# Professional replica of pattern recognition for textual data extraction

Candidate name: VechalapuAlekya Email: *alekya257@gmail.com*
Internal guide name: P.Pavithra Email: panthangi.*pavithra@gmail.com*
Vishakha institute of engineering and technology

## Abstract:

The discoveries on mining and to achieve the perfect and effectual pattern are still open challenge. But so many existing discoveries do not follow the offset of the prefix duplications before/after clustering. Here in this, we proposed how to get large and best patterns from huge data sets. For that we take a new approach after cleaning the documents called PNR and DYNPRO (prefix with indexing). In DYNPRO the documents are need to be compare Asynchronously with each other for extract the best patterns with less complexity. Then we have to generate invert matrix for retrieved document. So these two algorithms results consequently will provide the best feasible data patterns/ prefix wise. The Experiment results show best and significant improvement in searching performance over different documents over the datasets. Lastly, we evaluate the impact of different interference and models on the various categories of data mining.

**Key Terms***: Text classification, Invert Matrix, Pattern extraction, Document categorization.*

## Introduction:

Mining is the process of discovering and retrieve optimal and meaningful knowledge in a data set. It has been successfullyapplied to many real-life problems, for instance,web personalization, network intrusion detection, and customizedmarketing. Recent advances in computational scienceshave led to the application of data mining to variousdomains. As an area combiningideas from database systems, machine learning, andstatistical learning, data mining has been successfully appliedto many application domains.

In previous days, a number of data mining techniques have been proposed in order to performdifferent knowledge tasks. These techniques include associationrule mining, frequent item set mining, sequentialpattern mining, maximum pattern mining, and closedpattern mining. Most of them are proposed for the purposeof developing efficient mining algorithms to find different Patterns within a reasonable and acceptable time frame.With a large number of patterns generated by using datamining approaches, how to effectively use and update thesepatterns is still an open research issue.

Text mining is the discovery of interesting knowledge intext documents. It is a challenging issue to find accurate knowledge in text documents to help users tofind what they want. For that, in this paper, we propose Efficacious pattern mining which was taken as challenge to the previous techniques to retrieve best and optimal patterns as results. The information on Documents needs to be categorized as patterns in datasets. The user required documents are extracted in the best way. For every user given query the documents are need to be filtered and associated to give best patterned results according to the requirements. These results show knowledge and understandable patterns.

## Related work:

There are positive and negative document forming and this is very tedious, because once the all the documents scanned and framed these types, the combinations/clustering will be done on two deferent types like positive and negative. So the categories will be repeated in the initial process itself. This issue grows and shows much impact for further pattern discovery and double repetitive things will occur.

The documents are getting compared with the regular way (inner/outer loop) which is of normal looping, which leads to self/repetitive comparisons. Which is time consuming? If the documents size is 6 i.e. the comparisons would be 6^2, which is 36.

Patterns will be discovered using normal approach with clustering and non indexing methodology, which is very time consuming process for search engines. For the current sequences, extensions would be more tedious and extra complications for further discoveries.

## Effectual pattern frame:

our proposed algorithm of DYNPRO. In this approach we will not compare then after cleaning the documents each document is compared with other document without any repetition to get the possible combinations in NR possible non repetitive things).

The Documents which are already existed in datasets are needed to be filtered according to Document which is already compared with source and destination.

For instance, as shown in figure.1, the Documents are compared with each other. Here the documents are categorized as positive documents and Negative Documents. The motivation will be prepared as until they have to be clean. If *D1* is compared with *D2*, again *D2* will not be compared with *D1* for next iterations. And also whenever the documents are getting compared source and destinations will be marked for the small and big documents. This approach simplifies for next process of combining the patterns for the fine prefix patterns. If *D1* is compared with *D6* and if founds that *D6* is bigger than *D1* the combination in the matrix will be marked as **Pn(← +)**. This is the part of PNR approach.

DYNPRO is proposed algorithm for text mining to be in order of prefixes. The threshold will be maintained though out the document before clustering to achieve the best possible sequence.

**PNR for 6 available documents/categorization:**

| | D1 | D2 | D3 | D4 | D5 | D6 |
|---|---|---|---|---|---|---|
| D1 | $ | P1(→ -) | P2(→ +) | P3(→ +) | P4(→ -) | P5(→ +) |
| D2 | (P1← +) | $ | (P6→ +) | (P7→ -) | (P8→ +) | (P9→ -) |
| D3 | (P2← -) | (P6← -) | $ | P10(→ -) | (P11→ +) | (P12→ +) |
| D4 | (P3← -) | (P7← +) | P10(← +) | $ | P13(→ -) | P14(→ -) |
| D5 | (P4← +) | (P8← -) | P11(← -) | P13(← +) | $ | P15(→ -) |
| D6 | (P5← -) | (P9← +) | P12(← -) | P14(← +) | P15(← +) | $ |

**Picture 1**

## DYNPRO Method:

Once the documents scanned by new approach, documents will be categorized and will be

cleaned in an effective way like stop words removal, marking for most frequent categories, and segmentations will be framed. These segments will be framed up with the set of documents/category wise.

**DYNPRO Algorithm Process:**

1. We need to consider number of Datasets with number of different documents. Each document contains data with different patterns and also with different sizes.
2. We have to apply DynPro Algorithm for each data set.
3. The DynPro algorithm is used to compare the Documents to calculate time complexity for getting best document.
4. Compare Documents in every Dataset by taking the documents asynchronously.
5. The same document comparison will become taken as offset.
6. The different document comparison will give best document with time complexity.
7. Apply step4, 5, 6 to all datasets for which were we taken.
8. Then take time complexity from all datasets. And get best time complexity (low) from them.
9. Generate graph for Time complexity of all datasets.
10. Then get large, best and less time complexity Document based on the Graph.
11. Take that large document with less time complexity among datasets.

**Terminology:**

1) $\sum_0^{n-1} D \rightarrow$ Total documents

2) $\sum_0^i P_d \rightarrow$ Total positive documents

3) $\sum_0^j P_n \rightarrow$ Total negative documents

4) \$$\phi \rightarrow$ $\begin{cases} \text{neutralization (nullify)} \\ \text{No operation} \end{cases}$

5) MI ← insert matrix

6) Size (d) ← size of the document

7) Mark(d) ← marking the document

8) $\sum_0^{n-1} M$ ← Marked documents

9) $\sum_0^k C_m$ ←documents to be clustered

10) $F_{cd}$← final clustered document

**Algorithm:**

DYNPRO: Marking

**Input:** Documents of Datasets.

**Output:** patterns set

1. Initialization

2.      Count ←0

3.      $P_d$← 0

4.      $P_n$← 0

5.      for each d (document) in D

6.      Loop start:

7.            If size $(d_i)$ < size $(d_{i+1})$

8.            $P_n$←$d_i$

9.            Count++;

10           M ← Mark $(d_i)$

11.      Else if size $(d_i)$> size $(d_{i+1})$

12.            $P_d$← $d_i$

13.            Count ++

14         $M \leftarrow$ Mark $(d_i)$

15.    End if

16.     Count$\leftarrow$ 0 for each d in M

17.     $MI \leftarrow P_n \Omega P_d$

18.    Count++; // count increment

19. End for

## Document Clustering:

After cleaning the documents by DynPro Algorithm, the compared documents which are optimized as best results to show patterns are clustered. This technique is for close the labeled documents at a place for show results.

We simulate for the availability of 6 documents (assumption), the picture1 shows the combinations to compare for further DYNPRO:Clustering

1. Count$\leftarrow$ 0
2. For each k in M
3. Temp  $\leftarrow$ null
4. Temp $\leftarrow$ cluster (temp,$d_k$)
5. Count ++
6. End loop
7. $F_{cd} \leftarrow$ temp;

## InvertMatrix Generation:

We simulate for the availability of 6 documents (assumption), the picture1 shows the combinations to compare for further combinational documents. The first step of this process (P) of PNR is to get the weight of the source document (**D1**) to destination document (**D3**). If source is having the higher weight than destination, then the result in matrix box is **P2( →+)** . So by getting this type of results, it is very easy for further clustering process/ patterns by using initial words like prefixes.Dynpro discover invert matrix for best available

combinational documents. The first step of this process (P) of PNR is to get the weight of the source document (**D1**) to destination document (**D3**). If source is having the higher weight than destination, then the result in matrix box is **P2( →+)** . So by getting this type of results, it is very easy for further clustering process/ patterns by using initial words like prefixes.

Dynpro discover invert matrix for best available combinations for available prefixes. After invert matrix generation, by using a fixed threshold, the feasible/effective pattern/prefix will be generated. This approach is low cost solution for prefix mining after PNR. In Invert matrix the associations will be find out for optimal text patternsretrievals.

combinations for available prefixes. After invert matrix generation, by using a fixed threshold, the feasible/effective pattern/prefix will be generated. This approach is low cost solution for prefix mining after PNR.

In PNR, Invert matrix generation will be done by utilizing the documents which are cleaned by DYNPRO mechanism.The Invert matrix is used for evaluates various patterns of particular Document.The text data with non duplication patterns of each document are applied with this mechanism. The rows and columns will show with text data which was existed in documents. Column data will be indexed with numbers randomly. According to the text data by indexing we can find out the association patterns for required query which is given by the user for getting best search. The repeated terms in a document are discovering with the frequencythat's what we get in invert matrix calculation. Based on results the invert matrix shows efficacious and optimal patterns.

There will not be any dataset formations. So the best possible matrices will be generated to achieve the best possible patterns among the documents. The DYNPRO approach is the best effective pattern for search engines.

After that the sequences will be clustered to get the final patters. So the result will be maintaining the offsets and also best possible and fine patterns.

**Example:**

One four three five two three five four two one four two one three six four one three two five six three two

| | One | two | three | four | five | six | seven | eight | nine | ten |
|---|---|---|---|---|---|---|---|---|---|---|
| One – 1 | 5,1 | 5,2 | 3,3 | 2,4 | | | | | | |
| Two – 4 | 2,14,2 | $ | $ | | | | | | | |
| Three -2 | 2,1$$ | 6,4 | | | | | | | | |
| Four -5 | 4,14,2 | $3,4 | | | | | | | | |
| Five -6 | $$ | 5,3 | 2,4 | | | | | | | |
| Six -3 | 6,1 | $ | 3,3 | 4,4 | | | | | | |

Figure.2

**Algorithm for Invert Matrix:**

Dp= paragraph in document.

Dp(t)←terms in document

Step1:

Take the input as a document consisting of the stop words(the document which was positive ).

Initialization:

Dp(t)←0

I←0(index)

T←0(terms)

Step2: loop starts //take the next term of the term in document

i←0

Next (ti)←dp(t)

Ai←next (ti)

i++;

End loop

step3: loop starts //take the index values as randomly

j←0 to n

x←random(j);

j++;

End loop

step4: loop starts // for frequency value generation

Count←0

f←0

step5:  if // condition starts

Index (ti) ==index (ti)

Count++

f←count

Return f;

End loop

step6: if //condition starts

(ti,ai)∩dp(t)=Ǿ;

Then

[(ti,x),ti]→$

End if

Step7:loop starts // generate inverse matrix for rows

i←0

i<=size (t),

Increment i upto the size will compete.

Step8:          loop starts // for column terms of matrix

J←(ti,x)←0

j<=size (t)

Increment j upto the size will

compete

Step9: compare ( [(ti,x),ti])then

(ti,ai)→[Index (ai),f]

[(ti,x),ti]←[Index (ai),f];


Compare [(ti,x),ti+1]then

(ti,ai+1)→[Index (ai+1),f]

[(ti,x),ti+1]←[Index (ai+1),f]

......

......

......

Compare [(ti,x),tn]then

(ti,an)→[Index (an), f]

[(ti,x),tn]←[Index (an), f]

Step10: endloop

Endloop

**Conclusion:**

The techniques included in association rule mining,frequent item set mining, sequential pattern mining, maximumpattern mining, and closed pattern mining are closely related with this proposed one. However,using these discovered knowledge (or patterns) in the fieldof text mining is difficult and ineffective. By this proposed algorithm in this paper we could reduce the complexity in extracting the data. In this research work, aneffective pattern discovery technique has been proposed toovercome the low-frequency and misinterpretation problemsfor text mining. The proposed technique uses twoprocesses, pattern deploying and pattern evolving, to refinethe discovered patterns in text documents with low complexity.

**References:**

Armstrong, D., Gosling, A., Weinman, J., &artaeu, T. (1997). The place of inter-rater reliability in qualitative research: An empirical study. *Sociology, 31*, 597-606.

Bales, R. F. (1950). *Interaction process analysis*. Cambridge, MA: Addison-Wesley.

Blalock, H. M. (1964). *Causal inferences in nonexperimental research*. Chapel Hill, NC: University of North Carolina Press.

Camillo, F., Tosi, M., & Traldi, T. (2005). *Semiometric approach, qualitative research and text mining techniques for modeling the material culture of happiness*. Berlin, Germany: Springer.

Carey, J., Morgan, M., & Oxtoby, M. (1996). Inter-coder agreement in analysis of responses *Handbook of processing*. New York, NY: Marcel Dekker.

Feist, J., & Feist, G. (2006). *Theories of personality*. Boston, MA: McGraw-Hill.

Feldman, R., & Sanger, J. (2007). *The text mining handbook: Advanced approaches in analyzing unstructured data*. Cambridge: Cambridge University Press.

Frantzi, K. (2006, April). *Author identification*. Paper presented at National Centre for e-

to open-ended interview questions: Examples from tuberculosis research. *Cultural Anthropology Methods*, *8*(3), 1-5.

Carnap, R. (1937). *The logical syntax of language*. London: Routledge & Kegan Paul Ltd.

Chen, N., Kinshuk, Wei, C. W., & Chen, H. (2008). Mining e-Learning domain concept map from academic articles. *Computers & Education, 50*, 1009-1021.

Chomsky, N. (1957). *Syntactic structures*. The Hague, The Netherlands: Mouton.

Cohen, A., & Hersh, W. (2005). A survey of current work in biomedical text mining. *Briefings in Bioinformatics, 6,* 57-71.

Consoli, D. (2009). Analyzing customer opinions with text mining algorithms. *AIP Conference Proceedings, 1148*, 857-860.

Dale, R., & Moisl, H. (2000). (Eds). *natural language*