

# Data Mining: Investment risk in the bank

Ali Abdolahi

Sama Technical and Vocational Training College  
Islamic Azad University, Mahshahr Branch  
Mahshahr, Iran

Mohammad Farzizadeh

Sama Technical and Vocational Training College  
Islamic Azad University, Mahshahr Branch  
Mahshahr, Iran

---

**Abstract:** This paper will discuss the technology and methods behind data mining, how data mining works, how it helps to improve national security, and how sustainable the technology is. Sustainability, with regard to data mining, refers to the impact on the quality of life. Quality of life refers to the preservation of human rights and the ability to feel secure. The ethics and the fallbacks regarding privacy will also be discussed in depth, including the benefits that accompany these fallbacks, and whether they outweigh the cons. Both technical and ethical articles will be used to highlight and discuss the potential, good and bad, and the controversy of data mining. Applications of data mining to security will also be proposed.

Data mining methods are expanding rapidly allowing for the mass collection of information. This mass amount of information is then used by many government agencies to identify threats, gain intelligence, and obtain a better understanding of enemy networks. However, the ability to collect this information from any computer draws into question whether or not data mining leads to a violation of the average citizen's privacy and has created a debate as to if data mining is ethically plausible.

**Keywords:** Classes, Clustering, Data Mining, Neighborhoods, Networks and Rules, Security, Sustainability

---

## 1. INTRODUCTION

Within the past 5 years, data mining has become more and more prevalent in the United States due to recent scandals and exposes on the topic. Simply put by Jason Frand, a professor of computer technology, data mining "is the process of analyzing data from different perspectives and summarizing it into useful information" [1]. Due to major advances in technology and the average person's growing dependence on technology, data mining now affects almost every citizen in the United States, whether he or she is aware of it or not. Within the past ten years, data mining has become an essential tool that government organizations, such as the National Security Agency and the Central Intelligence Agency, use to protect the country and gain intelligence on potential threats.

Inside data mining there are several techniques, both old and new, referred to as clustering, classification, neighborhoods, decision trees, and neural networks. Each of these techniques utilizes algorithms to find connections and trends in people's everyday computer recorded activity [2]. This allows the government to identify potential threats within the country and outside of the country. However, in order to effectively data mine the government needs mass amounts of information from every citizen's computer activity. This leads to the fear that if the U.S. government has access to every citizen's computer activity, the government could have access to personal files and documents, which many could argue violates the rights of citizens [3]. Owing to the large amount of people data mining impacts, it could have a major effect on the quality of life of future generations. This draws into question how sustainable, in terms of quality of life, data mining is. Quality of life pertains to the retention of human rights and the ability to feel secure. Despite this, data mining is a continuous innovation [1]. It is constantly growing and changing as the technology the world uses grows and changes. New techniques and uses are frequently discovered; however, its most useful application is security, despite the controversy it generates.

## 2. TYPES OF DATA MINING

The subject of data mining is full of many complex techniques. Within data mining are two main classifications: classical data mining and next generational data mining [2]. Clustering, classification, and trees are all considered classical techniques, while the neighborhoods technique is considered next generational. Many of these techniques are combined, or build off of one another, to more efficiently data mine. In order to understand the more complex techniques and their potential, one must first understand the basics, or the building blocks

### 2.1 Clustering

Clustering, to put it simply, is the grouping of like things [2]. It is a building block technique that is often needed to use more complex data mining methods. This method can be used to sort and group numbers, topics, key words and anything else one may find useful. Within clustering there are two subcategories, hierarchical and non-hierarchical.

A hierarchical cluster starts with the broadest topic and breaks that topic into smaller groups or clusters, like shown in Figure 1. Hierarchical clusters are easier to understand and allow the analyst to define how many clusters are created, unlike their non-hierarchical counterpart [2]. A non-hierarchical cluster does not create this hierarchy, rather it just creates various clusters of data. Within non-hierarchical clustering are two other subcategories of methods referred to as single-pass methods and reallocation methods [2]. Reallocation methods move data from one cluster to another in order to better organize data within the clusters. The single pass method simply runs through the data once which results in less specific clusters [2].

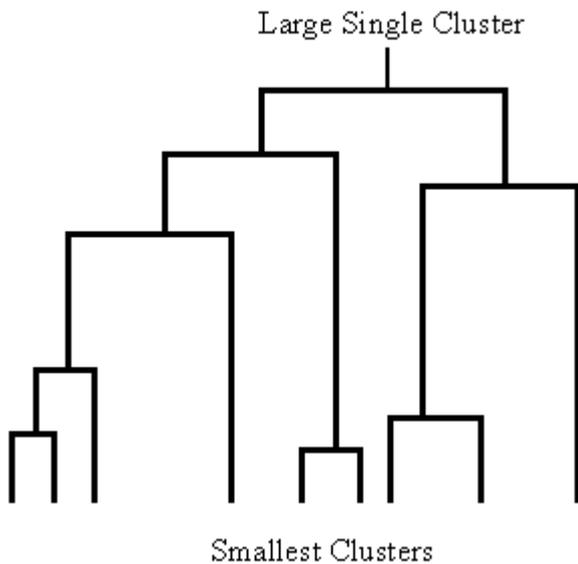


FIGURE 1 Above is a diagram of a hierarchy of clusters. The largest cluster is split up into smaller and smaller groups[2]

In order to choose which clustering method to use, one may consider how efficiently or timely he or she needs to organize the data and how specific the sorting needs to be. Non-hierarchical methods are quicker at sorting data than hierarchical methods, however, clusters are not broken down into more specific categories [2]. Data can be sorted even more quickly depending on which non-hierarchical method an analyst picks. Reallocating data requires more computer power and time than using a single-pass method. Therefore, specificity and efficiency are two factors to be taken into account when considering various methods of clustering. Specificity can be enhanced with classification, another method of data mining.

## 2.2 Classification

Classification, or classes, is a technique that stores data in order to locate different data in predetermined groups or classes, which will tell you more than the previously stored data alone [1]. This technique can be used in two ways, to build an archetype of an item, a product, or anything else one may find useful, or it can be used to add to other types of data mining such as trees and clustering [4]. For instance, attributes from different classes can be used to enhance clustering, specifically by using these attributes to find clusters [4].

In order to apply classification to security, a classification algorithm could sort through data and categorize potential threats based off people's age, criminal history, personal affiliations, etc. The algorithm would be categorizing based off of the known archetype of someone who, in the future, would pose a threat. This would allow organizations such as the Central Intelligence Agency or the National Security Agency to profile potential threats and possibly stop future attacks or security breaches. Classification and clustering are able to lead into another technique for data mining, neighborhoods.

### 2.2.1 Neighborhoods

According to Professor Jason Frand of UCLA, the neighborhood method "is a technique that classifies each record in a dataset based on the classes of the records" [1]. In short, the nearest neighborhood technique works by looking for similarities in data and making conclusions. Data that is similar to each other will have similar conclusions [2]. For example, if

you look at people who live in the same area, it is safe to conclude the residents make a similar income [2].

The neighborhood technique seems very much alike the previously presented methods in that it groups data based off of attributes that the data possesses. It is essentially a refined version of clustering. Nevertheless, it is more advanced than clustering in that the algorithm weights importance and can detect which information is more influential in coming to a conclusion [2]. Not only that, another difference between clustering and neighborhoods is that clustering is an "unsupervised learning" technique, whereas neighborhoods is "supervised learning" [2]. Unsupervised learning sorts data without a purpose, while supervised learning categorizes it for the purpose of performing a prediction [2].

With neighborhoods as a way to mine data for the purpose of domestic security, it can be used to group data from various crimes and make a conclusion about these crimes, based upon finding the "nearest neighbor" or most similar attribute. A conclusion could be made to show the crimes were committed by the same perpetrators which would give detectives an advantage when trying to find who committed them [2].

### 2.2.2 Decision Trees

Decision trees are tree-shaped structures that map decisions which generate rules for the classification of a dataset [1]. In short, a decision tree is a predictive model [2]. The goal of a decision tree is to provide a set of rules that can be applied to unclassified data to predict a certain outcome [1]. The two most common types of data trees are Classification and Regression Trees (CART) and Chi Square Automatic Interaction Detection (CHAID).

CART trees is an algorithm structured as a series of questions where the answers determine the next question [2]. A good metaphor stated in Alex Berson's book on data mining, the CART method is essentially "growing a forest and picking the best tree" [2]. On the other hand, CHAID is a decision tree that weights significance of data and asks questions accordingly. Using CHAID is basically a different way to determine the questions that are asked. The major advantage of CHAID over CART is the simplicity of the results and its ability to handle large sample sizes [2].

### 2.2.3 Neural Networks

An artificial neural network, much like the biological neural network within the human brain, detect patterns, make predictions and learn. Networks estimate conclusions that are dependent upon a large number of inputs, some of which are unknown [2]. Neural networks are far more advanced than any of the previously discussed techniques. This method represents the cutting edge of data mining technology and its applications are endless. Neural networks can be used alone, or they can be used as a supplement to clusters, neighborhoods, classification, and decision trees to overall enhance the quality of the results [2].

A neural net is composed of two parts: the node and the link. The node is very much like a neuron in a human brain. The link is most closely associated with the connections between neurons (the axons, dendrites and synapses) in the human brain. It is important to understand that although neural networks stem from the Artificial Intelligence community, neural networks are not a form of artificial intelligence. Networks can only perform brain-like activity [2].

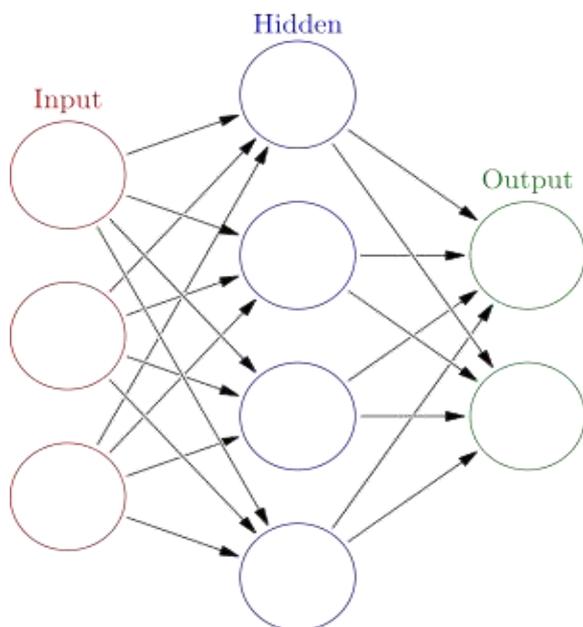


FIGURE 2 Above is a representation of a neural network; each circle is a node and each line is a link between the nodes and the outputs[2].

Although there are many different types of neural networks that have been created, the most used is the Kohonen feature maps. This algorithm is a simple version of neural networks, as it only has one input layer and one output layer. These levels compete amongst themselves to display the strongest conclusion [2]. Kohonen feature maps are extremely useful when combined with clustering. By making each output node a data cluster, each data byte would fall into only one cluster (the most common one). The other clusters that had less hits would still be shown, but they would be shown in the most likely to be the next best order. This is useful in that it allows the user to get a full analysis of all cluster categories [2].

In terms of security, neural networks are most commonly applied to clustering techniques. When networks are used for clustering, they are left in an unsupervised learning mode, or an autopilot-like setting where there is no user requesting specific types of conclusions. [2]. Therefore, the clusters the network creates are unbiased and may assist in finding patterns analysts would not have looked for otherwise. The clusters are formed by forcing the system to compress the data and create algorithms that create clusters that compete against one another for the information in the records it is sorting. This ensures the clusters overlap as little as possible, making the most useful and scientific analysis of the data [2].

#### 2.2.4 Rule Induction

Rule induction is one of the most popular forms of data mining in terms of knowledge discovery. It is also the process that most closely resembles what a person thinks when he or she hears the word data mining because rule induction quite literally mines for information. A rule induction program mines for a rule that is interesting. This rule would point out a pattern in data that would be near impossible for a human to find [2].

In a rule induction program, all possible patterns are systematically pulled out of the data and then put through accuracy and significance tests. The accuracy test tells the user how 'true' the patterns predictions are. The significance test

informs the user how widespread the pattern is. For example, if a pattern shows an outcome that is true 99% of the time, but it only applies to 20% of the data, this pattern is very accurate but low in significance. Since rule induction pulls so many patterns, the accuracy and significance test allows the program to rate which patterns may be of most use to the operator [2].

### 3. DATA MINING AND PRIVACY

In terms of privacy, data mining can be used to either protect it, or violate it. In order for data mining to be successful, the program needs to collect as much information as possible. To do this, many wish to cull mass amounts of information from computer users everywhere. This draws into question whether taking someone's online information is ethically plausible. It could also have a large effect on future generations and the durability of their future privacy rights. To attempt combat this ethical flaw, there is now research to see if there is a way to still have access to the information, without violating the privacy of the person whom the information came from.

#### 3.1 Protecting Privacy

Often, data contained in a database is personal, therefore people want to protect it. A database that is mined may contain peoples' phone numbers, addresses and credit card information; however, a database could contain anything. Data-mining can be used to protect that information but, mining a database can also lead to a breach in the security of this personal information. However, that information, if used correctly, could be very beneficial.

The two most popular types of privacy-preserving data mining algorithms are k-anonymity and l-diversity. As a whole, privacy preservation involves using algorithms that either protect private data from the mining process itself or censor the results of the mine [5]. Both k-anonymity and l-diversity are of the former. They both use generalization and suppression methods to prevent the sensitive data from ever being mined in the first place and preserve the anonymity of the individuals to whom the data which is mined belongs to [5].

Algorithms, proposed by Stergios G. Tsiafoulis and Vasilis C. Zorkadis in their conference paper on preserving privacy in data mining, use clustering techniques to help ease the major concerns over the threat to privacy data mining poses. Their algorithm combines the concepts of k-anonymity, l-diversity and clustering to maximize the anonymity of the people whose data they are using. The goal of their algorithm is to be able to make use of the data, but not violate people's privacy [5].

First, the algorithm organizes data sets into subsets based off similarity; then it creates a more relevant group of classes based off of the attributes they are mining for, such as spending information. Finally, it attempts to generalize the information in order to preserve the privacy of the data's owner while preventing data loss [5]. However, the program's efficiency in preserving privacy can be questionable, as it is not perfect. For many, this program is not enough, as their private information is not generalized enough and may still be able to be linked to themselves, despite the attempt at anonymization. However, the program is meant to keep the "disclosure possibility" and "information loss" negligible, as stated by Stergios Tsiafoulis and Vasilis Zorkadis [5].



connections a that human may not see [8]. These location based clusters are very useful in identifying a crime pattern or spree. This can be applied to serial killers, serial rapists or gang crimes due to the fact that these types of crimes often follow specific patterns [8].

Using the geographical based clusters, authorities are able to predict where a serial killer may strike next or where a gang's base is located. This is evidence of how data mining is beneficial and can help protect law-abiding citizens and prevent future crimes. However, due to security laws, information on many types of crimes, such as drugs and juvenile cases are more restricted as to who can access them [8]. This poses a problem for data analysts because clustering requires extensive information in order to make accurate conclusions. Since the data is limited, the conclusions have a larger margin for error, or may not have enough to make a proper conclusion. Also, analysts would need to be cautious with respect to where their data comes from and how it is mined in order to keep from breaking privacy laws. This once again brings up the issue of invading privacy with data mining. However, it could also be argued that incorporating data mining into crime investigation could help sustain a secure society in an increasingly dangerous world, improving the quality of life for future generations. Fortunately, data miners can work in accordance with the law to mine the proper crime data that is allowed for them to use.

## 6. COMPUTER ESPIONAGE

Today, one of the biggest threats to the United States is terrorism. Organizations like the National Security Agency, the Central Intelligence Agency and the Federal Bureau of Investigation are dedicated to preventing terrorist attacks and punishing the groups and individuals who commit them, from those deliberating within the country to those from foreign countries seeking to perform an act of terrorism. With the use of data mining, these government organizations are able to gather billions of pieces of data from phones, computers, Google, anything people use to communicate. With this data, operatives are able to make connections between people and help generate investigative leads [3]. To make these connections, analysts first tag all the data they have collected by sorting it into classes and clusters. From these classes and clusters, users analyze the data in an attempt to come to some conclusion about a possible attack, a possible suspect, or any potential threat to the country [3].

The tags that analysts collect operate by finding similarities and connections in video footage, audio tapes and phone records [3]. For example, data mining could flag a person who frequents terrorist websites or frequently search for words such as bomb-making and guns and put him or her on a watch list to help prevent a possible attack organized by this person [3].

The security advantages to data mining in this respect are evident. It helps the government prevent innocent citizens from violence and allows them to pursue criminals. However, when these organizations collect this data, they are not required to abide by the regular privacy laws and have access to information that usually requires a warrant [3]. No one monitors these organizations and how they are using this data. For many, this is a major breach in personal privacy and they believe it violates their constitutional rights. This is a major concern when it comes to the quality of life of future generations because they believe it is a violation of human rights. However, to others, the possible security benefits to the

preservation of safety outweighs the potential violation of rights.

## 7. ETHICAL CONSEQUENCES OF DATA MINING

As useful as data mining is, many believe it is unethical and possibly dangerous for anyone to have access to all this information. Data mining can be used both to protect and to harm. It is a very powerful tool and some believe it could be abused. If data mining becomes widespread, it will have a large effect, both good and bad, on the quality of life of future generations. Debate can be stirred up as to how well it protects personal information. For instance, if a government uses data mining to find potential domestic terrorists through peoples' internet search history, that data is perhaps not being anonymized if it can be linked to the person that becomes a potential threat. Despite this, for many, the security benefits outweigh the security risks, so long as we keep data mining in the right hands.

### 7.1 Security Benefits

The security advantages of data mining are clear. The most obvious benefit however, is data mining's potential in protecting the average citizen. According to best-selling author and respected journalist, James Bamford, data mining could be "useful for everything from predicting humanitarian crises to directing rescue efforts during natural disasters [by making it] available to the world" [8]. When used correctly, data mining can help protect important intelligence by identifying cyber-attacks. If this intelligence went unprotected it could seriously endanger many lives. Decision trees vastly improve upon the efficiency of these defense algorithms. The use of clustering can also assist law enforcement to help solve and prevent crimes, lowering the crime rates and creating a safer living environment. One of the greatest advantages and most useful practices of data mining is for the government to apply it to homeland security. By using clustering and classes, government agencies are able to identify possible terroristic threats and track down the people who have already committed these acts, protecting citizens from possible harm.

Data mining also has the potential to greatly improve the lives of future generations. It can improve quality of life due to its potential to help protect people and preserve their ability to feel safe. With the use of data mining, governments can help prevent attacks, gain valuable intelligence, and preserve homeland safety. Data mining also benefits ordinary police departments as it could assist in catching dangerous criminals, profiling gang activity, and predict future crime. It is sustainable in that it can preserve the quality of life of the next generation.

Despite the worry that data mining invades people's privacy, there are data mining programs that preserve the anonymity of the people the data is from or mine around sensitive data. This eliminates some risk since, even if sensitive data comes to light, it would be difficult to link it to the individual that it came from.

### 7.2 Privacy Drawbacks

Despite the usefulness of data mining in terms of national security, it poses a large threat to personal security. In order to be efficient and accurate, data mining requires an expansive

and all-encompassing data set. To get access to that amount of data, users collect data from everyone. This data could be phone records, personal videos and pictures, previous purchases, health records, and search histories. For many, this is a major invasion of privacy because, according to renowned author Elizabeth Svoboda, “some bit of seemingly harmless information that you post today could easily come back to haunt you years from now” [10]. The people who think data mining is unethical believe that the access to their personal information violates their constitutional rights and that any access to this information requires a warrant.

In terms of sustainability, many people believe that data mining is not sustainable because it will decrease the quality of life of future generations. This point of view is based on the belief that data mining violates human rights because it violates people’s basic desire for privacy. This violation of rights cancels out any benefits data mining might have. Many believe that the breach of human rights outweighs any positive effects on the quality of life that more advanced security could have. Opponents to data mining conclude that although data mining can improve domestic security, it breaches personal security along the way [10].

## 8. THE RESULTS: DO THE BENEFITS OUTWEIGHT THE RISKS?

After careful consideration of all the potential benefits and potential pitfalls of data mining, we can conclude that data mining is a useful tool which should be used by the government to enhance security. The main concern that is cited when people argue against data mining is the invasion of privacy through the gathering of everybody’s personal data. However, steps can be taken to anonymize the data so it cannot be connected to the individuals that it came from unless those individuals pose a threat. Collection of this data may be considered an invasion of privacy, but this is irrelevant if the data cannot be connected to the person. Mining anonymized data cannot be a violation of privacy, since the data has no face behind it.

Data mining can also be seen as a benefit from the sustainability point of view. It can greatly improve the quality of life for future citizens. The belief that data mining decreases quality of life by violating privacy rights is outweighed by the security benefits it provides. As social media and internet use becomes more popular, the value of privacy drops. Already the younger generations post all of their private information on the internet through social media. The value of privacy is

decreasing; however, the value of personal safety will never decrease. From catching criminals to protecting sensitive intelligence, the uses and advantages of data mining are worth the risk.

## 9. REFERENCES

- [1] J. Frand. (2010). “Data Mining: What is Data Mining?”. (online article). <http://www.anderson.ucla.edu/faculty/jason.frand/teacher/technologies/palace/datamining.htm>
- [2] A. Berson. (2012). “Building Data Mining Applications for CRM (Enterprise)”. McGraw-Hill Education. (Print book).
- [3] J. Pappalardo. (2013, Oct.). “NSA Data Mining: How It Works.” Popular Mechanics. (online article). DOI: 00324558
- [4] M. Brown. (2014). “Data mining techniques.” developerWorks. (online article). <http://www.ibm.com/developerworks/library/ba-data-mining-techniques/ba-data-mining-techniques-pdf>
- [5] G. Tsiafoulis, C. Zorkadis. (2012). “A neural-network clustering-based algorithm for privacy preserving data mining.” Computational Intelligence and Security (CIS). (online article). ISBN: 978-1-4244-9114-8. pp. 401-405
- [6] A. Bouguettaya, X. Yi, F. Rao, E. Bertino (2015). “Privacy-Preserving Association Rule Mining in Cloud Computing.” 10th ACM Symposium on Information, Computer and Communications Security. (online article). ISBN: 978-1-4503-3245-3.
- [7] M. Shree, J. Visumathi, P. Jayarin. (2016). “Identification of attacks using proficient data interested decision tree algorithm in data mining.” Advances in Intelligent Systems and Computing. (online article). DOI: 10.1007/978-81-322-2674-1\_60
- [8] S. Nath. (2016, Dec.). “Crime Pattern Detection Using Data Mining.” Web Intelligence and Intelligent Agent Technology. (online article). DOI: 10.1109/WI-IATW.2006.55
- [9] J. Bamford. (2015). “The Black-and-White Security Question.” Foreign Policy. (print article). pp.70-75
- [10] E. Svoboda. (2009). “Digital Exposure.” Discover. (print article). Vol. 30, Issue 10