

# Earthquake Prediction Using Data Mining

Radhika Kulkarni<sup>1</sup>, Rashmi Kulkarni<sup>2</sup>

<sup>1</sup>Cummins College of Engineering, Pune, BE ENT C

<sup>2</sup>P.E.S Modern College of Engineering, Pune, ME Computers

**Abstract:** *Data mining consists of evolving set of techniques that can be used to extract valuable information and knowledge from massive volumes of data. Data mining research & tools have focused on commercial sector applications. Only a few data mining research have focused on scientific data. This paper aims at further data mining study on scientific data. This paper highlights the data mining techniques applied to mine for surface changes over time (e.g. Earthquake rupture). The data mining techniques help researchers to predict the changes in the intensity of volcanoes. This paper uses predictive statistical models that can be applied to areas such as seismic activity, the spreading of fire. The basic problem in this class of systems is unobservable dynamics with respect to earthquakes. The space-time patterns associated with time, location and magnitude of the sudden events from the force threshold are observable. This paper highlights the observable space time earthquake patterns from unobservable dynamics using data mining techniques, pattern recognition and ensemble forecasting. Thus this paper gives insight on how data mining can be applied in finding the consequences of earthquakes and hence alerting the public.*

**Keywords:** data mining, seismic activity, earthquakes

## 1. Introduction

The field of data mining has evolved from its roots in databases, statistics, artificial intelligence, information theory and algorithms into a core set of techniques that have been applied to a range of problems. Computational simulation and data acquisition in scientific and engineering domains have made tremendous progress over the past two decades. A mix of advanced algorithms, exponentially increasing computing power and accurate sensing and measurement devices have resulted in more data repositories.

Advanced technologies in networks have enabled the communication of large volumes of data across the world. This results in a need of tools & Technologies for effectively analyzing the scientific data sets with the objective of interpreting the underlying physical phenomena. Data mining applications in geology and geophysics have achieved significant success in the areas as weather prediction, mineral prospecting, ecology, modeling etc and finally predicting the earthquakes from satellite maps.

An interesting aspect of many of these applications is that they combine both spatial and temporal aspects in the data and in the phenomena that is being mined. Data sets in these applications comes from both observations and simulation. Investigations on earthquake predictions are based on the assumption that all of the regional factors can be filtered out and general information about the earthquake precursory patterns can be extracted.

Feature extraction involves a pre selection process of various statistical properties of data and generation of a set of seismic parameters, which correspond to linearly independent coordinator in the feature space. The seismic parameters in the form of time series can be analyzed by using various pattern recognition techniques.

Statistical or pattern recognition methodology usually performs this extraction process. Thus this paper gives insight of mining the scientific data.

### Data mining - Definitions

- Data mining is defined as process of extraction of relevant data and hidden facts contained in databases and data warehouses.
- It refers to find out the new knowledge about an application domain using data on the domain usually stored in the databases. The application domain may be astrophysics, earth science or about solar system.

Data mining techniques support to identify nuggets of information and extracting this information in such a way that, this will support in decision making, prediction, forecasting and estimation.

### Data Mining Goals

- Bring together representatives of the data mining community and the domain science community so that they can understand the current capabilities and research objectives of each other communities related to data mining.
- Identify a set of research objectives from the domain science community that would be facilitated by current or anticipated data mining techniques.
- Identify a set of research objectives for the data mining community that could support the research objectives of the domain science community.

### Data Mining Models

Data mining is used to find patterns and relationships in data patterns. The relationships in data patterns can be analyzed via 2 types of models.

- 1) **Descriptive models:** Used to describe patterns and to create meaningful subgroups or clusters.
- 2) **Predictive models:** Used to forecast explicit values, based upon patterns in known results. Paper concentrate

on predictive model. In large databases data mining and knowledge discovery comes in two flavors:

### 1. Event based Mining:

#### • **Known events / known algorithms:**

Use existing physical models (descriptive models and algorithms) to locate known phenomena of interest either spatially or temporally within a large database.

#### • **Known events / unknown algorithms:**

Use pattern recognition and clustering properties of data to discover new observational (physical) relationships (algorithms) among known phenomena.

#### • **Unknown events / known algorithms:**

Use expected physical relationships (predictive models, Algorithms) among observational parameters of physical phenomena to predict the presence of previously unseen events within a large complex database. Paper's main line of action is on this type.

#### • **Unknown events / unknown algorithms:**

Use thresholds or trends to identify transient or otherwise unique events and therefore to discover new physical phenomena.

### 2. Relationship based Mining

• **Spatial Associations:** Identify events (e.g. astronomical objects) at the same location. (e.g. same region of the sky)

• **Temporal Associations:**

Identify events occurring during the same or related periods of time.

• **Coincidence Associations:**

Use clustering techniques to identify events that are co-located within a multi-dimensional parameter space.

• **User requirements for data mining in large scientific databases**

• **Cross identifications:** Refers to the classical problem of associating the source list in one database to the source list in another.

• **Cross correlation:** Refers to the search for correlations, tendencies, and trends between physical parameters in multidimensional data usually across databases. Main line of concentration of this paper is on this method.

• **Nearest neighbour identification :** Refers to the general application of clustering algorithms in multidimensional parameter space usually within a database.

• **Systematic data exploration:** Refers to the application of broad range of event based queries and relationship based queries to a database in making a serendipitous discovery of new objects or a new class .

### Data Mining Techniques

The various data mining techniques are

- 1) Statistics
- 2) Clustering
- 3) Visualization
- 4) Association
- 5) Classification & Prediction
- 6) Outlier analysis
- 7) Trend and evolution analysis

### 1) **Statistics:**

- a) Data cleansing i.e. the removal of erroneous or irrelevant data known as outliers.
- b) EDA Exploratory data analysis e.g. frequency counts histograms.
- c) Attribute redefinition e.g. bodies mass index.
- d) Data analysis is a measure of association and their relationships between attributes interestingness of rules, classification ,prediction etc.

### 2) **Visualization:**

Enhances EDA , make patterns visible in different views .

### 3) **Clustering(cluster analysis):**

- a) Clustering is a process of grouping similar data. The data which is are not part of clustering are called as outliers. How to cluster in different conditions,
- b) Class label is unknown: Group related data to form new classes, e.g., cluster houses to find distribution patterns
- c) Clustering based on the principle: maximizing the intra-class similarity and minimizing the interclass similarity
- d) It provides subgroups of population for further analysis or action –very important when dealing with large databases.

### 4) **Association (correlation and causality) :** Mining Classification and Prediction association rules finds the interesting correlation relationship among large databases .

- a) Finding models (functions) that describe and distinguish classes or concepts for future prediction e.g., classify countries based on climate, or classify cars based on gas mileage
- b) Presentation: decision-tree, classification rule, neural network
- c) Prediction: Predict some unknown or missing numerical values

### 5) **Outlier analysis**

Outlier: A data object that is irrelevant to general behaviour of the data ,it can be considered as an exception but is quite useful in fraud detection in rare events analysis

### 6) **Trend and evolution analysis**

- a) Trend and deviation: regression analysis
- b) Sequential pattern mining, periodicity analysis
- c) Similarity-based analysis
- d) Papers main intention is on clustering & visualization technique for predicting the earthquakes.

### Earthquake Prediction

- 1) Ground water levels
- 2) Chemical changes in Ground water
- 3) Radon Gas in Ground water wells.

### Ground Water Levels:-

Changing water levels in deep wells are recognized as precursor to earthquakes. The pre-seismic variations at observation wells are as follows.

- 1) A gradual lowering of water levels at a period of months or years.
- 2) An accelerated lowering of water levels in the last few months or weeks preceding the earthquake.

- 3) A rebound, where water levels begin to increase rapidly in the last few days or hours before the main shock.

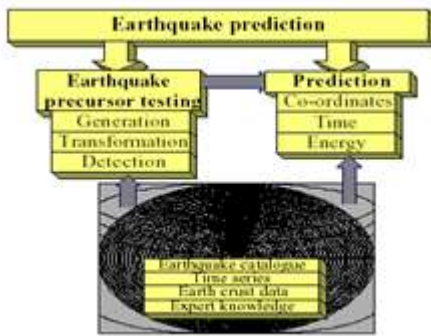
**Chemical Changes in Ground water**

- 1) The Chemical composition of ground water is affected by seismic events.
- 2) Researchers at the university of Tokyo tested the water after the earthquake occurred, the result of the study showed that the composition of water changed significantly in the period around earthquake area.
- 3) They observed that the chloride concentration is almost constant.
- 4) Levels of sulphate also showed a similar rise.

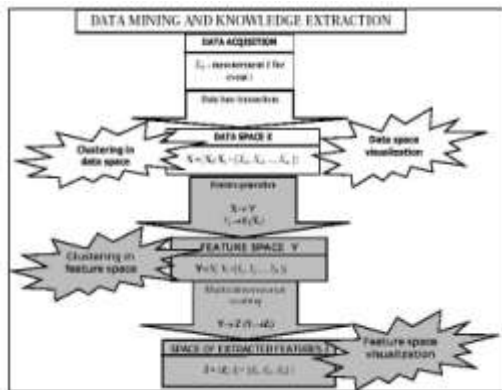
**Radon Gas in Ground water wells.**

- An increase level of radon gas in wells is a precursor of earthquakes recognized by research group.
- Although radon has relatively a short half life and is unlikely to seep the surface through rocks from the depths at which seismic is very soluble in water and can routinely be monitored in wells and springs often radon levels at such springs show reaction to seismic events and they are monitored for earthquake predictions..
- There is no effective solution to the problem.
- To solve this problem earthquake catalogs, geo-monitoring time series data about stationary seismotectonic properties of geological environment and expert knowledge and hypotheses.
- To solve this problem earthquake catalogs, geo-monitoring time series data about stationary seismotectonic properties of geological environment and expert knowledge and hypotheses about earthquake precursors

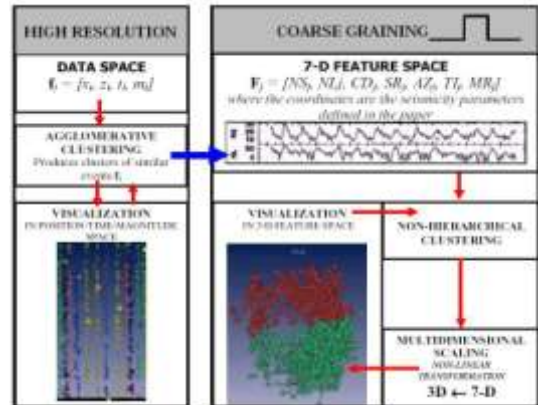
This proposes a multi-resolution approach, which combines local clustering techniques in the data space with a non-hierarchical clustering in the feature space. The raw data are represented by n-dimensional vector  $X_i$  of measurements  $X_k$ . The data space can be searched for patterns and can be visualized by using local or remote pattern recognition and by advanced visualization capabilities. The data space  $X$  is transformed to a new abstract space  $Y$  of vectors  $Y_j$ . The coordinates  $Y_l$  of these vectors represent nonlinear functions of measurements  $X_k$ , which are averaged in space and time in given space-time windows. This transformation allows for coarse graining of data (data quantization), amplification of their characteristic features and suppression of the noise and other random components. The new features  $Y_l$  form a N-dimensional feature space. We use multi-dimensional scaling procedures for visualizing the multi-dimensional events in 3D space. This transformation allows a visual inspection of the N-dimensional feature space. The visual analysis helps greatly in detecting subtle cluster structures which are not recognized by classical clustering techniques, selecting the best pattern detection procedure used for data clustering, classifying the anonymous data and formulating new hypothesis.



**Figure 1.1:** Earthquake Prediction



**Figure 1.2:** Prediction flow



**Figure 1.3:** Working

Clustering schemes Clustering analysis is a mathematical concept whose main role is to extract the most similar separated sets of objects according to a given similarity measure. This concept has been used for many years in pattern recognition. Depending on the data structures and goals of classification, different clustering schemes must be applied.

In our new approach we use two different classes of clustering algorithms for different resolutions. In data space we use agglomerative schemes, such as modified Mutual Nearest Neighbour algorithm (MNN). This type of clustering extracts the localized clusters in the high resolution data space. In the feature space we are searching for global clusters of time events comprising similar events from the whole time interval.

The non-hierarchical clustering algorithms are used mainly for extracting compact clusters by using global knowledge about the data structure. We use improved mean based schemes, such as a suite of moving schemes, which uses the k-means procedure and four strategies of its tuning by moving the data vectors between clusters to obtain a more precise location of the minimum of the goal function:



where  $z_j$  is the position of the center of mass of the cluster  $j$ , while  $x_i$  are the feature vectors closest to  $z_j$ . To find a global minimum of function  $J()$ , we repeat the clustering procedures at different initial conditions. Each new initial configuration is constructed in a special way from the previous results by using the methods. The cluster structure with the lowest  $J(w, n)$  minimum is selected.

### Hierarchical Clustering Methods

A hierarchical clustering method produces a classification in which small clusters of very similar molecules are nested within larger clusters of less closely-related molecules. Hierarchical *agglomerative* methods generate a classification in a bottom-up manner, by a series of agglomerations in which small clusters, initially containing individual molecules, are fused together to form progressively larger clusters. Hierarchical agglomerative methods are often characterized by the shape of the clusters they tend to find, as exemplified by the following range: single-link - tends to find long, straggly, chained clusters; Ward and group-average - tend to find globular clusters; complete-link - tends to find extremely compact clusters. Hierarchical *divisive* methods generate a classification in a top-down manner, by progressively sub-dividing the single cluster which represents an entire dataset. Monothetic (divisions based on just a single descriptor) hierarchical divisive methods are generally much faster in operation than the corresponding polythetic (divisions based on all descriptors) hierarchical divisive and hierarchical agglomerative methods, but tend to give poor results. One problem with these methods is how to choose which clusters or partitions to extract from the hierarchy because display of the complete hierarchy is not really appropriate for data sets of more than a few hundred compounds.

### Non Hierarchical Clustering Methods

A non-hierarchical method generates a classification by partitioning a dataset, giving a set of (generally) non-overlapping groups having no hierarchical relationships between them. A systematic evaluation of all possible partitions is quite infeasible, and many different heuristics have described to allow the identification of good, but possibly sub-optimal, partitions. Three of the main categories of non-hierarchical method are single-pass, relocation and nearest neighbour. Single-pass method (e.g. Leader) produce clusters that are dependent upon the order in which the compounds are processed, and so will not be considered further. Relocation methods, such as  $k$ -means, assign compounds to a user-defined number of seed clusters and then iteratively reassign compounds to produce the better clusters result. Such methods are prone to reaching local optimum rather than a global optimum, and it is generally not possible to determine when or where the global optimum solution has been reached. Nearest neighbour methods, such as the Jarvis-Patrick method, assign compounds to the same cluster as some number of their nearest neighbours. User-defined parameters determine how many nearest neighbours need to be considered, and the necessary level of similarity between nearest neighbour lists. Other non-hierarchical methods are generally inappropriate for use on large, high-dimensional datasets such as those used in chemical applications.

### Data mining Applications

- In Scientific discovery – super conductivity research, For Knowledge Acquisition.
- In Medicine – drug side effects, hospital cost analysis, genetic sequence analysis, prediction etc.
- In Engineering – automotive diagnostics expert systems, fault detection etc.,
- In Finance – stock market perdition, credit assessment, fraud detection etc.

## 2. Future Enhancements

The future of data mining lies in predictive analytics. The technology innovations in data mining since 2000 have been truly Darwinian and show promise of consolidating and stabilizing around predictive analytics. Nevertheless, the emerging market for predictive analytics has been sustained by professional services, service bureaus and profitable applications in verticals such as retail, consumer finance, telecommunications, travel and leisure, and related analytic applications. Predictive analytics have successfully proliferated into applications to support customer recommendations, customer value and churn management, campaign optimization, and fraud detection. On the product side, success stories in demand planning, just in time inventory and market basket optimization are a staple of predictive analytics. Predictive analytics should be used to get to know the customer, segment and predict customer behaviour and forecast product demand and related market dynamics. Finally, they are at different stages of growth in the life cycle of technology innovation.

## 3. Conclusion

The problem of earthquake prediction is based on data extraction of pre-cursory phenomena and it is highly challenging task various computational methods and tools are used for detection of pre-cursor by extracting general information from noisy data.

By using common frame work of clustering we are able to perform multi-resolutional analysis of seismic data starting from the raw data events described by their magnitude spatio-temporal data space. This new methodology can be also used for the analysis of the data from the geological phenomena e.g. We can apply this clustering method to volcanic eruptions.

## References

- [1] W.Dzwinel et al Non multidimensional scaling and visualization of earth quake cluster over space and feature space, nonlinear processes in geophysics 12[2005] pp1-12.
- [2] C.Lomnitz. Fundamentals of Earthquake prediction [1994]
- [3] B.Gutenberg & C.H. Richtro, Earthquake magnitude, intensity, energy & acceleration bulseism soc. Am 36, 105-145 [1996]
- [4] C.Brunk, J.Kelly & Rkohai “Mineset An integrate system for data access, Visual Data Mining &

- Analytical Data Mining”, proceeding of the 3<sup>rd</sup> conference on KDD 1997.
- [5] Andenberg M.R.Cluster Analysis for application, New York, Acedamic, Press 1973.
- [6] “Predicting the Earthquake using Bagging Method in Data Mining”, S.Sathiyabama, K.Thyagarajah, D. Ayyamuthukumar
- [7] “A Bagging Method using Decision Trees in the Role of Base Classifiers”, Kristína Machová, František Barčák, Peter Bednár
- [8] “Cluster Analysis, Data-Mining, Multi-dimensional Visualization of Earthquakes over Space, Time and Feature Space”, Witold Dzwiniel, David A. Yuen, Krzysztor Boryczko, Yehuda Ben-Zion, Shoichi Yoshioka, Takeo Ito
- [9] <http://cse.stanford.edu/class/sophomore-college/projects-00/neural-networks/Architecture/feedforward.html>
- [10] [www.dmreview.com](http://www.dmreview.com)
- [11] [www.aaai.org/Press/Books/kargupta2.php](http://www.aaai.org/Press/Books/kargupta2.php)
- [12] [www.forrester.com](http://www.forrester.com)
- [13] [www.ftiweb.com](http://www.ftiweb.com)