach is generalizeable to any region of the world, independent of the 42 style of seismicity or period of time. 43

#### 44 1 Introduction

That earthquakes form clusters in time and space, regardless of the tectonic context, is a 45 fundamental characteristic of earthquakes. The simplest classification of earthquakes breaks 46 down the total seismicity rate into two categories: background events that are generated 47 by long-term, large-scale tectonic forcings and aftershocks that are directly triggered by a 48 background event or another aftershock. Isolating the background seismicity is crucial for 49 a wide range of studies: from the monitoring of transient loading along faults (Marsan, 50 Prono, & Helmstetter, 2013; Marsan, Reverso, et al., 2013; Reverso et al., 2015), fluid 51 injections (Hainzl & Ogata, 2005; Bachmann et al., 2011; Kothari et al., 2020) for hazard 52 assessments where declustering is essential to remove the spatial bias introduced by clustered 53 seismic activity in non-declustered catalogs (Marzocchi & Taroni, 2014; Azak et al., 2018; 54 Galina et al., 2019; Taroni & Akinci, 2021). Aftershocks can be modeled with empirical 55 laws, such as Omori's law or the productivity law (Omori, 1894; Utsu, 1961), that predict 56

the rate and magnitude of aftershocks following some background event. The classification 57 of an earthquake as a mainshock or an aftershock then relies on the identification of the 58 relationship between a given event and the events that preceded it. If an earthquake is likely 59 to be triggered by a preceding earthquake, it can then be confidently labeled as an aftershock. 60 Several algorithms have been proposed to tackle this classification problem, also known 61 as declustering. Early declustering algorithms were mostly based on space-time windows 62 (Gardner & Knopoff, 1974; Reasenberg, 1985). More recently, declustering approaches have 63 made use of the Epidemic-Type Aftershock Sequence (ETAS) model (Ogata, 1988, 1998; 64 Zhuang et al., 2002) or space-time-magnitude metric describing the link between event-pair 65 (Baiesi & Paczuski, 2004; Zaliapin et al., 2008); we refer to (Van Stiphout et al., 2012) for 66 an exhaustive list of declustering algorithms. Each method possess its own parameters but 67 most of them rely on the estimation of the "reach" of each earthquakes to others. 68

Teng and Baker (2019) have recently compared several declustering methods: the space-69 time window based algorithms proposed by Gardner and Knopoff (1974) and by Reasenberg 70 (1985), the Nearest-Neighbor-Distance (NND) Algorithm (Baiesi & Paczuski, 2004; Zaliapin 71 et al., 2008) relying on a space-time-magnitude metric and stochastic declustering (Zhuang 72 et al., 2002). Their study suggests that both Reasenberg and NND declustering methods 73 are suitable for the declustering of the seismicity in particular to estimate the background 74 seismic rate useful for seismic hazard analysis. The NND algorithm presents non-negligible 75 advantages as it does not use tuning parameters other than characteristic features of the 76 observed seismicity, such as the b value or the fractal dimension. 77

However, as pointed out by Bayliss et al. (2019), there is room to improve NND-based 78 declustering. The frequency distribution of the NND metric is often observed as a bi-modal 79 distribution characterizing the background seismicity and the aftershocks. A threshold is 80 then used to split the distribution and classify the earthquakes into two categories, main-81 shocks or aftershocks, neglecting the overlap of both distributions. This implies that the 82 declustered background seismicity will contain clustered seismicity and vice-versa. Bayliss 83 et al. (2019) developed a probabilistic clustering framework using a Markov Chain Monte Carlo mixture modelling approach allowing the overlap of two Weibull function with the aim 85 to quantify the uncertainties in event-pair linkage. Zaliapin and Ben-Zion (2020) recently 86 introduced a modified version of the Nearest-Neighbor-Distance algorithm by introducing 3 87 additional steps. Their algorithm discriminates background events and aftershocks following 88 a random thinning approach which aims to remove events according to a space-dependent 89 threshold, which is estimated using randomized-reshuffled catalogs. In this work, we have 90 decided to revisit the problem of the declustering of earthquake catalogs with a machine-91 learning based take on the NND algorithm but with the aim to keep the simplicity of the 92 original algorithm. The declustering can be seen either as a clustering or a classification 93 problem; both utilize data features to describe an earthquake. Machine learning techniques 94 are now commonly used in seismology (Kong et al., 2019), especially for the detection and/or 95 classification of seismic waveforms with neural networks (Perol et al., 2018; Li et al., 2018; 96 Ross et al., 2018; Zhu & Beroza, 2019; Thomas et al., n.d.), driven by either supervised 97 (Rubin et al., 2012; Lara-Cueva et al., 2016) or unsupervised learning algorithms (Sevdoux 98 et al., 2020; Shi et al., 2021; Steinmann et al., 2021). Beyond waveform classification, qq many other applications have benefited from machine learning algorithms, including but 100 not limited to early warning systems (Kong et al., 2016), earthquake relocation (Trugman 101 & Shearer, 2017), predicting fault slip cycles (Rouet-Leduc et al., 2017), and forecasting 102 earthquake magnitudes (Asim et al., 2017; González et al., 2019; Hoque et al., 2020). How-103 ever, only few studies have used machine learning to classify earthquakes within catalogs 104 (Picozzi & Iaccarino, 2021), and they have not addressed the declustering problem. Here 105 we implement and compare different machine learning algorithms apply to the declus-106 tering of earthquake catalogs. Building on the nearest-neighbor algorithm, we perform a 107 thoughtful comparison between clustering and classification approaches apply to synthetic 108 catalogs with the aim to identify the most accurate algorithm. We apply our model to 2 109

real catalogs from Southern California and New Zealand and discuss the potential of our declustering model.

#### 112 **2** Method

113

140

#### 2.1 Nearest-Neighbor Distance algorithm

The NND represents a good compromise between computational efficiency and stability, relying on a generalizeable metric of earthquake catalogs. The NND does not depend directly on the location and magnitude of earthquakes, but is rather based on the computation of a space-time metric taking into account both the Gutenberg-Richter law (Gutenberg & Richter, 1955) and the Omori-Utsu law (Omori, 1894; Utsu et al., 1995). It consists of the estimation of the metric  $\eta_{ij}$  between the *j*th event and a preceding event *i*. We call the nearest-neighbor event of *j*, the event *i* that minimizes this distance:

$$\eta_{ij} = t_{ij} \times (r_{ij})^{d_f} \times 10^{-bm_i} \tag{1}$$

where  $t_{ij} = t_j - t_i$  is the inter-event time,  $r_{ij} = |r_i - r_j|$  the inter-event physical distance (epicenter to epicenter; we ignore depth here),  $m_i$  is the magnitude of the event iand  $d_f$  is the fractal dimension. Zaliapin et al. (2008) introduced a re-scaled time-difference  $T_{i,j}$  and spatial distance  $R_{i,j}$  for discriminating clustered and non-clustered events in a 2D visualisation, thus allowing to account for both time and space distributions:

$$\eta_{ij} = T_{ij} \times R_{ij}$$

$$T_{ij} = t_{ij} \times 10^{-\frac{1}{2}bm_i}$$

$$R_{ij} = (r_{ij})^{df} \times 10^{-\frac{1}{2}bm_i}$$
(2)

For simplification, we refer to  $T_{ij}$  as the rescaled time and  $R_{ij}$  as the rescaled distance.

We measured the NND  $\eta$  for the relocated seismicity of Southern California, between 127 1981 and 2019 (Hauksson et al., 2012), and the GeoNet catalog of New Zealand, between 128 2010 and 2020 (Figure 1a and b). We identify two populations: aftershocks have system-129 atically small rescaled times  $T_{i,j}$  and distances  $R_{i,j}$ , while the mainshocks exhibit a wider 130 distribution of rescaled time and distance (Figure 1c and d). The simplest way to dis-131 criminate between the two populations is by drawing a line at the local minimum of the  $\eta$ 132 distribution (Figure 1e and f); this vertical line to split the NND distribution becomes a di-133 agonal in the 2D representation of rescaled time and distance. Knowing that the background 134 distribution follows a Weibull distribution (Zaliapin et al., 2008), by cutting the  $\eta$  distribu-135 tion aggressively we neglect any overlap and thus a portion of mainshocks are considered as 136 mainshocks and vice versa for aftershocks. To better constrain this problem, we consider the 137 declustering of earthquake catalog as both a clustering problem and a classification problem 138 to compare both approaches. 139

#### 2.2 A Clustering or a Classification problem

To decluster an earthquake catalog is to distinguish the mainshocks from the after-141 shocks. In a clustering (unsupervised learning) approach, we would expect both popula-142 tions to separate without any a priori on the data set or on any sort of model that would 143 have generated the earthquakes in the catalog. This idealized approach ignores the diffi-144 culty of evaluating the quality of the prediction because there is no explicit loss function 145 that quantifies the quality of a given clustering. In a classification (supervised learning) 146 approach, one can evaluate how well a model performs during training and testing phases 147 by comparing the results with the synthetic labels. In our particular case of earthquake 148

declustering, training such a model can be achieved by using a Epidemic-Type Aftershock
 Sequence (ETAS) model to generate synthetic catalogs in an exercise of transfer learning.

In this work, we have compared the results of four well-known machine learning models. The unsupervised learning algorithms we tested here are the K-means algorithm (linear and deterministic) and a Gaussian-Mixture model (non-linear and probabilistic). The supervised learning models we tested rely on either a Support Vector Machine Classifier (linear and deterministic) or a Random Forest classifier (non-linear and probabilistic). We describe each of the four models below.

The K-means algorithm will separate the samples into k groups, here k = 2 to account for mainshocks and aftershocks, by minimizing the cluster variances, also called inertia,  $\sum_{i=0}^{N_k} (||x_i - \mu_k||^2).$ 

A Gaussian Mixture model (GMM) is a probabilistic model and considers the sample 160 distribution as a mixture of a finite number K of Gaussian distributions with unknown 161 parameters; each Gaussian is defined by a mean  $\mu_k$  and a variance  $\sigma_k$ . GMM is initialized 162 by applying the k-means algorithm in order to assign a label to each data point and a first 163 set of model parameters. Then a two-step technique called Expectation Maximization is 164 performed to estimate the mixture model parameters. First, the membership weights (or 165 probabilities)  $\phi_{ik}$  for each sample with the given set of model parameters are computed. 166 For each data point, the membership weight is defined as  $\sum_{k=1}^{K} \phi_{ik} = 1$ . In a second step, it infers the new set of model parameters for the given membership weights by maximizing 167 168 the log-likelihood function. This is an iterative process and eventually converges to a final 169 set of parameters. 170

A Support Vector Machine Classifier (SVC) is an unsupervised learning model that can 171 be used to solve classification or regression problems. We decided to use a linear kernel for 172 the SVC to have a supervised, linear and deterministic model to compare with K-means. 173 With SVC, a sample is considered as a k-dimensional vector, with k the number of features. 174 The linear SVC can be trained to separate the data set with k-1 dimensional hyperplane(s). 175 If the solution of this problem is not unique, the best hyperplane is the one that create the 176 largest separation or margin between the two populations; the "best" hyperplane maximizes 177 the distance from it to the nearest sample(s) on each side, also called the support vectors. 178

A Random Forest (RF) is an ensemble classifier represented by a forest of decision trees. 179 One main advantage of a RF is to prevent over-fitting of the training data set. The ensemble 180 of uncorrelated trees is generated following a method of Bootstrap Aggregating (also called 181 bagging): each tree is trained with a randomly selected (with replacement) sub-set of the 182 data. Until now, RFs have been most often used to classify events such as earthquakes, 183 landslides, gever or period of volcanic activation based on time-frequency features (Rubin 184 et al., 2012; Hibert et al., 2017; Maggi et al., 2017; Provost et al., 2017; Yuan et al., 2019; 185 Dempsey et al., 2020); an RF model has not yet been applied to seismic catalogs. During 186 the training of a tree, at each node, the best split is found either from all input features 187 or more generally from a random subset of it, thus contributing to the randomness which 188 reduces potential overfitting. From this set of features, the idea is to find the feature and the 189 associated threshold that allows the best splitting of the remaining data-set into two groups, 190 based on the provided data labels. The selected feature and the given threshold minimize 191 the class impurity of the two populations according to the Gini impurity index computed 192 from the known class of the samples. The Gini impurity at a split is the probability that a 193 randomly chosen sample would be incorrectly labeled. It is expressed for the node n: 194

$$I_G(n) = 1 - \sum_{i=1}^{K} p_i^2$$
(3)

This is done recursively until one sample remains in the leaves or when the trees reach a maximum depth. When all trees are built, the classification of a sample is done by feeding it to the ensemble of trees and comparing its features to the selected feature and its respective threshold at each node. At the end of each tree, in the leaf where the sample falls, a prediction is made. The class prediction is then obtained by a vote of all trees. The final probability to belong to each class is derived by the percentage of trees in the forest that voted for a given classification.

We also considered using Gradient Boosted Trees as a supervised learning model but knowing that the approach relies on the fit of the prediction residuals in a series of weak learners, i.e. shallow decision trees, we believe it would be difficult to use for the declustering of seismic catalogs without over-fitting the training data set. We also considered a neural network, but we judged the training and the estimation of the model's architecture parameters too demanding computationally and likely not adequate for this particular problem with a potential small number of features compared to a large number of samples.

#### 209 2.3 Synthetic catalogs for training

To train and later test the supervised learning models, the Support Vector Machine Classifier and the Random Forest, we need to build a labelled data set, i.e. synthetic catalogs where we know which earthquakes is a mainshock or an aftershock. To do so, we have followed the model of Epidemic-Type Aftershock Sequence (ETAS) (Ogata, 1988; Marsan, Reverso, et al., 2013; Marsan et al., 2017). The spatio-temporal distribution of earthquakes  $\lambda(x, y, t)$  is defined as the sum of the background seismicity  $\mu(x, y, t)$  and the aftershocks  $\nu(x, y, t)$ :

$$\lambda(x, y, t) = \mu(x, y, t) + \nu(x, y, t) \tag{4}$$

The mainshock distribution  $\mu(x, y, t)$  is characterized by a Poisson point-process: the events are independent from each other and the inter-event time follows the Poisson probability density function as:

$$f(\tau) = \frac{1}{T_0} \exp\left(-\frac{\tau}{T_0}\right) \quad \text{for } \tau \ge 0 \tag{5}$$

We have drawn coordinates from a 2D probability density function (PDF) to give a location for each mainshocks. This PDF can be uniform for a randomly distributed background seismicity or can be generated from an existing seismicity.= by

Following Marsan, Reverso, et al. (2013), the aftershock distribution  $\nu(x, y, t)$  can be expressed as a spatio-temporal distribution, which results from the product of the Omori-Utsu law, a productivity law and a power spatial density:

$$\nu_{(x,y,t)} = \sum_{i}^{N_{after}} \frac{K 10^{a(m_i - m_c)}}{(t - t_i + c)^p} \times \frac{(\gamma - 1)L_i^{\gamma - 1}}{2\pi \left((x - x_i)^2 + (y - y_i)^2 + L_i^2\right)^{(\gamma - 1)/2}} \tag{6}$$

 $K \text{ and } a \text{ are constant and part of the aftershock productivity of the } i-\text{th event which follow a Poisson law with the average, with K and a constants (Gu et al., 2013). } t_i \text{ is the occurrence time of the mainshock and } L_i \text{ is the characteristic length in kilometers (Utsu & Seki, 1955) } such as <math>L_i = L_0 \times 10^{(m_i - m_c)/2}$  with  $L_0 = 0.1$ km (Marsan, Reverso, et al., 2013; Reverso et al., 2016).

231

<sup>232</sup> We create the mainshock seismicity in three steps:

(1) To account for variability in the inter-event time, we have randomly selected an averaged
 number of events per day for each catalog ranging from 1 to 3. The total average number

of mainshocks  $\langle N \rangle$  is obtained by multiplying the average number of events per day by the duration of the catalog desired. Here the duration is set to 8000 days ( $\sim 22$  years) to account for long-term seismic interactions, which we will discuss later on. The actual number of mainshocks N is drawn from a Poisson-law with a mean equal to  $\langle N \rangle$ . The averaged inter-event time  $T_0$ , which is simply N divided by the actual duration of the catalog, is injected in equation 5 to draw N inter-event times. The origin times of the N mainshocks are simply obtained by computing the cumulative sum of the inter-event times.

(2) To account for variability in the inter-event distance, we have used a 2D-PDF obtained
from the declustered interface seismicity for Northern-Chile from Aden-Antoniów et al.
(2020). From this catalog, we built a earthquake density map, smoothed it with a Gaussian
filter and normalized it by the total event count. Using a 2D PDF that is not uniform and
generated from a subduction zone seismicity is allowing us to simulate a broader distribution
of distance between the events.

(3) To obtain the magnitude of each mainshock, we have randomly selected a b-value, 248 uniformly distributed between 0.8 and 1, and did the same for the completeness magnitude, 249 uniformly distributed from 2 to 3. The magnitudes are then drawn from the probability 250 density function following the Gutenberg-Richter law (Gutenberg & Richter, 1955) with 251 a maximum magnitude set to 7.5. We selected this maximum magnitude on the basis 252 that for larger magnitude, some catalogs with low b value presented aftershock sequences 253 lasting for decades continuously triggering M6+ events resulting in catalogs constituted of 254 a disproportionate number of aftershocks in comparison of mainshocks. 255

We have generated the aftershocks and the aftershocks of the aftershocks and so on 256 following the method presented in section 2.3. To generate each catalog, we have drawn 257 a set of parameters of equation 6 such as the p and c, corresponding to the parameters of 258 the Omori-Utsu law (Omori, 1894; Utsu, 1957; Utsu et al., 1995), the parameter  $\gamma$  for the 259 power spatial density function, K and a for the productivity law. The magnitude of each 260 aftershocks is drawn from the same PDF as the mainshocks of the a given catalog. An 261 example of these synthetic catalogs is shown Figure S1 and the range of each parameters 262 of the ETAS model are shown Table 1. Finally, we have labeled each events of the catalogs 263 with 0 if it is a mainshock or 1 if it is an aftershock. We have generated 200 synthetic 264 catalogs representing a total of 5,461.475 earthquakes, 68% of which are aftershocks. This 265 corresponds to an average of 27,307 events per catalog with a minimum of 9,382 and a 266 maximum of 386,054. 267

For each event j in each synthetic catalog, we conducted the search of its nearest-268 neighbor i which is not necessarily it's mainshock if j is an aftershock. We estimated the 269 fractal dimension  $d_f$  necessary to compute the nearest-neighbor-distance metric for each 270 catalog following the Minkowski-Bouligand approach, also called the box counting method. 271 For different size boxes that divide the area covered by the synthetic catalogs, we counted 272 the number of boxes actually containing earthquakes. The relation between the size of the 273 boxes and the number of non-empty boxes can be described by a power law with the fractal 274 dimension  $d_f$  as exponent. The  $d_f$  of the synthetic catalogs are on average 1.7 with minimal 275 variations. 276

For both supervised and unsupervised approaches, we need to identify the relevant 277 features that are useful to distinguish background events from aftershocks. To keep our 278 model generalizable, the features necessary to classify an earthquake should not rely on 279 absolute locations or origin times. The NND framework (equation 3) provides several relative 280 metrics well suited to use as data features. To predict the label of event j, we describe 281 the link with its nearest neighbor i with five features: (1) the rescaled time  $T_{ij}$  and (2) 282 the rescaled distance  $R_{ij}$  both described in equation 3; (3) the difference in magnitude 283  $\Delta m_{ij} = m_i - m_j$  between the event j and its nearest-neighbor i, which we would expect 284 to be high if j is an aftershock of a larger event; (4)  $N_p$  the number of siblings or events 285 that share the same nearest neighbor as event j, which if relatively high would suggest that 286

event j is one of many aftershocks following a background event; and (5)  $N_c$  the number of offspring of the event j, which if high would suggest that event j has generated many aftershocks. Associated together, these features provide information on the linkages between earthquakes so that we may distinguish between mainshocks and aftershocks.

Of the 200 synthetic catalogs we generated (see section 2.3), we used 100 catalogs to 291 train the supervised learning models; the remaining 100 serve to test and compare all the 292 different models on a known data set. Besides training a supervised learning model, one 293 should also estimate the best hyper-parameters of the model, i.e. the unique parameters 294 defining the architecture of the model. To do so we have conducted a stratified k-Fold 295 cross-validation (Mosteller & Tukey, 1968). This consists of first splitting the training data 296 set into k groups or folds, in a *stratified* fashion. In this context, stratified means that in 297 each fold, a similar percentage of samples of each label as the original data set is preserved, 298 this is particularly useful and recommended when the population to classify are not equally 299 represented in the data set. A new model is then trained k times and at each time a 300 different group is left out to be used later for validation purposes. The whole procedure is 301 repeated for each subset of potential hyperparameters, allowing a better estimation of the model performance and a more robust selection of the hyperparameters. The best model 303 hyperparameters are chosen by computing the average validation accuracy of the k folds. 304 One chose the number of folds with the goal that each fold should remain statistically 305 representative of the whole training data set; usually 5 or 10 folds are sufficient. 306

There are several hyper-parameters that control the architecture of the Random Forest, 307 including: (1) the number of trees; (2) the number of features that are randomly selected 308 at each split to minimize the Gini impurity, from 1 to total the number of features; (3) 309 the minimum number of training samples in a leaf. We do not consider here other param-310 eters such as the maximum depth of a tree or the minimum number of sample to allow 311 a split because they would be in competition with parameter (3). Our most accurate RF 312 model, with hyperparameters chosen after a Stratified 5-Fold Cross-Validation (approxi-313 mately 620,000 samples per fold), is composed of 100 trees, 2 features considered at each 314 split and a minimum number of sample in each leaf of 1. 315

The linear SVC is mainly controlled by the soft margin parameter C. It is a regular-316 ization parameter: a larger C will result in smaller acceptance margins, where the model 317 will allow fewer aftershocks to be on the mainshocks side of the hyper-plane; a lower C 318 will allow larger acceptance margins which in return will generate a simpler and less pre-319 cise hyper-plane, implying a lower training accuracy. Similar to the RF hyperparameters 320 estimation, we performed a 5-Fold cross-validation to look for the best parameter C with 321 the same training and testing data set and found that the value C = 0.1 gives the best 322 validation accuracy. 323

#### 324 **3 Results**

325

#### 3.1 Model Comparison

Fernández-Delgado et al. (2014) found that the RF is the best classification model 326 among many other available algorithms, including Support Vector Machine or Logistic Re-327 gression, but we prefer to conduct our own comparison tailored to the data set at hand. A 328 rather implicit hyper-parameter of this study is the number of feature to use. To investi-329 gate its impact on the accuracy of the different models, we have predicted the label of the 330 earthquakes in 100 test catalogs with a different number of features  $n_f$  for each model. If 331  $n_f = 1$ , only  $\eta_{ij}$ , the nearest-neighbor distance metric is used. For  $n_f = 2$ ,  $R_{ij}$  and  $T_{ij}$  are 332 used. Then if  $n_f = 3$ , the magnitude difference  $d_m$  is added to the list of features, then  $N_p$ 333 for  $n_f = 4$  and finally  $N_c$  for  $n_f = 5$ . We have trained both supervised learning models, 334 the Random Forest and the Support Vector Machine Classifier, beforehand on the training 335

data set with the corresponding number of features and each time their hyperparameters
 have been determined using a 5-Fold Cross-Validation.

We also compared the four models to the classical way of declustering with the NND, 338 which is to use a threshold  $\eta_0$  to split the nearest-neighbor distance distribution  $\eta_{ij}$  in 339 two. We estimated this threshold by rounding the membership probabilities given by the 340 GMM only using  $\eta_{ij}$  distribution.  $\eta_0$  is often defined as the local minimum or in its vicinity 341 between the two lobes of the  $\eta_{ij}$  distribution if the two populations are easily discernible 342 (e.g. Zaliapin & Ben-Zion, 2013; Gu et al., 2013; Davidsen et al., 2015; Shebalin & Narteau, 343 2017; Peresan & Gentili, 2018; Aden-Antoniów et al., 2020). To account for the overlap of 344 the two populations using the GMM, we drew a random number between 0 and 1 for each 345 event to test and if the probability to be an aftershocks is greater than this random number, 346 the event is labelled as aftershock; if not we label it as a mainshock. 347

An important consideration is most machine learning models benefit from a particular 348 pre-processing of the data. For most algorithms dealing with several features, such as the 349 K-means, GMM or the SVC, the different features need to be scaled or standardized, for 350 example by removing the mean and dividing by the standard deviation. Because most 351 models treat the features as an ensemble, they tend to weigh higher numerically large 352 features, regardless of the unit of the features, which is clearly problematic. Accordingly, 353 we have considered the logarithm of  $\eta_{ij}$ ,  $R_{ij}$  and  $T_{ij}$  because they span over a large range; 354 for example  $\eta_{ij}$  distribution is comprised between  $10^{-15}$  and  $10^3$  (Figure 1c and d); this is 355 sometimes called a kernel trick. We emphasize though that we did not have to do this for 356 the RF model. Because the RF randomly selects a single feature at each split it spares us 357 any extra step of scaling because the features are not directly compared to each other. The 358 comparison of the different models according to the number of features  $n_f$  on the testing 359 catalogs is shown in Figure 2a. For each number of features  $n_f$ , each testing catalog has 360 been declustered separately and we estimated the accuracy of the model's prediction. We 361 can then look at the testing accuracy confidence interval (32%-68%) for each model with 362 respect to the median (50%) and the mean (Figure 2a). 363

The classical approach consisting of splitting the  $\eta_{ij}$  into two groups with the threshold 364  $\eta_0$  performs well, with an average testing accuracy of 88%, in comparison to the unsupervised 365 learning models, except when 2 features are considered where the GMM performs slightly 366 better. We note that this is not too surprising as the large majority of events are easily 367 classified, such as aftershocks with small values of  $\eta_{ij}$ . The "tricky" part that is key is to 368 distinguish between aftershocks and mainshocks at the border between the two populations. 369 This is where we observe that the supervised learning models have an advantage over the 370 other approaches. Introducing this non-linear transformation of the nearest-neighbor dis-371 tance metrics clearly benefits the Linear Support Vector Machine Classifier, as it is the best 372 model when using one or two features with a maximum testing accuracy greater than 90%. 373 The additional features  $d_m$ ,  $N_p$  and  $N_c$  seem to only benefit the RF model which reaches 374 a maximum testing accuracy around 92% using all five features. We look into the feature 375 importance in the RF model as it is possible to measure which feature as been used the 376 most to do a split. We can have access to the uncertainty of the importance by considering 377 each tree or estimator separately. We can see Figure 2b that  $R_{ij}$  and  $T_{ij}$  constitute more 378 than 70% of the model. Each feature is important in this model as even  $N_c$ , which as a very 379 limited importance, is essential to explain more than 95% of the model. 380

Figure 3 shows the difference between the different models for the declustering of a 381 testing catalog as an example which can give us a lot of information about the differences 382 between the models (See other examples in Figures S2, S3 and S4). We selected the number 383 of features that gave the best testing accuracy for each model (Figure 2a): the best model 384 of GMM uses only  $T_{ij}$  and  $R_{ij}$  while the KM, the SVC and the RF show better results 385 using all the features. In Figure 3a, we can identify where the misclassified events are 386 located in the 2D rescaled-feature space. Not surprisingly, the classical approach shows a 387 clear boundary between the two populations of earthquakes and the resulting misclassified 388

events distribution correspond directly to the overlap between the two distributions. The 389 simple threshold on  $\eta_{ij}$  still was able to correctly classify almost 90% of the catalog that was 390 away from this boundary. For this catalog, the KM is the least accurate model with around 391 82% of accuracy; it has clearly missed large chunks of both aftershocks and mainshocks. 392 The GMM, on the contrary, identifies the region of overlap between the aftershock and the 393 mainshock distributions, however the accuracy is similar to that of the classical approach. 394 This means that if the GMM can handle the overlap as expected, it does not improve 395 significantly the quality of the overall prediction of the catalog. Both supervised learning 396 models show an improvement over the classical method and unsupervised learning models 397 with accuracies of around 93% and 98% respectively for the SVC and the RF. Because 398 we used a linear kernel for the SVC, it is not surprising to see we have a clear boundary 399 between the mainshock and aftershock distributions. It is clearly visible in the rescaled-400 feature space, especially because the additional features don't seem to improve the model 401 prediction capability as seen in Figure 2a. The RF, while greatly improving the accuracy 402 of this catalog, reduces the overlap uncertainties range and does not show a clear boundary 403 between the two populations. It in fact draws a region where it is difficult to access the 404 true nature of an earthquake, whether it is a long-term aftershock not directly "linked" to 405 its mainshock or a mainshock linked to a spatially distant earthquake but with a very small 406 separation in time. If we compare the  $\eta_{ij}$  distributions with the Weibull function of the 407 true mainshock distribution, we can also see where there are differences between the models (Figure 3b). The confusion matrices (Figure 3c) allow us to see that the RF model is as 409 good at labelling mainshocks as it is at labelling aftershocks, while the other models have 410 a tendency of favoring one population over the other. Only based on synthetic catalogs 411 built using an ETAS model, the RF seems to be a promising model to use for a real case 412 application as it is flexible, does not require a scaling of the features and is able to reproduce 413 with high fidelity the mainshock distribution. 414

415 416

#### 3.2 Application to Southern-California and New Zealand of the Random Forest Model

We have shown in the previous section that the Random Forest is our ideal machine 417 learning model for the declustering of seismic catalog, whose hyperparameters have been 418 optimized with a k-Fold Cross-Validation. We have also shown that the accuracy of the 419 model's predictions increases as we include additional data features: the rescaled distance 420  $R_{ij}$ , the rescaled time  $T_{ij}$ , the difference in magnitude  $d_m$  between the event considered and 421 its nearest neighbor,  $N_p$  and  $N_c$ . Here, we present the declustering with our RF model of 422 two actual, extensive seismic catalogs shown in Figure 1: the relocated Southern Californian 423 earthquake catalog (Hauksson et al., 2012) and the GeoNet catalog of New Zealand. 424

For Southern California, we have used the SCEDC catalog from 1981 to 2019 (Hauksson 425 et al., 2012), relocated with GrowClust (Trugman & Shearer, 2017). We estimated the cat-426 alog completeness at magnitude 2.3 with a b-value of 1.03 following the maximum curvature 427 method (Wiemer & Wyss, 2000) and maximum-likelihood method (Aki, 1965). The total 428 number of events whose magnitude is greater than 2.3 is 72,911. We estimated the fractal 429 dimension of this catalog to be approximately 1.6 using the box counting method, similar 430 to our synthetic catalogs and to what was used by Corral (2003), Zaliapin and Ben-Zion 431 (2013) and Moradpour et al. (2014). The result of declustering with the Random Forest 432 model is shown Figure 4. As expected from our model testing, the RF model is capable 433 of resolving an overlap of the aftershock and the mainshock distributions (Figure 4a and 434 b). Both the 2D distribution Rij- $T_{ij}$  and the distribution of the NND  $\eta$  (Figure 4a and 435 b) show that the separation of the aftershocks from the background seismicity is similar to 436 437 what we observed with our synthetic catalogs. We minimized the L2-norm to fit both the aftershocks and mainshocks distributions to two Weibull functions as suggested by Bayliss 438 et al. (2019). When we look at the prediction of the model and how the probability of being 439 an aftershock is distributed, we observe how the model struggled to predict the label of 440 the events. We show in Figure 4b the distribution of the probability to be an aftershock. 441

We note that it is more difficult for the model to label an event as a mainshock, because 442 the probability distribution is relatively flat from 0 to 50%; the model is more "decisive" 443 in determining an aftershock, given that more aftershocks are clearly labeled as aftershocks 444 (probabilities greater than  $\sim 75\%$ ). By looking at the classification of the catalog through 445 time, Figure 4c, we clearly identify the largest aftershock sequences associated with major 446 regional earthquakes (vertical black lines). We can also see how the mainshock seismicity 447 rate is changing through the years, although these variations could be related to changes in 448 network configuration or detection methods. Overall though, this background rate remains 449 remarkably constant. This suggests that our model is able to effectively distinguish between 450 real mainshocks and aftershocks. 451

We took a closer look at the aftershock sequences of four major earthquakes that oc-452 curred in Southern California: the M7.3 Landers earthquake in 1992 (Figure 4e), the M7.1 453 Hector Mine earthquake of 1999 (Figure 4f), the M7.2 Baja California in 2010 (Figure 4g), 454 and finally the recent 2019 M7.1 Ridgecrest earthquake (Figure 4h). We analyzed the seis-455 micity 180 days prior to and following these earthquakes, separating the seismicity into 456 inside and outside groups, depending on whether each earthquake was further or closer 457 than 350 km from the mainshock epicenter. We observe that we are able to well isolate the 458 aftershock sequence that resembles a characteristic Omori-type sequence, leaving a fairly 459 constant background sequence. We remark that the background or the mainshock seismic-460 ity rate slightly decreases in the days following the Landers, the Hector Mine and Ridgecrest 461 earthquakes (Figure 4f and h) only to recover after a week or two. We speculate that this is 462 the case is due to the fact that network analysts focus on the aftershock sequence for several 463 days and weeks following the mainshock, creating this artificial rate decrease in background 464 seismicity. Inside the 350km perimeter, the decrease of the background seismicity rate could also be explained by the fact that there are few mainshocks distinguishable (with the data 466 features we have used) from the ongoing an aftershock sequence. 467

We also applied our random forest model to the GeoNet New Zealand catalog; we 468 used the seismicity recorded by GeoNet from 2010 to 2020. Similar to Southern California, 469 we estimated the completeness magnitude to be 2.6 with the maximum curvature method 470 (Wiemer & Wyss, 2000), associated with a b value of 0.89. The resulting catalog is com-471 posed of 64,200 earthquakes. This catalog comprises a various type of tectonic and volcanic 472 settings originating from New Zealand volcanoes, the Hikurangi subduction zone where the 473 474 Pacific plate subducts beneath the North Island (Australian plate), and more shallow activity related to crustal intraplate faulting in the South Island along the Alpine fault and 475 the transpressional transition from the Alpine fault to the Hikurangi subduction zone. We 476 estimated the fractal dimension of the catalog at around 1.9 from the box counting method; this is slightly higher than our estimate for the Southern California catalog, but we suggest 478 it is reasonable given the presence of intermediate depth seismicity related to the Hikurangi 479 subduction zone and the larger study region considered. The results of our declustering 480 are shown Figure 5. Comparably to Southern California, the declustering of the GeoNet 481 catalog is in agreement with our expectations as we were able to fit two Weibull functions 482 to the mainshock and the aftershock  $\eta_{ij}$  distributions (Figure 5b). We also observe a similar 483 distribution of probability to be an aftershock (Figure 5c) with Southern California, where 484 aftershocks are more decisively classified. When looking at the seismicity rates for both mainshocks and aftershocks, we clearly identified the aftershock sequences of large earth-486 quakes which occurred in the last 10 years over New Zealand (Figure 5d). We can see that 487 the seismicity rate was slightly higher between 2010 and 2012, the background rate remains 488 constant afterwards. We focus again on four major aftershock sequences to evaluate the 489 quality of the declustering: the M7.1 Darfield earthquake of 2010 (Figure 5e), the M6.5 490 Cook Strait earthquake which occurred in 2013 (Figure 5f), the M7.1 East Coast event of 491 2016 (Figure 5g) preceding the complex M7.8 Kaikoura earthquake of 2016 (Figure 4h). For 492 all of these sequences, we observe that these major events did not affect the seismicity out-493 side a radius of 350km from the mainshock; away from these sequences the seismic activity 494 was calm as shown by the constant rate of mainshocks and aftershocks. We can also observe 495

that these major events had no incidence on the background seismicity inside their area as it does not show significant deviations from a constant rate.

We can speculate that our results could be different depending on how a catalog was 498 built. If events are detected visually and picked manually by an operator, the magnitude 499 completeness of a catalog could vary significantly in space and time during major after-500 shock sequences as more events would be added to the catalog close to these sequences. 501 By presenting these two applications of our model to real seismicity in very distinct tec-502 tonic contexts, we have demonstrated the model robustness and its increased capability of 503 classifying mainshocks and aftershocks compared to the classical approach of the Nearest-504 Neighbor-Distance approach, especially at the critical boundary between aftershocks and 505 mainshocks in the 2D NND space. 506

#### <sup>507</sup> 4 Discussion and conclusions

In this work we present a new transfer learning approach to solve one of the classical 508 problems in modern seismology, the declustering of aftershocks. As each seismicity cata-509 log is different, we built a generalizeable model that is capable of declustering regardless 510 of tectonic context. We based our model on one of the most used algorithms to decluster 511 earthquake catalogs, the nearest-neighbor distance (NND) approach. This approach relies 512 on the computation of a metric  $\eta_{ij}$  between one event and the events that precede it to 513 identify its nearest neighbor, interpreted here as the mainshock that triggered the event in 514 question. Recent work has focused on thresholding the NND  $\eta_{ij}$  to define the boundary 515 between mainshocks and aftershocks, but we have shown that this neglects in most catalogs 516 a significant overlap of the two populations. We tackled this problem with a data-science 517 and machine learning approach, testing different supervised or unsupervised learning models 518 to identify which model is best suited to decluster. We use features issued from the NND 519 algorithm, differences in magnitude, the number of times the nearest neighbor has been 520 selected as such, and the number of potential offspring. Our model of choice is a Random 521 Forest trained on synthetic catalogs built following an ETAS model with a different set of 522 parameters. The RF is the model with the best testing accuracy on synthetic catalogs. 523 The RF model does not require any preprocessing, standardization or normalization of the 524 features and it is scalable; regardless of the size of the catalog to be declustered, the model 525 classifies each earthquake one by one rather than the distribution of all earthquakes. The 526 hyperparameters of the model have been selected following a 5-Fold Cross-Validation. We 527 applied our declustering model to two real catalogs: the relocated Southern California cata-528 log from 1981 to 2019 and New Zealand GeoNet catalog from 2010 to 2020. We have shown 529 that our model significantly improves the quality of the declustering in comparison of the 530 classical use the NND metrics. Additionally, our approach is simple and fast to train, espe-531 cially in comparison to Monte Carlo Markov Chains inversions that can be computationally 532 demanding (Bayliss et al., 2019) or stochastic declustering (e.g. Zhuang et al., 2002) that 533 does not make use of relative NND metrics and relies on many more parameters. 534

One has to keep in mind that these results reflect the performance of the model trained 535 on a set of synthetic catalogs with a certain duration and fixed ETAS parameters for each 536 catalog and stationary background seismic rates. Hainzl et al. (2016) found that aftershock 537 sequences only last for about 100 days for moderate mainshocks (M = 6) and to a few 538 years for larger events  $(M \ge 7)$ . The estimated duration of an aftershock sequence is 539 impacted by the background activity rate, which can potentially hide the trailing edge of 540 a postseismic sequence as the process rate slows down with time. We have tested our 541 models on shorter catalogs (approximately 500 days long) and noticed that the relative 542 performance remains similar (see Figure S5 and S6). However, models trained on shorter 543 catalogs do not reach similar accuracy when applied to longer catalogs. This is not surprising 544 as long-term aftershocks could eventually be considered as mainshocks which means a larger 545 nearest-neighbor-distance metric with their real mainshocks than what the model has been 546 trained on. Indeed, another earthquake could have occurred closer in time and/or space but 547

eventually not sufficiently close for these long-term aftershocks to be labelled as aftershocks
 anymore. If this problem is difficult to evaluate, better prediction performances are actually
 achieved if we don't force the nearest neighbor of the aftershocks to be their real mainshock,
 thus somewhat challenging the model.

Another aspect to keep in mind is that we used only five features to build the model 552 presented here. It is of course possible to add more features by considering not only the 553 nearest-neighbor but the N-nearest-neighbors, providing information to the model to classify 554 an earthquake based on its relationship with many associated events, not just its single 555 nearest neighbor. We can also consider using the nearest neighbor of the nearest neighbor, 556 creating a family (Zaliapin & Ben-Zion, 2013). This would be a likely extension of our 557 approach, as one aftershock can trigger another, and so on. Finally we can also imagine using 558 N Random Forests to classify N neighbors based on a single set of features. Prelimary tests 559 with our training data set, however, did not improve the accuracy of our model significantly, 560 but this could represent a potential improvement for future work. 561

Regarding the comparison between the Southern California and the New Zealand cat-562 alogs, we could observe a significant difference between the shapes of their  $\eta_{ij}$  distribu-563 tions (Figure 1e and f). The GeoNet catalog exhibits a higher number of mainshocks than 564 aftershocks, respectively 57% against 43% (Figure 5), while the Southern California cata-565 log is 70% aftershocks (Figure 4a). This difference can be explained by comparing their 566 Nearest-Neighbor Distance distributions (see Figures 4b and 5b). Southern California's  $\eta_{ij}$ 567 distribution shows the two expected lobes associated with mainshocks and aftershocks; the 568 mainshock lobe of the  $\eta_{ij}$  distribution is less pronounced for New Zealand. The NZ catalog, 569 while being shorter, has a higher background seismic rate according to our declustering 570 model (see Figures 5d and 5d). We investigated this significant difference by comparing the 571  $\eta_{ij}$  distributions of the mainshock and aftershock seismicity as a function of depth (Figure 572 6a). We observe a rapid decrease in seismicity with depth of the Southern California seis-573 micity that reaches approximately a maximum of 30-40 km depth, while NZ earthquakes can 574 still be observed until 300 km depth. This reflects the difference in tectonics between the two 575 regions and the type of seismicity observed in both catalogs: crustal seismicity in Southern 576 California and a seismicity dominated by subduction earthquakes in New Zealand, whether 577 located on the upper plates, the plate interface or at intermediate depths. We see for both 578 regions that there are more aftershocks than mainshocks at shallow depth, above 25km. To 579 understand how much deep mainshocks account for the observed NND distribution for NZ, 580 we segmented the  $\eta_{ij}$  distribution for different depth ranges (Figure 6b). The seismicity 581 comprised between 0 and 25 km exhibits a bimodal NND distribution and includes most 582 of the aftershocks within the entire catalog. We see that the seismicity deeper than 50 km 583 corresponds mostly to mainshocks. A similar behavior has been observed for Northern Chile 584 by Aden-Antoniów et al. (2020), where no aftershocks were observed at intermediate depth 585  $(z \ge 70 \text{km})$ ; this results in a single-lobe NND distribution. The limited aftershock produc-586 tivity at intermediate depths has been documented for different subduction zones in the past 587 (e.g. Frohlich, 1987; Wiens et al., 1994). Only large intermediate-depths earthquakes seem 588 to be able to trigger significant aftershock sequences but it still depends on the subduction 589 zone (Frohlich, 1989). Temperature, geometry and hydration of the subducting slab seem to 590 be the main parameters to control the aftershock productivity at these depths (e.g. Wiens, 591 2001; Ye et al., 2017; Cabrera et al., 2021). 592

To conclude, it is important to keep in mind that this model has been trained on 593 synthetic data in an exercise of transfer learning. The synthetic data is controlled by range 594 of reasonable ETAS parameters (see Table 1), but neither do these parameters or even the 595 ETAS model reflect the entire spectrum of possibilities within earthquake catalogs. If this 596 model might not be applicable to very particular seismicity i.e. mining related seismicity or 597 seismicity with high earthquake rate variations, we suggest that it will be useful in a large 598 majority of study regions. This work also provides a route to go beyond the classical use of 599 the NND as the framework would be easy to reproduce. 600

#### 601 Acknowledgments

F.A-A would like to thank Sid Kothari from the University of Western Ontario and Eric Beaucé from the Massachusetts Institute of Technology for interesting and constructive discussions and comments. We also thank everyone who has contributed to the SCEDC and GeoNet catalogs. We would like to thank the editor X and two anonymous reviewers for their comments and suggestions that greatly improve the manuscript. The codes allowing the use of our best Random Forest model for the declustering of seismic catalog created for this study are available here.

#### 609 References

626

627

644

645

- Aden-Antoniów, F., Satriano, C., Bernard, P., Poiata, N., Aissaoui, E.-M., Vilotte, J. P., & Frank, W. (2020). Statistical analysis of the preparatory phase of the m w
   8.1 iquique earthquake, chile. Journal of Geophysical Research: Solid Earth, 125(6),
   e2019JB019337.
- Aki, K. (1965). Maximum likelihood estimate of b in the formula log n= a-bm and its confidence limits. *Bull. Earthq. Res. Inst., Tokyo Univ., 43*, 237–239.
- Asim, K., Martínez-Álvarez, F., Basit, A., & Iqbal, T. (2017). Earthquake magnitude
   prediction in hindukush region using machine learning techniques. Natural Hazards, 85(1), 471–486.
- Azak, T. E., Kalafat, D., Şeşetyan, K., & Demircioğlu, M. (2018). Effects of seismic declustering on seismic hazard assessment: a sensitivity study using the turkish earthquake catalogue. Bulletin of Earthquake Engineering, 16(8), 3339–3366.
- Bachmann, C. E., Wiemer, S., Woessner, J., & Hainzl, S. (2011). Statistical analysis of
   the induced basel 2006 earthquake sequence: introducing a probability-based moni toring approach for enhanced geothermal systems. *Geophysical Journal International*,
   *186*(2), 793–807.
  - Baiesi, M., & Paczuski, M. (2004). Scale-free networks of earthquakes and aftershocks. *Physical review E*, 69(6), 066106.
- Bayliss, K., Naylor, M., & Main, I. G. (2019). Probabilistic identification of earthquake
   clusters using rescaled nearest neighbour distance networks. *Geophysical Journal International*, 217(1), 487–503.
- Cabrera, L., Ruiz, S., Poli, P., Contreras-Reyes, E., Osses, A., & Mancini, R. (2021). North ern chile intermediate-depth earthquakes controlled by plate hydration. *Geophysical Journal International*, 226(1), 78–90.
- Corral, A. (2003). Local distributions and rate fluctuations in a unified scaling law for earthquakes. *Physical Review E*, 68(3), 035102.
- Davidsen, J., Gu, C., & Baiesi, M. (2015). Generalized omori-utsu law for aftershock
   sequences in southern california. *Geophysical Journal International*, 201(2), 965–978.
- Dempsey, D., Cronin, S. J., Mei, S., & Kempa-Liehr, A. W. (2020). Automatic precursor recognition and real-time forecasting of sudden explosive volcanic eruptions at whakaari, new zealand. *Nature communications*, 11(1), 1–8.
- Fernández-Delgado, M., Cernadas, E., Barro, S., & Amorim, D. (2014). Do we need
   hundreds of classifiers to solve real world classification problems? The journal of
   machine learning research, 15(1), 3133–3181.
  - Frohlich, C. (1987). Aftershocks and temporal clustering of deep earthquakes. Journal of Geophysical Research: Solid Earth, 92(B13), 13944–13956.
- <sup>646</sup> Frohlich, C. (1989). The nature of deep-focus earthquakes. Annual Review of Earth and <sup>647</sup> Planetary Sciences, 17(1), 227–254.
- Galina, N., Bykova, V., Vakarchuk, R., & Tatevosian, R. (2019). Effect of earthquake catalog
   declustering on seismic hazard assessment. Seismic Instruments, 55(1), 59–69.
- Gardner, J., & Knopoff, L. (1974). Is the sequence of earthquakes in southern california, with
   aftershocks removed, poissonian? Bulletin of the Seismological Society of America,
   64 (5), 1363–1367.

- González, J., Yu, W., & Telesca, L. (2019). Earthquake magnitude prediction using recurrent
   neural networks. In *Multidisciplinary digital publishing institute proceedings* (Vol. 24,
   p. 22).
- Gu, C., Schumann, A. Y., Baiesi, M., & Davidsen, J. (2013). Triggering cascades and
   statistical properties of aftershocks. *Journal of Geophysical Research: Solid Earth*,
   118(8), 4278–4295.
- Gutenberg, B., & Richter, C. (1955). Magnitude and energy of earthquakes. *Nature*, 176(4486), 795–795.
- Hainzl, S., Christophersen, A., Rhoades, D., & Harte, D. (2016). Statistical estimation
   of the duration of aftershock sequences. *Geophysical Journal International*, 205(2),
   1180–1189.
- Hainzl, S., & Ogata, Y. (2005). Detecting fluid signals in seismicity data through statistical
   earthquake modeling. Journal of Geophysical Research: Solid Earth, 110(B5).
- Hauksson, E., Yang, W., & Shearer, P. M. (2012). Waveform relocated earthquake catalog
   for southern california (1981 to june 2011). Bulletin of the Seismological Society of
   America, 102(5), 2239–2244.
- Hibert, C., Provost, F., Malet, J.-P., Maggi, A., Stumpf, A., & Ferrazzini, V. (2017).
  Automatic identification of rockfalls and volcano-tectonic earthquakes at the piton
  de la fournaise volcano using a random forest algorithm. Journal of Volcanology and
  Geothermal Research, 340, 130–142.

673

674

675

676

677

- Hoque, A., Raj, J., Saha, A., & Bhattacharya, P. (2020). Earthquake magnitude prediction using machine learning technique. In *International conference on computational intelligence, security and internet of things* (pp. 37–53).
- Kong, Q., Allen, R. M., & Schreier, L. (2016). Myshake: Initial observations from a global smartphone seismic network. *Geophysical Research Letters*, 43(18), 9588–9594.
- Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., & Gerstoft, P. (2019).
   Machine learning in seismology: Turning data into insights. *Seismological Research Letters*, 90(1), 3–14.
- Kothari, S., Shcherbakov, R., & Atkinson, G. (2020). Statistical modeling and characteri zation of induced seismicity within the western canada sedimentary basin. Journal of
   *Geophysical Research: Solid Earth*, 125(12), e2020JB020606.
- Lara-Cueva, R., Carrera, E. V., Morejon, J. F., & Benitez, D. (2016). Comparative analysis of automated classifiers applied to volcano event identification. In 2016 ieee colombian conference on communications and computing (colcom) (pp. 1–6).
- Li, Z., Meier, M.-A., Hauksson, E., Zhan, Z., & Andrews, J. (2018). Machine learning seismic wave discrimination: Application to earthquake early warning. *Geophysical Research Letters*, 45(10), 4773–4779.
- Maggi, A., Ferrazzini, V., Hibert, C., Beauducel, F., Boissier, P., & Amemoutou, A. (2017).
   Implementation of a multistation approach for automated event classification at piton
   de la fournaise volcano. Seismological Research Letters, 88(3), 878–891.
- Marsan, D., Bouchon, M., Gardonio, B., Perfettini, H., Socquet, A., & Enescu, B. (2017).
   Change in seismicity along the japan trench, 1990–2011, and its relationship with seismic coupling. *Journal of Geophysical Research: Solid Earth*, 122(6), 4645–4659.
- Marsan, D., Prono, E., & Helmstetter, A. (2013). Monitoring aseismic forcing in fault zones
   using earthquake time series. *Bull. Seismol. Soc. Am.*. doi: 10.1785/0120110304
- Marsan, D., Reverso, T., Helmstetter, A., & Enescu, B. (2013). Slow slip and aseismic deformation episodes associated with the subducting Pacific plate offshore Japan, revealed by changes in seismicity. J. Geophys. Res. E Planets. doi: 10.1002/jgrb.50323
- Marzocchi, W., & Taroni, M. (2014). Some thoughts on declustering in probabilistic seismic hazard analysis. Bulletin of the Seismological Society of America, 104 (4), 1838–1845.
- Moradpour, J., Hainzl, S., & Davidsen, J. (2014). Nontrivial decay of aftershock density
   with distance in southern california. Journal of Geophysical Research: Solid Earth,
   119(7), 5518–5535.
- Mosteller, F., & Tukey, J. W. (1968). Data analysis, including statistics. *Handbook of social* psychology, 2, 80–203.

- Ogata, Y. (1988). Statistical models for earthquake occurrences and residual analysis for point processes. *Journal of the American Statistical association*, 83(401), 9–27.
- Ogata, Y. (1998). Space-time point-process models for earthquake occurrences. Annals of the Institute of Statistical Mathematics, 50(2), 379–402.
- <sup>712</sup> Omori, F. (1894). Investigation of aftershocks. *Rep. Earthquake Inv. Comm*, 2, 103–139.
- Peresan, A., & Gentili, S. (2018). Seismic clusters analysis in northeastern italy by the
   nearest-neighbor approach. *Physics of the Earth and Planetary Interiors*, 274, 87–104.
- Perol, T., Gharbi, M., & Denolle, M. (2018). Convolutional neural network for earthquake
   detection and location. *Science Advances*, 4(2), e1700578.
- Picozzi, M., & Iaccarino, A. G. (2021). Forecasting the preparatory phase of induced earthquakes by recurrent neural network. *Forecasting*, 3(1), 17–36.
- Provost, F., Hibert, C., & Malet, J.-P. (2017). Automatic classification of endogenous
   landslide seismicity using the random forest supervised classifier. *Geophysical Research Letters*, 44(1), 113–120.
- Reasenberg, P. (1985). Second-order moment of central california seismicity, 1969–1982.
   Journal of Geophysical Research: Solid Earth, 90(B7), 5479–5495.
- Reverso, T., Marsan, D., & Helmstetter, A. (2015). Detection and characterization of
   transient forcing episodes affecting earthquake activity in the aleutian arc system.
   *Earth and Planetary Science Letters*, 412, 25–34.
- Reverso, T., Marsan, D., Helmstetter, A., & Enescu, B. (2016). Background seismicity in
   boso peninsula, japan: Long-term acceleration, and relationship with slow slip events.
   *Geophysical Research Letters*, 43(11), 5671–5679.
- Ross, Z. E., Meier, M.-A., Hauksson, E., & Heaton, T. H. (2018). Generalized seismic phase detection with deep learning. Bulletin of the Seismological Society of America, 108(5A), 2894–2901.
- Rouet-Leduc, B., Hulbert, C., Lubbers, N., Barros, K., Humphreys, C. J., & Johnson,
   P. A. (2017). Machine learning predicts laboratory earthquakes. *Geophysical Research Letters*, 44 (18), 9276–9282.
- Rubin, M. J., Camp, T., Van Herwijnen, A., & Schweizer, J. (2012). Automatically detecting
   avalanche events in passive seismic data. In 2012 11th international conference on
   machine learning and applications (Vol. 1, pp. 13–20).
- Seydoux, L., Balestriero, R., Poli, P., De Hoop, M., Campillo, M., & Baraniuk, R. (2020).
   Clustering earthquake signals and background noises in continuous seismic data with
   unsupervised deep learning. *Nature communications*, 11(1), 1–12.
- Shebalin, P., & Narteau, C. (2017). Depth dependent stress revealed by aftershocks. *Nature communications*, 8(1), 1–8.
- Shi, P., Seydoux, L., & Poli, P. (2021). Unsupervised learning of seismic wavefield features:
   clustering continuous array seismic data during the 2009 l'aquila earthquake. Journal
   of Geophysical Research: Solid Earth, 126(1), e2020JB020506.
- Steinmann, R., Seydoux, L., Beaucé, E., & Campillo, M. (2021). Hierarchical exploration
   of continuous seismograms with unsupervised learning. *Earth and Space Science Open* Archive ESSOAr.
- Taroni, M., & Akinci, A. (2021). Good practices in psha: declustering, b-value estima tion, foreshocks and aftershocks inclusion; a case study in italy. *Geophysical Journal International*, 224(2), 1174–1187.
- Teng, G., & Baker, J. W. (2019). Seismicity declustering and hazard analysis of the
   oklahoma-kansas region. Bulletin of the Seismological Society of America, 109(6),
   2356-2366.
- Thomas, A., Inbal, A., Searcy, J., Shelly, D., & Bürgmann, R. (n.d.). Identification of
   low-frequency earthquakes on the san andreas fault with deep learning. *Geophysical Research Letters*, e2021GL093157.
- Trugman, D. T., & Shearer, P. M. (2017). Growclust: A hierarchical clustering algorithm
   for relative earthquake relocation, with application to the spanish springs and sheldon,
   nevada, earthquake sequences. Seismological Research Letters, 88(2A), 379–391.

parameter	$\min$	max
$\overline{m_c}$	2.0	3.0
b	0.8	1.0
$< \tau >$	1	3
c	$10^{-8}$	1
p	1.0	1.3
K	0.12	0.18
a	0.8	1.05
$\gamma$	1.5	2.5

 Table 1. Parameter ranges for the modeling of synthetic catalogs. These values were used as boundaries to draw, for each catalog, the ETAS parameter from a uniform distribution.

- Utsu, T. (1957). Magnitudes of earthquakes and occurrence of their aftershocks. Zisin, Ser. 2, 10, 35–45.
- Utsu, T. (1961). A statistical study on the occurrence of aftershocks. *Geophys. Mag.*, 30, 521–605.
- <sup>767</sup> Utsu, T., Ogata, Y., et al. (1995). The centenary of the omori formula for a decay law of <sup>768</sup> aftershock activity. *Journal of Physics of the Earth*, 43(1), 1–33.
- Utsu, T., & Seki, A. (1955). A relation between the area of the aftershock region and the radius of the sensibility circle. Zisin,  $\Im(34)$ , 1955.
- Van Stiphout, T., Zhuang, J., Marsan, D., Stiphout, V., Zhuang, J., & Marsan, D. (2012).
   Theme V-Models and Techniques for Analyzing Seismicity Seismicity Declustering.
   CORSSA. doi: 10.5078/corssa
- Wiemer, S., & Wyss, M. (2000). Minimum magnitude of completeness in earthquake
   catalogs: Examples from alaska, the western united states, and japan. Bulletin of the
   Seismological Society of America, 90(4), 859–869.
- Wiens, D. A. (2001). Seismological constraints on the mechanism of deep earthquakes:
   Temperature dependence of deep earthquake source properties. *Physics of the Earth* and Planetary Interiors, 127(1-4), 145–163.
- Wiens, D. A., McGuire, J. J., Shore, P. J., Bevis, M. G., Draunidalo, K., Prasad, G., & Helu, S. P. (1994). A deep earthquake aftershock sequence and implications for the rupture mechanism of deep earthquakes. *Nature*, 372(6506), 540–543.

783

784

785

786

787

788

789

- Ye, L., Lay, T., Bai, Y., Cheung, K. F., & Kanamori, H. (2017). The 2017 mw 8.2 chiapas, mexico, earthquake: Energetic slab detachment. *Geophysical Research Letters*, 44 (23).
- Yuan, B., Tan, Y. J., Mudunuru, M. K., Marcillo, O. E., Delorey, A. A., Roberts, P. M., ... others (2019). Using machine learning to discern eruption in noisy environments: A case study using co2-driven cold-water geyser in chimayó, new mexico. Seismological Research Letters, 90(2A), 591–603.
- Zaliapin, I., & Ben-Zion, Y. (2013). Earthquake clusters in southern California I: Identification and stability. J. Geophys. Res. Solid Earth. doi: 10.1002/jgrb.50179
- Zaliapin, I., & Ben-Zion, Y. (2020). Earthquake declustering using the nearest-neighbor
   approach in space-time-magnitude domain. Journal of Geophysical Research: Solid
   Earth, 125(4), e2018JB017120.
- Zaliapin, I., Gabrielov, A., Keilis-Borok, V., & Wong, H. (2008). Clustering analysis of seismicity and aftershock identification. *Physical review letters*, 101(1), 018501.
- Zhu, W., & Beroza, G. C. (2019). Phasenet: a deep-neural-network-based seismic arrivaltime picking method. *Geophysical Journal International*, 216(1), 261–273.
- Zhuang, J., Ogata, Y., & Vere-Jones, D. (2002). Stochastic declustering of space-time
   earthquake occurrences. Journal of the American Statistical Association, 97(458),
   369–380.



Figure 1. Examples of real catalogs. The seismic catalogs used in this study are the Southern California catalog (a) from 1981 to 2019 (Hauksson et al., 2012) and New Zealand (b) from 2010 to 2020 from GeoNet. (c) and (d) show the respective 2D distribution of the Nearest-Neighbor distance (NND) as a function of the rescaled time  $T_{ij}$  and distance  $R_{ij}$  while (e) and (f) show their 1-D NND metric distribution. The dashed orange line represents the threshold that would be used for a classical approach of the NND for the declustering of earthquake catalogs.



**Figure 2.** Machine Learning algorithms comparison and feature importance for the Random Forest model. (a) We declustered 100 synthetic catalogs using different number of features with different models and estimated the testing accuracies. The shaded area represent the 38% - 68% inter-quartiles range of the accuracy distributions, the solid line represent the median of this distributions and the dashed line is their mean. (b) We estimated the features importance in the Random Forest using all 5 features. The black lines show the relative importance of these features in a descending order for the model while the black line is the cumulative. The vertical bars shows the standard variation of the features importance. This was computed from the distribution of each feature importance across all trees of the forest.



Figure 3. Example of the declustering of a synthetic catalog with different machine learning algorithms. (a) For the different algorithms, we show the distribution of the nearest neighbor rescaled distance  $R_{ij}$  and time  $T_{ij}$  where the color corresponds to the mis-labelled events whether there are originally mainshocks (purple) or aftershocks (blue). The accuracy of the model is displayed in the title of each subplot. (b)  $\eta_{ij}$  distribution obtained by applying each model. The orange corresponds to the inferred mainshock population while yellow is for the aftershock population. The black curve shows the real background distribution. (c) shows the different confusion matrices where we show the number of correctly identified samples in the colored cells and mis-labelled events with their corresponding proportion regarding their population in the white cells.



**Figure 4.** Our best Random Forest model for the declustering of earthquake catalog applied to Southern California seismicity between 1981 and 2019. (a) shows the nearest-neighbor rescaled distance  $R_{ij}$  and time  $T_{ij}$  distribution. The color indicate if the events have been classified as mainshocks (orange) or aftershocks (yellow). The dashed grey line shows where the classical approach of the NND whould have separated the two populations. (b) shows the obtained the two stacked  $\eta_{ij}$  distributions. The dashed black lines corresponds to the fit of a Weibull function to both distributions while the black line shows the resulting sum and fit to the overall  $\eta_{ij}$  distribution. (c) displays the probability density function of the probability to be an aftershock, obtained from the data. (d) exhibits the stacked histograms on the earthquake count per window of 30 days. The black vertical lines mark the date of large earthquake that occurred during the period covered by the catalog. (e) - (h) show the normalized cumulative number of mainshocks (dashed curves) and aftershocks (plain curves) for 4 remarkable sequences of aftershocks. The color indicates whether these events are located at distance larger or lower than 350km.



Figure 5. Our best Random Forest model for the declustering of earthquake catalog applied to New Zealand seismicity between 2010 and 2020. Similar to Figure 4.



Figure 6. Mainshock/Aftershock depth distributions in relation to the Nearest-Neighbor Distance distribution for New Zealand. a) The histograms represent the distribution of the depth of the background seismicity and aftershocks for the Southern California and New Zealand. b) The NND  $\eta_{ij}$  distribution for the New Zealand catalog is represented by the colored histograms for different depth range. The green histogram represents the complete NND  $\eta_{ij}$  distribution for Southern California for comparison purposes.

# Supporting Information for "Transfer Learning to Build a Scalable Model for the Declustering of Earthquake Catalogs"

F. Aden-Antoniów<sup>1</sup>, W. B. Frank<sup>1</sup>, L. Seydoux<sup>2</sup>

<sup>1</sup>Department of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA

 $^{2}$ Équipe Ondes et Structures, ISTerre, Université Grenoble-Alpes, UMR 5375, Gières, France

## Contents of this file

1. Figures S1 to S7

## Figures complementing information reported in the main text



:

**Figure S1.** Example of a short synthetic catalog generated for this study. (left) Spatial distribution of the synthetic mainshocks (orange) and aftershocks (yellow). The black contours represent the 2D probability function use to generate the location of these earthquakes. (right) Cumulative number of the mainshocks and the aftershocks.



Figure S2. Example of the declustering of a synthetic catalog with different machine learning algorithms. (a) For the different algorithms, we show the distribution of the nearest neighbor rescaled distance  $R_{ij}$  and time  $T_{ij}$  where the color corresponds to the mis-labelled events whether there are originally mainshocks (purple) or aftershocks (blue). The accuracy of the model is displayed in the title of each subplot. (b)  $\eta_{ij}$  distribution obtained by applying each model. The orange corresponds to the inferred mainshock population while yellow is for the aftershock population. The black curve shows the real background distribution. (c) shows the different confusion matrices where we show the number of correctly identified samples in the colored cells and mis-labelled events with their corresponding proportion regarding their population in the white cells.

September 16, 2021, 1:56am





:

Figure S3. Other example of the declustering of a synthetic catalog with different machine learning algorithms. Same as S2.



**Figure S4.** Other example of the declustering of a synthetic catalog with different machine learning algorithms. Same as S2.



**Figure S5.** Machine Learning algorithms comparison with different catalog duration for training and testing phases. (a) We declustered 100 synthetic catalogs of an approximated duration of 1.5 years using different number of features with different machine learning models (KM: K-Means, GMM: Gaussian Mixture Model, LSVC: Linear Support Vector Classifier, RF: Random Forest) and estimated the testing accuracies. The Linear Support Vector Classifier and the Random Forest have been trained on synthetic catalogs with a duration of approximately 10 years. The shaded area represent the 38% - 68% inter-quartiles range of the accuracy distributions, the solid line represent the median of this distributions and the dashed line is their mean.

September 16, 2021, 1:56am



Figure S6. Example of the declustering of a short synthetic catalog (1.5 year) with different machine learning algorithms. (a) For the different algorithms, we show the distribution of the nearest neighbor rescaled distance  $R_{ij}$  and time  $T_{ij}$  where the color corresponds to the mislabelled events whether there are originally mainshocks (purple) or aftershocks (blue). The accuracy of the model is displayed in the title of each subplot. (b)  $\eta_{ij}$  distribution obtained by applying each model. The orange corresponds to the inferred mainshock population while yellow is for the aftershock population. The black curve shows the real background distribution. (c) shows the different confusion matrices where we show the number of correctly identified samples in the colored cells and mis-labelled events with their corresponding proportion regarding their population in the white cells.

September 16, 2021, 1:56am