# DOCUMENT CLASSIFICATION AND SUMMARIZER USING PROBABILISTIC CLASSIFIER.

Vinaya Kulkarni

Shruti Rothe                      Shivani Devgirikar

Rasika Deshpande                  Sneha Sultane

**Abstract:**

In this paper, we propose a straightforward component extraction algorithm that can accomplish high archive order precision with regards to advancement driven subjects.Document classification has become an important field of research due to the increase of unstructured text documents available in digital form.Programmed content synopsis turns into a significant method for finding applicable data correctly in huge content in a brief span with little endeavors.In the proposed system,we take PDF as input,classify it according to it"s domain and summarize it using power of Natural Language Processing and Machine Learning.Here we are providing options for selecting required line of summary,the prediction result will send to the user email with summary and classified domain.Having generated summary help us to tell whether to deep dive in detail or not.The main motive of project to save time required in document classification and understanding it.

**Keyword**: Naïve Bayes (NB), Machine Learning, Natural Language processing, TF- IDF(Term-Frequency inverse document frequency)

**Introduction:**

Document classification is that the work of uncertain documents into classes supported their content. There area unit several classification strategies for documents. Classification is outlined as categorizing document into one in every of a set variety of predefined categories with one document happiness to just one category. The document is entered as input during this system, then the system can do preprocessing steps. Main words area unit solely retrieved. That means summarize words and words area unit matched with all the information needed for every field hold on within the info.

Then count the necessary words exploitation term frequency (TF) in feature choice. Finally, calculate the accuracy by exploitation holdout technique. This paper is comparing the results of naïve Bayes classifier techniques. Highlighted 5 algorithms that area unit applied to the text classification and expressed comparative study on differing types of classifiers with their enhancements.

This project uses machine learning techniques for classification and summarizers of documents. Machine Learning enables systems to recognize patterns based on existing algorithms and data sets and to develop adequate solution concepts. Machine Learning undoubtedly helps people to work more creatively and efficiently. Therefore, in Machine Learning, artificial knowledge is generated based on experience. In this project, by classifying and summarize a document, one or more categories are assigned to a document, making it easier to manage and sort. This is especially useful for publishers, news sites, blogs or anyone who deals with a lot of content.

**Related work:**

Nowadays finding the papers related to a particular domain is very difficult until we read the whole paper. Classification and summarize is a machine learning technique that assigns categories to a collection of data to aid in more accurate predictions and analysis. The document is entered as input in this system, and then the system will do preprocessing steps. The main words are only retrieved. The main words are matched with all the data required for each field stored in the database and then count the important words using Term

Frequency in feature selection. The document is entered as input in this system, and then the system will do preprocessing steps. The main words are only retrieved. The main words are matched with all the data required for each field stored in the database. Summarize that document by using NLP algorithms And then count the important words using Term Frequency (TF) in feature selection. Finally, calculate the accuracy by using the holdout method. This section provides the characteristics and descriptions of the data sets used for performing our experiments. In this system, the user can know any field of document and calculate the probability of the words of the document. Only words including more and more in a document are to display what kind of field. But if the count of words is the same, it needs to make accuracy so that the program displays what kind of field. This classification and technique are based on Bayes theorem with an assumption of independence between predictors. NLP is used for summarization techniques. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

Text and document classification processes are often used in areas such as sentiment analysis, text summarization, etc. The author Mehmet Baygin has performed document classification using the Naive Bayes approach.

Clustering techniques can be applied only on structured data [9]. So unstructured data need to be converted to structured data. But while converting unstructured data into structured data the algorithm efficiency decreases. So to increase the efficiency For this, we are going to use the TF-IDF approach.

Text classification has become more important due to the growth of big data with which we could obtain huge data daily. It has many applications like information retrieval, spam detection, language identification, sentiment analysis and plays a major role in natural language processing as well

Text classification and classifiers: a survey of text classification[4] its process and also the overview of classifiers. They also tried to compare some existing classifiers. The existing classification methods are compared and contrasted based on various parameters. They also found that the performance of the classification algorithm is greatly affected by the quality of the data source. A decision tree is an important method for both induction research and data mining, which is mainly used for model classification and prediction[6]. ID3 and C4.5 algorithm is the most widely used algorithm in the decision tree. We aim to implement the algorithms in a very time and space-effective manner and throughput and response time for the application will be promoted as the performance measures. Our project aims to implement these algorithms and graphically compare the complexities and efficiencies of these algorithms.

## System Architecture:

The system architecture describes the overall flow of the system. This system is useful for the early classification of the post. The user who will use this system needs to first register into the system. The details will be stored in the database. Now the user will enter the details like name ,email address, etc. After registration, the user will log in to the system using the login.JSP page. Firstly user upload the document then System will ask line of summary want after selecting required line of summary,the prediction result will send to the user email with summary and classified domain.

The algorithm used in the system is Core NLP for text mining. For classification, the Naïve Bayes algorithm is used.
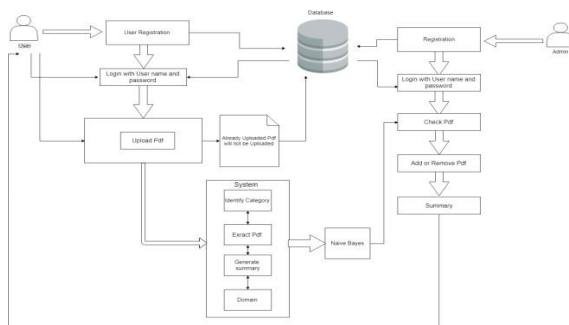


Fig. System Architecture

**Algorithms:**

In Naïve Bayes classifier, we predicate the result, depending upon the trained dataset.In this Naïve Bayes classifier is used to classify the values. Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity. Naive Bayes is known to outperform even highly sophisticated classification methods. In Naïve Bayes classifier, we predicate the result ,depending upon trained dataset. In this Naïve Bayes classifier is used to classify the values.

Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

Formula is as given below:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

-Finding the probability of A given that B has already occurred.

**Natural language processing:**

Tokenization: The process of converting a text into tokens

Stemming: Stemming is a rudimentary rule- based process of stripping the suffixes ("ing", "ly", "es", "s" etc) from a word.

Stop word removal: Language stop words (commonly used

Words of a language – is, am, the, of, in, etc.), URLs or links, social media entities (mentions, hash tags), punctuations and industry-specific words. This step deals with the removal of all types of noisy entities present in the text

Entity Extraction: Entities are defined as the most important chunks of a sentence – noun phrases, verb phrases or both. Entity

Detection algorithms are generally ensemble models of rule-based parsing, dictionary lookups, post tagging, and dependency parsing. The applicability of entity detection can be seen in automated chatbots, content analyzers, and consumer insight

**Algorithm for Document Summarization :**

Step1:Text  Extraction  from  PDF- document

Step 2:Cleaning Extracted data

$$TF(t) = \frac{Number\ of\ times\ term\ t\ appears\ in\ a\ document}{Total\ number\ of\ terms\ in\ the\ document}$$

Step 3:Count Vectorizer Step 4:Result Visualization

**Conclusion:**

In this system, we have developed an application by which we can identify the domain of our document. Due to the use of our system, we can easily be finding the domain of our document which is useful business or education purposes. This paper expressed the extraction of fields document related to IT research paper. It applied the naive Bayes algorithms to classify documents automatically. This classifier gives a correct and accurate result.

**Reference:**

[1] Aggarwal, C. C., & Zhai, C. (2012). A survey of text classification algorithms. In mining text data (pp. 163-222). Springer US.

[2] Mamoun, R., & Ahmed, M. A. (2014). A Comparative Study on Different Types of Approaches to the Arabic text classification. In Proceedings of the 1st International Conference of Recent Trends in Information and (Vol. 2, No. 3) A.Helen Victoria, M.Vijayalakshmi.An Outcome-based Comparative study of different Text Classification Algorithm.
Volume 118 No. 22 2018, 1871-1877

[3] Ahmed, H. A ., & Esrra, H. A. A (2017). Comparative Study of Five Text Classification Algorithms with their Improvements International Journal of Applied Engineering Research, 12(14),4309-4319

[4] Korde, V., & Mahender, C. N. (2012). Text classification and classifiers: A survey. International Journal of Artificial Intelligence & Applications, 3(2), 85.

[5] Badgujar, M. G. V., & Sawant, K. (2016). Improved C4. 5 Decision Tree Classifier Algorithms for Analysis of Data Mining Application. International Journal, 1(8).

[6] Amna Rahman, Usman Qamar(2016). A Bayesian Classifiers based Combination Model for Automatic Text Classification. International Journal, 1(6).

[7] Petre, R. (2015). Enhancing Forecasting Performance of Naive-Bayes Classifiers with Discretization Techniques. Database Systems Journal, 6(2), 24-30.

[8] Mehdi Allahyari, Seyedamin Pouriyeh, A Brief of Text Mining: Classification, Clustering, and Extraction Techniques. KDD Bigdas, August 2017, Halifax, Canada.

[9] Mehdi Allahyari and Krys Kohut. 2016. Discovering Coherent Topics with Entity Topic Models. In Web Intelligence (WI), 2016 IEEE/WIC/ACM International Conference on. IEEE, 26–33.