

# Text Summarization of Large Document

**Vikky M.Chaple**

**Kirti D. Singh**

**Sonali K. Borkute**

**Anjali H. Rai**

U.G. student, Department of computer science and engineering, Ballarpur Institute of Technology,  
Ballarpur, Maharashtra, India

**Prof. Aarti Vaidya(project guide)**

Assistant Professor, Department of computer science and engineering, Ballarpur Institute of Technology,  
Ballarpur, Maharashtra, India

## Abstract:

The main motive of project is that to classify the documents into two modules first one is to get an exact models describing vital information and the second one is to predict future information. At the present stage of the development of information technologies, in the whole world. The tasks related to the computer search for information acquire special urgency, and in particular special attention is attached to the definition of the meaning of textual documents.

The amount of data has increased exponentially over the past few years. Most of the text documents are unstructured and not organized and therefore user-facing additional difficult to search out his data desires. It creates a challenge

for information retrieval for ranking documents according to the user relevance one amongst.

Text classification is that the method of assigning a document to at least one or additional target classes supported its contents. coaching and classification area unit performed exploitation Naïve Bayes(NB) classifiers. This approach will facilitate the user to induce the foremost relevant documents at the start. The experiments on TREC assortment show that the proposed technique is giving higher results compared to the baselines and also the linguistics ranking ESA.

This technique summarizes text on some areas like "Data Management", "Software Engineering", "information Mining", "Artificial Intelligent", "Image processing", "Data Base

Management System", "Big Data " and so on. The main aim is to make proper decision based on the document summary.

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

### Introduction:

Due to the Huge domains in the IT industry, it's very difficult to find classify the domains from a huge number of PDFs. We are going to generate a short Description of that paper so that we can save the time of classification of Papers and find out their relevance.

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 work area is to easily retrieve summarized document. That means summarize words and words that are presents in the document are matched.

This paper is comparing the results of naïve Bayes classifier technique, then count the

necessary words exploitation term frequency (TF) in feature choice. Finally, calculate the accuracy by exploitation holdout technique. Highlighted 5 algorithms that are used for text classificationandexpressed comparative study on different types of classifiers.

We are going to generate a short Description of that paper so that we can save the time to find out their relevance. Ranking of queries is one of the fundamental problems in information retrieval. Ranking relevance is that the work of uncertain documents into classes supported their content. There area unit several classification and ranking strategies for documents. The purpose of the ranking algorithm is to retrieve from the collection of documents the most relevant onesRankingdirectly affects retrieval quality. For each document, a score iscomputed which reflects his relevance concerning the givenquery. Then, documents are ranked according to this score andreturned as a result.

Semantic analysis is the process of relating syntactic structures, from the levels of phrases, clauses, sentences and paragraphs to the level of the writing as a whole, to their language-independent meanings.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 the peoples who use lots of contents such as new sites, publisher etc.

**Related work:**

Nowadays finding the papers related to a particular domain is very difficult until we read the whole paper whether the paper is of Machine Learning, Data Mining, Image Processing, etc. The main motivation of this project to classify the field of document. We are going to generate a short Description of that paper so that we can save the time of classification of Papers and find their relevance. To Classify the Domain of paper from the number of pdf and find their relevance. To Generate a Short Summary of that paper. We are going to Classify the domain according to their content and divide it into groups.

The document is entered as input during this system, then the system can do preprocessing steps. Main work area is to easily retrieve summarized document. That means summarize words and words that are presents in the document 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 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. In this system, the user can know any field of document and calculate the probability of the words of the document. 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. For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as 'Naive'. 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.

**1. Classification of text documents based on naïve Bayes using n-gram features**

**Description:** 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

## 2. Document clustering: TF-IDF approach

**Description:** TF: Term Frequency, which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear much more times in long documents than shorter ones. Thus, the term frequency is often divided by the document length (aka. the total number of terms in the document) as a way of normalization:  $TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$ .

IDF: Inverse Document Frequency, which measures how important a term is. While computing TF, all terms are considered equally important. However it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:  $IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$ .

## 3. An outcome-based comparative study of different text classification algorithms

**Description:** Due to the Huge domains in the IT industry, it's very difficult to find classify the domains from a huge number of PDFs. We are going to generate a short Description of that paper so that we can save the time of classification of Papers and find out their relevance. 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

## 4. Text classification and classifiers: a survey

**Description:** In this paper, A.Helen Victoria and M.Vijayalakshmi, has introduced 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

## 5. Upgradedc4.5 conclusionhierarchy classifier procedure for investigation of data excavatingsolicitation

**Description:** C4.5 is an algorithm used to generate a decision tree. C4.5 is an extension of ID3 algorithm. The decision trees generated by C4.5 can be used for classification, and for this reason, C4.5 is often referred to as a statistical classifier. ID3 and C4.5 algorithm is the most widely used algorithm in the decision tree. We aim to instrument the algorithms in a very time and space-effective routine and quantity and reply time for the presentation will be endorsed as the presentation procedures.

### **Motivation:**

At the present stage of the development of information technologies, in the whole world. the tasks related to the computer search for information acquire special urgency, and in particular special attention is attached to the definition of the meaning of textual documents.

Nowadays finding the papers related to a particular domain and find the frequency relevance of that paper. The main motivation of this project to classify the field of document and ranking of that paper. To generate a short Description of that paper so that we can save the time of classification of Papers and find their relevance frequency.

### **Problem Statement**

To Classify the Domain of paper from the number of pdf and find their relevance of

frequency. To Generate a Short keyword Summary of that paper and find their ranking. To find frequency relevance according to their content and also find their domain

### **Overview:**

Nowadays finding the papers related to a particular domain and ranking is very difficult until we read the whole paper. Classification and finding ranking relevance is a machine learning technique that assigns categories to a collection of data to aid in more accurate predictions and analysis.

Document Classification and ranking of that paper is a problem in information science. This system helps to solve this problem. It will be helpful for students, faculties, researchers in categorizing the papers into different domains according to their content. It makes it easier for them to find the relevant information at the right time and for filtering and routing documents directly to users.

### **Literature Survey:**

- 1) **Paper Name:-** Classification of text documents based on naïve Bayes using n-gram features

**Author Name:-** Mehmet Begin

**Year:-**2018

**Paper Description:**-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.

**Paper Limitations:**-Need to 2-gram, 3-gram and 4-gram features of all documents were extracted and the test procedures were performed separately for each feature.

2) **Paper Name:**-An outcome-based comparative study of different text classification algorithms

**Author Name:**-1.A.Helen Victoria,  
2.M.Vijayalakshmi

**Year:**-2018

**Paper Description:**-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

**Paper Limitations:**-The Decision trees require the features that have to be checked in a specific

order, which limits their ability to exploit features that are relatively independent of each other.

### **Problem Statement:**

To enhance the accuracy of document clusters incrementally by using a supervised document clustering technique that incorporates domain identification and its summarization.

Nowadays, classifying set and finding a ranking of documents (i.e. IEEE papers) into its respective domains is difficult, until we read the whole paper. So, to make our task easy we are developing this web application. This system will identify the domain of the document by using the TF-IDF algorithm and classify them as per their domain with the help of the Naïve Bayes algorithm.

### **Objectives:**

To identify the domain of the given document.

To generate a summary of keywords from the given document.

To finding ranking relevance of the given document.

### **Module:-**

❖ **User:**

- Registration

- Login
- Give Pdf Data
- Classify Domain
- View Description
- Ranking of that document
- Logout
- ❖ **Admin:**
  - View all Pdf
  - Remove File

## PROJECT REQUIREMENT:

### Hardware Requirements:

1. Processor – Intel Core2Duo,
2. Speed – 2.4 GHz (min)
3. RAM - 2 GB (min)
4. Hard Disk - 50 GB (min)

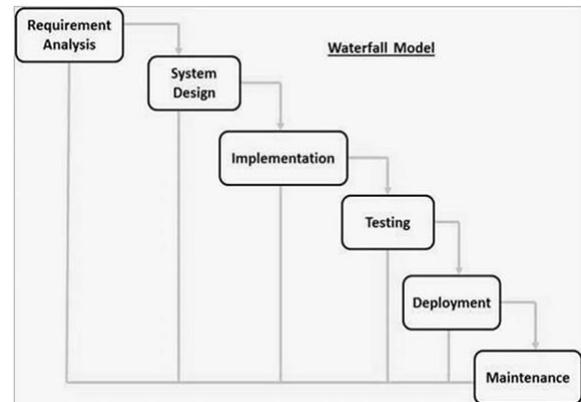
### Software Requirements:

1. Operating System: Windows 7, 8,10
2. Frontend : HTML, JSP, CSS
3. Backend: JAVA , JSP ,Servlet ,JDBC
4. Database: MySQL
5. JDK 1.8
6. TOMCAT Server

### Analysis Models:

SDLC Model to be applied

The waterfall model is a classical model used in the system development life cycle to create a system with a linear and sequential



approach. It is termed as a waterfall because the model develops systematically from one phase to another in a downward fashion. This model is divided into different phases and the output of one phase is used as the input of the next phase. Every phase has to be completed before the next phase starts and there is no overlapping of the phases.

1. Requirement Gathering- All possible requirements are captured in product requirement documents.

2. Analysis Read - the requirement and based on analysis define the schemas, models and business rules.

3. System Design -- Based on analysis design the software architecture.

4. Implementation Development of the software in the small units with functional testing.

5. Integration and Testing Integrating of each unit developed in the previous phase and post-integration test the entire system for any faults.

6. Deployment of system - Make the product live on the production environment after all functional and nonfunctional testing completed.
7. Maintenance Fixing issues and release a new version with the issue patches as required.

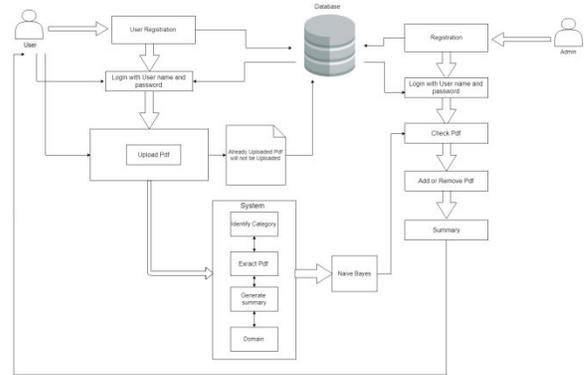


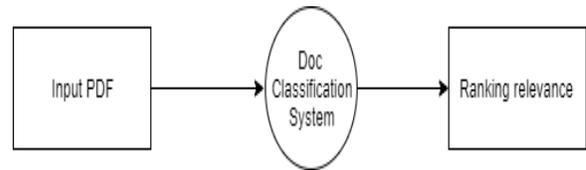
Fig. System Architecture

**System Architecture:**

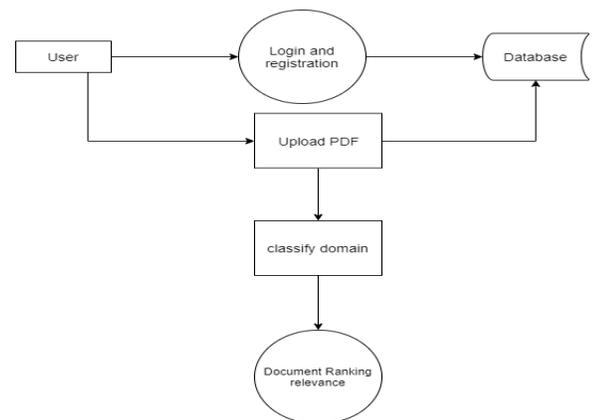
The system architecture describes the overall flow of the system. This system is useful for the early classification of the post. The user who will use this system needs to first register into the system. The details will be stored in the database. After registration, the user will log in to the system using the login.JSP page. Now the user will enter the details like age, gender, etc which is mentioned in the form which we are using. The algorithm used in the system is Core NLP for text mining. For classification, the Naïve Bayes algorithm is used. The prediction result is shown on the Notification page.

**DFD'S:**

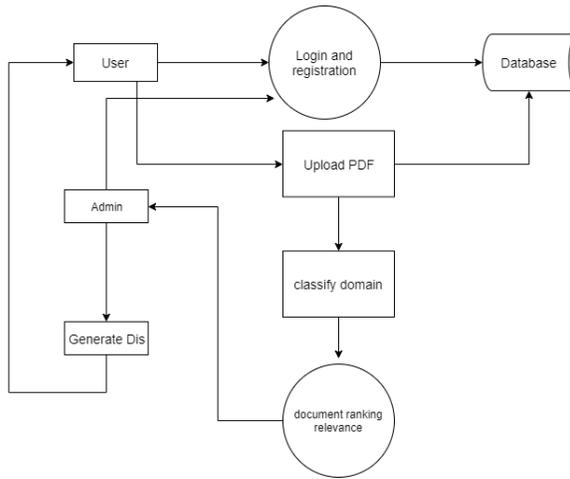
- **DFD 0 :-**



- **DFD 1:-**



• DFD 2:-



UML Diagrams:

• Class Diagram:

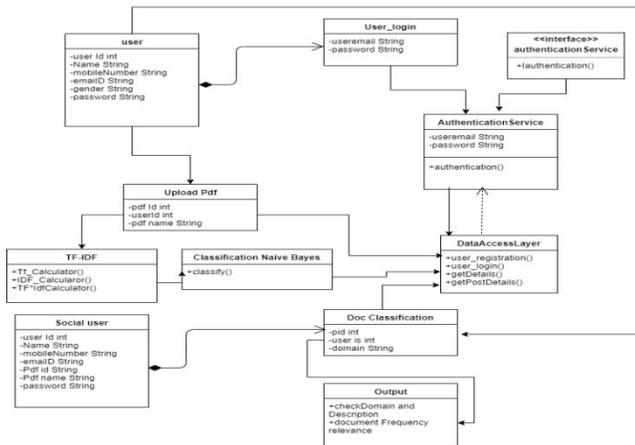


Figure: Class Diagram

• Use- case Diagram:-

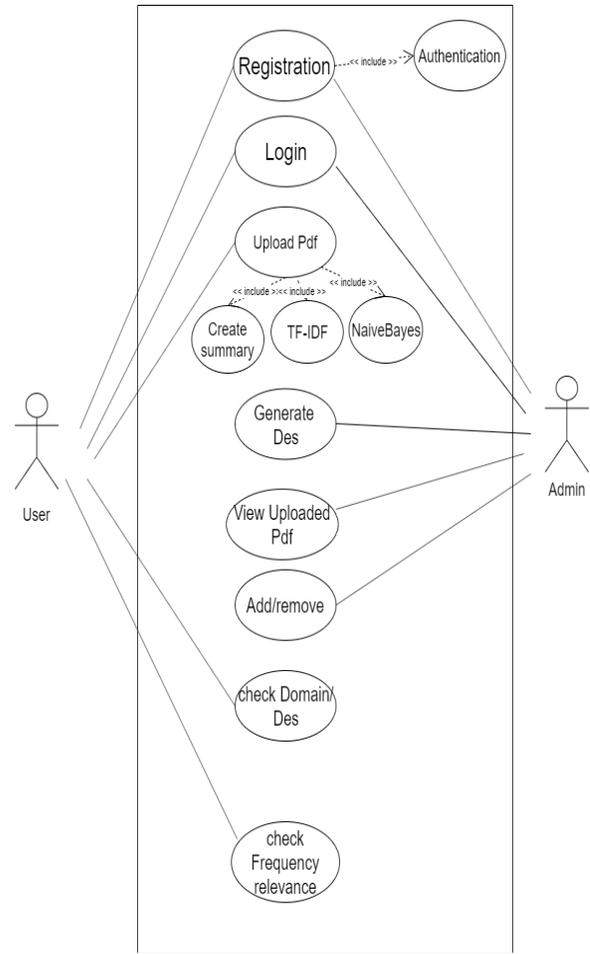


Figure: Use-Case Diagram

• Sequence Diagram:-

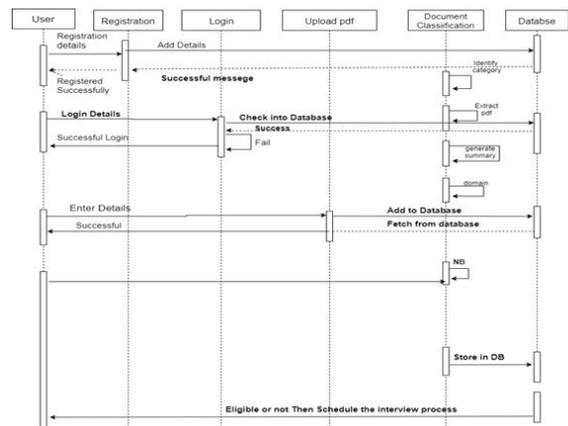
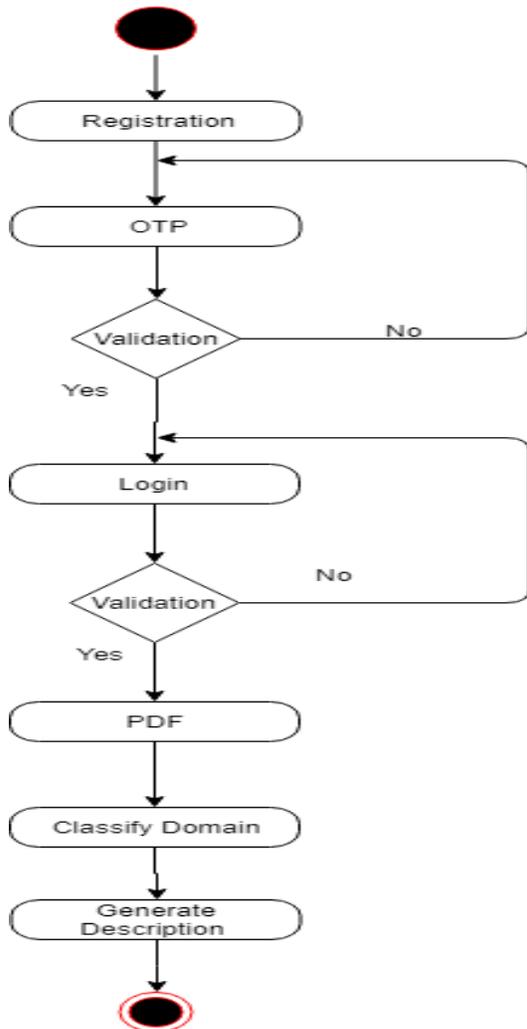


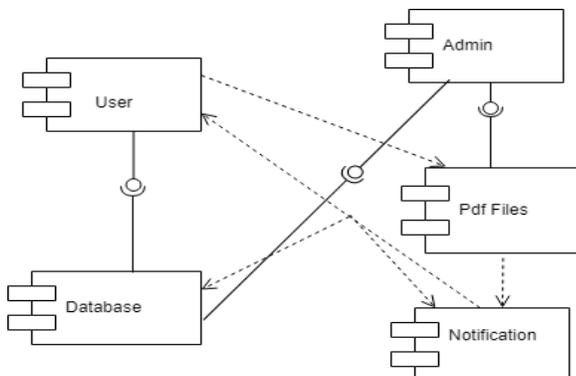
Figure: Sequence Diagram

• **Activity Diagram:-**



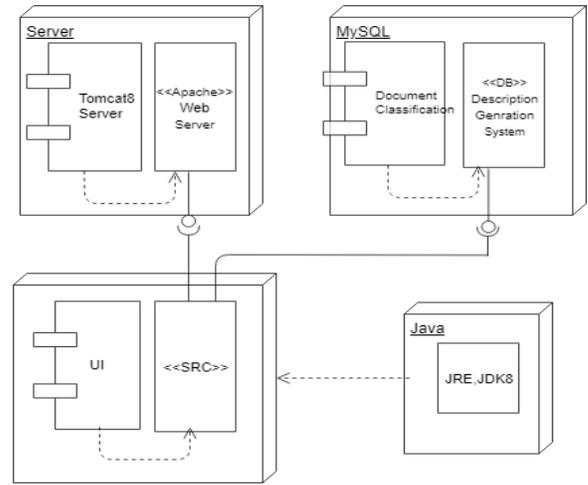
**Figure: Activity Diagram**

• **Component Diagram:-**



**Figure: Component Diagram**

• **Deployment Diagram:-**



**Figure: Deployment Diagram**

**Implementation:**

**Algorithms**

**NLP:** NLP is a way for computers to analyze, understand, and derive meaning from human language in a smart and useful way. By utilizing NLP, developers can organize and structure knowledge to perform tasks such as automatic summarization, translation, named entity recognition, relationship extraction, sentiment analysis, speech recognition, and topic segmentation.

NLP is used to analyze text, allowing machines to understand how human speaks. This human-computer interaction enables real-world applications like automatic text summarization, sentiment analysis, topic extraction, named entity recognition, parts-of-

speech tagging, relationship extraction, stemming, and more. NLP is commonly used for text mining, machine translation, and automated question answering.

- **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, hashtags), 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 the

automated chatbots, content analyzers, and consumer insight.

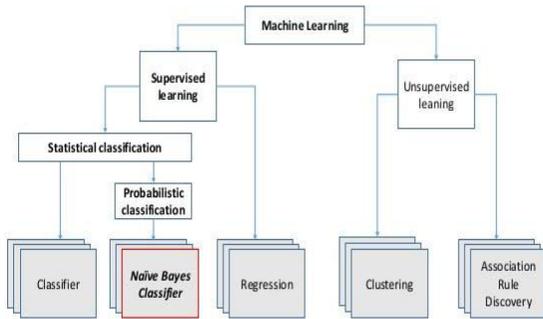
- **Sentiment:** NLP finds the sentiment from the statement whether the sentence is positive, negative or Very negative.

### Naïve Bayes Algorithm:

- 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. For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as ‘Naive’.
- 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.

### What is Naïve Bayes Classifier



Bayes theorem provides a way of calculating the posterior probability,  $P(c|x)$ , from  $P(c)$ ,  $P(x)$ , and  $P(x|c)$ . Naive Bayes classifier assumes that the effect of the value of a predictor ( $x$ ) on a given class ( $c$ ) is independent of the values of other predictors. This assumption is called class conditional independence.

$P(c|x)$  is the posterior probability of class (target) given predictor (attribute).

$P(c)$  is the prior probability of class.

$P(x|c)$  is the likelihood which is the probability of predictor given class.

$P(x)$  is the prior probability of predictor.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability  
Posterior Probability
Predictor Prior Probability

$$P(c|\mathcal{X}) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Problem: Players will play if the weather is sunny. Is this statement is correct

We can solve it using the above-discussed method of posterior probability.

$$P(\text{Yes} | \text{Sunny}) = P(\text{Sunny} | \text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$$

Here we have  $P(\text{Sunny} | \text{Yes}) = 3/9 = 0.33$ ,  $P(\text{Sunny}) = 5/14 = 0.36$ ,  $P(\text{Yes}) = 9/14 = 0.64$

Now,  $P(\text{Yes} | \text{Sunny}) = 0.33 * 0.64 / 0.36 = 0.60$ , which has higher probability

### TF-IDF Algorithm:

TF: Term Frequency, which measures how frequently a term occurs in a document. Since every document is different in length, It is possible that a term would appear much more time in long documents than shorter ones. Thus, the term frequency is often divided by the document length (aka. the total number of terms in the document) as a way of normalization:

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

IDF: Inverse Document Frequency, which measures how important a term is. While

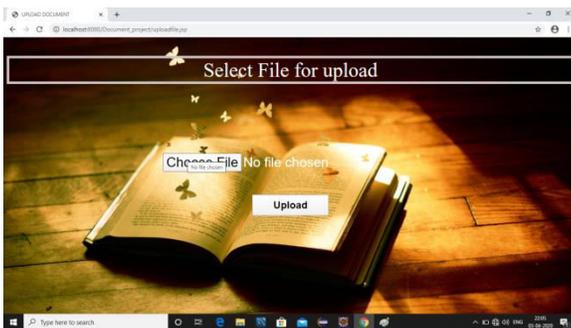
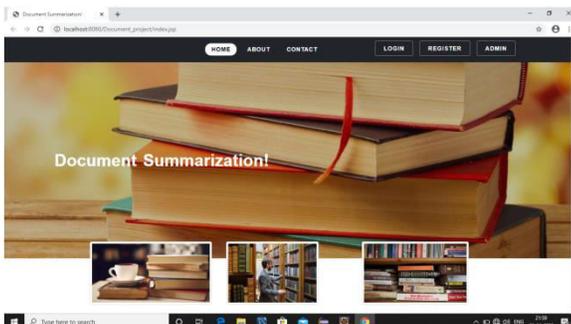
computing TF, all terms are considered equally important.

However, it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance.

Thus we need to weigh down the frequent terms while scaling up the rare ones, by computing the following:

**Formula:**  $IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it}).$

**Screenshot:**



**Conclusion:**

In this system, we have developed an application by which we can identify the domain of our document. And also relevance. 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.

**Future Scope:**

To Classify the Domain of paper from the number of pdf and also find out frequency relevanceToGenerate a Short Summary of that paper. We are going to classify the domain according to its content and divide it into groups. In the future, we can classify the Data of Image.

**Reference:**

[1]. Mehmet Baygin 2018 "classification of text document based on naive bayes using N-grams Features

[2].PrafullaBafna,DhanyaPramod,AnaghaVaidya(2016)"Document clustering:TF-IDF Approach"

[3]A.Helen Victoria, M.Vijayalakshmi.An Outcome based Comparative study of different Text Classification Algorithm. Volume 118 No. 22 2018, 1871-187

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

[5]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),4309-431.

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

[7] AmnaRahman, UsmanQamar(2016). A Bayesian Classifiers based GroupingPrototypical for Unconscious Text Arrangement.International Journal, 1(6).

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

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

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