

RESEARCH ARTICLE

Seismic Event Clustering in Mainland Portugal: DBSCAN Approach

André Martins Brito^a

^a Faculdade de Ciências e Tecnologia, Universidade Nova de Lisboa

ARTICLE HISTORY

Compiled May 9, 2023

Abstract

This template is for authors who are preparing a manuscript for a Taylor & Francis journal using the L^AT_EX document preparation system and the **interact** class file, which is available via selected journals' home pages on the Taylor & Francis website.

KEYWORDS

Sections; lists; figures; tables; mathematics; fonts; references; appendices

1. Introduction

The Portuguese Mainland is located relatively close to a geological complex region. Just south of the Iberian Peninsula, of which Mainland Portugal is a part of, is the boundary between the Eurasian and African Plates. This boundary extends along the Azores-Gibraltar Fracture Zone, starting at the Azores Archipelago, passing through the Gibraltar strait and continuing on until the Mediterranean Sea. For its part, the Azores Archipelago is located where this Fractures Zone meets the American Plate, called the Triple Junction of the Mid-Atlantic Rift, a very seismic intensive region which as been studied by many authors, (e.g. [9]). Due to the complexity of the Africa/Eurasia plate boundary, this region is prone to generate powerful earthquakes that, given its proximity to Mainland Portugal, can produce major material damages in land and, thus, represents a risk to the Portuguese Population. The infamous Lisbon earthquake of 1755, with an estimated magnitude of 8.5-8.7 on the Richter scale, which devastated the city of Lisbon and several towns in the Algarve is the recurrent example given in most of the literature of the seismic risk inherent to the Portuguese Mainland.

[7] present a geophysical and geological overview of the geological structures of the Iberian Peninsula and surrounding maritime area. The authors give particular focus to the western margin of Portugal stating that, most of the seismic activity recorded in the Portuguese Mainland originates at the Africa/Eurasia plate boundary or close to it. Banco do Gorringe is a region of this boundary that is pointed as especially active and the main seismic source. Furthermore, the authors assert that there is moderate intraplate seismicity, especially at the Tore Seamount located west of Lisbon, and also identify several active onshore faults in Mainland Portugal.

[12] propose several seismic zones for Mainland Portugal and surrounding maritime area based on their classification as SCR or ACR, meaning Stable Continental Region

and Active Continental Region respectively, and on seismic criteria. The definition of seismic zones allowed the computation of a and b – *values* of the Gutenberg-Richter formula for each zone. In the context of seismic zonation, [9] define seven zones in the Azores Archipelago area based on seismicity criteria such as, the number of seismic events in a year and the magnitude of these events. The purpose of defining different seismic zones is to establish surface areas with homogeneous seismic characteristics, thus, allowing for similar risk evaluation.

The definition of seismic zones usually results in polygonal areas that follow straight and regular shapes. However, the complexity and interaction between geological structures, such as faults and interplate boundaries, tend to produce seismic events that when mapped by their epicentres create agglomerates of irregular shapes and sizes. Quantitative methods, such as clustering methods, may assist in the process of seismic zone definition since their purpose of application is to group objects by identifying similar characteristics and underlying patterns, [5]. In order to demonstrate the advantages of applying quantitative methods to the definition of seismic zones, [14] used a hierarchical clustering algorithm to group a grid of area units in Iran. The authors used a dataset composed of 25 distinct variables representing geological and geophysical parameters, measured in a grid of 1°longitude per 1°latitude constrained to the country’s borders. The usefulness of this method is highlighted stating that, the method is able to reveal tectonic evolution trends of a region and, additionally, allows for a continuous tectonic classification and assists in determining the interaction between the own variables used for classification. Moreover, [1] highlight once again the usefulness of clustering algorithms by developing a fuzzy clustering algorithm and applying it to a seismic catalog in Iran, emphasising its ability to recognise patterns.

With the purpose of categorising seismic data, [8] have implemented a density based algorithm, the DBSCAN, in data sourced from South African mine. In this work, the authors demonstrated that the algorithm is able to identify clusters of precursor events that precede bigger seismic events. The DBSCAN algorithm was initially presented by [2] and enables the identification of irregular shaped clusters, without the need of previously specifying the number of clusters to be created, thus, the resulting number of clusters is influenced only by the data. Besides, the algorithm also classifies some data points as outliers or noise. [6] applied the DBSCAN algorithm to seismic data, consisting of seismic event’s geographical location and depth, to identify seismic zones in India. The authors compared the resulting clusters to another seismic zonation used by an Indian governmental entity and concluded their similarity.

By identifying the necessity of the existence of algorithms that take into account spatial seismic data, the complex relationship between geological structures, namely the interaction between faults, and the that allow the identification of irregular and elongated shaped clusters, [3] present a two-phased clustering method to apply on seismic data, where the first phase is a modification of the DBSCAN algorithm based on the notion of “seismic mass” and the second phase is an agglomerating phase dropping time information. This two-phased method was referred to as the “SM-DBSCAN”, where “SM” stands for “seismic mass”, and was applied to data sourced from the Hellenic Arc.

The identification of seismic zones continues to be a relevant topic, mainly in regions with considerable seismic activity, as is the case of Portugal, and the application of quantitative methods has shown to be highly versatile, providing important results in a continuous fashion. Having said that, the aim of this work is to obtain seismic zones for Mainland Portugal by applying the DBSCAN algorithm to seismic data, considering a modified distance function, and confront the results with known geological structures

in the region and other seismic zonations defined by other authors.

This work will have the following structure: firstly, the dataset used is to be presented and there will be brief exploratory analysis made, secondly the DBSCAN algorithm is going to be presented as well as the proposed distance index. Subsequently the algorithm will be applied,... the results and a discussion shown, then the results will be confronted with known geological structures and compared to other seismic zonations.

2. Dataset

The dataset used in this work was provided by — , it contains information regarding seismic events in Mainland Portugal and its maritime surroundings. Given the complexity of the boundary between the Eurasian and African plates and the influence this structure has on the seismic activity felt by the Portuguese population [7], it is important include the maximum number of events that originate in the western margin of the country. Therefore, records with Longitude between 6° and 16° West and with Latitude between 34° and 44° North are included in the dataset. The dataset also contains records starting from 1900 to 1992. For each seismic event, one has information about its epicentre's geographical location, longitude and latitude, its magnitude on the Richter scale and the time of its occurrence in decimal years. However, there is a large number of events with missing magnitude recording. Table 1 summarises the dataset's characteristics.

Table 1. Dataset's characteristics

Dataset's characteristics	
Time period (in years)	1900 - 1992
Longitude interval considered	6°W - 16°W
Latitude interval considered	34°N - 44°N
Total number of records	4132
Total number of records with magnitude recording	1856
Total number of records with missing magnitude recording	2276

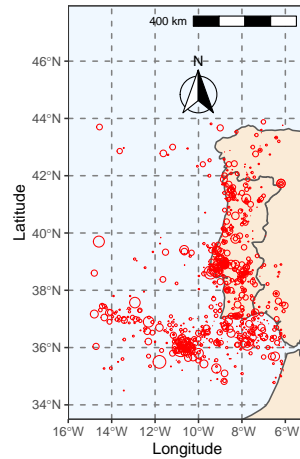


Figure 1. Seismic events in Mainland Portugal with magnitude recording, since 1900 until 1992.

3. Exploratory Analysis of the Seismic Data

The R software was used to perform the analysis presented in this section, along with a set of packages for the construction of the maps. The packages used are listed in the references.

Considering the exploratory analysis of the seismic data for the Azores Archipelago presented in [9] and the dataset used in this work, there can be defined the following variables: annual seismicity, time and size.

3.1. Annual seismicity

Annual seismicity can be defined as the number of seismic events (earthquakes) that occurred in the timeframe of a year. Figure 2 shows the annual seismicity when all data is considered and the annual seismicity when only records with magnitude information are considered.

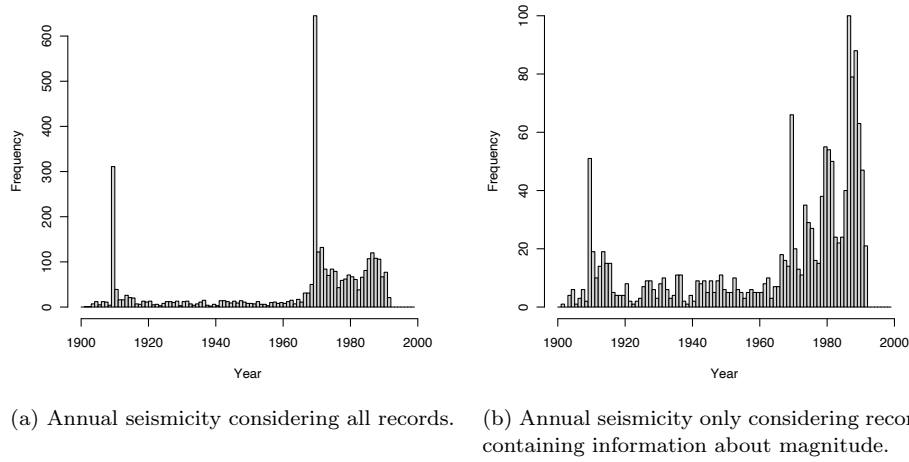


Figure 2. Annual seismicity.

Looking at figure 2 (a), there are two significant spikes in annual seismicity that, comparing with figure 2 (b), seem to result from an agglomerate of recordings without information about magnitude. In figure 2 (b) the frequency scale reduces by six-fold in comparison to (a), and there cease to be any significant spikes, nevertheless the annual seismicity continues to be fairly heterogeneous and there is a clear increase in the number of seismic events per year from mid-1960s onward. Regardless, from now on, all analysis and subsequent clustering in section 5 will focus on the data with information about magnitude.

Figure 3 displays the distribution of the number of earthquakes per year with recorded magnitude and table 2 summarises the statistical characteristics of the variable. It shows some variability, ranging from fewer than 5 events per year to 100 events in others.

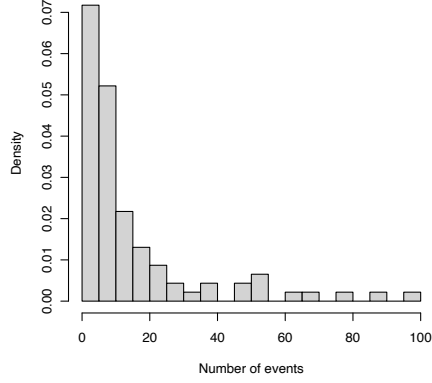


Figure 3. Histogram of the number of seismic events per year.

Table 2. Statistics of the number of seismic events per year

Statistic	Value
Mean	15.717
Standard deviation	19.998
Skewness	2.244
Maximum	100
Minimum	0
Quantile	
0.1	2
0.25	4.75
0.5	8
0.75	16.5
0.9	46.3

3.2. Time variable

Let DT be the time difference between two consecutive seismic events. This variable shows how often there is an earthquake in mainland Portugal and surrounding maritime area. The records have a time stamp measured as a decimal year in order to simplify an analytical analysis, thus, DT will also be measured in decimal years. Figure 4 shows the histogram of DT and table 3 summarises the statistical characteristics of the DT . On average, there is an earthquake every 0.062 decimal years, approximately 22 days, in the Portuguese Mainland and surrounding coast.

3.3. Size variable

The size variable relates to the magnitude of the seismic events. Let S be the magnitude, in the Richter scale, of a seismic event that occurred in the Portuguese Mainland or surrounding coast. Figure 5 displays the histogram of S and table 4 summarises the statistical characteristics of the S .

The S variable has an approximately normal distribution, averaging 3.5 magnitude on the Richter scale.

The variables defined in this section, as shown by [9], can be used as a measure of dissimilation between seismic zones when computed within defined sections. These

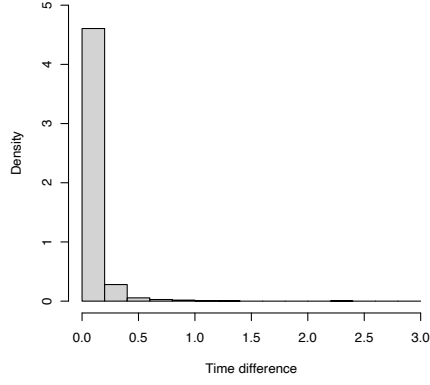


Figure 4. Histogram of DT .

Table 3. Statistics of DT

Statistic	Value
Mean	0.062
Standard deviation	0.132
Skewness	6.599
Maximum	2.294
Minimum	0
Quantile	
0.1	0.001
0.25	0.005
0.5	0.0196
0.75	0.063
0.9	0.153

sections may well be represented by clusters of events. Later in this work, the same variables will be used to characterise clusters, resulting from the application of the DB-SCAN algorithm, and a comparison between them will be made in order to determine their dissimilarity or not.

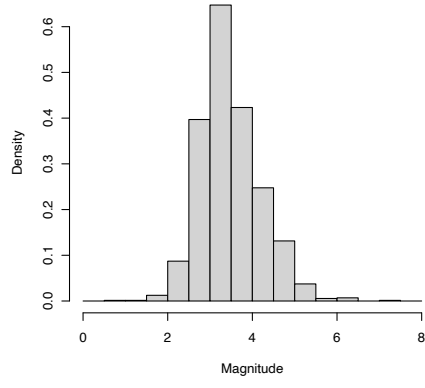


Figure 5. Histogram of S .

Table 4. Statistics of S

Statistic	Value
Mean	3.510
Standard deviation	0.706
Skewness	0.581
Maximum	7.3
Minimum	0.76
Quantile	
0.1	2.7
0.25	3.01
0.5	3.4
0.75	3.928
0.9	4.5

4. Methods

In this section the DBSCAN algorithm will be presented, as well as, all the definitions needed to understand it. Next, the proposed distance index is presented along with the reasoning behind it. Furthermore, seismic data was simulated to better display the intended results from applying the distance index.

4.1. DBSCAN

The method that is going to be used in this work in order to cluster seismic events is the DBSCAN algorithm. DBSCAN stands for Density Based Spatial Clustering of Applications with Noise and it was first presented by [2]. The authors' intention was to develop an efficient clustering algorithm that required little input from the user and could find clusters of irregular shapes. The DBSCAN algorithm relies on the notion of "density" to cluster objects. Two objects are clustered together if it is possible to reach one another without leaving a defined "dense" zone. In order to understand the underlying mechanics of the algorithm, one needs to comprehend several definitions firstly presented by [2]. In this work, the definitions presented are adapted from [2] and [4].

From now on, D is the set of points to be clustered.

Definition 4.1. A ϵ -neighbourhood of a given point $p \in D$, denoted by $N_\epsilon(p)$, is defined by,

$$N_\epsilon(p) = \{q \in D : \text{dist}(p, q) \leq \epsilon\} \quad (1)$$

where dist is any distance function.

The cardinality of a ϵ -neighbourhood of a point $p \in D$ determines the notion of density. For a given point, one can define the minimal number of points (minPts) that should be in the ϵ -neighbourhood in order for it to be called dense.

Definition 4.2. A point $p \in D$ is classified, given ϵ and minPts , as,

- a *core point*, if $N_\epsilon(p)$ is dense, i.e., $|N_\epsilon(p)| \geq \text{minPts}$,
- a *border point*, if p is not a *core point*, however it belongs to the neighbourhood of a *core point* $q \in D$, i.e., $p \in N_\epsilon(q)$,
- a *noise point*, otherwise.

The application of the DBSCAN algorithm also results in the identification of the so called *noise points*, a point that doesn't belong to any dense area.

Definition 4.3. A given point $q \in D$ is *directly density-reachable* from a point $p \in D$, given ϵ and minPts if,

- (1) $p \in N_\epsilon(q)$
- (2) $|N_\epsilon(q)| \geq \text{minPts}$.

Where 2 is the requirement for p to be a *core point*. Meaning that q belongs to a dense neighbourhood.

Definition 4.4. A given point $p \in D$ is *density-reachable* from a point $q \in D$, given ϵ and minPts , if there is a chain of points p_1, \dots, p_n , $p_1 = q$ and $p_n = p$, such that p_{i+1}

is *directly density-reachable* from p_i .

Definition 4.5. A given point $p \in D$ is *density-connected* to a point $q \in d$, given ϵ and $MinPts$, if there is a point $t \in D$ such that, p and q are *density-reachable* from a point $o \in D$.

Definition 4.6. A *density-based cluster* C is a non-empty subset of D that satisfies the following conditions,

- (1) Maximality: If $p \in C$ and q is *density-reachable* from p , then $q \in C$.
- (2) Connectivity: $\forall p, q \in C$, p is *density-connected* to q .

The following algorithm describes DBSCAN and was adapted from [13].

Let D be the set of points to be clustered, D_s the set of points already classified and S the set of points classified as *noise point*, given ϵ and $minPts$.

Algorithm 1: DBSCAN

Result: List \mathcal{L} with the cluster assignment to each point $p \in D$ through

```

    ordered pairs  $(i, p)$ 
    {Each cluster is identified by the index  $i$ };
     $i = 0$ ;
     $D_s = \emptyset$ ;
     $\mathcal{L} = \{\}$ ;
    for each point  $x \in D$  do
        if  $x \notin D_s$  then
             $D_s = D_s \cup \{x\}$ ;
            if  $|N_\epsilon(x)| < minPts$  then
                 $S = S \cup \{x\}$ ;
            else
                 $i = i + 1$ ;
                for each point  $p \in N_\epsilon(x)$  do
                     $\mathcal{L} = \mathcal{L} + \{(i, p)\}$ ;
                end
                for each point  $y \in N_\epsilon(x)$  e  $y \notin D_s$  do
                     $D_s = D_s \cup \{y\}$ ;
                    if  $|N_\epsilon(y)| \geq minPts$  then
                        for each point  $p \in N_\epsilon(y)$  do
                             $\mathcal{L} = \mathcal{L} + \{(i, p)\}$ ;
                            if  $p \in S$  then
                                 $S = S - \{y\}$ ;
                            end
                        end
                    end
                end
            end
        end
    end
end

```

The application of the DBSCAN relies solely on the input of ϵ , $minPts$ and a defined distance function. One can manipulate what a dense neighbourhood is by setting dif-

ferent values for ϵ and *minPts* or even by changing the distance function. Through the distance function, one is able to “encourage” the joint grouping of two given points by decreasing the distance between them or “discourage” their joint grouping by increasing the distance between them. This opens up a new dimension within DBSCAN, where by manipulating the distance function one can also manipulate the clustering, besides adjusting ϵ or *minPts*.

In the next subsection, a modified distance function for the clustering of seismic data will be presented, as well as its objective and the reasoning behind it.

4.2. Distance Index for Seismic Data Clustering

Considering that a seismic event i is characterised by (t_i, m_i, Lo_i, La_j) , one defines the distance index as:

$$\text{dist}(i, j) = k_t(t_i - t_j)^2 + (1 - k_s \max\{m_i, m_j\}) \times \text{Hav}((Lo_i, La_j), (Lo_j, La_j)) \quad (2)$$

where,

- t_i is the time of occurrence in decimal year;
- m_i is the magnitude;
- (Lo_i, La_j) are the geographical coordinates of the epicenter;

of a seismic event i . And,

- k_t is a scalar $k_t \geq 0$;
- k_s is a scalar $(0 \leq k_s < \frac{1}{M})$;
- Hav is the Haversine formula;

So as to the distance to remain positive, M corresponds the maximum magnitude of any given seismic event ever recorded. The Haversine Formula [11] determines the distance between two points on a sphere given their longitudes and latitudes.

As previously mentioned, by manipulating the distance between points, one is able to “encourage” the joint grouping of two given objects by decreasing the distance between them or “discourage” their joint grouping by increasing the distance between them. In the context of seismic clustering, the magnitude of an event, due to the interaction of geological structures, influences the occurrence of other seismic events and this influence is directly proportional to the magnitude. The stronger the event, the bigger its influence in surrounding structures. As stated by [3], the “importance” of a seismic event depends on its magnitude and, thus, magnitude should be taken into account when clustering. Furthermore, [10] has shown through statistical methods that, the occurrence of an earthquake depends on the location of preceding earthquakes. In order to capture this influence through the **dist** index and reflect it in the clustering, two events will be “virtually” closer, the bigger the magnitude of either events. In other words, the distance between two events is proportionately inverse to the magnitude of either events. This refers to the $(1 - k_s \max\{m_i, m_j\})$ part of the calculation, where k_s determines how much influence will magnitude have on the distance index value. On the other hand, [10] also showed that, there is a clear reduction of the dependence between events over time, therefore, the time difference between events should also be taken into account considering that, the influence of a seismic event over another lessens as time passes. To capture this effect, the distance between two events will be proportional to the time difference between occurrences, meaning events will be

“virtually” further apart the larger the time difference between them. This relates to the $k_t(t_i - t_j)^2$ part of the calculation, where k_t determines how much influence will the time difference have on the distance. Additionally, all other effects aside, the distance will also depend on the Haversine distance between the epicenter’s geographical location of two seismic events. By including the time of occurrence of the seismic events, the clustering will be in a three dimensional space which, if we consider only the longitude latitude plane and project the clusters, will result in overlapping clusters. The first stage of the method presented in [3] also results in spatially overlapping cluster. The authors state that this could help identify closely located geological structures, like faults, whose seismic activity could be separated in time or even identify the seismic activity of one geological structure in different time periods. The authors then propose a second stage with an agglomerate method that uses a single linkage procedure to fuse the clusters and project them onto the space plane in order to define seismic zones. The same procedure could be applied in this case.

So as to illustrate the expected results of the proposed modified distance index, a simulation of seismic data was made. The simulation was based on the seismic data for Mainland Portugal. Intervals for the longitude and latitude were defined and the events within that region were selected. Then, epicenters’ geographical coordinates were sampled from the selected events, with longitude and latitude sampled one by one in order not to have duplicated epicenters. Magnitude for the simulated events was also sampled from the selected event’s magnitude. The time of occurrence for the simulated seismic events was done by sampling for three different normal distributions with averages 1920, 1950 and 1980 respectively, and all with standard deviation 6. There were simulated 180 seismic events. Figure 6 displays the simulated data in a three dimensional scatter plot, where the x and y axis are longitude and latitude, respectively, and z axle represents the time. The points’ diameter is larger the bigger its magnitude.

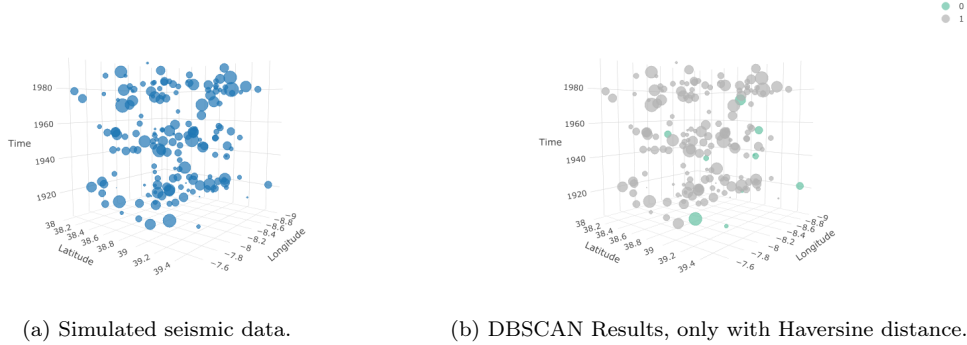


Figure 6. Simulation Results.

Figure 6(a) clearly displays, at least, three distinct cluster in the time. By applying the BDSCAN algorithm with $\epsilon = 20$ and $minPts = 5$ to the simulated data, considering only the Haversine distance between the events’ epicenters, the results, displayed in figure 6(b), show 1 cluster was obtained, which corresponds to the points displayed in grey, while the green points represent noise. There were identified 10 seismic events as noise.

Figure 7 displays the results from the DBSCAN algorithm now using the proposed distance index to determine the distances between the seismic events. There were

obtained three clusters clearly separated in time, with some parameter tuning. In this particular case, it was considered $k_t = 0.5$ and it definitely has the expected effect, clusters separated in time. On the other hand, these results show only 7 seismic events were identified as noise compared to the 10 obtained when only the Haversine distance was used. This is clearly an effect of the k_s scalar, which in this case was considered $k_s = 0.05$. The parameters for the DBSCAN algorithm used were $\epsilon = \frac{1}{1-k_t} \times 20$ and $minPts = 5$, where the unit measure for ϵ is the kilometre. The increase in the distance index generated by the time difference between two seismic events needed to be accounted for when applying the DBSCAN algorithm, especially when defining the ϵ -neighbourhood. Trough experimentation, applying the $\frac{1}{1-k_t}$ factor to the previously defined ϵ value showed the best results for clustering when using the distance index.

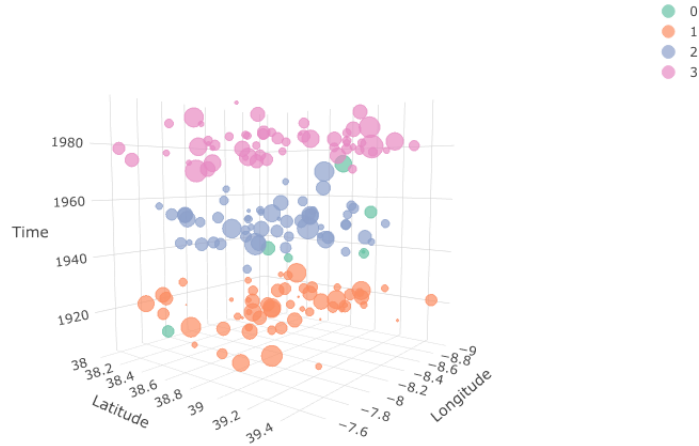


Figure 7. Results with the distance index on simulated data.

The proposed distance index, with simulated data, shows that intended behaviour was achieved, separated clusters in time and joint grouping of events, that would otherwise be separated, through the influence of magnitude.

In the next section, the DBSCAN algorithm will be applied to the seismic data of Mainland Portugal, considering only the Haversine distance, and the results will be compared to a seismic zonation proposed by another author. From the clustering results, there will be an exploratory analysis made considering some of the clusters, taking into account the seismic variables defined previously in this work, in order to determine their dissimilarity. Then, the DBSCAN algorithm, in conjunction with the proposed distance index, is applied to the Portuguese Mainland’s seismic data and the results displayed.

5. Results & Discussion

The DBSCAN algorithm used to obtain the results presented in this and in the previous sections was implemented by [4] for the R software.

5.1. Seismic Zones

In this first subsection, the DBSCAN algorithm was applied considering only the spatial distance between seismic events, in other words, only the Haversine distance was considered to determine the distance between two distinct seismic events, meaning $k_t = 0, k_s = 0$. There were obtained 29 clusters, along with the 484 records considered noise. For these results, it was used $\epsilon = 15$ and $minPts = 5$. In this case ϵ is measured in kilometres.

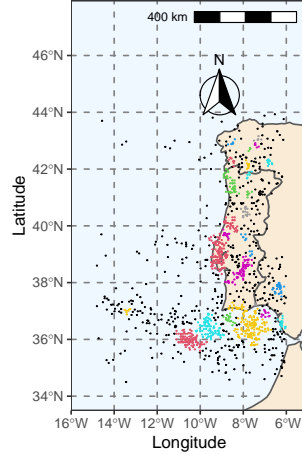


Figure 8. Seismic event clusters for Mainland Portugal.

Figure 8 displays the clusters plotted on Mainland Portugal’s map, while figure 9 shows the zonations proposed by [12]. There is a clear match between the cluster results and the zonation presented. In figure 9, the Active Continental Regions (ACR) are displayed with a striped pattern while the Stable Continental Regions (SCR) are transparent. ACR have more seismic activity than SCR, this seismic activity results in bigger denser clusters in these regions. In addition, the clusters that overlap Lisbon Metropolitan Area and the above coast fit plainly with the zone identified as S05 in figure 9. Furthermore, the purple cluster located close to the middle of Portugal coincides with a known geological fault, designated as the Messejana Fault. DBSCAN is clearly able to identify irregular shaped clusters despite the complex shape of geological structures like the Messejana Fault. Besides this, there’s another relevant result that needs to be taken into account when using the DBSCAN algorithm, the noise points. No clusters are formed outside of the considered seismic zones meaning, all seismic events outside the identified zones are marked as noise. When a neighbourhood of a given point doesn’t comply with the definition of a dense zone, that point is considered noise. Translating this in the context of seismic events, if in a given zone there is a smaller number of seismic events, these seismic events will be considered noise and the zone will be considered a background zone. This shows an advantage of the way DBSCAN works and proves to be very useful in the case of seismic clustering.

The very good fit between clusters and the considered seismic zones obtained, suggest that DBSCAN algorithm is a practical and fast tool for seismic zonation. If applied continuously, considering the most recent records of seismic events, the resulting clusters could provide a clear picture of how seismic behaviour is evolving in certain areas, moreover, if this evolution happens to be geographical, the resulting clusters may include new areas indicating its growth or translation. This could be used to amend or

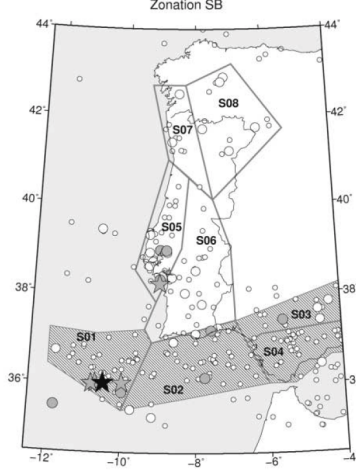


Figure 9. Seismic zonation for Mainland Portugal, adapted from [12].

propose changes to already defined seismic zones. For instances, in the results presented in this work, the south of Portugal is contained in a cluster that is included in the seismic zone identified as S02 in figure 9. This could suggest that more of the area of southern Portugal should be included in S02.

Now, with the clustering results, there will be an exploratory analysis made considering some of the clusters, taking into account the seismic variables previously defined: annual seismicity, time and size. This analysis may indicate that, the clusters differ between each other, in regards to the seismic variables. This would mean that each zone overlapped by a cluster has, inherently, a characteristic seismic behaviour. To do this analysis, the clusters composed by more than 50 events were selected in order to have a representative set of events for each one. Figure 10 displays the clusters, and their respective identification, that are composed by more than 50 events and will be object of the analysis.

Table 5 presents the values for each seismic variable for each of the five clusters selected. There doesn't seem to be very clear differences between the clusters in regards to the seismic variables. Cluster C_1 shows to have more seismic events per year in average, while the C_{17} cluster displays the largest variance in annual seismicity. Furthermore, considering the time between seismic events, C_1 , C_{17} and C_{20} have approximately, on average, a seismic event every 5 months whereas C_5 has a seismic event every 11 months, on average. With respect to the magnitude of the seismic events, the values for the variable Size are very similar between clusters. They all average a magnitude of approximately 3.5 on the Richter scale, however the C_{17} cluster has the events with the largest magnitude.

5.2. Clustering results with the proposed distance index

In this subsection, the DBSCAN is applied to seismic data for Mainland Portugal, where the distances between seismic events are given by the distance index proposed. A set of values for k_t and k_s was used to understand the effect of these values on the clustering. Furthermore, the parameters of the DBSCAN algorithm used were the same for each combination of k_t and k_s , those being $\epsilon = \frac{1}{1-k_t} \times 15$ and $minPts = 5$. Note

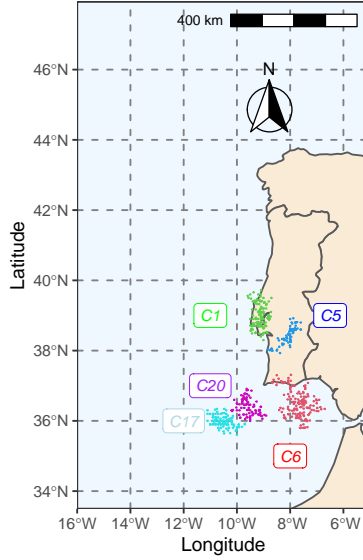


Figure 10. Seismic event clusters for Mainland Portugal, with more than 50 events.

that the adjustment factor $\frac{1}{1-k_t}$ was also taken into account. Table 6 summarises the results for each combination of values. Firstly, considering the number of noise points, it can be seen that increasing both k_t and k_s leads to a decrease in the number of noise points, being k_s the factor that has the most influence in this aspect. Secondly, focusing now on the number of resulting clusters, these too are reduced when k_s and k_t are increased, again k_s having more influence than k_t in the reduction of the number of clusters.

When $k_t = 0.25$ and $k_s = 0.1$, the number of clusters generated is the equal to the number of clusters obtained when only the Haversine distance between was considered, however the number of noise events decreased by two hundreds, result of the k_s factor. Despite the number of clusters being equal, the proposed distance index does not fail to separate events in the time dimension. Figure 11 shows the clusters that resulted from the DBSCAN algorithm with the proposed distance index with parameters $k_t = 0.25$, $k_s = 0.1$, $\epsilon = \frac{1}{1-k_t} \times 15$ and $minPts = 5$, without the events considered noise. There are clear geographically overlapping clusters, meaning the seismic events originated from the same geological structures are grouped together but in different time frames.

Figure 12 displays the clusters in three dimensional space, in order to better present the time separated clusters that overlap geographically. This allows for the identification of distinct sets of activity over time from the same geological structures, which could be used to study how often this activity is expected to happen. Moreover, the k_s factor aggregates seismic events to clusters that otherwise would be considered as noise, hence the reduction in the number of events classified as noise. This could help relate earthquakes with epicentres in background zones to bigger events that originated elsewhere. Comparing these clusters to the previously presented seismic zones (figure 9), the earthquakes generated in the zones classified as Active Continental Regions by [12] are clustered together, fitting the defined ACR very well, aggregating almost all of the seismic events of these regions into the same cluster. This suggests that, all of the zones considered to be ACR influence each other trough large magnitude earthquakes.

Table 5. Statistics of variables Annual Seismicity, Time (DT) and Size (S) for each cluster.

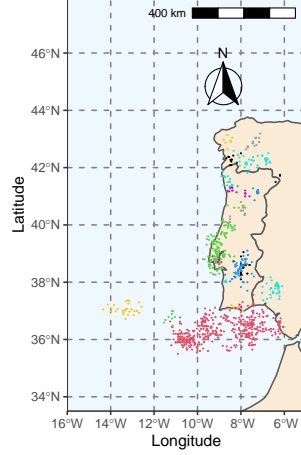
Variable Statistics	C_1	C_5	C_6	C_{17}	C_{20}
Annual Seismicity					
Mean	2.011	1.033	1.728	1.217	0.957
Standard deviation	1.386	3.935	2.794	5.302	2.158
Skewness	6.021	1.582	1.976	8.025	3.313
Maximum	34	6	12	49	12
Minimum	0	0	0	0	0
Quantile					
0.1	0	0	0	0	0
0.25	0	0	0	0	0
0.5	1	1	1	0	0
0.75	2	1	1.25	0	1
0.9	5	3	7	3	3
Time					
Mean	0.475	0.911	0.565	0.442	0.431
Standard deviation	0.724	1.215	1.123	1.84	0.775
Skewness	2.263	2.894	3.741	6.586	3.483
Maximum	3.654	8.264	8.383	14.026	4.699
Minimum	0	0	0	0	0
Quantile					
0.1	0.003	0.019	0.012	0	0.024
0.25	0.020	0.151	0.049	0.002	0.051
0.5	0.157	0.489	0.147	0.039	0.177
0.75	0.651	1.343	0.452	0.227	0.392
0.9	1.346	2.64	1.825	0.552	0.977
Size					
Mean	3.631	3.508	3.48	3.853	3.432
Standard deviation	0.849	0.68	0.623	0.677	0.493
Skewness	0.439	0.309	1.034	1.566	0.997
Maximum	6.25	5.34	6.2	7.3	4.8
Minimum	1.99	2.03	2.38	2.7	2.6
Quantile					
0.1	2.474	2.728	2.872	3.1	2.921
0.25	3.1	3.065	3	3.4	3.1
0.5	3.54	3.4	3.3	3.7	3.3
0.75	4.24	3.925	3.8	4.2	3.585
0.9	4.696	4.506	4.492	4.69	4.23

6. Conclusion

In this work, in a first phase the DBSCAN algorithm is applied to a catalog of seismic events that occurred in the Portuguese Mainland and surrounding maritime area, taking only into account the geographic distance, more specifically the Haversine distance, between the epicentres of seismic events. The resulting clusters are used to identify seismic zones for the area under consideration which are then compared to other seismic zonation already defined in the literature. In this case, the seismic zonation used for comparison was proposed by [12]. The resulting clusters reveal a very good fit with the zonation considered. Furthermore, no seismic event clusters were identified outside of the seismic zones considered, as the events outside these zones were considered noise by the DBSCAN algorithm which, in geological terms, translates to have had originated in so-called background zones. Background zones include small numbers of events when compared to other areas. These results validate the DBSCAN algorithm as a useful analytical tool for identifying seismic zones. With little information, con-

Table 6. Results from DBSCAN with the proposed distance index for different values of k_t and k_s .

	$k_t = 0.20$		$k_t = 0.25$	
	Number of clusters	Number of noise events	Number of clusters	Number of noise events
$k_s = 0.05$	40	641	35	628
$k_s = 0.1$	32	456	29	450

**Figure 11.** Seismic event clusters for Mainland Portugal, with proposed distance index ($k_t = 0.25$ and $k_s = 0.1$).

sidering only the epicentre’s longitude and latitude, and some parameter tuning, the algorithm provides good results when compared to seismic zones defined by geological experts. The notion of noise, inherent to the DBSCAN algorithm, shows to have a translation in a geological dimension, converting into background zones. Given the resulting clusters, the ones composed by more than 50 seismic events were object of a exploratory analysis considering the seismic variables: annual seismicity, time and size. This analyses didn’t show differences between the clusters by observing the statistics calculated, hence, future work should focus on this exploratory analysis in order to determine the dissimilarity between the clusters through statistical methods.

In a second phase, a distance index to be applied along with the DBSCAN algorithm is proposed. By manipulating the distance between two given points, one can encourage the joint clustering of those points. In a seismic zoning context, given that the size of an event influences other events and that this influences fades over time, one wants to determine how this influence can affect the clustering of events. This was the reasoning behind the proposed distance index. A test was made with simulated data in order to determine if the intended behaviour was achieved. The DBSCAN algorithm was applied considering the distances between seismic events calculated with the distance index proposed and, with some parameter tuning, the results show exactly what was intended, a separation in time of three geographically overlapping clusters. Subsequently, the distance index was used to cluster the seismic events for Mainland Portugal. A set of experimentation are made to determine the influence of k_t and k_s on the cluster formation. The increase of both factors determines the decrease in the number of resulting clusters and the decrease in the number of events classified as noise. The application of the distance index in the seismic catalog for Mainland Portugal has clearly interesting results, as it enables the identification of distinct clusters in time that overlap geographically. Further work should focus more on the impact of both

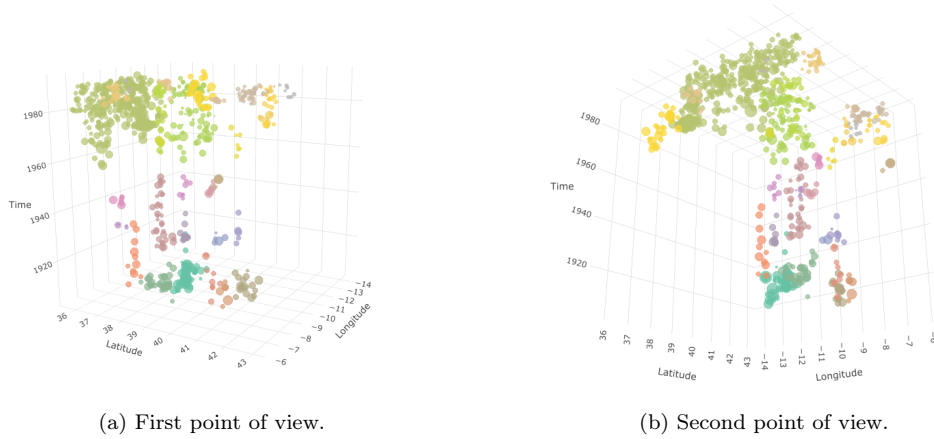


Figure 12. Results for Mainland Portugal’s seismic data, DBSCAN with the proposed distance index.

factors k_t and k_t on clustering results and what could be done with the results achieved by combining the distance index proposed with the DBSCAN algorithm, namely an analysis of how often should a cluster of seismic events be identified provided they originate from a specific geological structure.

Furthermore, it should be noted that the disregard of seismic events that were missing information about magnitude crippled the results. There was information that was wasted despite doubts about its quality. Although, the reasoning behind the proposed distance index relies heavily on the information about a seismic event’s magnitude, thus, a choice was made to ignore records with missing magnitude information.

References

- [1] A. Ansari, A. Noorzad, and H. Zafarani, *Clustering analysis of the seismic catalog of Iran*, Computers and Geosciences 35 (2009), pp. 475–486.
- [2] M. Ester, H.P. Kriegel, J. Sander, and X. Xu, *A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise*, Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (1996), pp. 226–231.
- [3] G. Georgoulas, A. Konstantaras, E. Katsifarakis, C.D. Stylios, E. Maravelakis, and G.J. Vachtsevanos, *"Seismic-mass" density-based algorithm for spatio-temporal clustering*, Expert Systems with Applications 40 (2013), pp. 4183–4189.
- [4] M. Hahsler, M. Piekenbrock, and D. Doran, *Dbscan: Fast density-based clustering with R*, Journal of Statistical Software 91 (2019).
- [5] A.K. Jain, *Data clustering: 50 years beyond K-means*, Pattern Recognition Letters 31 (2010), pp. 651–666. Available at <http://dx.doi.org/10.1016/j.patrec.2009.09.011>.
- [6] N.A. Karri, M. Yousuf Ansari, and A. Pathak, *Identification of seismic zones of India using DBSCAN*, 2018 International Conference on Computing, Power and Communication Technologies, GUCON 2018 (2019), pp. 65–69.
- [7] L. Pinheiro, R. Wilson, R. Pena dos Reis, R. Whitmarsh, and A. Ribeiro, *The western Iberia Margin: a geophysical and geological overview*, Proceedings of the Ocean Drilling Program, 149 Scientific Results 149 (1996).
- [8] D. Rebuli and S. Kohler, *Using clustering algorithms to assist short-term seismic hazard analysis in deep South African mines*, Proceedings of the Seventh International Conference on Deep and High Stress Mining (2014), pp. 699–708.
- [9] M.C. Rodrigues and C.S. Oliveira, *Seismic zones for Azores based on statistical criteria*, Natural Hazards and Earth System Sciences 13 (2013), pp. 2337–2351.
- [10] M.C.M. Rodrigues and C.S. Oliveira, *Considering spatial memory to estimate seismic risk: the case of the Azores Archipelago*, GEM - International Journal on Geomathematics 11 (2020), pp. 1–15. Available at <https://doi.org/10.1007/s13137-020-00152-0>.
- [11] R.W. Sinnott, *Virtues of the Haversine*, Sky and Telescope 68 (1984), p. 158.
- [12] S.P. Vilanova and J.F. Fonseca, *Probabilistic seismic-hazard assessment for Portugal*, Bulletin of the Seismological Society of America 97 (2007), pp. 1702–1717.
- [13] P. Viswanath and R. Pinkesh, *L-DBSCAN: A fast hybrid density based clustering method*, Proceedings - International Conference on Pattern Recognition 1 (2006), pp. 912–915.
- [14] A. Zamani and N. Hashemi, *Computer-based self-organized tectonic zoning: A tentative pattern recognition for Iran*, Computers and Geosciences 30 (2004), pp. 705–718.