

## Impressive Order Invention in Pattern Evolution for Text Mining

M.Revathi

*Department of Computer Science and Engineering  
Madanapalle, Andhra Pradesh, India*

P.Raja Rajeswari

*Department of Computer Science and Engineering  
Assistant Professor, Department of Computer Science and Engineering, Madanapalle,  
Andhra Pradesh, India*

Dr.D.Vasumathi

*Associate Professor, Department of Computer Science and Engineering  
Hyderabad, Andhra Pradesh, India*

**Abstract - We Provide an effective pattern discovery technique, which first calculates discovered specificities of patterns and then evaluates term weights according to the distribution of terms in the discovered patterns rather than the distribution in documents for solving the misinterpretation problem. It also considers the influence of patterns from the negative training examples to find ambiguous (noisy) patterns and try to reduce their influence for the low-frequency problem. The process of updating ambiguous patterns can be referred as pattern evolution. The proposed approach can improve the accuracy of evaluating term weights because discovered patterns are more specific than whole documents. We also conduct numerous experiments on the latest data collection, Reuters Corpus Volume 1 (RCV1) and Text Retrieval Conference (TREC) filtering topics, to evaluate the proposed technique. The results show that the proposed technique outperforms up-to-date data mining-based methods, concept-based models and the state-of-the-art term based methods. We propose Techniques include association rule mining, frequent item set mining, Maximum pattern mining, and closed pattern mining. There is a plethora of text mining and visualization tools available on the market to facilitate the innovative process in uncovering “hidden nuggets” of information about emerging technologies. By using discovered knowledge in the field of text mining is difficult ineffective. Long pattern with high specificity lack in support i.e. low frequency problem we know that all frequent short patterns are useful. In This paper we proposed an effective pattern discovery technique has been introduced to overcome the low-frequency and misinterpretation problems for text mining. We propose two techniques processes, pattern deploying and pattern evolving, the proposed model outperforms not only other pure data mining-based methods and the concept-based model also term based state of art models such as BM25 and SVM-based models.**

**Keywords: Text mining, Pattern taxonomy, pattern mining, information filtering**

### I. INTRODUCTION

Text analysis referred as text mining. The way to make qualitative or unstructured data usable by a computer. Qualitative data is descriptive data that cannot be measured in numbers and often includes qualities of appearance like color, texture, and textual description. Quantitative data is numerical, structured data that can be measured. However, there is often slippage between qualitative and quantitative categories. For example, a photograph might traditionally be considered “qualitative data” but when you break it down to the level of pixels, which can be measured. Various tools are based on standard analysis techniques and mainly differentiate in their capabilities to use different data sources and visualize in different ways tables, maps, graphs, and matrices. To define experiences at Bristol-Myers Squibb (BMS) is learning about and evaluating selected tools. It provides the impressions of the capabilities of tools. It is not our intention to cover every available tool on the market nor is it to cover every aspect of the tools and to provide new insights into the capabilities, data sources, results, our perceived strengths and potential limitations for each of the tools. Given a set of transactions  $D$ , the problem of mining association rules is to generate all association rules that support and condense greater than the user-specified minimum support and minimum condense. For example,  $D$  could be a data file, a relational table, or the result of a relational expression. An algorithm for finding all association rules, referred to as the AIS algorithm, another algorithm for this task, called the SETM algorithm, we present two new algorithms, Apriori and Apriori Tid, that differ fundamentally from these algorithms.

## II. RELATED WORK

The Analysis Group at BMS established a project to evaluate various text mining and data visualization tools. Gathering background information and brainstorming, identifying the potential tools to evaluate, and conducting on-site demonstrations. The project includes piloting a few of the select tools and identifying potential clients groups and real-life case studies for using the tools. Capabilities, the data sources permitted, the results generated, our perceived strengths and potential limitations. The “type” includes whether the tool is software designed for text mining/visualization or a patent database content provider, or both. The capabilities evaluated include tools performing keyword, statistical and/or linguistic analysis. Keyword analysis refers to extracting nouns or noun phrases in text without understanding their meaning or relationships (i.e. out of context). Statistical analysis refers to word frequency-based analysis or counting the number of times a word appears in the text. Linguistic analysis refers to using a trained agent to do natural language processing or semantic analysis Data sources.

The specific text mining data sources include:

- (1) unstructured text (such as full-text documents and emails).
- (2) structured text (such as database records from STN\_ or Pub Med), and
- (3) Hybrid content (such as patents where the front page information is structured but the remaining text is not).

The output from each tool is typically generated in lists of documents, tables, charts, graphs, or maps. Perceived strengths Based on vendor demonstrations, we highlight key features that we perceived as strong points for each of the tools. Potential limitations Based on vendor demonstrations, we list certain features that we perceived as lacking or inadequate. These potential limitations are our own personal views.

*A set of paragraphs*

Paragraph	Terms
dp1	t1,t2
dp2	t3,t4,t6
dp3	t3,t4,t5,t6
dp4	t3,t4,t5,t6
dp5	t1,t2,t6,t7
dp6	t1,t2,t6,t7

A two staged model that used both term based methods and pattern based methods was introduced to improve the performance of information filtering. . A text episode is defined as a pair  $\alpha = (V, <)$  where  $V$  is a collection of feature vectors, and  $<$  is a partial order on  $V$ . Given a text sequence  $S$ , a text episode  $\alpha = (V, <)$  occurs within  $S$  if there is a way of satisfying the feature vectors in  $S$  using the feature vectors in  $V$  so that the partial order  $<$  is respected. The feature vectors of  $V$  can be found within  $S$  in an order that satisfies the partial order  $<$ .

## III. LITERATURE SURVEY

Literature survey is the most important step in software development process. Before developing the tool it is necessary to determine the time factor, economy n company strength. Once these things are satisfied, then next steps are to determine which operating system and language can be used for developing the tool. Once the programmers start building the tool the programmers need lot of external support. This support can be obtained from senior programmers, from book or from websites. Before building the system the above consideration are taken into account for developing the proposed system.

### 3.1. Pattern Taxonomy

Discovering all association rules can be decomposed into two sub problems

1. Find all sets of items that have transaction support above minimum support. The support for an item set is the number of transactions that contain the item set. Item sets with minimum support are called large item sets, and all others small Item sets. The Algorithms, Apriori and Apriori Tid, for solving problem.
2. Use the large item sets to generate the desired rules. Here is a straightforward algorithm for this task. For every large itemset  $l$  all non-empty subsets of  $l$ . For every such subset  $a$ , output a rule of the form  $a \Rightarrow (l \square a)$  if the ratio of

support(l) to support(a) is at least min conf. We need to consider all subsets of l to generate rules with multiple consequents. Due to lack of space, we show the relative performance of the proposed Apriori and AprioriTid algorithms against the AIS and SETM algorithms. To make the paper self-contained, we include an overview of the AIS and SETM algorithms in this section. We describe how the Apriori and AprioriTid algorithms can be combined into a hybrid algorithm, Apriori Hybrid, and demonstrate the scale up properties of this algorithm. We conclude by pointing out some related open problems.

### 3.2 Pattern Mining

We assume that all documents are split into paragraphs. So a given document d yields a set of paragraphs. Let D be a training set of documents, which consists of a set of positive documents and a set of negative documents. To derive this method patterns in text documents for information filtering systems simplifies the process. It is defined as Let  $p_1 + p_2 = \{(t, x_1 + x_2) / (t, x_1) \in p_1, (t, x_2) \in p_2\} \cup \{(t, x) / t, x \in p_1 \cup p_2, \text{not}((t, -) \in p_1 \cap p_2)\}$

Where  $\_$  is the wildcard that matches any number.

In order to use the semantic information in the pattern taxonomy to improve the performance of closed patterns in text mining, we need to interpret discovered patterns by summarizing them as d-patterns (see the definition below) in order to accurately evaluate term weights (supports). The rationale behind this motivation is that d-patterns include more semantic meaning than terms that are selected based on a term-based technique (e.g.,  $tf \cdot idf$ ).

### 3.3 INNER PATTERN EVOLUTION :

we 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. The proposed model includes two phases: the training phase and the testing phase. In the training phase, the proposed model first calls Algorithm PTM ( $D_p, \text{min sup}$ ) to find d-patterns in positive documents ( $D_p$ ) based on a min sup, and evaluates term supports by deploying d-patterns to terms. It also calls Algorithm IPEvolving ( $D_p, D\_ , DP, \_$ ) to revise term supports using noise negative documents in D, based on an experimental coefficient  $\_$ . In the testing phase, it evaluates weights for all incoming documents. The incoming documents can be sorted based on these weights.

## IV. EVALUATION AND DISCUSSION

Reuters text collection is used to evaluate the proposed approach. Term stemming and stopword removal techniques are used in the prior stage of text preprocessing. Several common measures are then applied for performance evaluation and our results are compared with the state-of-art approaches in data mining, concept-based, and term-based methods. Baseline Models we choose three classes of models as the baseline models. The first class includes several data mining-based methods that we have introduced. In the following, we introduce other two classes: the concept-based model and term-based methods.

1. Concept-Based Model.
2. Term-Based Methods.

### 4.1 Concept-Based Model

This model represents both sentence and document levels. This model uses a verb-argument structure which splits a sentence into verb.

Ex: "John hits the ball", where hits is the verb, and "John" or "the ball" are the arguments of "hits". Arguments are assigned as subjects or objects.

For a document d,  $tf(c)$  is the number of occurrences of concept c in d; and  $ctf(c)$  is called the conceptual term frequency of concept c in a sentence s, which is the number of occurrences of concept c in the verb-argument structure of sentences.

### 4.2 Term-Based Methods.

There are many classic term-based approaches. The Rocchio algorithm which defines widely adopted information retrieval built text representation of a training set using a centroid  $\hat{c}$  as follows

$$c = \alpha \sum_{d \in D^+} |d| - \beta \sum_{d \in D^-} |d|$$

Where  $\alpha$  and  $\beta$  are empirical parameters:  $D^+$  and  $D^-$  are the sets of positive and negative documents, respectively:  $d$  denotes a document.

The input design is the link between the information system and the user. It comprises the developing specification and procedures for data preparation and those steps are necessary to put transaction data in to a usable form for processing can be achieved by inspecting the computer to read data from a written or printed document or it can occur by having people keying the data directly into the system. The design of input focuses on controlling the amount of input required, controlling the errors, avoiding delay, avoiding extra steps and keeping the process simple. The input is designed in such a way so that it provides security and ease of use with retaining the privacy. Input Design considered the following things:

- What data should be given as input?
- How the data should be arranged or coded?
- The dialog to guide the operating personnel in providing input.
- Methods for preparing input validations and steps to follow when error occur.

## V. OBJECTIVES OF RESULTS

1. Input Design is the process of converting a user-oriented description of the input into a computer-based system. This design is important to avoid errors in the data input process and show the correct direction to the management for getting correct information from the computerized system.

2. It is achieved by creating user-friendly screens for the data entry to handle large volume of data. The goal of designing input is to make data entry easier and to be free from errors. The data entry screen is designed in such a way that all the data manipulates can be performed. It also provides record viewing facilities.

3. When the data is entered it will check for its validity. Data can be entered with the help of screens. Appropriate messages are provided as when needed so that the user will not be in maize of instant. Thus the objective of input design is to create an input layout that is easy to follow.

A quality output is one, which meets the requirements of the end user and presents the information clearly. In any system results of processing are communicated to the users and to other system through outputs. In output design it is determined how the information is to be displaced for immediate need and also the hard copy output. It is the most important and direct source information to the user. Efficient and intelligent output design improves the system's relationship to help user decision-making.

1. Designing computer output should proceed in an organized, well thought out manner; the right output must be developed while ensuring that each output element is designed so that people will find the system can use easily and effectively. When analysis design computer output, they should Identify the specific output that is needed to meet the requirements.

2. Select methods for presenting information.

3. Create document, report, or other formats that contain information produced by the system.

The output form of an information system should accomplish one or more of the following objectives.

- ❖ Convey information about past activities, current status or projections of the
- ❖ Future.
- ❖ Signal important events, opportunities, problems, or warnings.
- ❖ Trigger an action.
- ❖ Confirm an action.

## VI. IMPLEMENTATION

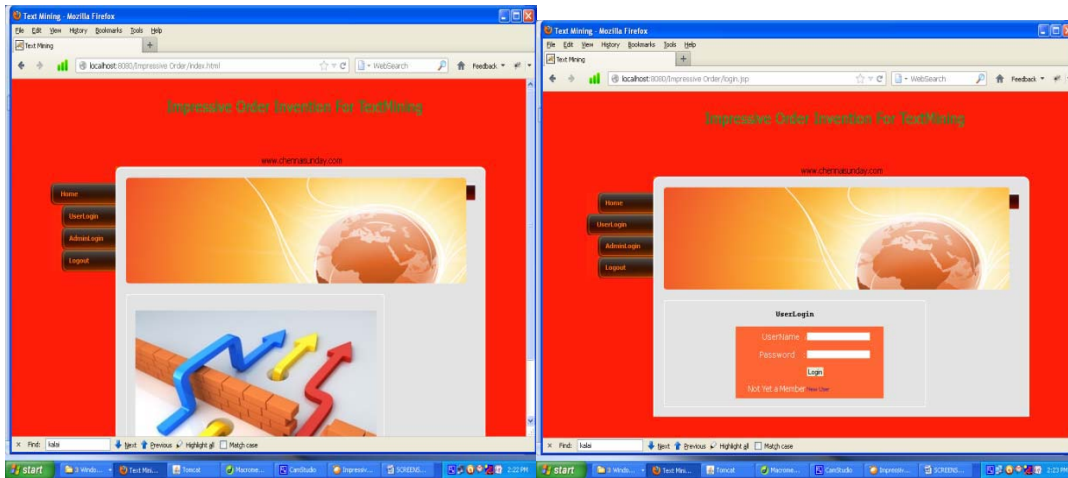
Implementation is the stage of the project when the theoretical design is turned out into a working system. Thus it can be considered to be the most critical stage in achieving a successful new system and in giving the user, confidence that the new system will work and be effective. The implementation stage involves careful planning, investigation of the existing system and it's constraints on implementation, designing of methods to achieve changeover and evaluation of changeover methods.

VII. RESULTS

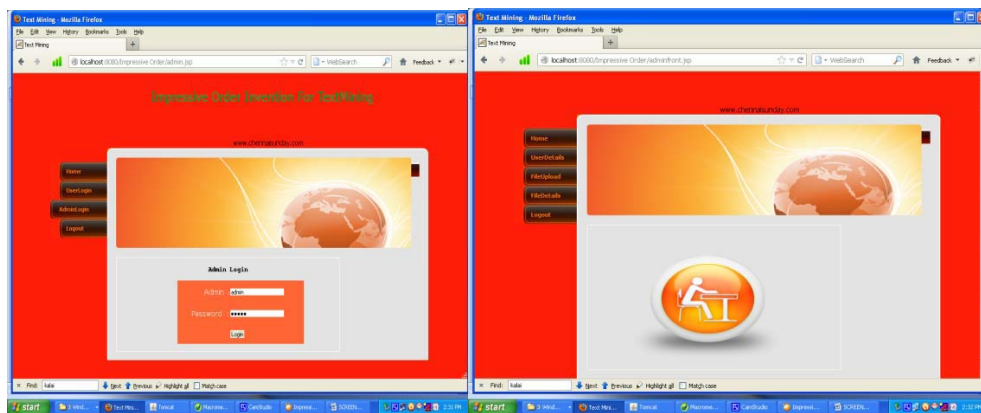
The results of overall comparisons are presented where PTM performs not only the pattern- mining methods but also term-based methods and also performs CBM pattern matching and CBM in five measures.

Method	Top-20	b/p	MAP	$F_{\beta=1}$	<i>IAP</i>
SVM	0.477	0.409	0.408	0.421	0.434
BM25	0.434	0.399	0.401	0.410	0.422
TFIDF	0.321	0.321	0.322	0.355	0.348
CBM	0.448	0.409	0.415	0.423	0.440
CBM pattern matching	0.329	0.282	0.283	0.320	0.311

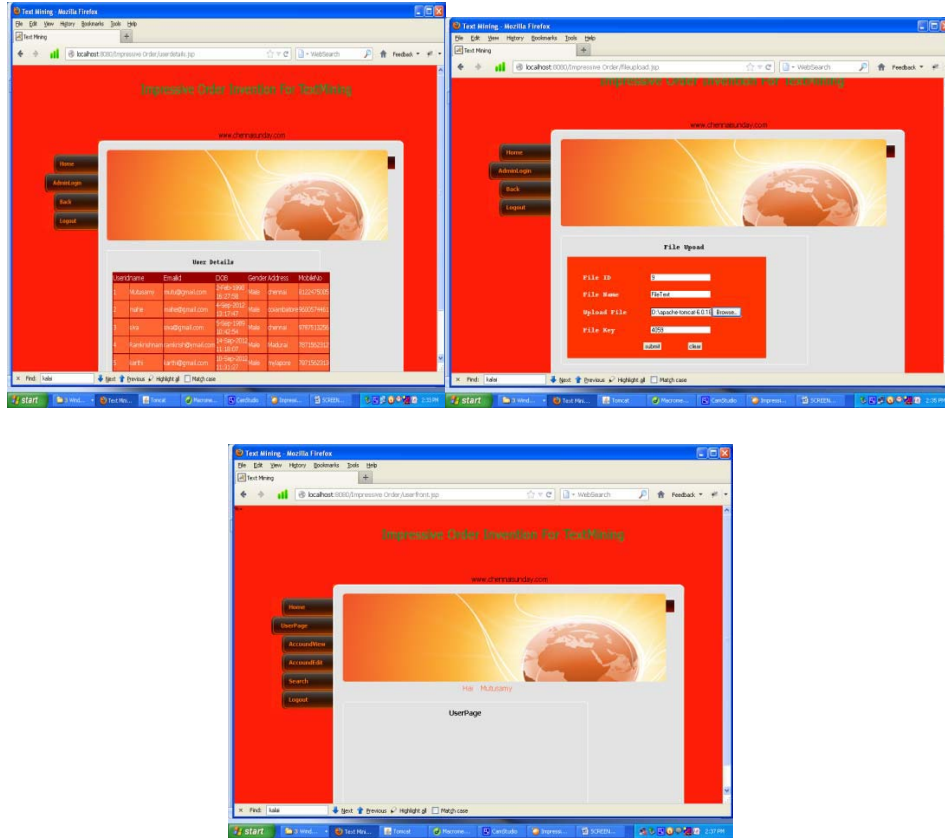
7.1 Experimental Results



The snap shot defines the home page and enters the login page.



In login page we enter the key word to search among many documents where we can evolve effective pattern by mining text data.



Designing of methods over changing where the admin login page was displayed and the entering the user details searching indent goes through that and finally search results of simulation both pre and post processing roles and weighting results are displayed.

### VIII. CONCLUSION

Various algorithms have been proposed for text documents to mining frequent patterns. But how to efficiently find these patterns is still an issue in text mining domain. Traditionally, texts have been analyzed by using various information retrieval related methods, such as full-text analysis, and natural language processing. However, only few examples of data mining in text, particularly in full text, are available. In this paper we present a framework for text mining using descriptive phrase extraction. The framework follows the general knowledge discovery process, thus containing steps from preprocessing to the utilization of the results. We apply generalized episodes and episode rules data mining method. We introduce a weighting scheme that helps in pruning out redundant or non-descriptive phrases. Simulation results relived that episodes and episode rules produced discriminate between documents. Both pre and post processing have essential roles in pruning and weighting the results.

### REFERENCES

- [1] Phrase-based Document Categorization Cornelis H.A. Koster<sup>1</sup>, Jean G. Beney<sup>2</sup>, Suzan Verberne<sup>1</sup>, and Merijn Vogel<sup>1</sup> <sup>1</sup> Computing Science Institute ICIS, Univ. of Nijmegen.
- [2] Towards Modernised and Web-Specific Stoplists for Web Document Analysis Mark P Sinka *University of Reading, Reading, UK*.m.p.sinka@reading.ac.uk.
- [3] Text mining and visualization tools – Impressions of emerging capabilities YunYun Yang \*, Lucy Akers, Thomas Klose, Cynthia Barcelon Yang.
- [4] Text categorization: : AA A Survey KJERSTI Aas and Line Eikvil.
- [5] Modern information Retrieval Ricardo Baeza-yates Addison Wesley ACM New York.
- [6] Descriptive Phrase Extraction in Text mining ,Alekhya.v, B.Govinda laxmi, Ap .
- [7] Nicola canceledda, Nicolocesa Bianchi, Alexcivinokourov, Kernel methods for document filtering.
- [8] Yun Yun Yand,lucky Akers,Thomas Klose,Cynthia Barcelon Yang, Text mining and Visualization tools.
- [9] Rakesh Agarwal Ramakrishnan srikanth, Fast algorithm for mining association
- [10] Effective pattern Discovery for Text Mining,Ning Zhong,Yuefeng Li,and sheng-Tang wu