

Text Categorization and Feature Extraction Using Bayesian Classification

Snehlata, Mr. Ramesh Loar

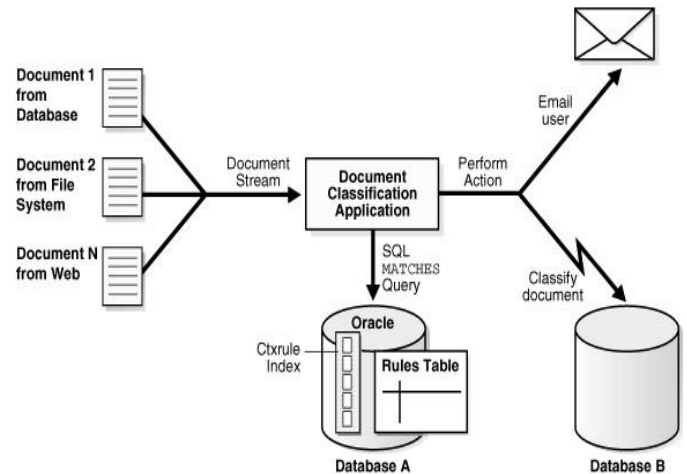
Department of Computer Science and Engineering, Rao Pahlad Singh Group of Institutions, Balana, Mohindergarh

Abstract- In this paper, a novel approach is proposed for extract eminent features for classifier. Instead of traditional feature selection techniques used for text document classification. We introduce a new model based on probability and over all class frequency of term. We applied this new technique to extract features from training text documents to generate training set for machine learning. Using these machine learning training set to automatic classify documents into corresponding class labels and improve the classification accuracy. The results on these proposed feature selection method illustrates that the proposed method performs much better than traditional methods.

Keywords- Data mining, Text mining, Text mining framework, Text mining techniques, Bayesian.

I. INTRODUCTION

Information Mining Data Mining insinuates use for removing or mining data from a great deal of data.[1] Data Mining is an of technique discovering potential, supportive, assurance, novel, interesting and as of now cloud case from immense measure of data. With the use fitting estimation we can find relevant information [1]. Data mining is moreover called "taking in disclosure from data".(KDD) There are various diverse terms like data mining, for instance, data extraction, data burrowing, data pale history. The information and learning increment can be used as a piece of market examination, blackmail revelation, creation control and consistent data examination [1]. Content Mining and Text Mining Frame Work. Content mining is one kind of data mining framework. The method use for expelling or mining gaining from the substance record. Content mining discover the already cloud information removing it thus from different source.[2] Text mining resembles data mining. Regardless, the data mining overseeing structure data and substance mining overseeing unstructured or semi structure data. Like email, content record et cetera in a substance mining major goal is to discover the officially cloud information. Additionally, the issue is that the result isn't correlated to customers require. In a substance mining the collection of report from various unmistakable sources. Social affair information is basic anyway fining pertinent information on ask for is troublesome. Content mining technique or substance mining plot work start with the social affair of file from different source. Content mining contraptions recoup a document and perform preprocessing on it.



II. TEXT MINING TECHNIQUES

Innovation like data extraction, bunching, synopsis, order and representation are utilized as a part of content mining outline work or process. Here in taking after area we talk about the content mining techniques.[2]

Information Extraction

Information extraction is basic walk for PC to separate unstructured substance and its relationship. This methodology is done by case organizing is used to look for pre portray course of action of substance. IE is join ID, sentence division. This techniques is to a great degree significant for generous substance file. Various testing in electronic information is as regular tongue planning and IE deal with this issue change content record into structure mastermind.

A. Clustering

Batching is unsupervised system. Batching technique used to hoard similar files anyway it contrasts from classification, in this reports are gathered. This system relies upon apportioning practically identical substance into same bundle. Each group contain different equivalent reports.

B. Summarization

On account of far reaching measure of data we need to consolidate the data from the amount of report .which pack the

data without change significance of substance, and the length of data. Also, make summary from the social affair of report. Accordingly whole record set is supplanted by the summary. Summary is valuable for the customer read short summary of record as opposed to extended reports.

C. Visualization

In content mining recognition improve the straightforwardness to discover the information. Social occasion of record or a single document content pennant used to indicate report and shading used. This procedure give speedier player and legitimate information. Which find or mine the case from social event of record. Its use different shading, relationship discrete et cetera.

D. Categorization

Course of action resembles content request [4] Categorization is a directed framework since it relies upon information yield cases to gathering. Content classifier is used to classification of the substance chronicle into pre portray class. Furthermore, pre describe class is dole out in perspective of substance record content. A normal substance arrange handle include preprocessing, requesting, estimations reductions and course of action. The target of the request is to get ready classifier on the commence of known and cloud representation are masterminded normally. To orchestrate the substance number of substance arrange techniques used which we will discuss as a piece of the going with section (3)

III. PROPOSED METHODOLOGY

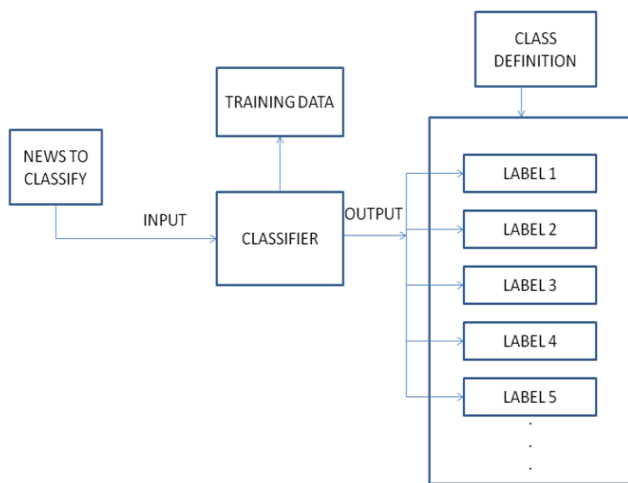


Fig: Proposed System

From the above figure, you can see the system architecture clearly. System consist of new to be classified as its input and the output will be the label to which the news probably belongs to. Classifier is the main module of the system which is the implementation of the naïve Bayes algorithm. It uses the training data as its input and classifies the input documents.

Training data consist of large number of documents preprocessed i.e. term frequency and document frequency is calculated. Using this data the input file is classified.

IV. RESULT ANALYSIS

Comparison of Proposed Method with Text Classifier using Decision Tree Overall Comparative Presentation

Technique	(%)Training Data	(%)Accuracy
Association Rule		
Based Decision Tree	50	80
Proposed		
Algorithm	76	87

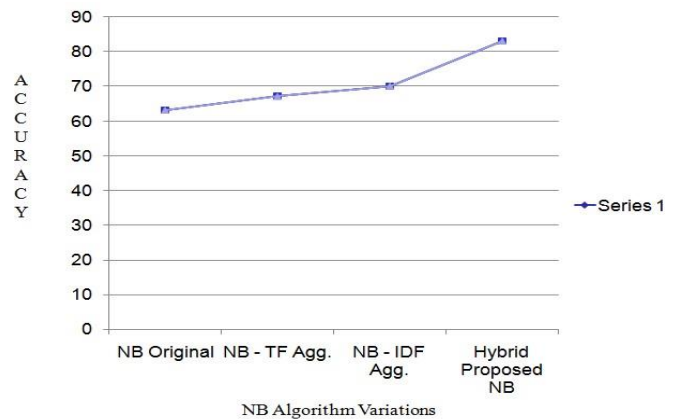


Figure: Accuracy Found for Different Variations

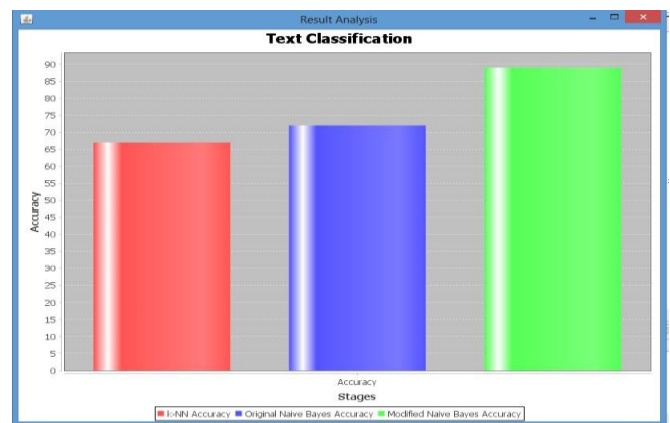


Figure: Result Analysis Graph



Figure: Graph of text reduction

Above figure shows result analysis graph. Text file are appropriately classified using this application. It shows the different stages depend upon accuracy.

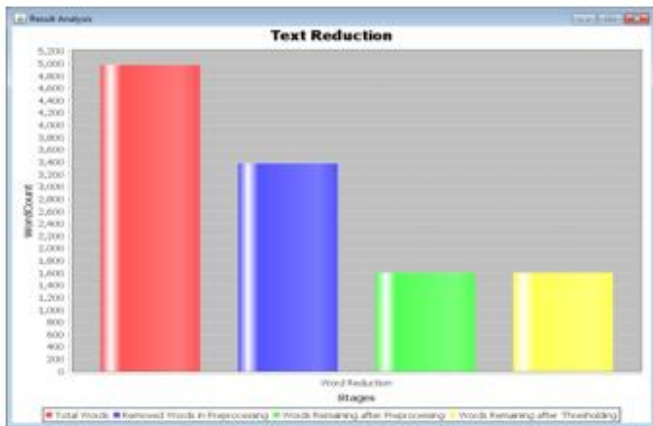


Figure: Accuracy classifications for NB variation

V. CONCLUSION

The Text Classification using analytical approach project proposed a design of the application that can effectively classify text files into appropriate folder depending upon the theme of the file, using the training data to model the classifier. This application automates the text classification process otherwise would take long time doing manually the same task. Text file are appropriately classified using this application. This application allows you to select the test data, training data. In the future, a similar concept can be used for different purposes like arrange your computer, classify various documents with various applications and analyze them.

VI. REFERENCES

- [1]. Jiawei Han and Micheline Kamber "Data Mining Concepts And Techniques" Morgan kaufman publishers, San Francisco, Elsevier, 2011, pp. 285351
- [2]. M.Sukanyal, S.Biruntha2 "Techniques on Text Mining" International Conference on Advanced Communication Control and Computing Technologies, IEEE-2012
- [3]. Sonali Vijay Gaikwad, Archana Chaugule, Pramod Patil "Text Mining Methods and Techniques" International Journal of Computer Applications (0975 – 8887) Volume 85 – No 17, January 2014
- [4]. Nidhi, Vishal Gupta "Recent Trends in Text Classification Techniques" International Journal of Computer Applications (0975 – 8887) Volume 35– No.6, December 2011
- [5]. S. Subbaiah "Extracting Knowledge using Probabilistic Classifier for Text Mining" Proceedings of the 2013 International Conference on Pattern Recognition, Informatics and Mobile Engineering, February 21-22, IEEE-2013
- [6]. M. Janaki Meena , K. R. Chandran "Naive Bayes Text Classification with Positive Features Selected by Statistical Method" ©2009 IEEE vaishali Bhujade, N.J.Janwe "knowledge discovery in text mining techniques using association rule extraction" International Conference on Computational Intelligence and Communication Systems, IEEE2011
- [7]. Zhou Faguo, Zhang Fan "Research on Short Text Classification Algorithm Based on Statistics and Rules" 2010 Third International Symposium on Electronic Commerce and Security © 2010 IEEE
- [8]. Shuzlina Abdul-Rahman, Sofianita Mutalib, Nur Amira Khanafi, Azliza Mohd Ali "Exploring Feature Selection and Support Vector Machine in Text Categorization" 16th International Conference on Computational Science and Engineering, IEEE-2013
- [9]. Xianfei Zhang, Bicheng Li, Xianzhu Sun "A kNearest Neighbor Text Classification algorithm Based on Fuzzy Integral" Sixth International Conference on Natural Computation, IEEE-2010