# Study on Pattern Revealing in Text Mining

Prof. Thirumahal Rajkumar, Prof. Gayatri Dantal, Mr. Balmukund Dubey
I.T Engineering, Computer Engineering, Computer Engineering
Mumbai University, Mumbai University, Mumbai University
Mumbai India, Mumbai India, Mumbai India
balmukunddubey1989@gmail.com

**Abstract**: Today most of the data available in the digital form. In the past lots of people using the phrases related to the hypothesis to the document of the information and topic should be performing better result of the terms. In this paper we use the more important point include in the data mining method and this method give the better improve the effectiveness or performance of the patterns also in that we implementing the pattern detection method these are solved the problem of the term based method and give the good result and which is helpful in information retrieval system. This paper present the effective pattern discovery technique which include the process of the pattern deployed and pattern evolving and it will give the better improvement to the effectiveness.

**Keywords**: Text Mining, Pattern Mining, Pattern Taxonomy, Pattern Evolution.

## INTRODUCTION

The growth of the digital data made available in the last few years. The very important point is used in the Information discovery in the data mining technique and data processing has a good deal of the with associate with the close need for turning such that knowledge or basic idea which can be useful  into helpful data and knowledge [1]. There are contain the large number of the several applications like business management and market analysis and research and, it will profit by the or the user and employment of the information and information extracted from large amount of the database and  outsized quantity of data. That useful in the Information discovery will be viewed or look like a because that method of nontrivial is to be extraction of data from massive or huge amount of the databases, information that is local conferred within the knowledge discovery of the method, in the last part unknown and probably helpful for users. Data mining is to be containing the very important absolutely necessary step within the method of knowledge discovery in databases [4]. There are very important parts which contain the many types of data mining techniques are which represent the used sequential pattern mining and closed sequential pattern etc. this two pattern mining technique are very useful in the data mining. The very important concept is to be a  Data mining and they can easy to be allow to the process of retrieving or accessing to the interesting or relevant data knowledge from the huge amount of the database or storage area like database [2]. In this proposed system, contain the very important part is to be discovery of patterns will be done efficient through the very important part is to be pattern evolution and another important part is to be pattern deploying technique [3]. And this technique useful in the system and it will not only find the useful patterns but also efficiently use the and update them and to be find the requirement of the relevant data and important requirement of the interesting information from the database. Always very difficult to solved the problem of low frequency and misinterpretation. In that the system is supposed to develop the concept of the knowledge discovery model and this model handle the problem and they can efficiently use and easy to update and understand the patterns [1].

## LITERATURE SURVEY

### Study on Phrase Based Approach

In these phrases based method contains the huge amount of data and occurs the number of the problem in the term. In this we assume that phrase-based method give better result to the term based method [1][3], that why always phrases may carry more semantic meaning like contain the huge amount of information. We know that phrases are always handles the small ambiguous and large types of discriminative or information than individual terms, there are some reasons for the discouraging performance include:
1) In that Phrases always contain the inferior statistical number of the properties to terms,
2) In this Phrase based always problem arias on the low frequency in the operation handling.
3) the most big problem arises of the Phrase based is the huge or large numbers of redundant or repeating and noisy or bugs phrases [1]

### 1. Feature Selection and Feature Extraction for Text Categorization

In this we study the effect of selecting problem varying numbers and related kinds of information features for use in the predicting large number of the category of the membership was investigated on the large number of the Reuters and MUC-3 which

contain the number of the text categorization with the data sets [4]. And Good categorization performance the good achieved using a large number of the statistical classifier and require. the proportional assignments of the strategy and method of the feature selection and feature extraction. In that important optimal feature set size for word- based indexing was found to be surprisingly low between the 10 to 15 features contain the large training sets or huge data. And extraction of large number of the new text features by syntactic analysis and generating the feature clustering was investigated on the number of the Reuters data set are to be contain in the data [2]. The Syntactic indexing phrases are very important things which use in the feature, and contain the clusters of this number of the phrases, and also contain the clusters of words and these are always present on the data were all found and provide the less or small effective representations than of the individual words or single word[5]. In that we use the indexing language which is very useful and used to represent the number of texts influences how they can handle the problem in the easily and effectively a text categorization system can be built, also which can be builds whether the system is built by human engineering and statistical training, or a combination of these two main part. The important things that are to be simplest indexing languages are very impotent which are formed by treating each word as a requirement of the feature. These are very important words which contain the same amount of the properties, such as synonymy and polysemy that make them a less than ideal indexing language. Synonymy which share the same meaning and polysemy which share the two or more different meaning that contains [5].

**Study on Pattern Based Approach**

**1. Identifying Comparative Sentences in Text Documents**

In this paper we study the number of problem is to be identifying in the sentences in the large number of text documents. In that contain the problem is related to the text document arises quite opinion and identification of the sentence or classification. In this sentiment classification studies the large number of problem of classifying a number of documents or a sentence based on the opinion of the author [3]. An the large number of  important application area of the sentiment and opinion identification is business intelligence as a large number of the  product manufacturer always needs to the consumers' related require opinions on the requirement of the user products or component. The differentials on the different method or other hand can be related to the subjective or objective requirement. Furthermore, a comparison is not only concerned with an object in isolation. Instead, it also compares the object with number of the others product. In this pattern based approach tells the identifying is the comparative sentences is also very useful in practice and understanding of the problem because direct comparisons are always very important thing perhaps to the one of the most convincing way to the author ways of evaluation, which may  be more powerful  than opinions on the each related  individual object [3][2].

**Study on Keyword Based Approach**

In this Keyword based contain the bag of words scheme is a typical nothing but keyword based representation in the area of information retrieval. The main disadvantage of the keyword based is that relationship among word cannot be reflected. And another problem are coming that considering single word as a feature is the semantic ambiguity which contain the two part Synonyms and homonym [3].
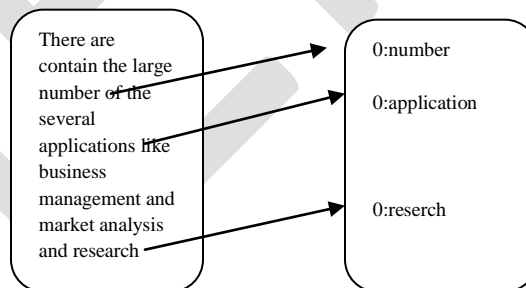


Figure 1: Bag of word

## PROPOSED SYSTEM

1. In this problem we use the An effective pattern discovery technique is to discovered the problem.
2. And the next is to be Evaluates of the patterns and then evaluates term weights according to the number of the distribution of terms in the form of the discovered patterns.
3. Solves Missing value or element of the Problem
4. In this we Considers the influence of the patterns from the negative document of the training for examples to understand find ambiguous or noisy patterns and handle the problem and tries to reduce the influence for the related of the low-frequency problem.
5. In that process of updating number of the related noise or ambiguous and this can provide to the patterns can be referred as pattern evolution.

6. The proposed approach can improve the accuracy of evaluating term weights because discovered patterns are more specific than whole documents.

7. The important point of the In training phase the discovery patterns in positive documents based on a min_sup are found, and another part contain the evaluates term supports by deploying patterns to terms

8. In this Testing Phase to check the revise term supports and containing the document using noise negative documents in D based on an example of the algorithms coefficient

9. The number of the incoming documents then can be sorted in the based on these weights or order in the weight.

Advantages of proposed system:

1. In this proposed method is used to improve the accuracy of the evaluating term weights method.

2. The another part discovered patterns are more useful in the specific than contain the whole documents.

3. Avoid the problem of phrase-based method to the using the very important point is to be pattern-based method.

4. The very useful Pattern mining techniques can be used to find list of the text patterns.
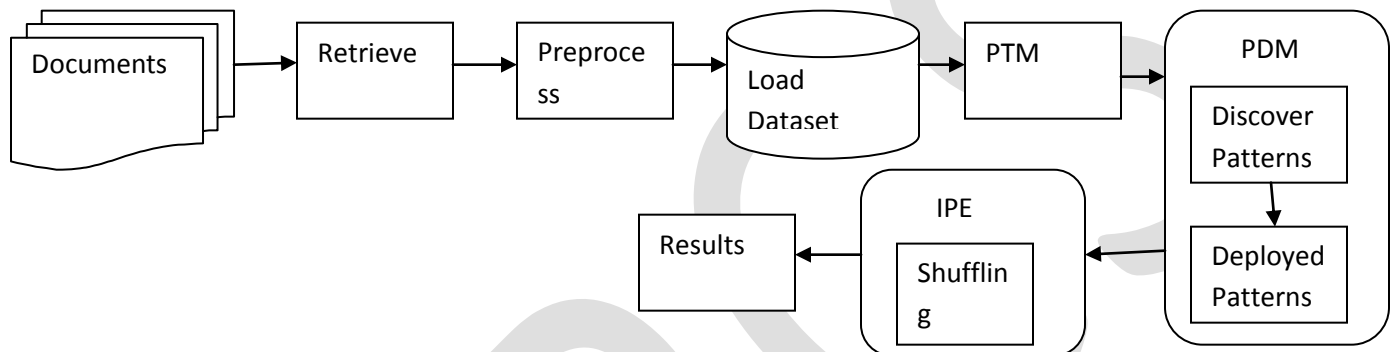


Figure 2: System Architecture

1. Loading document: In this above System Architecture, to load the number of all documents. And this Document is pass to the next step. And the user to retrieve number of the documents. This document is given to next process. That process is name call as preprocessing [3][4].

2. Text Preprocessing: once a completed the load document the next step is that Text Preprocessing. In that retrieved number of document preprocessing is done in module. In this text Preprocessing contain the two types of process is doing one is the stop words removal and second is that text stemming .The meaning of that Stop words are words which are filtered out the prior to, or after, processing of natural language data. And another meaning of the Stemming is the process for reducing inflected or sometimes derived words to their stem base or root form. It generally a written word forms [3][4].

3. Pattern taxonomy process: the third part of the system is the Pattern taxonomy model in this module; the number of documents is split into set of the paragraphs. Each paragraph is considered to be number of document. In each document, the set of terms are extracted. The terms which represent can be extracted from set of positive documents [3][4].

4. Pattern deploying: The second last step is that pattern discovery discovered patterns are summarized. The Deploying with Relevance Method algorithm is used to discover all patterns in positive documents are composed. The term supports are calculated by all terms in Deploying with Relevance Method. Term support means weight of the term is evaluated [3][4].

5. Pattern evolving: the final step is the pattern evolving In this module used to identify the noisy or bugs patterns in number of documents. Sometimes system which give the falsely identified negative document as a positive. So that noise is occurred in positive document. The noised pattern named as offender. If partial conflict offender contains in positive documents [3][4].

## USING ALGORITHMS

There list of algorithms are use in the pattern taxonomy

1) SPMining [5]

2) Deploying with Relevance Method

3) Evolving of pattern

4) Shuffling

**Algorithm 1.** SPMining (PL, min_sup) [5]

**Input:** A set of nTerms frequent sequential patterns.PL; Minimum support (min_sup).

**Output:** A set of frequent sequential patterns (SP).

**Method:**

1)  SP ←SP – {Pa ∈SP |∃Pb ∈PL} such that_len (Pa) = len( Pb) -1 ∧Pa ⊂Pb ∧suppa(Pa)= suppa(Pb) } /* pruning */

2)  SP ←SP ∪PL     /* add found patterns */

3)  PL′ ← {∅}    /* PL′: set of (n+1) Terms frequent sequential patterns */

4)  foreach pattern $p$ in PL do begin

5)  Generate p-projected database PD

6)  foreach frequent term $t$ in PD do begin

7)  P′ ←p ⋈t   /* P′: set of (n+1) Terms sequential candidates */

8)   If supp (P′) ≥ min_sup then

9)     PL′ ←PL′ ∪ P′

10)    End if

11)  End for

12) End for

13) If |PL′ | = 0 then

14)   Return   /* no more patterns found */

15) Else

16)   call SPMining(PL′, min_sup)

17) End if

18) Output frequent sequential patterns in SP

In this Sequential Pattern Mining Algorithms we use the pruning scheme for meaning that removing the non-closed related pattern during the process of Sequential Pattern discovery. And the main concept is to find the a set of frequent Sequential Pattern in this algorithms [5].

**Algorithm 2.**Deploying with Relevance Method [1]

**Input:** a set of positive documents ($D^+$), minimum Support (min_sup);

**Output:** New set of terms, a set of vectors (Δ).

**Method:**

1)Δ← ∅

2) Foreach document d in ($D^+$), do begin

3) Extract 1Term frequent pattern PL from d

4) SP=SPMining (PL, min_sup) // call Algorithm SPMining

5)  $\vec{d} \leftarrow \emptyset$

6) Foreach pattern p in SP do begin

7) $\vec{d} \leftarrow \vec{d} \oplus p$ // p' is the expanded form of p

8) End for

9) $\Delta \leftarrow \Delta \cup \{\vec{d}\}$ p

10) End for

     In this Deploying with Relevance Method algorithms the main process of it pattern Deploying occurs in the above algorithms and return the output is to the set of vectors [1].


**Algorithm 3**. Evolving of pattern ($\Omega, D^+, D^-$) [1]

**Input:** A list of deployed patterns $\Omega$; a list of positive and negative documents, $D^+$ and $D^-$

**Output:** A set of term weight pairs $\vec{d}$

**Method:**

1.$\vec{d} \leftarrow \emptyset$ // It give minimum threshold Value

2. $\tau \leftarrow$ Threshold ($D^-$)

3. $foreach\ negative\ documents$ "nd "$in\ D^-$ do begin

4. If Threshold ({nd}) $> \tau$ then

5.$\Delta p \leftarrow \{dp \in \Omega|$ termset (dp)$\cap$ nd $\neq \emptyset\}$

6. Shuffling ($nd, \Delta_p$)  // call shuffling Algorithms

7. End if

8. foreach deployed  pattern dp in $\Omega$ do begin

9. $\vec{d} \leftarrow \vec{d} \oplus p$

10. End for

11. End for

     In this Evolving of pattern algorithms we take the input A list of deployed patterns $\Omega$; a list of positive and negative documents, $D^+$ and $D^-$ and it will return a set of term weight pairs $\vec{d}$ also in that we call the Shuffling algorithms.

**Algorithm 4.** Shuffling ($nd, \Delta_p$) [5]

**Input:**  negative document $nd$and a list deployed patterns $\Delta p$.

**Output:** updated deployed patterns.

**Method:**

1: foreach deployed patterns $dp$ in $\Delta_p$ **do begin**

**2:** if $termset(dp) \subseteq nd$ then   // complete conflict offender

3:   $\Omega = \Omega - \{dp\}$

4: else    // partial conflict offender

5:   $offering = (1 - \frac{1}{\mu}) \times \sum_{t \in termset\ (dp)} \{t.weight | t \in nd\}$

6:   base $= \sum_{t \in termset\ (dp)} \{t.weight | t \notin nd\}$

7:   foreach term $t$ in $termset(dp)$ do begin

8:     if $t \in nd$ then //shrink offender weight

9      $t.weight = (\frac{1}{\mu}) \times t.weight$

10:      else   // shuffle weight

11:        $t.weight = t.weight \times (1 + offering \div base)$

12:      end if

13:     end for

14:    end if

15:  end for

The main purpose of the Shuffling Algorithms is to the assign the weight Distribution to the Terms in the part of the deployed pattern and it will perform the different operation in these Algorithms for the each type of offenders [5].

## CONCLUSION

In this paper, the main issue regarding the pattern based user item information to rank the item features approach is low frequency and misinterpretation. In order to enable an effective clustering process, the word frequencies need to be normalized in terms of their relative frequency of presence in the document and over the entire collection This presents research pattern taxonomy model which includes pattern evolving and deploying method helps in the updating of useful pattern efficiently and the two issues can be solved. It helps in finding the useful information to the user. The inner pattern evolution outperforms the pattern deploying method.

### REFERENCES:

[1] Ning Zhong, Yuefeng Li "Effective pattern discovery in text mining"*IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*, VOL., NO.1, January 2012.

[2] S-T. Wu, Y. Li, and Y. Xu. "An effective deploying algorithm for using pattern-taxonomy". *In Proceedings of the 7th international Conference information Integration and Web-based Applications & Services*(1iWASO5), pages 1013-1022, 2005

[3] Rajuta Taware , Prof. Sanchika A. Bajpai International Journal of Advanced Research in Computer Engineering & Technology (IJARCET) Volume 3, Issue 2, February 2014

[4] Miss Dipti S.Charjan, Prof. Mukesh A.Pund international Journal of Engineering Trends and Technology (IJETT)-Volume 4 Issue 10-Oct 2013. ISSN: 2231-2803

[5] M. Zaki, "SPADE: An efficient algorithm for mining frequent sequences," Machine Learning vol. 40, 2001, pp. 31-60