

Improving the Generation of Earthquake Risk Models Using Evolutionary Algorithms tempered by Domain Knowledge

Yuri Lavinas
University of Brasilia
Department department
yclavinas@gmail.com

Marcelo Ladeira
University of Brasilia
Department department
mladeira@unb.br

Claus Aranha
University of Tsukuba
Graduate School of SIE
caranha@cs.tsukuba.ac.jp

Abstract—Earthquake Risk Models describe the risk of occurrence of seismic events on a given area based on information such as past earthquakes in nearby regions and the seismic properties of the area under study. These models can be used to help to better understand earthquakes, their patterns and their mechanisms.

In previous work, we showed that Genetic Algorithms (GA) could generate risk models with the same degree of precision as the Relative Intensity (RI) method, which is considered a benchmark for this problem. However, a few shortcomings were also defined in that approach: (1) The representation of the model in the Genetic Algorithm was too sparse, (2) Domain knowledge was not used to create the model, and (3) The relationship between foreshocks and aftershocks were not taken into account.

In this work, we try to address these three concerns. We propose a new representation of a seismic risk model to be used as the genome of the Genetic Algorithm. We introduces a hybrid model that incorporates seismic theories about earthquake distribution (such as the Omori-Utsu formula). And we introduce two methods to filter the earthquake catalog in order to remove earthquakes that are likely to be aftershocks before generating the risk model.

We examine each of these changes through simulations using the catalog of Japanese earthquakes between 2000 and 2010. Our results show that (XXX)(YYY). These results allow us to draw recommendations for future development in this field.

1. Introduction

Earthquakes may cause soil rupture or movement, tsunamis and more. They may cause great losses and that can be explicit by some examples, such as the earthquakes in Tōhoku (2011), a 9.0 M_w earthquake [?], and considered the most powerful earthquake to ever hit Japan, and in Nepal(2015) with moment magnitude (M_w) of 7.8 and considered the largest since 1934 [?]. To be able to minimize the consequences of these events, we look to create forecast earthquake occurrences models. Hence the characteristics that most influence the earthquakes events may vary both in time and place, these methods should be to adapt their

behavior to be able to forecast earthquakes events which reflects well the reality.

This project aims to improve the GAModel [?], a statistical method of analysis of earthquakes risk using the Genetic Algorithm technique (GA). Two ideas are proposed for this. The first one, is to change the candidate solution representation. By that, we objective to make the GAModel more specialized, focusing only on areas on which earthquakes happened already in a near past. It is expect that it would lead to a faster convergence, once the amount of parameters is smaller and consequently, the search space gets smaller.

Formulated on this idea, we propose the ReducedGAModel. Its genome only has information of areas that already had occurrences in the past. This helps the method to converge gets faster, by minimizing the number of parameters the method has to deal with.

The other idea is based on the assumption that earthquakes cluster in both space and time, and we want incorporate in the Genetic Algorithm technique (GA) some geophysical knowledge. It is a hybridisation of the models generated with GA some empirical laws, such as the modified Omori law. First, the background intensity (the independent earthquakes or mainshocks), which is a function of the space, is forecast using the GA. Then, we use some empirical laws to obtain the dependent earthquakes (aftershocks) for a specific time interval.

The Emp-GAModel is a method proposed that incorporates some geophysical knowledge. It is a hybridization of the models generated by the GAModel with the these empirical laws, see Section 3.

Finally, there is the Emp-ReducedGAModel. This method is a combination of the two ideas. Therefore, it also performs a hybridization of models with the group of empirical law. Though, for this method, the models are generated by the ReducedGAModel.

The forecast models produced by those methods and the ones produced by the GAModel were all analyzed by their log-likelihood values calculated as suggested by the Regional Earthquake Likelihood Model (RELM) [?].

For developing the methods and to be able to compare them we used the earthquake catalog from the Japanese Meteorological Agency (JMA), using event data from 2005 to 2010.

This paper is organized as: in Section 2 reviews applications of Evolutionary Computation in the context of seismology research. The next Section, Section 3, we give a details of each of the forecast proposed covering the Collaboratory for the Study of Earthquake Predictability (CSEP) framework and the empirical laws. In Section ??, we give the description of the tests proposed in [?]. After that, in ??, we define the target areas used for the experiment and the data from the JMA; we clarify the design followed during the experiments and how we compared the forecast models derived from our methods. Finally, we show the results and conclude this work in 3.8.4 and 4.

2. Evolutionary Computation for Earthquake Risk Analysis

In this section we will briefly discuss some reports of the application of Evolutionary Computation and related method for Earthquake Risk Analysis.

The usage of Evolutionary Computation in the field of earthquake risk models is somewhat sporadic. Zhang and Wang [?] used Genetic Algorithms to fine tune an Artificial Neural Network (ANN) and use this system to produce a forecast model. Zhou and Zu [?] also proposed a combination of ANN and EC, but their system only forecasts the magnitude parameter of earthquakes. Sadat, in the paper [?], follows the idea of Zhou and Zu, aiming to predict the magnitude of the earthquakes in North Iran, but in this case, they used ANN and GA.

Nicknam et al. [?] simulated some components from a seismogram station and predicted seismograms for other stations. They combined the Empirical Green's Function (EGF) with GA. The EGF method is used to synthesise acceleration time histories and the GA approach is developed to optimise the seismological model. They found that this method obtained good agreement with the observed data, but are not sure that results are free from uncertainties.

Kennett and Sambridge [?] used GA and associated teleseisms procedures to determine the Fault Model parameters of an earthquake. By doing so, they demonstrated that non-linear inversion can be achieved for teleseismic problems without any calculation of waves travel times.

Some seismological models were developed aiming to estimate parameter values by using Evolutionary Computation. For example, Evolutionary Computation was used to estimate the Peak Ground Acceleration (PGA) of seismically active areas [?], [?], [?], [?]. The works done by Kerh et al. [?], [?] are basely a combination of ANN and GA to estimate or predict PGA in Taiwan. These work are based on the benefits of mixing both techniques. They state that the usage of a purely ANN method to estimate PGA may fall into a local minimum and that can be avoid by combining ANN with GA, hence GA is a good method to find global optimums.

Ramos and Vázquez [?] used Genetic Algorithms to decide the location of sensing stations. In this work they achieved, in general, better results with the GA method

when compared with the Seismic Alert System (SAS) method and a greedy algorithm method.

Ramos et al. work is a important work because it helps the population to avoid bigger disasters caused by earthquakes by increasing the time response of the Seismic-Sense Stations. It has some similar feature as the one present in this document: it uses GA to prevent earthquake disasters and tries to locate targets in a given area (though the targets of this work is sensing stations and ours works target is the earthquakes location) and it proposes a methodology to do a GA parameter setting to find which combination of values for the GA parameters achieve higher results. It is interesting to state that once a solution places a station in an area that is not possible to have sensors, this possible solution suffers some penalties.

Saeidian et al. [?] also based on the same idea of locating sensing stations. Their work differs from the work of Ramos et al. because it makes a comparison in performance between the GA and Bees Algorithm (BA) to decide which of those techniques would perform better when choosing the location of sensing stations. He found out that the GA was faster than the BA.

Huda and Santosa [?] published a paper in which the goal is to find, via Genetic Algorithm, the speed of the waves P and S in the mantle and in the earth crust. P waves are indicated as the first fault found in seismological data and S waves are the changes caused in the phase of a P wave [?]. This research aims to obtain a structure of the Japanese underground and geographically focuses in the same region as our work, though it uses data from two kinds of waves which are not available to us.

3. The Forecast Models Using Genetic Algorithm

All forecast models proposed in this paper are based in the Collaboratory for the Study of Earthquake Predictability (CSEP) framework.

Each individual has its own representation of the framework based on different perceptions of what are the best aspects of the framework.

For all methods, the population is trained on earthquake event data for a training period, which is anterior to the target test period. After completing the evaluation limit, the best individual is chosen to be the final forecast.

3.1. 1-year Forecast Models

The CSEP framework, a forecast model uses a gridded rate forecast [?], one common format in the literature. In this format, a geographical region is divided in sections, bins, during a start date and an end date. The forecast will estimate the risk of earthquake occurrence in this target region, during the target time interval. For this study we considered the target time interval of one year [?].

Large and independent earthquakes, also known as mainshocks, are followed by a wave of others earthquakes, the

aftershocks [?]. Hence there is no physical measurement to identify mainshocks and its aftershocks [?], we divided the forecast models in two groups: the ones that only forecasts mainshocks, using only GA techniques, and those that forecast both mainshocks and aftershocks using both GA techniques and empirical laws, such as the modified Omori law. These laws are used to derive the aftershocks from a synthetic data of mainshocks.

Both classes forecast earthquakes with magnitude greater than 3.0 for every scenario proposed, with a binning of 0.1, here named as cells to avoid conflicts with the location bin. That results in magnitude cells of [3.0, 3.1), [3.1, 3.2), until [9.9, 10).

3.2. Mainshock Models

There are two mainshock models. The ones generated by the GAModel and the ones generated by the ReducedGAModel.

The GAModel is considered as one method to generate space-rate-time forecasts. It also could be described as:

$$\Lambda(t, x, y, M | \Upsilon_t) = \mu(x, y) \quad (1)$$

where you can denote the number of earthquakes forecast in all bins as $\Lambda(t, x, y)$ [?] given that Υ_t is the earthquake observation data up to time t .

The ReducedGAModel, which represents the idea of changing the candidates solution representation (see Section 1), is a method with a similar description of the GAModel. The difference is that, in the ReducedGAModel each possible solution represents only a fraction of the forecast where we expect to find specific risk areas.

3.2.1. GAModel. The GAModel is completely based on the framework suggested by the CSEP. In it, one forecast is defined as a region in a specific time interval and is divided in bins. Each bin represents a geographical interval. The whole target area of study is covered by a group of these bins where each bin has an earthquake forecast value. This groups of bin represent the $\mu(x, y)$, the background intensity [?]. In the GAModel, each possible solution is represented as an entire forecast model.

Genome Representation. In the GAModel each individual represents an entire forecast model. Each gene of the individual is a real value, corresponding to one bin in the desired model. The values are sampled from the interval [0, 1). These real values are converted to a integer forecast, we use the same modification of the Poisson deviates extraction algorithm 3.2.1 used in [?]. In the algorithm 3.2.1 x is the real value that will be converted and μ is the mean of the earthquakes observations in the real data.

The genome is a real valued array X , where each element corresponds to one bin in the desired model (the number of bins n is defined by the problem). Each element $x_i \in X$ takes a value from [0, 1). In the initial population, these values are sampled from a uniform distribution and they are randomly generated. For more details of the genome representation, please refer to [?].

Algorithm 1 Obtain a Poisson deviate from a $[0, 1)$ value [H]

```

 $L \leftarrow \exp(-\mu), k \leftarrow 1, prob \leftarrow 1 * x$ 
repeat
  increment  $k$ 
   $prob \leftarrow prob * x$ 
until  $prob > L$ 
return  $k$ 
while  $prob > L$  do
   $k \leftarrow k + 1$ 
   $prob \leftarrow prob * x$ 
end while
return  $k$ 

```

To clarify how the GAModel works, we use the same example as the one used in [?]. The Kanto region, one of the four areas used in both studies, is divided into 2025 bins (a grid of 45x45 squares). Each bin has an area of approximately 25km². The GAModel then calculates an expected number of earthquakes for every bin on a determined time interval, so the GA searches for good values in 2025 bins.

Fitness Function. The GAModel, as all the other methods, uses the log-likelihood value, as the fitness function. The fittest individual among all the others, is preserved in the next generation, to make the solution of one generation as good as the its last generation. The bins, a gene of the genome representation, b_n , define the set β and n is the size of the set β :

$$\beta := b_1, b_2, \dots, b_n, n = |\beta|. \quad (2)$$

The probability values of the model j , expressed by the symbol Λ , is made of expectations λ_i^j by bin b_i . The vector is define as:

$$\Lambda^j = (\lambda_1^j, \lambda_2^j, \dots, \lambda_i^j); \lambda_i^j := \lambda_i^j(b_i), b_i \in \beta \quad (3)$$

The vector of earthquake quantity expectations is defined as: earthquake by time. The Ω vector is composed by observations ω_i per bin b_i , as the Λ vector:

$$\Omega = (\omega_1, \omega_2, \dots, \omega_i); \omega_i = \omega_i(b_i), b_i \in \beta \quad (4)$$

The calculation of the log-likelihood value for the ω_i observation with a given expectation λ is defined as:

$$L(\omega_i | \lambda_i^j) = -\lambda_i^j + \omega_i \log \lambda_i^j - \log \omega_i! \quad (5)$$

The joint probability is the product of the likelihood of each bin, so the logarithm $L(\Omega | \Lambda^j)$ is the sum of for $L(\omega_i | \lambda_i^j)$ every bin b_i :

$$\begin{aligned}
 L^j &= L(\Omega | \Lambda^j) = \sum_{i=1}^n L(\omega_i | \lambda_i^j) \\
 &= \sum_{i=1}^n -\lambda_i^j + \omega_i \log \lambda_i^j - \log \omega_i!
 \end{aligned} \quad (6)$$

The fitness function is a coded version of the equation 6. It uses the probabilities of the bins of each individual of model for the λ values.

Evolutionary Operators. The GAModel use a combination of operators made available by the Distributed Evolutionary Algorithms in Python (DEAP) [?]. We used the One Point Crossover for the crossover operator, the Polynomial Bounded Mutation for the mutation operator and for selection, we used Tournament selection and Elitism. The parameters are described in the Table 1.

Table 1. PARAMETERS USED IN GAMODEL AND EMP-GAMODEL

| | |
|-------------------------------|--------------------------|
| Population Size | 500 |
| Generation Number | 100 |
| Elite Size | 1 |
| Tournament Size | 3 |
| Crossover Chance | 0.9 |
| Mutation Chance (individual) | 0.1 |
| Polynomial Bounded parameters | eta = 1, low = 0, up = 1 |

The parameters of the Polynomial Bounded mutation function are:

- 1) eta = 1. Crowding degree of the mutation. A high eta will produce a mutant resembling its parent, while a small eta will produce a solution much more different;
- 2) low = 0. The lower bound of the search space;
- 3) up = 1. The upper bound of the search space.

The chance of applying both mutation operator function and crossover operator function takes into account only their chance of occurrence. This means that it may be the case that one of them or both are not applied.

3.2.2. ReducedGAModel. The GAModel defines a expected number of earthquakes for every single bin in the target region. That could lead to exhaustive and, sometimes worthless, searches. That is caused by the number of bins in the forecast and also because in some bins there are no earthquake occurrences in the observation data. That means that the GAModel has a lot of parameters and may of its bins have null values (values equal to 0). To avoid such unnecessary task we proposed the ReducedGAModel.

With this method, we aim to minimise the search space and the quantity of parameters the GA has to deal with. For that we changed the individual representation. The individuals in the ReducedGAModel only define expected number of earthquakes in bins that already had some occurrence in the past, giving a direction to where the GA should search. That helps the ReducedGAModel in the search for better solutions and it makes the convergence faster once the space search is smaller.

The ReducedGAModel has a similar description of the GAModel. The difference is that, in the ReducedGAModel, each possible solution represents only a fraction of the forecast where we expect to find specific risk areas. To do so, this method obtain the position of past occurrences. Then

it calculates some expected number of earthquakes only for the bins related to those positions. These positions may vary during the evolution of the method, including positions that never had earthquake events before. That is important to add some variation to the method.

Genome Representation. The genome representation in the ReducedGAModel is a simplified version of the genome of the GAModel. For the ReducedGAModel, the genome is a list of ordered pairs. The first element of the pair are the coordinates of a bin in the model. The second element of the pair is a number that indicates an earthquake occurrence estimate for this bin.

To calculate the size of the individual we use the real data from the worse 5 years and create a list of every bin that had events in it, even if only once.

In the ReducedGAModel, each individual is a list of a sub-region of the forecast model. This list initially refers to bins where earthquake events happened in the past. During the develop of the ReducedGAModel, the list may refer to positions that never had occurrences before. Each element of the list, a gene, also contains one real value between [0,1). In the initial population, these values are sampled from a uniform distribution and they are randomly generated. When needed, every real value is converted to a integer forecast by the same Algorithm, as in the GAModel.

To generate the forecast model we need to do an intermediate step. We map every location from the list with a bin in the forecast model.

The genome size is usually smaller than the one used in the GAModel and the Emp-GAModel, once the amount of sub-regions where earthquakes with magnitude above 3.0 happened for any given area is smaller then the total number of genes of the individual.

To exemplify, we use a similar example as the one in the Section 3.2.1. Lets consider that there are 10 bins with occurrences in Kanto in the last 5 years. It will make the GA start searching for good values for only those 10 bins, leaving the other 2015 bins empty, representing zero occurrence. It is important to highlight that in the worst case, it will make the same amount of searches as the GAModel. The final forecast model will maintain the amount of bins with occurrence, but the number of events for every bin and their location may change.

Fitness Function. The fitness function is the same as in the GAModel, 3.2.1. Here is also important to generate the forecast model by applying the map function on the individual.

Evolutionary Operators. All operators in the ReducedGAModel are the same as the operators of the GAModel, except the mutation function. We use a simple mutation operator which samples entirely two new values, both sampled from uniform distributions. The first, is a new real value from [0,1) and the second one, a new integer value from [0,x), where x is the maximum amount of bins a model can have in the target region. For the parameters see Table 1.

As in the GAModel, the chance of applying both mutation operator function and crossover operator are indepen-

dent and they may or may not be used.

3.3. Mainshock+Aftershock Models

The mainshock and aftershock methods are a two-step methods. The first step is as defined for the mainshocks methods, therefore, we first use GA techniques to obtain a synthetic mainshock data. The second step is to use seismological empirical equations to obtain the aftershocks from the mainshocks.

Hence earthquakes cluster in space and inspired by the space-time epidemic-type aftershock sequence (ETAS), we proposed two methods, called Emp-GAModel and Emp-ReducedGAModel. They represent the idea of associating the GA with seismological empirical equations. They are described as:

$$\Lambda(t, x, y, M | \Upsilon_t) = \mu(x, y) J(M) \quad (7)$$

That can be expanded to:

$$\Lambda(t, x, y | \Upsilon_t) = \mu(x, y) + \sum_{t_i \in t} K(M_i) g(t - t_i) P(x, y) \quad (8)$$

methods use $\mu(x, y)$ as defined for mainshock methods 3.2. It is calculated as an expected number of earthquakes for every bin in the target region, given that Υ_t is the earthquake observation data up to time t .

3.3.1. Empirical Equations. The Omori law, $g(t)$, which is considered one empirical formula of great success [?] [?] [?], is a power law that relates the earthquake occurrence and its magnitude with the decay of aftershocks activity with time. For this approach we used the probability density function (PDF) form of the modified Omori law [?]:

$$g(t) = \frac{(p-1)}{c(1 + \frac{t}{c})^{-(p)}} \quad (9)$$

The variable p is a index of this equation and the variable c is a constant, given in days. In the paper [?], Utsu summarise most of the studies in Japan and described the range for these variables. For p the range is between 0.9 and 1.4 and for c 0.003 and 0.3 days. These values were based on the Davidon-Fletcher-Powell optimisation procedure and used in ETAS [?]. Also there is the variable t that is the time limit to when a mainshock may influence the cause a aftershock.

Based on paper [?], we set the values of 1.3 for p and 0.003 for c for our experiments. We set the time interval t between a mainshock and its aftershocks at one month. In the paper, it says that if the t value is too short, the number of aftershocks is too small, but if it is too big, we may also consider background activity and suggest the use of a 30 days period.

For $K(M_i)$, the total amount of triggered events, we count aftershocks within a given area, A , using the following

formula, where M_c is the magnitude threshold, with $M_c = 3.0$:

$$K(M_i) = A \exp([\alpha(M_i - M_c)]) \quad (10)$$

In the paper [?], it states that α should be equal to the inverse of the magnitude of an event, or $magnitude^{-1}$. To obtain A , the following equation from [?], was used:

$$A = e^{(1.02M-4)} \quad (11)$$

With the $K(M_i)$ and $g(t)$, the PDF Omori, equations it is possible to calculate the total number of earthquakes. For that we must sum the product of the equations, varying t :

$$\sum_{t_i \in t} K(M_i) g(t - t_i) \quad (12)$$

This result will lead to a number of aftershocks related to a single mainshock. Then, we can use the $P(x, y)$ equation to distribute the aftershock to the bins near the mainshocks position. $P(x, y)$ calculates the position of the aftershocks with base on the origin of the mainshock. It is a simple space distributing function, that allocates the aftershocks in one of the following positions: upper, lower, left or right. It runs for a number of steps, getting further from the origin at each step or as when there are no more events to be allocated. $P(x, y)$ can be split into 4 equations, one for each position:

$$\begin{aligned} model[x + y] &= (aftershocks - [model[x] - 2 * x]) / 4; \\ model[x - y] &= (aftershocks - [model[x] - 2 * x]) / 4; \\ model[x - y * row] &= (aftershocks - [model[x] - 2 * x]) / 4; \\ model[x + y * row] &= (aftershocks - [model[x] - 2 * x]) / 4 \end{aligned}$$

and lastly, the $J(M)$ is obtained by using the function etasim, from the SAPP R package [?] that simulates magnitude by Gutenberg-Richter's Law.

3.3.2. Emp-GAModel. The Emp-GAModel is a specialisation of the GAModel. This is achieved by the use of empirical equations after the forecast is provided. This means that the its genome representation are the same as in the GAModel.

Genome Representation. The genome representation is the same as in the GAModel, Section 3.2.1.

Fitness Function. The fitness function is the same as in the GAModel, Section 3.2.1, and the ReducedGAModel.

Evolutionary Operators. The Emp-GAModel use the same combination of operators that the GAModel. For more explanation, please see the Section 3.2.1.

Emp-ReducedGAModel. The Emp-ReducedGAModel is a specialisation of the ReducedGAModel. This is achieved by the use of empirical equations after the forecast is provided. This means that the its genome representation are the same of ReducedGAModel.

Genome Representation. The genome representation is the same as in the ReducedGAModel, Section 3.2.2.

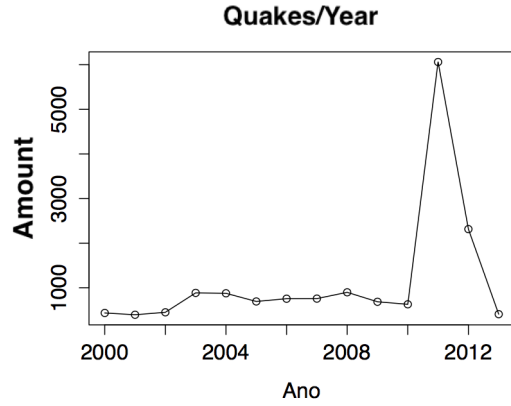


Figure 1. Amount of earthquake by year.

Fitness Function. The fitness function is the same as for all methods, Section 3.2.1. Here is also important to generate the forecast model by applying the map function on the individual as in the last Section, 3.2.2.

Evolutionary Operators. The Emp-ReducedGAModel use the same combination of operators that the ReducedGAModel. For more explanation, please see 3.2.1.

3.4. Experimental Data

Here we describe the earthquake catalogue, how we used it and the regions in Japan selected for the experiments.

We also preprocessed the catalogue. We wanted to analyse how earthquakes characteristic changed with the magnitude and the depth. Also we explain briefly how we classified the mainshocks and aftershocks.

3.4.1. Earthquake data. The goal of this research is to find existing patterns in the occurrence of earthquakes. For that it is essential to access trustful data and to explore its details. From the *Japan Meteorological Agency* web page we obtained earthquake data about earthquakes in Japan. In this data there are information about earthquakes that happened in or nearby Japan, with the variables: time of the occurrence, magnitude, latitude and longitude and epicentre depth, for the years of 2000 to 2013.

During the preprocessing phase, we discovered a higher number of occurrences of earthquakes during the year of 2011, when a 9.0 M_w earthquake happened, see Section ???. This earthquake triggered too many after called aftershocks in all Japan. It is considered that big earthquakes may cause others earthquakes [?]. In Figure 1 it is possible to visualise a great number of earthquakes for the year of 2011. Because of this abnormal behaviour and because we decided to focus on more stable occurrences, we limited the training base to earthquakes until 2010.

Based on the statement done before and considering that we want earthquakes that follow more stable patterns, we

selected the ones that happened in land areas or very shallow sea areas, with maximum depth of 100km.

3.4.2. Regions. For the experiments, the data was changed into slices for every year. Each slice is as follows: if the base contains data about a time interval of 10 years, it will be split in 10 slices.

We also selected some sub-areas in Japan to better extract and understand earthquakes characteristics and patterns. Those areas are Kanto, Kansai, Touhoku and East Japan. The Figure 2 shows how we defined them.

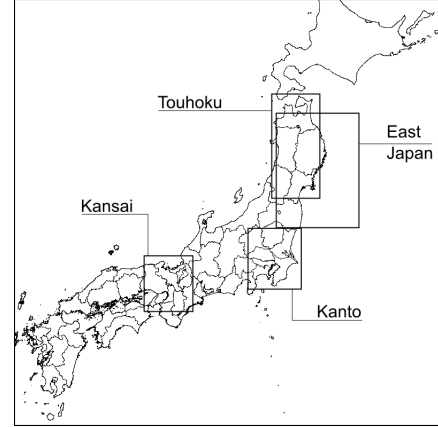


Figure 2. Japan and the areas used in this studied.

They are described as follows:

Kanto. Kanto is the region around Tokyo. It is an area with high seismologic activity during the years we studied. Its coordinates are 34.8 North, 138.8 West, with 2025 bins. Each bin covers an area of approximately 25km².

Kansai. Kansai is the region that includes Kyoto, Osaka and many others historical cities. In this area, rather than Kanto area, there is a small seismic activity. Its coordinates are 34 North, 134.5 West, with 1600 bins. Each bin covers an area of approximately 25km².

Touhoku. Touhoku is the region in the North of the main Japanese island. It has some clusters of seismic activities during the years we studied. Its coordinates are 37.8 North, 139.8 West, with 800 bins. Each bin covers an area of approximately 100km².

East Japan. Is the region that is related with the east coast of Japan. It is the most different area, because it has earthquakes that happened both in land or in the sea. It was in this region that the 9.0 M_w earthquake happened. Its coordinates are 37 North, 140 West, with 1600 bins. Each bin covers an area of approximately 100km².

3.4.3. Depth Histogram of Earthquakes. The patterns of earthquakes are dependent of the epicentre. We wanted to explore the relation between the depth of the earthquakes and how would our models behave on those situations.

In Figure 3, it is possible to understand that most of the earthquakes happened with depths smaller or equal to 100 km. The earthquakes deeper than 100 km are fewer

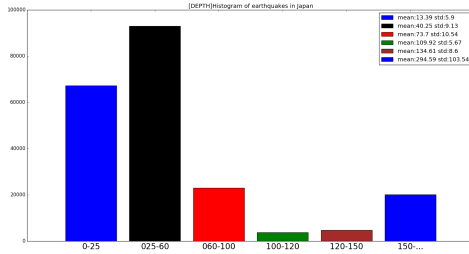


Figure 3. Depth Histogram of earthquakes.

and more distant, as it is in the same Figure.

The reason we decided to groups as: earthquakes with depth until 25 km, until 60 km or until 100 km. This is because shallow earthquakes are considered to be more independent earthquakes [?].

3.4.4. Mainshocks and Aftershocks - Clustering. In the Section ??, we explained that we have two kinds of models, the ones that only consider aftershocks and those that consider both mainshocks and aftershocks. Therefore, it is needed to isolate, to classify the earthquakes into one of these two groups.

The question is how should it be done. The simplest way, is to select earthquakes with magnitude above 3.0 in the Richter Scale and then to consider those as the mainshocks. The distribution of earthquakes after this selection is exemplified in the Figure 4.

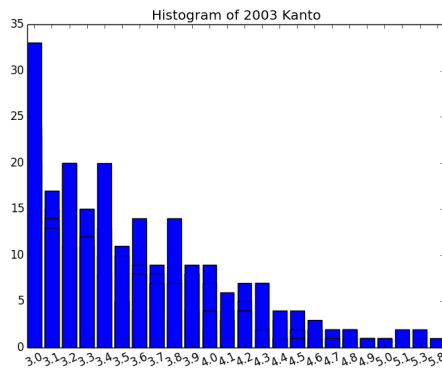


Figure 4. Histogram of earthquakes stronger than 3.0 in the Richter Scale in Kanto

The problem with this simple idea is: if a big mainshock happens and it triggers some aftershocks with magnitude higher than 3.0 in the Richter Scale it would be considered as a mainshock. To avoid this problem we used two methods proposed in the literature: Window Methods and the Single-Link Cluster. For more information about these methods, see reference [?].

3.5. Experimental Design

The first experiment was made to compare the all the models proposed with each other and to discover which method would achieve higher log-likelihood values. We created new scenarios, applying the methods for all regions and for the years of 2005-2010. We also used 3 kinds of catalogues: the JMA and the declustered catalogues form the Window method and the SLC method. Then, we compared the means of the models log-likelihood values using the ANOVA test. If a group of variables considered for the ANOVA test showed no statistically significant difference, we applied the Paired Student t-test, in the case all groups showed statistically significant difference, the Tukey HSD methodology analysis was used.

We also made a magnitude experiment. This experiment was done to explore the influence of the magnitude in all models generated. We split them into slices composed of earthquakes that have magnitude in a given magnitude interval. The we calculated the log-likelihood of these slices and applied the ANOVA test to compared these sliced-models.

3.5.1. The Mainshock Models and Mainshock with Aftershock Method Experiments. Here we describe the catalogues and the evolutionary operators used for the experiments. Then we specify the models comparison.

The catalogues. The data used was from JMA catalogue, with the minimum magnitude of 3.0 and the two declustered catalogues, obtained from the methods explained in the Section 3.4.4. The models that use these catalogues have in the word Window appended at their names, for the methods that used the Window declustering, or SLC, for the methods that used the Single Link Cluster.

3.5.2. Models Comparison. For this new experiment, we used even more scenarios (space/time regions) than the others. Each scenario contains the earthquakes for the regions of Kanto, Kansai, Touhoku and East Japan for a given year (2005-2010). We wanted to explore if there exists any influence in the performance of the models that are caused by the depth of an earthquake. So, the scenarios are also composed by introducing a three groups depths thresholds. They are: of earthquakes with depth smaller than 25km, or between 0km and 60km or even between 0km and 100km.

These methods are stochastic methods and hence are variations of the GAModel, we decided to maintain the number of repetitions without redoing the Power of the Student t-test.

3.6. Statistical Analysis of the Results

The goal is to discover if there is any variation between the methods and which are the most influential variables. For achieve that, we will use the ANOVA test.

In the ANOVA test, of variance of one specific variable with 95% confidence level, with “p-value” < 0.05 it means

that there exists a statistical significant evidence that the variables variance are different from.

There are some tests hypothesis for this experiment that we want to analyse. They all can be generalised as follows:

$$\begin{cases} H_0 : \text{The population means are equal.} \\ H_1 : \text{The population means are different.} \end{cases}$$

Then, if there is no statistical significant difference between the means, we apply the Tukey HSD. We apply it on the results obtained from the ANOVA test to specify which groups differ. Tukey's methodology analysis shows the means of a case with the means of every other case. Doing so, it identifies differences between means :

$$\begin{cases} \mu_a - \mu_b, \text{ where } \mu_a \text{ is the mean of the first group} \\ \mu_b \text{ is the mean of the second group.} \end{cases}$$

In the case where statistical significant difference exists, we explore this by pairing the measures observations of two groups [?].

That is:

$$\begin{cases} H_0 : \mu = 0, \text{ the difference between observations is 0.} \\ H_1 : \mu \neq 0, \text{ difference between observations is not 0.} \end{cases}$$

3.6.1. Results from The Mainshock Models Mainshock with Aftershock Models Experiment. An one-way between subjects ANOVA was conducted to compare the effects of the models, the years and regions on the log-likelihood value. In this study there are the models: *ReducedGAModel*, *Emp-ReducedGAModelSLC*, *Emp-GAModel*, *Emp-ReducedGAModel*, *GAModelWindow*, *ReducedGAModelWindow*, *GAModelSLC*, *ReducedGAModelSLC*, *Emp-GAModelWindow*, *Emp-ReducedGAModelWindow*, *GAModel* and *Emp-GAModelSLC*.

Based on the results of this first test, it is evident that all variables are significantly different. The results of the experiments are in the Table 2. For all, the confidence level is set to 95% .

Because we found statistically significant result, we applied a Post-hoc comparisons using the Tukey HSD analysis methodology. It compared each condition with all others. For example, it compares the values from the *GAModel* with the *GAModelWindow*. It indicated that the models *Emp-ReducedGAModelWindow*, *GAModelWindow*, *ReducedGAModelWindow*, *GAModelSLC*, *ReducedGAModelSLC* achieve statistically better or equal results in terms of log-likelihood when compared with the other models and when compared with themselves, they are statistically similar. The results of the experiments are in the Table 3. For all, again, the confidence level is set to 95% .

Therefore, to confirm that statistically the models are similar, we applied the conducted the ANOVA test considering only with the models indicated by the Tukey HSD, see Figure 5.

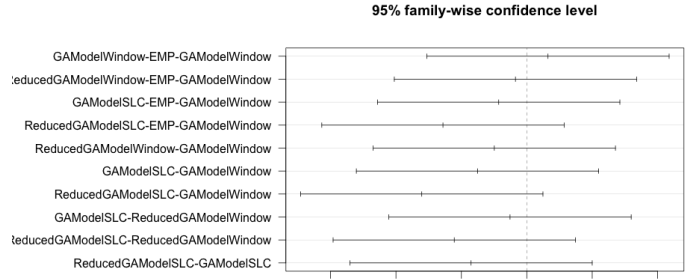


Figure 5. Intervals of Confidence 95% of differences between the Mainshock Models Mainshock and the Aftershock Models, taken two by two.

This time, we found statistically significant difference only for the year and region condition. To show that the models results are not statistically different from each other, we applied a pairing analysis.

From the pairing analysis, we decided to use the *ReducedGAModelSLC* as the representative method of this study. That is because, in most cases when its values were compared, it showed a little better performance in the means of the log-likelihood values. For the results, see the Table 4.

In this Table, the column labelled $\mu_a - \mu_b$ shows the result of paired difference between the models referred in the "Models Compared" column. The p-value shows the significance value of the paired *Student's t-test* for the null hypothesis "The paired difference of the means of the models is equal".

3.7. The Models Examples And The Real Data

The Figure ?? shows a model from the *GAModel* method for the year 2005 in East Japan. The next Figure, ?? shows a model from the *ReducedGAModel* ?? method for the year 2005 in East Japan.

All Figures, ?? ?? ?? ??, indicate a low earthquake intensity as white while the more intensity areas, are shown in red. They are, in order, the data visualisation for the model from: the *GAModel*, the *ReducedGAModel*, the *Emp-GAModel* and the *Emp-ReducedGAModel* for East Japan in 2005. The Figure ?? represents the earthquake occurrences in the same region and year.

3.8. Magnitude Experiment

In this experiment, we focused on studying how the magnitude of an earthquake affects the model quality, because the patterns of the earthquakes are depended of its magnitude. We wanted to explore the relation between the magnitude of the earthquakes and how would the models behave on those situations.

| Variable | Degrees of Freedom | Sum Sq | Mean Sq | F Value | Pr(>F) |
|----------|--------------------|-----------|-----------|---------|--------|
| Model | 15 | 149303768 | 9953585 | 63.72 | <2e-16 |
| Year | 5 | 414016420 | 82803284 | 530.06 | <2e-16 |
| Region | 3 | 869821655 | 289940552 | 1856.02 | <2e-16 |

Table 2. ANOVA TEST RESULTS VALUES - MAINSHOCK MODELS MAINSHOCK AND AFTERSHOCK MODELS.

| Variable | Degrees of Freedom | Sum Sq | Mean Sq | F Value | Pr(>F) |
|----------|--------------------|-----------|----------|---------|--------|
| Model | 4 | 884882 | 221220 | 1.604 | 0.171 |
| Year | 5 | 150297410 | 30059482 | 217.955 | <2e-16 |
| Region | 3 | 234225270 | 78075090 | 566.107 | <2e-16 |

Table 3. ANOVA TEST RESULTS VALUES - EMP-REDUCEDGAMODELWINDOW, GAMODELWINDOW, REDUCEDGAMODELWINDOW, GAMODELSLC, REDUCEDGAMODELSLC.

| Region | Models Compared | Mean of $\mu_a - \mu_b$ | p-value |
|------------|---|-------------------------|-------------------|
| Kansai | EMP-GAModelWindow - GAModelWindow | 38.67553 | 3.304e-05 |
| | EMP-GAModelWindow - ReducedGAModelWindow | 4.272185 | 0.2607 |
| | EMP-GAModelWindow - GAModelSLC | 112.0424 | 1.122e-05 |
| | EMP-GAModelWindow - ReducedGAModelSLC | -1.787262 | 0.5673 |
| | GAModelWindow - ReducedGAModelWindow | -34.40335 | 0.000963 |
| | GAModelWindow - GAModelSLC | 73.36687 | 9.065e-06 |
| | GAModelWindow - ReducedGAModelSLC | -40.46279 | 6.32e-05 |
| | ReducedGAModelWindow - GAModelSLC | 107.7702 | 2.632e-05 |
| | ReducedGAModelWindow - ReducedGAModelSLC | -6.059447 | 0.2982 |
| | GAModelSLC - ReducedGAModelSLC | -113.8297 | 1.2e-05 |
| Touhoku | EMP-GAModelWindow - GAModelWindow | 3.34556 | 0.546 |
| | EMP-GAModelWindow - ReducedGAModelWindow | 81.60965 | 5.225e-07 |
| | EMP-GAModelWindow - GAModelSLC | 63.02216 | 0.01971 |
| | EMP-GAModelWindow - ReducedGAModelSLC | -62.70586 | 0.007075 |
| | GAModelWindow - ReducedGAModelWindow | 78.26409 | 2.938e-05 |
| | GAModelWindow - GAModelSLC | 59.6766 | 0.04829 |
| | GAModelWindow - ReducedGAModelSLC | -66.05142 | 0.001231 |
| | ReducedGAModelWindow - GAModelSLC | -18.58749 | 0.3443 |
| | ReducedGAModelWindow - ReducedGAModelSLC | -144.3155 | 0.000214 |
| | GAModelSLC - ReducedGAModelSLC | -125.728 | 0.01216 |
| East Japan | EMP-GAModelWindow - GAModelWindow | 1.872764 | 0.9539 |
| | EMP-GAModelWindow - ReducedGAModelWindow | 194.4944 | 1.834e-06 |
| | EMP-GAModelWindow - GAModelSLC | 189.1155 | 0.0003456 |
| | EMP-GAModelWindow - ReducedGAModelSLC | -274.9858 | 4.961e-05 |
| | GAModelWindow - ReducedGAModelWindow | 192.6217 | 0.003738 |
| | GAModelWindow - GAModelSLC | 187.2428 | 9.495e-06 |
| | GAModelWindow - ReducedGAModelSLC | -276.8586 | 4.636e-05 |
| | ReducedGAModelWindow - GAModelSLC | -5.378912 | 0.8576 |
| | ReducedGAModelWindow - ReducedGAModelSLC | -469.4803 | 1.446e-05 |
| | GAModelSLC - ReducedGAModelSLC | -464.1014 | 2.38e-06 |
| Kanto | EMP-GAModelWindow - GAModelWindow | 57.95612 | p-value = 0.00138 |
| | EMP-GAModelWindow - ReducedGAModelWindow | 79.60781 | 3.441e-05 |
| | EMP-GAModelWindow - GAModelSLC | 274.3114 | 5.717e-06 |
| | EMP-GAModelWindow - ReducedGAModelSLC | -96.61803 | 6.22e-07 |
| | GAModelWindow - ReducedGAModelWindow | 21.65169 | 0.1105 |
| | GAModelWindow - GAModelSLC | 216.3553 | 2.302e-07 |
| | GAModelWindow - ReducedGAModelSLC | -154.5741 | 1.741e-05 |
| | ReducedGAModelWindow - GAModelSLC | 194.7036 | 3.678e-05 |
| | ReducedGAModelWindow - ReducedGAModelSLC | -176.2258 | 4.337e-06 |
| | GAModelSLC - ReducedGAModelSLC | -370.9294 | 1.942e-06 |

Table 4. PAIRED EXPERIMENT RESULT.

For that, we created magnitude intervals, where each interval is named as a slice. A slice is an closed interval of 1.0 degree starting from 3.0 degrees of magnitude, see Section 3.5.1, and ending in 10.0 degrees. For example, $[3.0 - 4.0]$ or $[7.0 - 8.0]$ are two different slices. For each model, we selected only the earthquakes that belong to a slice. Then, we calculate the log-likelihood value.

3.8.1. Magnitude Study. From the results already obtain and showed in the section 3.6.1, when selected the models with earthquakes with depth smaller or equal to 25 km and then we split the models in magnitude intervals, as defined in 3.8.

After that, we compared those split models against themselves. Based on the results of this test, it is evident that all variables are still significantly different. The results of the experiments are in the Table 5. For all, as before, we choose

| Variable | Degrees of Freedom | Sum Sq | Mean Sq | F Value | Pr(>F) |
|-----------|--------------------|-----------|-----------|---------|--------|
| Model | 5 | 2.360e+09 | 4.720e+08 | 2828 | <2e-16 |
| Year | 3 | 4.624e+09 | 1.541e+09 | 9234 | <2e-16 |
| Magnitude | 7 | 3.726e+09 | 5.322e+08 | 3189 | <2e-16 |

Table 5. ANOVA TEST RESULTS VALUES - MAGNITUDE STUDY.

the confidence level to be 95%.

We found statistically significant result and, as before, we applied the Tukey HSD test. The results are shown in Figure ?? and the *NULL* field was used as the model with all magnitude intervals (the complete model).

It indicated that the interval $[3.0 - 4.0]$ always performed, in terms of log-likelihood values, worse than all other intervals. this phenomenon also happens in the interval $[4.0 - 5.0]$, though in this case, the difference is not as big as the last one. The other intervals show no significant difference. From the results found, we decided to chose only earthquakes with magnitude higher than 4.0 as our threshold value.

3.8.2. Catalogues and Models . This experiment used the same catalogues used in the previous experiment 3.5.1.

The models also are the models from the last experiment. We also created the new models, considering the slices and add them to the comparison. That lead to a comparison with the models from the last experiment and the models sliced.

3.8.3. Statistical Analysis. The goal is to discover if the magnitude influences any variation in the methods and how it does.

For this experiment, we followed the same design from the Section 3.6.

3.8.4. Results.

4. Conclusions

Acknowledgments

The authors would like to thank Bogdan M. Enescu, from the department of Earth Evolution Sciences in the university of Tsukuba for his useful comments.

We would also like to thank the Japan Meteorological Agency for the earthquake catalog used in this study.