

A STUDY OF THE KEY INFLUENCERS IN BIOMEDICAL EARLY DETECTION, DIAGNOSIS AND MANAGEMENT

Anoushka Gupta

ABSTRACT

The finding of disease is a troublesome work that needs to do in a precise way. Content mining bargains an incredible occupation in this field. A colossal mass of information is accessible in the biomedical field, utilizing this information we can determine numerous infections by content mining procedures in a productive way. Content mining strategies are utilized to recover helpful learning from huge information. The point of this paper is to audit a few document mining strategies utilized in the biomedical field. This study is useful to choose the best book digging strategy for biomedical information. In this paper, the characterization strategy is utilized to contemplate the biomedical content digging for diagnosing illnesses. In the field of biomedical, characterization should be possible based on quiet illness example to isolate the patients into a high hazard or generally safe. The grouping procedures have two techniques they are Binary contains two classes and staggered contains multiple classes. The characterization technique is broadly utilized in biomedical content mining. In this paper, diverse characterization strategies can be applied to sort the content they are SVM (Support Vector Machine) NN (Neural Network), K-NN (K-Nearest Neighbor), Bayesian Method and DT (Decision Tree). In this paper, diverse grouping procedures were studied and their benefits and restrictions have been talked about. The different arrangement strategies were applied in therapeutic information where valuable examples and learning were removed. The significant errand is that to choose the reasonable information and characterization technique for sickness determination. The target is that how the order techniques are applied in the biomedical applications and to choose which strategy is appropriate and proficient for the conclusion of a specific ailment. The fundamental preferred position of the overview is that it is very well applied to any sort of dataset. For future improvement, we will normalize our proposed strategy on utilizing some significant chest infections datasets and estimated execution as far as preparing time and exact determination.

1. INTRODUCTION

Biomedical research refers to the study of the medical issue and problems using biological methodologies, including basic medical research and clinical medical research¹. Biomedical text mining can make easier to begin the process of discovery and to integrate the data present in the biomedical literature. Bioinformatics translator has been prominent with integrating biological and clinical data. The main view of this research is focused on biomedical text mining, aimed at correlating diseases and molecular entities². Data mining is used to point out the hidden information in biomedical data and correct differentiate pathological from normal data. It can be used to extract hidden features of a group of patients and state of diseases that can aid in automated decision making. Data mining offers a clear examination in the field of biomedical³.

2. METHODS

2.1 Text Mining

Text Mining⁴ is the task of discovering unknown information that it may be new or previous information, extracting automatically from various text documents.

2.2 Text Mining Techniques

The text mining techniques are IE (Information Extraction), IR (Information Retrieval), Categorization or Classification, Topic Tracking, Clustering, Summarization and Concept Linkage. Classification includes SVM, K-NN, NN, Bayesian Method and DT. Clustering includes Partition Method, DB (Density Based) Clustering and Hierarchical clustering. These models are described in the following sections.

2.2.1 IE

IE (Information Extraction)⁵ is that structured information is automatically extracted from unstructured/ semi structured documents. This is mostly done by NLP.

2.2.2 IR (Information Retrieval)

IR is nothing but finding and extracting information in documents. The documents may be web documents contain text or image⁶.

2.2.3 Summarization

Summarizing of text in a compressed form of its input, which specifies human Consumption. The document may be in clustered or single.

2.2.4 Classification

It is one of the managed systems. It is the undertaking of looking through a model that clarifies and unmistakable information classes or ideas. These models are gotten from the examination of prepared information. The model is utilized to guess the items' class name that is unknown⁸. Various arrangement systems can be applied to classify the content, for example, SVM (Support Vector Machine), K-NN (K-Nearest Neighbour), NN (Neural Network), Bayesian Method and DT (Decision Tree).

2.2.4.1 K-NN (K-Nearest Neighbour)

K-NN (K-Nearest Neighbour) is a technique of correlation learning, which is comparing the training tuples with the given test tuple, which are similar to it⁹.

2.2.4.2 DT (Decision Tree)

DT (Decision tree) is like a tree structure, in which test attributes as internal node, test outcomes as branch nodes and class labels as leaf nodes. Root node is the topmost node in the tree. Decision tree has been used in operations research to find the conditional probabilities¹⁰.

2.2.4.3 SVM (Support Vector Machine)

SVM (A Support Vector Machine) converts the training data, where it finds a hyper plane using support vectors. The hyper plane splits the data by class.

2.2.4.4 NN (Neural Network)

NN (Neural Network) is a huge number of individual neurons similar to processing nodes and a huge number of weights between these nodes¹⁰.

2.2.4.5 Bayesian Method

In classification and probabilistic learning, Bayes theorem played an important role. Prior Knowledge and observed data is combined by the Probabilistic model. Naïve Bayes classification is one of the simplest Bayesian Algorithms. It has two phases they are learning phase and test phase¹¹.

2.2.5 Clustering

Clustering is the task of partitioning a dataset objects into subsets. Every subset is called a cluster; the objects are similar to one another in one cluster and dissimilar in another cluster. Clustering Methods Are Partition Method, DB (Density Based) Clustering and Hierarchical clustering. These models are described in the following sections¹².

2.2.5.1 Hierarchical Clustering

This clustering is made by multiple levels. It can be categorized as either divisive or agglomerative, based on the decay is formed¹³.

2.2.5.2 Partition Method

In this method dataset objects are partitioned into several clusters, Formally, given set S , number of objects as N and number of clusters as M to form a partitioning algorithm classifies the objects into M partitions ($M \leq N$), where

2.2.5.3 DB (Density Based) Clustering

DB method is the process of partitioning a dataset objects into multiple or a hierarchy clusters. DB clustering can be stretched into subspace from full space clustering¹².

2.2.6 Topic Tracking

Topic tracking is used to track the user views, to find the changes in the IAI (Information Area of Interest) and regularly produces a summarized report of changes, this reveals the emerging topic in the particular information area¹⁴.

2.2.7 Concept Linkage

Concept linkage is a technique to find the corresponding documents which share the same concepts¹⁵. Concept linkage is mainly used to promote browsing.

2.2.8 Question Answering

The question answering method is used to ask questions in World Wide Web and then it gets the related answer. This technique has allowed in many websites¹⁶.

3. RELATED WORKS AND DISCUSSION

Microarray data were analysed by various classification method such as SVM, Decision tree, Bagging, Boosting and Random Forest. The data set obtained from Kent Ridge were comparatively analysed by 10-fold cross validation approach. Among all classification methods random forest shows accuracy result¹⁷. The early warning system of chronic disease was promoted by KNN and Linear Discriminate Analysis(LDA). The connection between the heart disease and hypertension were analysed and minimized the complication occurrences of the disease by constructing an early warning system¹⁸. Patient those who are having chronic disease are classified by their actions, using universal hybrid decision tree. They extended their research work to get more accuracy by classifying various activities of patients¹⁹. SVM classification method is used to classify many diseases. For diagnosing diseases, the combination of both SVM and K means clustering were applied to microarray data²⁰. ANN was used to examine chest diseases, comparative analysis also done for chest diseases that was conducted by probabilistic neural network, generalized regression and multilayer neural networks²¹. Bayesian method has an outstanding performance in diagnosis of psychiatric disease. The dataset of psychiatric patient was taken from Lugo municipal hospital²². Genetic support vector machines (GSVM) was performed better analysis for heart valve diseases. GSVM classifies the ultrasound signal of heart valve and also extracts important features. In this work the automatic system examines heart valve diseases from 215 samples. After evaluation of samples the result was effectively find the Doppler heart sounds²³. The comparison of different classification methods is shown in Table 1.

Author	Algorithm	Working Mode	Advantages	Limitations
Hu et al ¹⁷ .	Random Forest	i) Construct with many trees. ii) After each tree is built, all of the data are run down the tree. ii) Proximities are computed for each pair of cases.	i) It is unexcelled in accuracy among current algorithms. ii) It runs efficiently on large databases. iii) It can handle thousands number of input variables without variable deletion	i) Random forests have been observed to over fit for some datasets with noisy Classification/ Regression tasks. ii) It is not reliable for categorical variables with different number of levels.
Jen et al ¹⁸ .	K-Nearest Neighbor	i) Find out the unidentified data point using the previously known data points (nearest neighbor).	i) It is easy to implement. ii) Training is done in a faster manner.	i) Testing is slow. ii) It requires a large storage area. iii) Sensitive to noise.

Chief et al ¹⁹ .	Decision Tree(DT)	<ul style="list-style-type: none"> i) Search based on the topic or previously viewed by the user. ii) The topic is conjectured by the interest of the user. 	<ul style="list-style-type: none"> i) Simple to understand and interpret. ii) There are no requirements of domain knowledge in the construction of decision trees. iii). It reduces the ambiguity of complicated decisions and assigns exact values to outcomes of various actions. iv).Performs well with large datasets 	<ul style="list-style-type: none"> i)It is restricted to one output attribute. ii) It generates categorical output. iii) It is an unstable classifier i.e. performance of classifier is depend upon the type of dataset.
Soliman et al ²⁰ .	Support Vector Machine(SVM)	<ul style="list-style-type: none"> i) First select the main sentences and paragraphs. ii) Join them into an abstract form. iii) Perceive main concepts of the text. iv) Convey main concepts in natural language. 	<ul style="list-style-type: none"> i) Better Accuracy as compare to other classifier. ii) Easily handle by complicated nonlinear data points. iii) Over fitting problem is not as much as other methods. 	<ul style="list-style-type: none"> i) Computational is very expensive. ii) The main drawback is to choose a right kernel function. For each dataset different kernel function shows various results.
Er et al ²¹ .	Neural Network	<ul style="list-style-type: none"> i) Used to perform non-linear statistical modeling ii) Uses gradient descent method. iii) Based on neurons. iv) Having multiple interconnected processing elements known as neurons. 	<ul style="list-style-type: none"> i) Easily find the complex relationships between dependent and independent variables. ii) Able to handle noisy data. 	<ul style="list-style-type: none"> i) Local minima. ii) Over-fitting. iii) The processing of ANN network is difficult to interpret and require high processing time if there are large neural networks.
Curia et al ²² .	Bayesian Method	<ul style="list-style-type: none"> i)Based on Bayes theory. ii) Concerted on prior, posterior and discrete probability distributions of data items. 	<ul style="list-style-type: none"> i) It makes the computation process easier. ii) Have better speed and accuracy for large datasets. 	<ul style="list-style-type: none"> i) It does not give exact results in some cases where there exists dependency among variables.
E. Avci et al ²³ .	Genetic support vector machines (GSVM)	<ul style="list-style-type: none"> i) Promote solutions to optimization problems using techniques inspired by natural evolution, such as inheritance, mutation, selection and crossover. 	<ul style="list-style-type: none"> i) Prediction accuracy is generally high. ii) Robust, works when training examples contain errors. iii) Fast evaluation of the learned target function. 	<ul style="list-style-type: none"> i) Long training time. ii) Difficult to understand the learned function (weights). iii) Not easy to incorporate domain knowledge.

4. CONCLUSION

The different content mining strategies in the biomedical field were dissected in this overview. Their benefits and restrictions have been talked about. The motivation behind the examination is how the order techniques are applied in the biomedical field and to choose a strategy that is appropriate for diagnosing a specific ailment. As per the presentation, the characterization

strategy is well reasonable for diagnosing sicknesses. The order should be possible based on understanding malady example to isolate the patients into a high hazard or okay. Ordering patients related data were analysed that prompts great outcome.