# Survey on Research Paper Classification based on TF-IDF and Stemming Technique using Classification Algorithm

Kshitija G. Deshmukh
PG Student, Department of Computer Engineering
PES Modern college of Engineering
Pune, India

S. A. Itkar
Professor , Departement of Computer Engineering
PES Modern college of Engineering
Pune, India

## ABSTRACT
Classification System is the system of categorizing objects into classes or into groups of classes. It is used in many wide applications including text classification, web page classification, image classification, research paper classification etc. Various research papers are published online and offline. Text classification and class prediction is important for paper classification to reduce the feature size and to speed up the learning process of classifiers. Text classification is a growing interest within the research of text mining. This paper presents a survey on classification algorithm and stemming technique used for Text classification.

## Keywords
Text Classification, Stemming Technique, classes

## 1. INTRODUCTION
Number of research papers have been published online and offline with the increasing advance of computer technologies, which makes it difficult for users to search and categorize their research papers for a specific domain [1]. Therefore, it is desired that these large numbers of research papers are classified with similar domain so that users can find their interesting research papers easily [1].

The classification into predefine classes done in 3 ways, unsupervised, supervised and deep learning ways. Unsupervised learning is a machine learning technique, wherever ought not to supervise the model. Unsupervised learning formula ranked bunch [10],k-Means clustering[3].Supervised learning is that the machine learning task of learning a operate that maps an input to an output supporte example input output pairs. Supervised learning algorithm formula call tree [11],Naïve Bayes algorithm [12], SVM algorithm [13] and deep learning algorithms are Convolutional Neural Network (CNN) [14], Recurrent Neural Networks (RNNs) [4], HAN [15].

Stemming is the method of conflating the variant forms of a word into a common representation, the stem. For example, the words: "presentation", "presented", "presenting" could all common representation of "present". This is used procedure in text processing for information retrieval based on the assumption that move the question with term presenting implies an interest in documents containing the words presentation and presented [6].

The remaining paper is standardized follows. Related work of text classification is done in section 2.Text

Classification process in section 3.Section 4 comparative analysis of classification algorithm. Section 5 conclude the study and explain challenges.

## 2. RELATED WORK
### 2.1 Realated Work
The K-means algorithm is applied to classify the total papers into research papers with similar subjects, using the Term frequency inverse document frequency (TF-IDF) values of every paper [1] [5].

An artificial datasets such as news 20, Reuters, email, and analysis papers on completely different topics. Term Frequency-Inverse Document Frequency formula is employed along with fuzzy K means and hierarchical formula [2].TF-IDF method and framework for text classification. The framework allows classification in line with varied parameter, measure and analysis of results [3] [4].

The main of porter stemmer uses suffix denudation in English language. This stemmer could be a linear step stemmer. Specifically, it's 5 steps applying rules inside every step [11]. Inside every step, if a suffix rule matched to a word, then the conditions connected to its rule are tested on what would be the ensuing stem, if that suffix was removed, within the means outlined by the rule[9][10].

## 3. TEXT CLASSIFICATION
Text classification is the method of distribution tags or classes to text in line with its content. This section explains Classification Process for text using pre-processing, stemming technique, TF-IDF, a classification algorithm.

### 3.1 Data Preprocessing
Data processing methods are required in order to extract useful knowledge. The process of extracting information and knowledge from unstructured text documents is possible by using text mining [16].

### 3.1.1 Removal of stop word and symbol
Text classification preprocessor is processes words by removing Symbols removal, Stop words removal. The symbols are removed in pre-processing step and a stop word list is a list of commonly repeated features which appears in every abstract. The common features such as it, he, she and conjunctions such as and, or, but etc. are to be removed because they do not have effect on the categorization process [17].

The main motive behind the remove of stop words is to increase the execution speed and the accuracy. Stop words are typically one set of words. It means that completely different

for various varieties of application. As an example, a stop word list will contain

- Determiners: the, a, an, another

- Coordinating conjunctions: for, an, nor, but, or, yet, so

- Prepositions: in, under, towards, before

## 3.2 Stemming Algorithm

In text classification, stemming is that the method of reducing the words into root kind effectively.

### 3.2.1 Affix Removal Algorithm

In text classification affix removal stemming is to remove the endings of the word hold on first n letters, i.e. to truncate a word up to $n^{th}$ character and remove the remaining [18]. Affix removal stemmers remove the suffixes or prefixes kind of the terms leaving the stem. One of the examples of the affix removal stemmer is one that removes the plurals kind the terms. Some set of rules for a stemmer are as follows [6].

1) If a word ends in "ies" however not "eies" or "aies" Then "ies" -> "y"

2) If a word ends in "es" however not "aes", or "ees" or "oes" Then "es" -> "e"

3) If a word ends in "s" however not "us" or "ss" Then "s"

-> "NULL"

### 3.2.2 Table Lookup method

In text classification table lookup stemming one task to do stemming is to store a table of all index terms and their stems. Terms from the queries and indexes could then be stemmed via a lookup table by using b-trees or hash tables [6]. There are problems with this approach. The primary is that therefore making these lookup tables author want to extensively work on a language. There will be some probability that these tables could miss out some exceptional cases. Another drawback is that the storage overhead for such a table [20].

### 3.2.3 Successor Variety Stemming

Successor variety stemmers are based on the structural linguistics which determines the term and morpheme boundaries based on the distribution of phonemes. Successor variety type of a string is that the variety of characters that follow it in words in some body of text. As an example, consider a body of text consisting of following words. Able, ape, beatable, read, readable, reading, reads, red, ripe [6]. When text classification successor variety stemming given a set of a word's morphological variants, a potential stem may be derived heuristically, using a skillful analysis of prefix frequency and prefix length among the variants. E. g., the longest common prefix of the words \connection", \connect", \connectivity", \connecting" is \connect", that is also the stem. [21].

### 3.2.4 Porter Stemming Algorithm

It is an affix removal stemming algorithm. The context sensitive suffix removal algorithm is a stemmer. It is used of all the stemmers and implementations in various languages are available [6]. The main idea of porter stemmer uses suffix removing in English Language. This stemmer consists of 5 or 6 steps depend upon the method that is used to give the final stem. Original algorithm consists of only five steps. Every step applied the rules, and conditions also involved. If the rule is properly accepted, the suffixes are automatically removed

according to the condition, and then next step performed. The resultant stem end at rules and conditions [22].

### 3.2.5 Lovins Stemming Algorithm

It is an affix removal stemming algorithm. Lovins stemmer is to remove suffix from the term. This algorithm involved list of two hundred ninety four suffices 29 conditions and 35 transformation rules which have been used for longest match principal. The word is recoded using totally different table after the ending is removed [22]. It is a one pass, context sensitive stemmer that removes endings based on the longest-match principle [6].

## 3.3 TF-IDF Model

TF-IDF technique is used which eliminates the most common terms and extracts only most relevant terms from the corpus[7].Term frequency-inverse document frequency (TF-IDF) is a numerical statistic method that allows the determination of weight for every word in each document. The process is often used in natural language processing or in information retrieval and text mining [8].

### 3.3.1 Term Frequency (TF)

The method computes the number of repetitions of a word in the document [8].Suppose collection of text documents and want to rank that document is most relevant to the question, "the brown cow" an easy thanks to begin out is by eliminating documents that don't contain all 3 words "the", "brown", and "cow", however this still leaves several documents.

$$TFij = \frac{ni,j}{\sum_k nkj} \qquad (1)$$

Where,

ni,j :The number of occurrences of word i in document j

and $\sum_k$

nk,j : The total number of occurrences of words in document j

### 3.3.2 Inverse Document Frequency

While the TF means the number of occurrences of each word in a document, the inverse DF means how many times each word appears in the collection of documents. Inverse DF is calculated by dividing the total number of documents by the number of documents that contain a specific word. It is defined as [1]

$$DFij = \frac{|D|}{|\{d\epsilon D : t\epsilon d\}|} \qquad (2)$$

Where,

|D|: The total number of documents

|{d ∈ D: t ∈ d}| : The number of documents that keyword t

occurs.

### 3.3.3 TF-IDF

TF-IDF is used to convert a document into structured format [7].

TF-IDF process determines the relative frequency of terms in a particular document through an inverse proportion of the term over the total document corpus [8]. TF-IDF is an efficient and simple for matching words to documents that are relevant to the query [9].It is defined as:

$$TF\text{-}IDF=TF * IDF \qquad (3)$$

## 3.4 Classification Algorithm

### 3.4.1 Supervised Algorithm

Supervised learning algorithms are decision tree [11], Naive Bayes algorithm [12], and SVM algorithm [13]

### 3.4.1.1 Decision Tree

The Decision Tree is a classification method in data mining. As the classical algorithmic program of the choice tree ID3, C4.5, C5.0 algorithms have the benefits of high classifying speed, robust brain and straightforward construction.When using it to classify, there does exists the problem of falling to choose attribute which have more values, and observing attributes which have less values[11].

### 3.4.1.2 Naive Bayes

Naive Bayes is a machine learning algorithm whose classification efficiency is proven in applications like document categorization and e-mail spam filtering. This classifier learns through a document classification algorithm, and is based on an easy usage of the Bayes' rule [12].

$$P(c|d) = \frac{P(c|d)P(c)}{P(d)} \qquad (4)$$

### 3.4.1.3. SVM

SVMs are sets of related supervised learning process used for classification and regression. They belong to a family of generalized linear classification. SVM minimizes the classification error and maximize the geometric margin. Therefore SVM called Maximum Margin Classifiers. SVM is based on the Structural risk minimization .SVM map input vector to a higher dimensional space where a maximal separating hyper plane is built. [13].

### 3.4.2. Unsupervised Algorithm

Unsupervised learning algorithm Hierarchical clustering [10], k-Means clustering [23]

### 3.4.2.1 Hierarchical Clustering

A hierarchical clustering algorithm is a hierarchical merging or splitting based on a given data set. The production process of a set of nested clusters organized as a hierarchical tree. There are two main types of Hierarchical clustering algorithm, agglomerative and divisive hierarchical clustering. The agglomerative algorithm starts with the down points and every document is an individual cluster or class at first. At every step, the algorithm merges the closest pair of classes to a parent class until only one class or k classes left. The divisive algorithm starts with a up class or k classes. At every step, the farthest one point is divided from its parent class until each class contains a point [10].

### 3.4.2.2 K-means

K-means is one of the easiest clustering algorithms to group data, which aims to partition the samples into k sets with minimizing cluster error. K-means methodology to capture many cluster centroids for every class, and then select the high frequency words in centroids because the text features for categorization. The words extracted by k-means not only can represent every class clustering well, but also own high quality for semantic expression. [23].

### 3.4.3 Deep Learning Algorithm

Deep learning algorithms are convolutional neural network (CNN) [14], recurrent neural networks (RNN) [2], and HAN [15].

### 3.4.3.1 Convolutional Neural Network (CNN)

Convolutional Neural Networks, originally invented for computer vision, have been shown to achieve strong performance on text classification tasks as well as other traditional Natural Language Processing tasks, even when considering relatively simple one-layer models [14].

### 3.4.3.2.Hierarchical Attention Networks (HAN)

It consists of many parts: a word sequence encoder, a word-level attention layer, a sentence encoder and a sentence-level attention layer. Hierarchical Attention Network (HAN) that is designed to capture two basic insights about document structure [15]. The HAN classification model has two characteristics: first, it has two levels of attention mechanisms applied at the word and sentence level, enabling it to attend deferentially to greater and less important content once constructing the document representation. Second it's a hierarchical structure which mirrors the hierarchical structure of documents [15].

## 3.4.Recurrent Neural Networks (RNN)

Recurrent Neural Networks are one of the common Neural Networks used in Natural Language Processing. The idea behind a RNN is to consider the sequence of the input. Predict the next term (word) in the sentence need to remember what term appeared in the previous time step. These neural networks are known as Recurrent because this step is carried out for each input. As these neural networks consider the previous term during predicting, it acts as a memory storage unit that stores it for a short period of time. [2]. A perennial neural network is also a method of a sequence of arbitrary length by recursively applying a transition technique to its internal hidden state vector ht of the input sequence.The activation of the hidden state ht at time-step t is computed as a method f of the current input symbol xt and the previous hidden state ht−1 [24].

$$ht = \sum_{f(ht-1,Xt)}^{0} \quad \begin{matrix} t = o \\ otherwise \end{matrix} \qquad (5)$$

RNN model, which only exhibits a time complexity O(n).RNN model analyzes a text word by word, stores the semantics of all the previous text in a fixed-sized hidden layer [25].

## 4. COMPARTIVE STUDY

This paper compared the different algorithms and produces some results as follows, K-mean Classification - Higher time and space complexity stores all the instance, Noisy features degrades the classification accuracy, Naive Bayes classification - Independence assumption of features, Decision Tree Classification - Noise handling is bad, RNN - Ability to better capture the contextual information [2].

This section will provide the comparative study of k means, Naïve Bayes, RNN and Affix Removal method algorithm

| Sr. No. | Parameter | K means | Naïve Bayes | RNN |
|---------|-----------|---------|-------------|-----|
| 1 | Type of Algorithm | Unsupervised | Supervised | Deep learning |
| 2 | Basic | K- Mean clustering is to partition n observations into k cluster in which each observation | Naïve bayes classifier is a simple probabilistic classifiers based on | A recurrent neural network (RNN ) is is a category of artificial neural networks wher |

| | | | | | |
|---|---|---|---|---|---|
| | | belong to the cluster with the nearest mean. . | applying B ayes' theorem | ever connections between nodes form a directed graph | |
| 3 | Classific ation | Classification on the basis of k cluster | Classificati on on basis of Probability | Classification on basis of hidden neural network layer | |
| 4 | Learning method | It good learning method Without supervision | It better learning method as compare to K means because It working Under supervision | It better learning method as compare to Naïve Bayes because It working on deep learner method | |

**Fig 1: Comparative Study of Classification Algorithm**

| Sr. N o. | Paramet er | Lovins Stemmer | Porter Stemmer |
|---|---|---|---|
| 1 | Type of Algorith m | Affix Removal | Affix Removal |
| 2 | Basic | The Lovins stemmer removes the longest suffix from a word. Once the ending is removed, the word is recoded emplo ying a different table that produces varied changes to convert these stem into valid words. | Combination of smaller and less complicated suffixes. It's steps, and among every step, rules area unit applied till one of them passes the conditions. If a rule is accepted, the suffix removed consequently, and also the next step is performed |
| 3 | Advanta ge | Fast – single pass algorithm | Produces the best output as compared to other stemmers. |
| 4 | issue | Not all suffixes available. | Time consuming |

**Fig 2: Comparative Study of Affix Removal Method Algorithm**

## 5. CONCLUSION

The intension of this paper is to survey on research paper classification. Authors studies in space conclude that reduce the high dimensionality of feature space using porter stemming technique. Classification systems can classify research papers in advance by both of keywords and stemming with the support of high performance computing techniques.

There are various types of leaning algorithm like K-means, hierarchical clustering ,naïve bayes, decision tree, SVM and there are various types of stemming algorithm like affix removal, lookup table, successor variety stemming. More research is required in classification algorithm and stemming technique.

## 6. ACKNOWLEDMENT

## 7. REFERENCES

[1] Sang‑Woon Kim and Joon‑Min Gil, Research paper classification systems based on TF‑IDF and LDA schemes, Kim and Gil Hum.Cent. Comput. Inf. Sci. (2019).

[2] D. Yogeshwaran1, Dr. N. Yuvaraj ,Text Classification using Recurrent Neural Network in Quora, International Research Journal of Engineering and Technology (IRJET) Volume: 06 Issue: 02 Feb 2019.

[3] Fatiha Barigou ,Impact of Instance Selection on kNN-Based Text Categorization, Journal Information Process System, Vol.14, No.2, pp.418~434, April 2018.

[4] Pema Gurung and Rupali Wagh ,A study on Topic Identification using K means clustering algorithm: Big vs. Small Documents ,Advances in Computational Sciences and Technology ISSN 0973- 6107 ume 10, Number 2 (2017) pp. 221-233.

[5] Ms. Anjali Ganesh Jivani, A Comparative Study of Stemming lgorithms" Int. J. Comp. Tech. Appl. IJCTA NOV-DEC 2011, Vol 2 (6), 1930-1938

[6] Vairaprakash Gurusamy, S.Kannan, K.Nandhini,Performance Analysis: Stemming Algorithm for the English Language ,IJSRD International Journal for Scientific Research & Development Vol. 5, Issue 05, 2017 ISSN (online): 2321-0613

[7] Prafulla Bafna, Dhanya Pramod, Anagha Vaidya, "Document Clustering: TF-IDF approach", IEEE 2016.

[8] Bruno Trstenjak,Sasa Mikac, Dzenana Donko, "KNN with TF-IDF Based Framework for Text Categorization", Procedia Engineering 69,science Direct ( 2014 ) 1356 – 1364

[9] Juan Ramos,Using TF-IDF to Determine Word Relevance in Document Queries

[10] Anping Zeng, Yongping Huang," A Text Classification Algorithm Based on Rocchio and Hierarchical Clustering", D.-S. Huang et al. (Eds.): ICIC 2011, LNCS 6838, pp. 432–439, 2011.Springer-Verlag Berlin Heidelberg 2011.

[11] Mr. Brijain R Patel, Mr. Kushik K Rana,A Survey on Decision Tree Algorithm For Classification", 2014 IJEDR Volume 2, Issue 1 ISSN: 2321-9939.

[12] Dalibor Buzic, Jasminka Dobsa,Lyrics Classification using Naive Bayes, May 2018.

[13] Durgesh K. Srivastava, lekha Bhambhu, Data Classification Using Support Vector Machine, Journal of Theoretical and Applied Information Technology · February 2010.

[14] Alon Jacovi, Oren Sar Shalom, Yoav Goldberg, Understanding Convolutional Neural Networks for Text Classification, Proceedings of the 2018 EMNLP Workshop Black box NLP: Analysing and Interpreting Neural Networks for NLP, November 1, 2018.

[15] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, Eduard Hovy,"Hierarchical Attention NetworksforDocumentClassification,hovyg@cs.cmu.edu

[16] Hanumanthappa M and Narayana Swamy M,Language Independent Categorization of Documents Based on the Domain", Advances in Natural and Applied Sciences, 9(6) Special 2015.

[17] Jashanjot Kaur, Preetpal Kaur Buttar,A Systematic Review on Stopword Removal Algorithms", International Journal on Future Revolution in Computer Science & Communication Engineering April 2018 ISSN: 2454-4248 Volume: 4 Issue: 4

[18] Sandeep R. Sirsat, Dr. Vinay Chavan, Dr. Hemant S. Mahalle,Strength and Accuracy Analysis of Affix Removal Stemming Algorithms, Sandeep R. Sirsat et al, / (IJCSIT) International Journal of Computer Science and Information Technologies Aug 2015, Vol. 4 (2) , 2015, 265 – 269

[19] Mrs. R. Jayanthi , Ms. C. Jeevitha,"An Approach for Effective Text Pre-Processing Using Improved Porters Stemming Algorithm, IJISET - International Journal of Innovative Science, Engineering & Technology, Vol. 2 Issue 7, July 2015.

[20] Deepika Sharma, Stemming Algorithms: A Comparative Study and their Analysis, International Journal of Applied Information Systems (IJAIS) – ISSN : 2249-0868 Foundation of Computer Science FCS, New York, USA Volume 4– No.3, September 2012.

[21] Benno Stein and Martin Potthast,Putting Successor Variety Stemming to Work, Advances in Data Analysis Selected Papers from the 30th Annual Conference of the German Classi_cation Society (GfKl) Berlin, ISBN 978-3-540-70980-0, pp. 367-374, c Springer 2007.

[22] S.P.Ruba Rani, B.Ramesh, M.Anusha, Dr.J.G.R.Sathiaseelan,Evaluation of Stemming Techniques for Text Classification, International Journal of Computer Science and Mobile Computing, Vol.4 Issue.3, March- 2015, pg. 165-171.

[23] Xiaofei Zhou, Yue Hu, Li Guo,Text Categorization Based on Clustering Feature Selection, 1877-0509 © 2014 Published by Elsevier.

[24] Pengfei Liu ,Xipeng Qiu, Xuanjing Huang,Recurrent Neural Network for Text Classification with Multi-Task Learning, Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI-16).

[25] Siwei Lai, Liheng Xu, Kang Liu, Jun Zhao,Recurrent Convolutional Neural Networks for Text Classification, Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence.