# WEB DOCUMENT MINING USING INCREMENTAL APPROACH OF NEURAL NETWORK IN BPP

**[1]A.P.Tawdar, [2]M.S.Bewoor, [3]Prof.Dr.S.H.Patil**

M.Tech Department of Computer Engineering, Associate Professor, Department of Computer Engineering, Professor, Department of Computer Engineering

Bharati Vidyapeeth University College of Engineering Pune, India

*Abstract— Text and Document classification has gained a tremendous importance due to huge amount of information stored in databases and web user's dependency on the search engine's to retrieve relevant information. Text Classification also known as Text Categorization is the task of automatically classifying a set of text documents into different categories from a predefined set. Text Classification helps in quick and relevant information retrieval from relational databases, documents, text, multimedia files, and World Wide Web. The applications are wide and not limited to only text summarization, search engines, document clustering and spam filtering. For Information Retrieval (IR) and Machine Learning (ML), TC uses several tools and has received much attention in the last decades. In this paper, first classifies the text documents using MLP based machine learning approach (BPP) and then return the most relevant documents. And also describes a proposed back propagation neural network classifier that performs cross validation for original Neural Network. In order to optimize the classification accuracy, training time. Proposed web content mining methodology in the exploration with the aid of BPP. The main objective of this investigation is web document extraction and utilizing different grouping algorithm. This work extricates the data from the web URL.*

*Index Terms— Back Propagation Algorithm, Neural Network, Information Retrieval, Information Extraction, Clustering, Steaming, Stop word, Feature Extraction.*

## I. INTRODUCTION

Classification and prediction are two modes of data analysis and data cannot be used without analysis. Classification is used to extract models describing important data classes and prediction model gives the future data trend. Information Retrieval (IR) is the science of searching for information within databases, and the World Wide Web. The breakthrough of the Internet and web search engines have urged scientists and large firms to create very large scale retrieval systems to keep pace with the exponential growth of online data. In IR system is user first submits a query which is executed over the retrieval system. The latter, consults a database of document collection and returns the matching document. In general, in order to correctly classify unseen documents, it is necessary to train it with some pre-classified documents from each category that is tanning dataset, in such a way that the classifier is then able to understand the model it has learned from the pre-classified documents and use that model to correctly classify the unseen documents.

Extraction of net data could be an important method for data integration .Web pages may give the same or analogous information. This makes the addition of information fascinating task. The deep Web contents are accessed by queries submitted to net databases and also for retrieved data. Web pages are data records. The distinctive Web pages are made progressively and hard to list by routine crawler based web crawlers, in particular Google and Yahoo. In this paper, depict this kind of exceptional Web pages as deep Web pages. A noteworthy issue of online web crawlers is that the unit results are a whole Web document. Human exertion is obliged to inspect each of the returned sections to separate exact information. Automatic information extraction frameworks can automate the task of successfully recognizing the pertinent content sections inside of the document. Information Extraction (IE) is concerned with extracting pertinent information from a gathering of archives. It includes methods and algorithms extract knowledge from distinctive data repositories such as transactional databases, data warehouses, text files, WWW and then forwards. Most of the online resources square measure as machine-readable text Mark-up Language (HTML) documents, that square measure seen by net browsers. In this manner, the need for automated, adaptable Web Information Retrieval (IR) tools that that extract information and data from the online pages and transfer into a significance and valuable structures for more investigation will turn into an extraordinary need.

## II. RELATED WORKS:

Tak-Lam Wong and Wai Lam [1] developed a new attribute discovery via Bayesian learning approach which can automatically adapt the information extraction patterns learned previously in a source web site to new unseen web sites and discover new attributes together with semantic labels. Extensive experiments from more than 30 real time web sites are in three different domains were conducted and the results exhibit that the framework achieves a very promising performance.

Rajendra Kumar Roul [2] has proposed web document clustering using data mining. This paper studies some clustering methods relevant to the clustering document collections and, in consequence, web data. This method of cluster analysis seems to be relevant in approaching the cluster web data. The graph clustering is also described in its methods to contribute significantly in clustering web data. Based on previously presented information, the core section provides an overview approaches to clustering in the web environment.

Jiang Su et al [3] proposed data classification using semi-supervised multi-modal Naive Bayes. It presents Semi-supervised Frequency Estimate (SFE), a novel semi-supervised parameter learning method for MNB. They first point out that EM's objective function, Maximizing Marginal Log Likelihood (MLL), is quite different from the goal of classification learning, i.e. maximizing conditional log likelihood (CLL). Then propose SFE that uses the estimates of word probabilities obtained from unlabelled data, and class conditional probability given a word, learned from labeled data, to learn parameters of an MNB model.

Amit Ganatra [4] has proposed initial classification through back propagation algorithm. This paper says initial classification using genetic and neural network algorithm. Performing weight adjustment in order to minimize the Mean Square Error between obtained output and desired output is the main goal of this hybrid algorithm. For reducing the search space of Genetic algorithm it is better to apply back propagation algorithm first. Hence the problem of local minima is solved. For the purpose of accelerating neural network training the

proposed algorithm exploits the optimization advantages of GA. BP algorithm is sensitive to initial parameters and GA is not. As compared to the GA, BP algorithm has high convergence speed.

Yan Liu [5] proposed a novel deep learning model for query-oriented multi documents summarization. Accordingly, the empirical validation on three standard datasets, the results not only show the distinguishing extraction ability of QODE but also clearly demonstrate our intention to provide human-like multi document summarization for nature language processing.

Saduf, Mohd Arif Wani [6] proposed the comparative study of learning in neural network.

Citra Ramadhenal [7] has proposed classification based on error rate. In order to minimize the Mean Square Error Hybrid algorithms are used to perform weight adjustment. First create a model by running the algorithm on the training data. Then test the model to identify class of new data for a class label. Then for classification this data is given to the Back propagation algorithm. After applying Back propagation algorithm, for weight adjustment genetic algorithm is applied. The developed model can then be applied to classify the unknown tuples from the given database and this information may be used by decision maker to make useful decision. But it provides less efficiency.

Daniel Soudry [8] proposed data classification in neural network using discrete continuous weight. In intrusion detection and classification using back propagation neural network approach were followed. It first collects the data set then the data is pre-processed. BPNN classifier is built for detection and classification of events. In BPNN classifier, first design network and set parameters then initialize weights with random values; finally calculate the actual output from the input. Finally, the Results showed are, it classifies instances into several attack types with low detection rate.

## III. PROPOSED SYSTEM:

A web document is similar in concept to a web page. Every web document has its individual URI. Note that a Web document is not the same as a file: a single web document can be accessible in various arrangements and dialects, and a single document, for instance a PHP script, may be in charge of creating a substantial number of web documents with different URIs. A Web document is characterized as something that has a URI and can return representations of the identified asset in response of HTTP requests.

The usual web content extraction methods concentrate only on extracting the content without checking whether the content is relevant or not. The proposed approaches include diverse algorithm in a comparative manner to evaluate the performance measures of the web contents in an efficient manner. For extraction purpose it offers web URL or web documents.
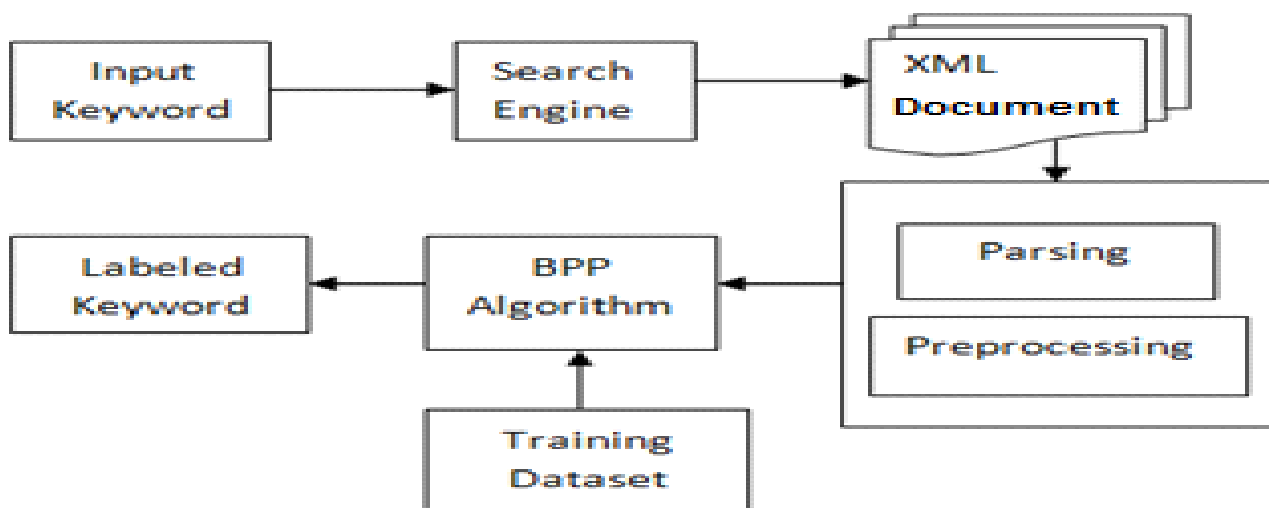


**Figure1. Proposed System.**

- **Input:** Get Unknown Text Document of Searched keyword
- **Output:** Labeled to the Text Document

  Class: Sports, Education, Technology etc

  **Modules**
- **XML Document Extraction**
- **Parsing**
- **Preprocessing**
- **Classification**

**1) Parsing:**

Jsoup is a java html parser is used. It is a java library that is used to parse HTML document. Jsoup provides api to extract and manipulate data from URL or HTML file. It uses DOM, CSS and Jquery-like methods for extracting and manipulating file.

The parser will make every attempt to create a clean parse from the HTML you provide, regardless of whether the HTML is well-formed or not. It handles:

• unclosed tags (e.g. <p>Lorem <p>Ipsum parses to <p>Lorem</p> <p>Ipsum</p>)

• Implicit tags (e.g. a naked <td>Table data</td> is wrapped into a <table><tr><td>...)

• reliably creating the document structure (html containing a head and body, and only appropriate elements within the head)

The parse (String html, String baseURI) method parses the input HTML into a new Document. The base URI argument is used to resolve relative URLs into absolute URLs, and should be set to the URL where the document was fetched from. If that's not applicable, or if you know the HTML has a base element, you can use the parse (String html) method.

**2) Preprocessing:**

It include two modules:

**a)** Stop Word Removal

**b)** Stemming

**Stop Word Removal:**

Sometimes, some extremely common words which would appear to be of little value in helping select documents matching a user need are excluded from the vocabulary entirely. These words are called *stop words*. The general strategy for determining a stop list is to sort the terms by *collection frequency*, and then to take the most frequent terms, often hand-filtered for their semantic content relative to the domain of the documents being indexed, as a *stop list*, the members of which are then discarded during indexing.

**Stemming**:

In this method words shorter than n are kept as it is. The chances of over stemming increases when the word length is small.

Rules in porter stemming algorithm are separated into five distinct steps:

**1)** Gets rid of plurals and -ed or -ing. eg-> caress ponies -> ponities -> ti caress -> caress cats -> cat

**2)** Turns terminal y to i when there is another vowel in the stem. eg happy->happi

**3)** Maps double suffices to single ones. so -ization ( = -ize plus -ation) maps to -ize etc.

**4)** Deals with -ic-, -full, -ness etc. similar strategy to step3.

**5)** Takes off -ant, -ence etc.

## 1) Implementation of Back Propagation Algorithm

The back-propagation algorithm consists of the following steps: 1. Initialization: At first the algorithm has to be initialized considering no prior information is known and picking the synaptic weights and thresholds from a uniform distribution. The type of activation function is sigmoid. 2. Presentations by Training Examples: The network has to be presented by epochs of training examples to perform forward and backward computations. 3. Forward Computation: Let us consider, the input vector to the layer of sensory nodes is x(n) and the desired response vector is d(n) which is in the output layer of computation nodes. In forward computation, the network's local fields and function signals are computed by proceeding forward through the network by layer by layer basis. Implementation of Back Propagation Algorithm the back-propagation algorithm consists of the following steps: 1. Initialization: At first the algorithm has to be initialized considering no prior information is known and picking the synaptic weights and thresholds from a uniform distribution. The type of activation function is sigmoid. 2. Presentations by Training Examples: The network has to be presented by epochs of training examples to perform forward and backward computations. 3. Forward Computation: Let us consider, the input vector to the layer of sensory nodes is x(n) and the desired response vector is d(n) which is in the output layer of computation nodes. In forward computation, the network's local fields and function signals are computed by proceeding forward through the network by layer by layer basis. If sigmoid function is used, the output signal is obtained by the equation below:

$$= \varphi_j(V_j(n)) \qquad (1)$$

If l=1 which means the j neuron is in the first hidden layer then we get,

$$= X_j(n) \qquad (2)$$

Here, $X_j(n)$ is the jth element of the input vector x(n). Let, L is the depth of network. If the neuron j is in the output layer that means l= L then

$$= O_j(n) \qquad (3)$$

So the error signal will be

$$e_j(n) = d_j(n) - o_j(n) \qquad (4)$$

Here, $d_j(n)$ is the jth element of the vector of desired response d(n).

5. Iteration: Finally the forward and backward computations have to be iterated until the chosen stopping criterion is met. As the number of iterations increases the momentum and learning-rate parameters are adjusted by decreasing the values.
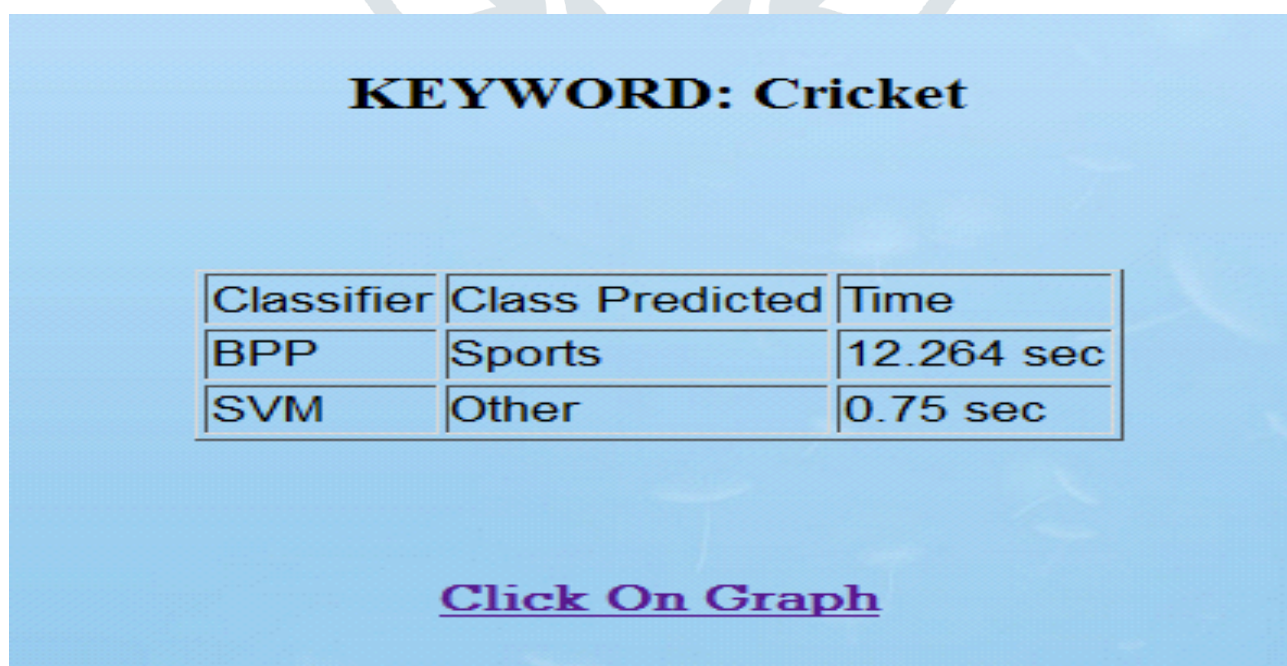
## IV. Result Analysis:



**KEYWORD: Cricket**

| Classifier | Class Predicted | Time |
|---|---|---|
| BPP | Sports | 12.264 sec |
| SVM | Other | 0.75 sec |

**Click On Graph**

**Figure2. Different Classification Comparison**

This table shows comparative analysis between Back Propagation Algorithm and Support Vector Machine. As soon as dataset increased, this will take less time than SVM for classification.
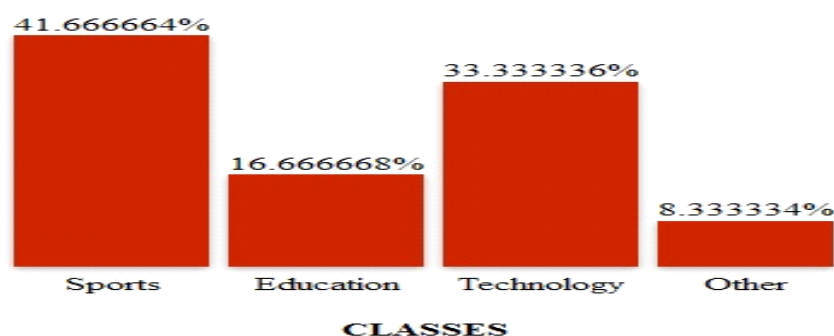
GENERATED GRAPH



**Figure3. Graph Generated**

The graph shows that the weight assigned for each keyword in data set. When there is new search is happened then that search keyword will match with training dataset that have implemented. When keyword not matches with the training dataset then that keyword will categorized as other.

## V. CONCLISION:

First, utilized our developed text mining algorithms, including text mining techniques based on classification of data in several data collections. After that, employ exiting neural network to deal with measure the training time for five data sets.

BPN is a very popular algorithm in the applications pattern matching, character recognition etc., Here this algorithm is discussed with reference to the Text categorization problem. This algorithm is yet to be implemented for this problem.

## VI. FUTURE WORK:

Implementation of this algorithm is considered as part of the future work. Apart from BPP there are other algorithms found in neural network systems. Also as in this system some fixed labeling is considered instead of this system can automatically generate new class as new keyword is found. Training dataset should train as new keyword found.

## VII. REFERENCES:

[1] Tak-Lam Wong and Wai Lam, "Learning to Adapt Web Information Extraction Knowledge and Discovering New Attributes via a Bayesian Approach", IEEE Transactions on Knowledge and Data Engineering, Vol. 22, No. 4, pp. 523-536, 2010.

[2] R.K. Roul and S.K. Sahay, "An Effective Approach for Web Document Classification using the Concept of Association Analysis of Data Mining", International Journal of Computer Science and Engineering Technology, Vol. 3, No. 10, pp. 483-491, 2012.

[3] Jiang Su, Jelber Sayyad Shirab and Stan Matwin "Large Scale Text Classification using Semi-supervised Multinomial Naive Bayes", Proceedings of the 28th International Conference on Machine Learning, pp. 97-104, 2011.

[4] Amit Ganatra, Y. P. Kosta, Gaurang Panchal and Chintan Gajjar, "Initial Classification Through Back Neural Network Following Optimization Through GA to Evaluate the Fitness of an Algorithm", International Journal of Computer Science & Information Technology, Vol. 3, No. 1, pp. 98-116, 2011.

[5] Yan Liu, Sheng-hua Zhong, Wenjie Li, "Query-Oriented Multi-Document Summarization via Unsupervised Deep Learning", Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, pp. 1699-1705, 2012.

[6] Saduf and Mohd Arif Wani, "Comparative Study of Back Propagation Learning Algorithms for Neural Issue Networks", International Journal of Advanced Research in Computer Science and Software Engineering, Vol. 3, No. 12, pp. 1151-1156, 2013.

[7] Citra Ramadhena, Ashraf Osman Ibrahim and Sarina Sulaiman, "Weights Adjustment of Two-Term Back-Propagation Network Using Adaptive and Fixed Learning Methods", International Journal of Advances in Soft Computing and its Application, Vol. 5, No. 2, 2013.

[8] Daniel Soudry, Itay Hubara and Ron Meir, "Expectation Back propagation: Parameter-Free Training of Multilayer Neural Networks with Continuous or Discrete Weights", Advances in Neural Information Processing Systems, pp. 963-971, 2014.

[9] A. P. Tawdar, M. S. Bewoor and Prof. Dr. S. H. Patil, "BACK PROPOGATION BASED TRAINING ALGORITHEM FOR TEXT AND DOCUMENT MINING", International Journal of Engineering Technology and Applied Research Volume 1 Issue 2, July 2014 ISSN: 2347 – 9396