

Data Mining of the Large Dataset for Classification Based on Rule and Tree Based Classifiers: A Review

Renu Bala*

Assistant Professor in Computer Science, DAV College, Malout

Abstract – Data mining is outlined because the procedure of extracting info from huge sets of knowledge. Data processing is mining knowledge from data. Data processing is additionally employed in the fields of master card services and telecommunication to observe frauds. In fraud phone calls, it helps to search out the destination of the decision, length of the decision, time of the day or week, etc. It conjointly analyzes the patterns that deviate from expected norms. Within the method of knowledge mining varied forms of classifiers are used for call analysis method. During this paper varied approaches are mentioned that may be used for classification of various datasets. On the idea of rules, and trees varied classifiers are reviewed and their method of classification of knowledge has been mentioned during this paper.

Keywords: Data Mining, Decision Table, Decision Tree, SVM, Naïve Bayes.

-----X-----

1. INTRODUCTION

1.1 Data Mining

It is the method of taking hidden information from a good store of data. The information should be new, and one should be able to use it. Data processing has been outlined as “It is that the science of taking vital data from wide databases”. It’s one amongst the tasks within the method of information discovery from the info. Data processing is employed to find information out of information and gift the info in a very simple and understood ready type. It’s a method to look at giant amounts of information habitually collected. It’s a cooperative effort of humans and computers. Best results are achieved by reconciliation the information of human specialists in describing issues and goals with the search capabilities of computers. 2 goals of information mining are prediction and outline. Prediction tells US regarding the unknown worth of future variables.

1.2 Architecture for Data Mining

To best apply these advanced techniques, they need to be totally integrated with an information warehouse likewise as versatile interactive business analysis tools. several data processing tools presently operate outside of the warehouse, requiring additional steps for extracting, importing, and analyzing the information.

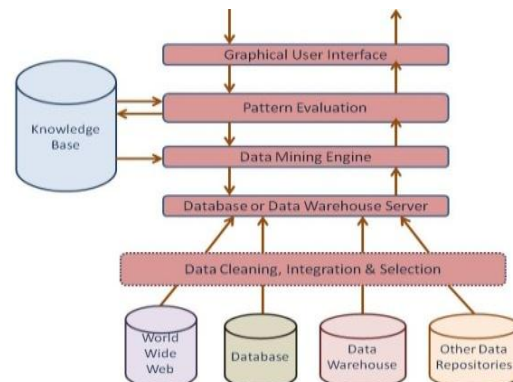


Figure 1.1 Integrated Data Mining Architecture

Furthermore, once new insights need operational implementation, integration with the warehouse simplifies the appliance of results from data processing. The ensuing analytic knowledge warehouse are often applied to enhance business processes throughout the organization, in areas like promotional campaign management, fraud detection, new product rollout, and so on. the best start line may be a knowledge warehouse containing a mix of internal knowledge pursuit all client contact plus external market knowledge regarding challenger activity. Background data on potential customers conjointly provides a superb basis for prospecting. This warehouse is often enforced in a very style of electronic information service systems: Sybase, Oracle, Redbrick, and so

on, and may be optimized for versatile and quick knowledge access. Associate OLAP (On-Line Analytical Processing) server permits a additional refined end-user business model to be applied once navigating the information warehouse. The multidimensional structures enable the user to research the information as they require to look at their business – summarizing by product, region, and alternative key views of their business. The information Mining Server should be integrated with the information warehouse and also the OLAP server to enter ROI-focused business analysis directly into this infrastructure. a sophisticated, process-centric information guide defines the information mining objectives for specific business problems like campaign management, prospecting, and promotion improvement. Integration with the information warehouse permits operational selections to be directly enforced and tracked because the warehouse grows with new selections and results, the organization will regularly mine the most effective practices and apply them to future selections.

1.3 Key Phases in Data Mining Process

1.3.1 Information association

The one in all the foremost acquainted and simple feature of this method is that here we tend to created association between 2 or a lot of things or usually of an equivalent sort to formulate specific pattern. Am fond of it is incredibly acknowledging etiological association between smoking and carcinoma. We've got to gather information involved with smoking habit details together with numbers of smoke per day, period of smoking, style of smoking either bidis, cigarettes, specific brands, way and age of patient etc.

1.3.2 Information classification

This is the second introduce this we will classify the collected info in step with our objectives like etiological factors, investigation purpose, drug treatment plans and results. as an example the etiological info collected from carcinoma patients will be classified on the idea of length of smoking habit, variety of exposure, variety of exposure, age of patient etc.

1.3.3 Pattern Sequencing

This is the next step in module preparation. The pattern sequencing can be prepared with the help of readymade software packages available in market.

1.3.4 Preparation of decision tree

This is final step of prediction system.

1.3.5 Implementation

This is directly involved with last step. You'll have choice either long run or short term processing. Every data processing system has their totally different objectives. Data process area unit loosely developed either as supervised run supervised learning. Supervised learning is that style of learning during which a coaching set is employed to find out model parameters however in unsupervised learning no coaching set is employed. This area unit loosely dived either classification or prediction primarily based pattern. Decision Trees and Neural Networks use classification algorithms whereas Regression, Association Rules and clump use prediction algorithms.

1.4 Data Mining Techniques

Data mining technique is connected with processing, characteristic patterns and trends in info. Or we are able to say that process merely suggests that assortment and processing knowledge in general manner by victimization based mostly programs and subsequent formation of illness prediction or patient management system aid. With the invention of data technology, currently recently it's even additional current. You'll perform data processing with relatively modest information systems and easy tools, together with making and writing your own, or victimization off the shelf code packages. Advanced data processing advantages from the past expertise and algorithms outlined with existing code and packages. this system is habitually use in sizable amount of industries like engineering, medicine, crime analysis, knowledgeable prediction, Web mining, and mobile computing, besides others utilize data processing.

2. REVIEW OF LITERATURE

Thuraisingham, B.et al [1]“Data Mining for Malicious Code Detection and Security Applications ”In this paper author need to mention that the process of sitting queries and attractive patterns from massive quantities of knowledge victimization pattern matching or another reasoning techniques. Data processing has several applications in security as well as for national security furthermore as for cyber security. Threats embody in national security assaultive buildings, destroying essential infrastructures like power grids and telecommunication systems. data processing techniques square measure being investigated to seek out United Nations agency the suspicious folks square measure and United Nations agency is capable of closing terrorist activities. Cyber security is involved protective the pc and network systems against corruption because of Trojan horses, worms and viruses. data processing is additionally being applied to produce solutions like intrusion detection and auditing.

Thuraisingham, B.et al [2]“Data mining for security applications” Author need to projected that the presentation can offer an outline of knowledge mining and security threats so discuss the applications of knowledge mining for cyber security and national security as well as in intrusion detection and statistics. Privacy concerns as well as a discussion of privacy protective data processing also will run.

Asghar, S.et al [3] “Automated data processing Techniques: A essential Literature Review ” during this paper author need to projected that data processing has emerged joined of the key analysis domain within the recent decades so as to extract implicit and helpful data. This data will be apprehended by humans simply. This data extraction was computed and evaluated manually victimization applied math techniques. After, semi-automated data processing techniques emerged due to the advancement within the technology. Such advancement was additionally within the sort of storage that will increase the strain of study. In such case, semi-automated techniques became in economical. Therefore automatic data processing techniques were introduced to synthesis data expeditiously. Consequently

RanaAlaa El-DeenAhmeda et al. [4] “Performance study of classification algorithms for client on-line searching attitudes and behavior victimization knowledge mining”, Author projected eleven data processing classification techniques that square measure relatively tested to seek out the most effective classifier appropriate client on-line searching attitudes and behavior in step with obtained dataset for giant agency of on-line searching. The results show that call table classifier and filtered classifier provide the very best accuracy and also the lowest accuracy is achieved by classification via clump and easy cart. Also, this paper provides a recommender system supported multidimensional language classifier serving to the client to seek out the merchandise he/she is finding out in some e-commerce websites. Recommender system learns from the knowledge concerning customers and merchandise and provides applicable personalised recommendations to customers to seek out the required merchandise.

PareshTanna et al. [5] “A performance comparison between classification techniques with CRM application” Author declared data exploration from the big set of knowledge generated as a results of the varied processing activities because of data processing solely. Frequent Pattern Mining is taken into account a awfully vital enterprise in data processing. Apriori approach applied to get frequent item set usually espouse candidate generation and pruning techniques for the satisfaction of the required objective. This paper shows however the various approaches win the target of frequent mining beside the complexities needed to perform the duty. This

paper demonstrates the employment of rail tool for association rule mining victimization Apriori formula.

Ila Padhiet al. [6] “Predicting Missing things in handcart victimization Associative Classification Mining” Author bestowed a way known as the “Combo Matrix” whose main diagonal parts represent the association among things and searching to the main diagonal parts, the client will choose what else the opposite things will be purchased with the presently contents of the handcart and additionally scale back the rule mining value. The association among things is shown through Graph. The frequent item sets square measure generated from the dance band Matrix. Then association rules square measure to be generated from the already generated frequent item sets. The association rules generated type the idea for prediction. The incoming item sets i.e. the contents of the handcart are going to be pictured by set of distinctive indexed numbers and also the association among things is generated through the dance band Matrix. Finally the expected things square measure urged to the client.

3. APPROACHES USED

Naive Bayes is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. It is not a single algorithm for training such classifiers, but a family of algorithms based on a common principle: all naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable. For example, a fruit may be considered to be an apple if it is red, round, and about 10 cm in diameter. A naive Bayes classifier considers each of these features to contribute independently to the probability that this fruit is an apple, regardless of any possible correlations between the color, roundness and diameter features.

SVM training algorithm builds a model that assigns new examples into one category or the other, making it a non-probabilistic binary linear classifier. An SVM model is a representation of the examples as points in space, mapped so that the examples of the separate categories are divided by a clear gap that is as wide as possible. New examples are then mapped into that same space and predicted to belong to a category based on which side of the gap they fall on.

Decision tables are a precise yet compact way to model complex rule sets and their corresponding actions. Decision tables, like flowcharts and if-then-else and switch-case statements, associate conditions with actions to perform, but in many cases do so in a more elegant way.

Decision tree learning uses a decision tree as a predictive model which maps observations about an item to conclusions about the item's target value. It is one of the predictive modeling approaches used in statistics, data mining and machine learning. Tree models where the target variable can take a finite set of values are called **classification trees**. In decision analysis, a decision tree can be used to visually and explicitly represent decisions and decision making. In data mining, a decision tree describes data but not decisions; rather the resulting classification tree can be an input for decision making. This page deals with decision trees in data mining.

4. CONCLUSION

Data mining is that the field of knowledge that process or data deposition that has been used for extraction of valuable information from the data supported numerous set of rules. Within the method of knowledge mining bunch, classification and attribute choice has been done. Attribute choice is employed for choice of best set of attributes that have minimum dependency on different attributes that area unit out there within the dataset. During this paper numerous classification approaches are mentioned. Classification has been in hot water prediction of assorted knowledge attributes so best set of the foundations are often extracted which will be used for extraction of best decision making process. On the idea of classification rule primarily based classifiers, tree based classifiers and chance based classifiers are reviewed. On the idea of those classifier one will say that rules based classifiers give higher potency for little scale datasets whereas tree based classifiers are often used for classification of the dataset that contain massive instances.

REFERENCES

1. Thuraisingham (2011). "Data Mining for Malicious Code Detection and Security Applications", 978-0-7695-4406-9, 4 – 5, IEEE.
2. Asghar, S. (2009). "Automated Data Mining Techniques: A Critical Literature Review" 978-0-7695-3595-1, 75 – 79, IEEE, 2009.
3. Rana Alaa El-Deen Ahmeda (2015). "Performance study of classification algorithms for consumer online shopping attitudes and behavior using data mining", Fifth International Conference on Communication Systems and Network Technologies, pp. 1344-1349.
4. Dalia Ahmed Refaat Mohamed (2015). "A performance comparison between classification techniques with CRM application", IEEE International Conference on AI Intelligent Systems, pp. 112–119.
5. Hossin, M. (2015). "A Review on Evaluation Metrics for Data Classification Evaluations", International Journal of Data Mining & Knowledge Management Process, pp. 1-6.
6. Nedaabdel Hamid (2015). "Emerging trends in associative classification data mining" International journal of electronics and electrical engineering, Feb 2015, pp. 56-62.
7. Shrey Bavisi A (2015). "A Comparative Study of Different Data Mining Algorithms", International Journal of Current Engineering and Technology, pp. 3248-3252.
8. Kamal R. (2014). "Adaptive Pointing Theory (APT) Artificial Neural Network", International Journal of Computer and Communication Engineering, pp. 212-215.
9. Meenakshi "Survey on Classification Methods using WEKA", International Journal of Computer Applications, 2014, pp. 16-19.
10. Mohammed Al-Maolegi (2014). "An Improved Apriori Algorithm For Association Rules", International Journal on Natural Language Computing (IJNLC), pp. 21-29.
11. Paresh Tanna (2014). "Using Apriori with WEKA for Frequent Pattern Mining", International Journal of Engineering Trends and Technology (IJETT), pp. 127-131.

Corresponding Author

Renu Bala*

Assistant Professor in Computer Science, DAV College, Malout

thakral.renu345@gmail.com