

# Effective Pattern Discovery for Text Mining By Pattern Deploying and Pattern Evolving

Shaikh Rijwan  
Department of computer  
JSPM Narhe  
Pune, India

Ankush Vyavhare  
Department of computer  
JSPM Narhe  
Pune, India

Mane Nandkumar  
Department of computer  
JSPM Narhe  
Pune, India

P.S.Patil  
Department of computer  
JSPM Narhe  
Pune, India

**Abstract**—Text Mining is one of technique of data mining to discover or extract useful information from textual data. The textual data is in unstructured format by using text mining we can organize this unstructured data into structure data. Text mining is used to extract the effective patterns from large amount of text data. The paper focuses on developing effective algorithm for discovering patterns from large documents. Since most existing text mining method based on term-based method, but they faces some problems. For overcome this problem pattern-based approach is used paper focus on developing effective pattern discovery technique which support process of pattern deploying and pattern evolving .It used to improve effectiveness of discovering interesting patterns.

**Keywords**— Text mining; text classification; pattern mining; pattern evolving; information filtering; pattern taxonomy; sequential patterns; closed patterns

## I. INTRODUCTION

Information from large databases, information which is gather from various sources in business from that extract information or discovery knowledge in databases to take decisions for business growth and management. There are number of data mining techniques have been proposed to perform different knowledge tasks. Techniques like association rule mining, frequent item set mining, sequential pattern mining, maximum pattern mining, and closed pattern mining. Pattern mining means finding of extracting a specific set of textual data from large amount of data or document. Sequential pattern mining is a topic of data mining in which is used to find statically relevant patterns between data example and values delivered in sequence. A frequent closed sequential pattern is a frequent sequential pattern such that it is not include in another sequential pattern for this BIDE+ algorithm is used. Text mining is the discovery of interesting knowledge in text documents. It is a challenging to find accurate knowledge in text documents from large text document which help user to find what they actually want. Initially, Information Retrieval (IR) provided many term-based methods to solve this this issue, such as Rocchio and probabilistic models [1], rough set models [2], BM25 and support vector machine (SVM) [3] based filtering models. The advantages of term-based methods is efficient computational performance, which have maximum time used over the last couple of decades from the IR. Term-

based methods faces the problems of polysemy and synonymy, where polysemy means a word has multiple meanings, and synonymy is multiple words having the same meaning.

In recent year, people use that (pattern) phrase-based approaches could perform better than the term-based ones, as phrases may carry more “semantics” like information. This hypothesis has not fared too well in the history of IR [4], [5], [6]. Although pattern are less ambiguous and more discriminative than individual terms, the likely reasons for the discouraging performance include: patterns have inferior statistical properties to terms, low frequency of occurrence, an there are large numbers of redundant and noisy phrases among them [6].In the presence of these setbacks, sequential patterns used in data mining community have turned out to be a promising alternative to pattern [7], [8] because sequential patterns enjoy good statistical properties like terms. To overcome the disadvantages of pattern-based approaches, pattern mining-based approaches pattern taxonomy models. (PTM) [8], [9] have been proposed, which adopted the concept of closed sequential patterns, and pruned non-closed patterns.

## II. LITERATURE SURVEY

In that paper they present a survey on the different pattern mining techniques to extract the patterns as per the user's requirement and need [10]. This paper concern with the concept of text categorization techniques to update and use the discovered patterns form the text mining. The pattern taxonomy model, pattern deploying and Inner Pattern Evolving to identify the patterns from the document .The most innovative and effective pattern discovery technique is the process of pattern taxonomy, pattern deploying and pattern evolving, to improve the effectiveness of using and updating discovered patterns for finding relevant and interesting information. Provides an effective pattern discovery to increase the effectiveness and also to enhance the discovery of patterns for identifying the relevant data. It is used for overcome the problems of misinterpretation and low frequency.

## III. PROBLEM STATEMENT

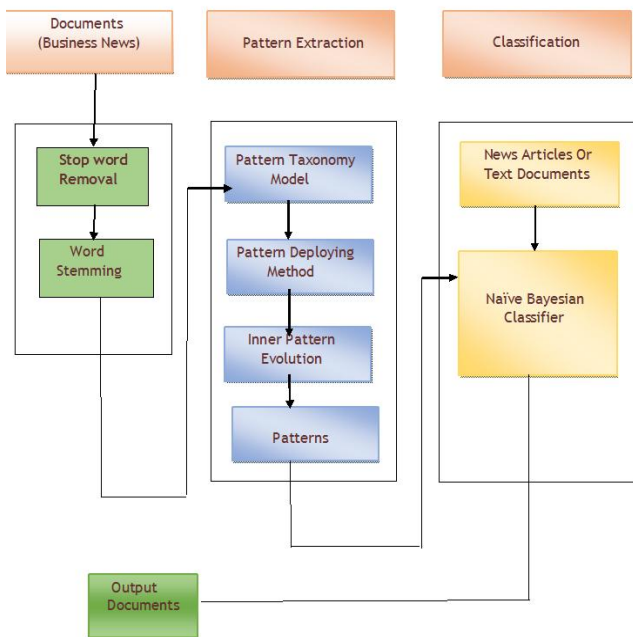
To discover patterns and then compute specificities of patterns for evaluating term weights as per their distribution in the discovered patterns form large amount of text document.

To improve the effectiveness by effectively using closed patterns in text mining. An effective pattern discovery technique has been proposed to overcome the low-frequency and misinterpretation problems for text mining. The proposed technique uses two processes, pattern deploying and pattern evolving, to refine the discovered patterns in text documents. The advantage of the concept-based model is that it can effectively discriminate between non important terms and meaningful terms which describe a sentence meaning. In this system we are giving high priority for long sequence in the evaluated pattern. In this system we are giving term weight based on occurrence of term in long pattern (sequence). The proposed model outperforms not only other pure data mining-based methods and the concept based model, but also term-based state-of-the-art models. The Proposed system will be the output of the project, It will include the discovering the patterns from the large amount of data and search for interesting patterns that user want.

IV. APPLICATION

- For Enhancing Web Search.
- For Dimensionality Reduction.
- For Mining Bibliographic Data.
- For Sentiment Classification.
- In business analysis.
- In development industry

V. SYSTEM ARCHITECTURE



VI. PROPOSED SYSTEM

In proposed scheme, the pattern-based approach is used in which, the discovered patterns are more specific than the whole documents. The patterns are a set of terms extracted from the

documents. In this work, we have focused on the discovering the patterns from the large amount of data and search for interesting patterns that user want. In proposed technique, the pattern taxonomy model is used to extract the pattern from the documents. We also propose a classification algorithm, called naive Bayesian classifier in which the extracted patterns are considered to be training set on text documents.

In architectural design it divide mainly into four module as follows.

- Text Pre-processing
- Pattern Extraction
- Classification(classification of text document)
- Output Text or Pattern.

The main objective of pre-processing is to obtain the key features or key terms from text documents and to enhance the relevancy between the word and document. It includes two techniques.

1) *Stop-Word Removal*: The most common words in any text document that does not provide any meaning of the documents. And hence these words are eliminated. Example of such words are 'the', 'in', 'a', 'an' etc.

2) *Word-Stemming*: Stemming is a technique that reduces the words into their root or stem. The hypothesis of stemming is that the word with the same stem or root describes the same or relatively close concepts in the text document. Example, agree, agreed, agreeing, agreement belong to root word 'agree'. In pattern extraction PTM is used for extraction, for effective pattern discovery technique used the process of pattern deploying and pattern evolving to improve effectiveness. After that naïve bayesian classifier is get used and from that final output document is generated.

VII. ALGORITHMS

A. *Pattern Taxonomy Model*

In this paper, we are considering that documents are splitting in paragraphs. So we consider  $d$  as document and it includes paragraphs. Consider  $D$  is training set of documents, which includes set of positive documents  $D^+$  and set of negative documents  $D^-$ . Consider  $T = \{t_1; t_2; \dots; t_m\}$  be a set of terms (or keywords) which can be extracted from the set of positive documents,  $D$ .

B. *Pattern Taxonomy*

Patterns can be structured into a taxonomy by using the is-a(or subset) relation .For example of table 1, where we have describe asset of paragraphs of a document  $d$  and found 10 frequent patterns in table 2.

TABLE I. A SET OF PARAGRAPHS

Paragraphs	Terms
dp1	t1t2
dp2	t3t4t6
dp3	t3t4t5t6
dp4	t3t4t5t6
dp5	t1t2t6t7
dp6	t1t2t6t7

TABLE II. FREQUENT PATTERNS AND COVERING SETS

Frequent Pattern	Covering Set
{ t3,t4,t6 }	{ dp2,dp3,dp4 }
{t3,t4}	{ dp2,dp3,dp4 }
{t3,t6}	{ dp2,dp3,dp4 }
{t4,t6}	{ dp2,dp3,dp4 }
{t3}	{ dp2,dp3,dp4 }
{t4}	{ dp2,dp3,dp4 }
{t1,t2}	{ dp1,dp5,dp6 }
{t1}	{ dp1,dp5,dp6 }
{t2}	{ dp1,dp5,dp6 }
{t6}	{ dp2,dp3,dp4,dp5,dp6 }

C. Pattern Deploying Method

In this we use the semantic information in the pattern taxonomy for better performance of closed patterns in text mining we need to interpret discover patterns by summarizing them as d- patterns for accurately evaluate turn weights.

Algorithm 1: Pattern Taxonomy Model Algorithm (D+,min\_sup)

Input- Take positive documents D+, minimum support, min\_sup.

Output- d-patterns DP, and supports of terms.

```

1 DP = ∅ ;
2 For each document d ∈ D+ do
3 let PS(d) be the set of paragraphs in d;
4 SP=SPMining(PS(d), min_sup);
5 d^ = ∅
6 foreach pattern pi ∈ SP do
7 p= {(t,1)|t ∈ pi };
8 d^ = d^ ⊕ p;
9 end
10 DP=DP ∪ {d^};
11 end
12 T={t|(t,f) ∈ p,p ∈ DP};
13 foreach term t ∈ T do
14 support(t) =0;
```

```

15 end
16 foreach d-patterns p ∈ DP do
17 foreach (t,w) ∈ β (p) do
18 support(t) = support(t) + w;
19 end
20 end
```

D. D-Pattern Mining Algorithm

To improve the efficiency of the pattern taxonomy mining, an algorithm, SPMining, was proposed in [50] to find all closed sequential patterns, which used the well-known Apriori property in order to reduce the searching space. . The main focus of this paper is the deploying process, which consists of the d-pattern discovery and term support evaluation.

E. Inner Pattern Evolution

In this section, they discuss how to reshuffle supports of terms within normal forms of d-patterns based on negative documents in the training set. The technique will be useful to reduce the side effects of noisy patterns because of the low-frequency problem. This technique is called inner pattern evolution here, because it only changes a pattern's term supports within the pattern. A threshold is usually used to classify documents into relevant or irrelevant categories. Using the d-patterns, the threshold can be defined naturally as follows.

$$\text{Threshold (DP)} = \min_{p \in DP} (\sum_{(t,w) \in \beta(p)} \text{sup} t)$$

Algorithm 2: IPEvolving

Input: A training set D =D+ D-; a set of d-pattern DP; and experimental coefficient

Output: A set of term-support pairs np.

```

1 np
2 threshold=Threshold (DP); //
3 foreach noise negative document nd
4 if weight (nd) ≥threshold then (nd) = {p};
5 NDP = { };
6 Shuffling (nd,);
7 foreach p
8 np
9 end
10 end
```

F. Future Scope

In last ten years many data mining techniques have been proposed for mining useful patterns of users wish. It includes association rule mining, sequential pattern mining, maximum and closed pattern mining. Discovered patterns used in text mining field is difficult and ineffective because useful long patterns with specificity and lack in support. In this we have focused on finding the patterns from large amount of data as per user's requirements. In n proposed system we are using

pattern taxonomy model for extracting the pattern from given documents. In future we are going to use pattern deploying and inner pattern evolution to minimize the noisy patterns and to calculate the terms weights. We propose to use the Naive Bayesian classification algorithm to classify the documents.

#### CONCLUSION

In this effective pattern discovery technique has been proposed to minimize the misinterpretation and low-frequency problems. Pattern deploying and Pattern evolving processes are used in proposed system to refine the discovered patterns in text documents. The results show the proposed model outperforms not only other pure data mining-based methods and the concept-based model, but also term-based state-of-the-art models.

#### REFERENCES

- [1] R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*. Addison Wesley, 1999.
- [2] Y. Li, C. Zhang, and J.R. Swan, "An Information Filtering Model on the Web and Its Application in Jobagent," *Knowledge-Based Systems*, vol. 13, no. 5, pp. 285-296, 2000.
- [3] S. Robertson and I. Soboroff, "The Trec 2002 Filtering Track Report," *TREC,2002*, trec.nist.gov/pubs/trec11/papers/OVER.FILTERING.ps.gz.
- [4] D.D. Lewis, "Feature Selection and Feature Extraction for Text Categorization," *Proc. Workshop Speech and Natural Language*, 212-217, 1992.
- [5] S. Scott and S. Matwin, "Feature Engineering for Text Classification," *Proc. 16th Int'l Conf. Machine Learning (ICML '99)*, pp. 379-388, 1999.
- [6] F. Sebastiani, "Machine Learning in Automated Text Categorization," *ACM Computing Surveys*, vol. 34, no. 1, pp. 1-47, 2002.
- [7] N. Jindal and B. Liu, "Identifying Comparative Sentences in Text Documents," *Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '06)*, pp. 244-251, 2006.
- [8] S.-T. Wu, Y. Li, and Y. Xu, "Deploying Approaches for Pattern Refinement in Text Mining," *Proc. IEEE Sixth Int'l Conf. Data Mining (ICDM '06)*, pp. 1157-1161, 2006.
- [9] S.-T. Wu, Y. Li, Y. Xu, B. Pham, and P. Chen, "Automatic Pattern-Taxonomy Extraction for Web Mining," *Proc. IEEE/WIC/ACM Int'l Conf. Web Intelligence (WI '04)*, pp. 242-248, 2004.
- [10] Bhushan.V and Ujwala patil, "A Comparative Study on Different Types of Effective in Text Mining: A Survey (IJCET), March 2013.

