

## Content Systemization using Naive Bayes Technique

Isha Pandya, Deepti Theng

*Computer Science and Engineering G. H. Rasoni College of Engineering Nagpur, India*  
*Computer Science and Engineering G. H. Rasoni College of Engineering Nagpur, India*

---

**Abstract-** *In today's scenario, the internet is growing day-by-day. And so the data and information is growing large. More and more users rely on search engines. It is the need of the hour to improve the speed and efficiency of the search engine. When a user searches a query he should be able to get most relevant results in very less time. Therefore there is need to do content systemization for efficient and faster search. One of the efficient ways to achieve is to use a supervised learning algorithm called Naive Bayes algorithm.*

**Keywords-** *Content systemization, supervised learning, Naive Bayes Technique*

---

### I. Introduction

World Wide Web is growing day by day and so the web pages are growing from thousands to millions. When user is searching for a query in search engine, he is expecting the correct and most relevant result. Because of this reason Content systemization has become important.

Content systemization is the process to group documents according to some predefined knowledge. Documents of same concept are grouped together. And non-related documents are in other category. For grouping of documents clustering techniques will be used.

Machine learning is an important concept used here. Machine learning is programming computers to optimize a performance criterion using example data or past experience. There are two things that need to be achieved in the learning process. First, the training needs to be done with an efficient algorithm to solve the optimized problem, and to store and process the data. Second, the trained model needs an efficient representation and algorithm for inference.

Naive Bayes technique is a supervised learning algorithm. These algorithms generate a function based on the training data set that maps inputs to predetermined labels or classes.

### II. Literature Survey

There are many algorithm like K-Nearest Neighbor, Neural networks, Support Vector Machine, decision tree. Many researchers have used different techniques for systemization. According to S.L. Ting.et.el, Naive Bayes is the best classifier used for content systemization. NB gives more accurate results than Support Vector Machine, NN and decision tree. NB gave 97% accuracy, SVM gave 96.6% accuracy, NN gave 93% accuracy, and Decision tree gave 91.1% accuracy. Time required is less as compared to other four algorithms. K-Nearest Neighbor is a lazy classifier and is bad if speed is important Naive Bayes is faster than K-NN. NB handles missing data while k-NN cannot.

According to Daniela Xhemali.et.el,[8] Accuracy of Naive bayes classifier is 95.20% while that of Decision tree is 94.85%. Precision of Naive bayes classifier is 99.37% while that of decision tree is 98.31%. Recall of Naive bayes classifier is 95.23% while that of decision tree is 95.90%. F-Measure of NB classifier is 97.26% while that of decision tree is 97.09%. Tokenization with preprocessing gives more accurate and fast result as compared to tokenization without preprocessing.

### III. Proposed system

The input of the system is URL. The contents of the URL is classified into predefined groups with the help of training data.

URL contents:

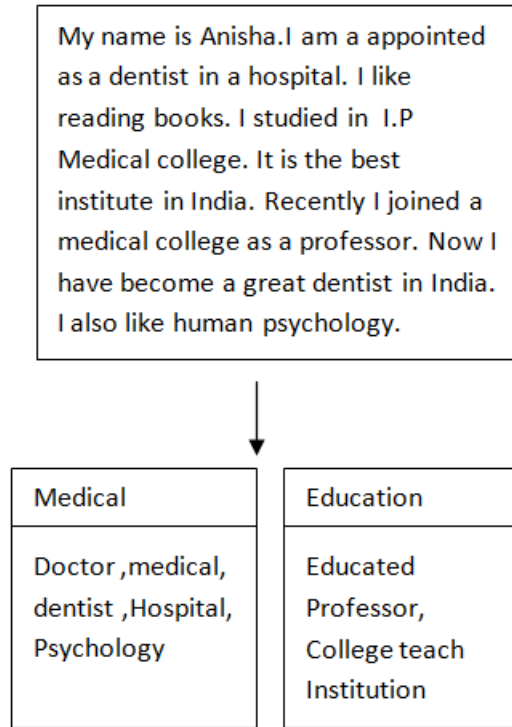


Fig. 1 Predefined categories

#### IV. Methodology

The input is supplied as a URL. The Contents of URL will undergo through following steps.

The first step is preprocessing. It includes Tokenization. Given a sequence of string as input, tokenization is the task of chopping the string into pieces called Tokens. It removes some characters like punctuation. Second step in preprocessing is removal of stop words. Some stop words are defined. And if in the contents of URL these words occur, they are removed because they have no significance in content systemization. Third step is to count frequency of words and sort them. Then preprocessed and tokenized text is obtained. This text is given as input to classifier. The classifier takes URL and output of training module as input and classifies it into a category. It uses Naïve Bayes theorem is based on concept of probability

1. Prior probability (A) = (Probability of URL occurring in class)

$$\frac{\text{no. of words in A}}{\text{total no. of words}}$$

2. Likelihood of X (new) given A =

$$\frac{\text{no. of words in A in the vicinity of X}}{\text{total no. of A}}$$

3. Posterior probability of X being in A,

**$\log(\text{Prior probability}) + \log(\text{likelihood of X being A})$**

where A= category and Highest probability value is the winner.

The intermediate step is to train the classifier. In training the classifier some categories are defined which contains some pre-defined words and serializes them to make learning data. These categories are given as input to trainer module. The trainer module finds the probability of each word. It serializes them and stores it in wordprob.srl which acts as local database. The output of training module is word probability.

## V. Result and discussion

The URL of a website is taken as input.



Fig. 2 Request URL for Searching

Tokenization is done to the contents of the website and removing punctuation marks are removed.

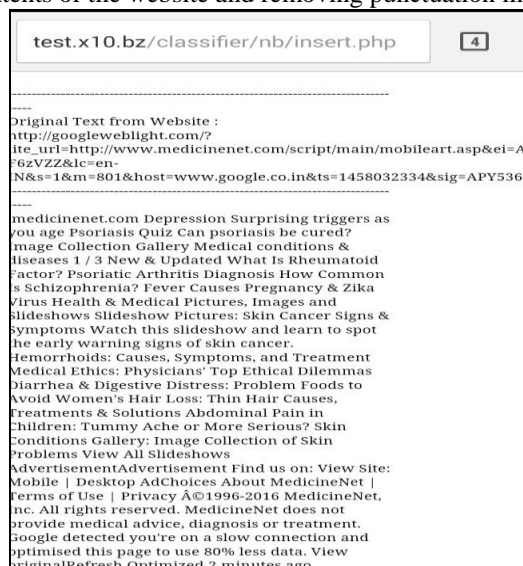


Fig. 3 Tokenization

Stop words are removed from the tokenized text.



Fig. 4 Removal of stop words

The words are sorted in the decreasing order of their frequency of occurrence.

```

Preprocessed and Tokenized text :
-----
Array ( [amp] => 7 [View] => 3 [Skin] => 3 [Causes]
=> 3 [MedicineNet] => 3 [Medical] => 3 [or] => 2 [All]
=> 2 [Image] => 2 [Collection] => 2 [Symptoms] => 2
[Slideshows] => 2 [Pictures] => 2 [Is] => 2 [Hair] => 2
[this] => 2 [Gallery] => 2
[AdvertisementAdvertisement] => 1 [Find] => 1
[Loss] => 1 [Women's] => 1 [Mobile] => 1 [Site] => 1
[us] => 1 [Conditions] => 1 [Children] => 1 [Tummy]
=> 1 [in] => 1 [Pain] => 1 [Abdominal] => 1 [Ache] =>
1 [More] => 1 [Solutions] => 1 [Problems] => 1
[Treatments] => 1 [Serious] => 1 [Desktop] => 1
[Thin] => 1 [Terms] => 1 [connection] => 1
[optimised] => 1 [slow] => 1 [you're] => 1 [Google]
=> 1 [detected] => 1 [page] => 1 [use] => 1 [minutes]
=> 1 [ago] => 1 [Optimized] => 1 [originalRefresh] =>
1 [less] => 1 [data] => 1 [treatment] => 1 [diagnosis]
=> 1 [-] => 1 [Inc] => 1 [Privacy] => 1 [Use] => 1
[About] => 1 [Avoid] => 1 [rights] => 1 [reserved] =>
1 [medical] => 1 [advice] => 1 [provide] => 1 [not] =>
1 [does] => 1 [AdChoices] => 1 [Ethical] => 1 [What]
=> 1 [Rheumatoid] => 1 [Updated] => 1 [New] => 1
[conditions] => 1 [diseases] => 1 [Factor] => 1
[Psoriatic] => 1 [Common] => 1 [Schizophrenia] => 1
[How] => 1 [Diagnosis] => 1 [Arthritis] => 1 [cured]
=> 1 [be] => 1 [Surprising] => 1 [triggers] => 1
[Depression] => 1 [com] => 1 [medicinenet] => 1 [as]
=> 1 [you] => 1 [Can] => 1 [psoriasis] => 1 [Quiz] => 1
[Psoriasis] => 1 [age] => 1 [Fever] => 1 [Pregnancy]
=> 1 [Ethics] => 1 [Physicians'] => 1 [Treatment] => 1
[Hemorrhoids] => 1 [skin] => 1 [cancer] => 1 [Top]
=> 1 [nbsp] => 1 [Distress] => 1 [Problem] => 1
[Digestive] => 1 [Diarrhea] => 1 [Dilemmas] => 1
[signs] => 1 [warning] => 1 [Images] => 1
[Slideshow] => 1 [Health] => 1 [Virus] => 1 [Zika] =>
1 [Cancer] => 1 [Signs] => 1 [spot] => 1 [early] => 1
[learn] => 1 [slideshow] => 1 [Watch] => 1 [Foods]
    
```

Fig 5. Preprocessed and tokenized text

Preprocessed and tokenized text is given input to the classifier.

```

-----
Input to classifier : amp, View, Skin, Causes,
MedicineNet, Medical, or, All, Image, Collection,
Symptoms, Slideshows, Pictures, Is, Hair, this,
Gallery, AdvertisementAdvertisement, Find, Loss,
Women's, Mobile, Site, us, Conditions, Children,
Tummy, in, Pain, Abdominal, Ache, More, Solutions,
Problems, Treatments, Serious, Desktop, Thin,
Terms, connection, optimised, slow, you're, Google,
detected, page, use, minutes, ago, Optimized,
originalRefresh, less, data, treatment, diagnosis, -,
Inc, Privacy, Use, About, Avoid, rights, reserved,
medical, advice, provide, not, does, AdChoices,
Ethical, What, Rheumatoid, Updated, New,
conditions, diseases, Factor, Psoriatic, Common,
Schizophrenia, How, Diagnosis, Arthritis, cured, be,
Surprising, triggers, Depression, com, medicinenet,
as, you, Can, psoriasis, Quiz, Psoriasis, age, Fever,
Pregnancy, Ethics, Physicians', Treatment,
Hemorrhoids, skin, cancer, Top, nbsp, Distress,
Problem, Digestive, Diarrhea, Dilemmas, signs,
warning, Images, Slideshow, Health, Virus, Zika,
Cancer, Signs, spot, early, learn, slideshow, Watch,
Foods
    
```

Fig 6. Input to classifier

```

-----
Classified as :
480.010309278350515480.010309278350515480.010309278350515480.01
581.65035086342engineering
450.010989010989011450.010989010989011450.010989010989011450.01
573.60439282275medical
medical
    
```

Fig. 7 Classified output

## VI. Future work

To further improve the search results and quality of classification we can apply Feature selection method. It makes training and applying a classifier more efficient by decreasing the size of the effective vocabulary. Furthermore preprocessing techniques like pruning, stemming can be applied in step 1 for

improving quality of search result. Semantic analysis can also be considered for more accurate classification. Semantic analysis refers to formal analysis of the meaning.

## VII. Conclusion

We succeeded in building NB classifier Tokenization with Preprocessing techniques are applied to improve the quality of systemization and search quality. Supervised learning technique called Naïve Bayes is applied for content systemization. The disadvantages of SVM, Neural network, Decision tree and K-nearest neighbor are overcome by Naïve bayes technique. The contents of URL obtained dynamically are classified accurately by this classifier. This gives high performance of search engines. We also presented some preprocessing techniques like tokenization, removal of stop words and frequency count of words. Preprocessing has huge impact on the performance of systemization. It improved quality of classification.

## References

- [1]. Muhammad Husnain Zafar<sup>1</sup> and Muhammad Ilyas<sup>2</sup> ,” A Clustering Based Study of Classification Algorithms ,” International Journal of Database Theory and Application Vol.8, No.1 (2015)
- [2]. Soumen Swarnakar<sup>1</sup>, Sangita Karmakar<sup>2</sup> , “ CONCEPT BASED CATEGORIZATION OF DOCUMENTS FOR SEARCH ENGINES “Volume: 04 Issue: 10 | Oct-2015
- [3]. Teena Rani, Ankush Goyal ,”Survey of Clustering Techniques for Information Retrieval in Data Mining ”, International Journal of Science, Engineering and Technology Research (IJSETR), Volume 4, Issue 4, April 2015
- [4]. Khushboo Thakkar, Urmila Shrawankar, “Test Model for Text Categorization and Text Summarization” , International Journal on Computer Science and Engineering (IJCSSE) ISSN :0975-3397 Vol. 3 No. 4 Apr 2011.
- [5]. Bhumika<sup>1</sup>, Prof Sukhjot Singh Sehra<sup>2</sup>, Prof Anand Nayyar<sup>3</sup>, “A REVIEW PAPER ON ALGORITHMS USED FOR TEXT CLASSIFICATION”, *International Journal of Application or Innovation in Engineering & Management (IJAEM)*, Volume 2, Issue 3, March 2013.
- [6]. JIANG Xiao-Yu, FAN Xiao-Zhong, Wang Zhi-Fei, Jia Ke-Liang, “Improving the Performance of Text Categorization using Automatic Summarization”, 978-0-7695-3562-3/09 \$25.00 © 2009 IEEE DOI 10.1109/ICCMS.2009.29
- [7]. S. Revathi , Dr.T.Nalini ,” Performance Comparison of Various Clustering Algorithm ”, International Journal of Advanced Research in Computer Science and Software Engineering , © 2013, IJARCSSE .
- [8]. Daniela XHEMALI, Christopher J. HINDE and Roger G. STONE,” Naïve Bayes vs. Decision Trees vs. Neural Networks in the Classification of Training Web Pages”, IJCSI International Journal of Computer Science Issues, Vol. 4, No. 1, 2009
- [9]. S.L. Ting, W.H. Ip, Albert H.C. Tsang Is Naïve Bayes a Good Classifier for Document Classification International Journal of Software Engineering and Its Applications Vol. 5, No. 3, July, 2011