

# Enhancement Of Pattern Matching In Data Mining And Comparative Analysis With Knuth-Morris-Pratt, Boyer-Moore Algorithm

Ankit Jangra<sup>1</sup>, Shilpa Nagpal<sup>1</sup>

<sup>1</sup>Department of CSE ,

Prannath Parnami Institute of Management & Technology Hisar, Haryana, India

shilpa.n.ppimt@ppu.edu.in

**Abstract:** The Purpose of data mining is extracting vital information from huge databases or the data warehouses. Many Data mining applications have used for commercial & scientific sides. This type of study emphatically discusses Data Mining applications into scientific side. Here Scientific data mining differentiates itself and explores that nature of datasets is various from present market concentrated data mining applications. Most people use pattern matching in some form. Search engines on Web use pattern matching to locate information of interest.

**Keyword:** information, applications, Scientific, Search, warehouses.

## I. INTRODUCTION

The huge amount of data is available in the Information Industry. There is no use of this data until it has to be converted in useful information. It is very important to analyze this big amount of data & to find the useful information from it.

Extraction of data is not just the procedure that are expected to play out; the information mining likewise incorporates other many procedures like Data Integration, Data Cleaning, Data Mining, Pattern Evaluation, Data Transformation, & Data Presentation.

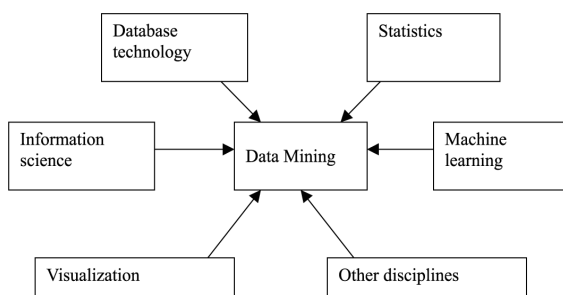


Fig 1 Classification Data mining

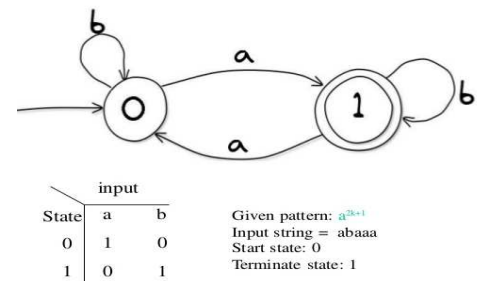


Figure 1: An automaton.

## Fig 2 Knuth-Morris-Pratt

### Pre-processing

Before the data mining procedure could be used step by step, a target data must be settled. As data mining may be exposed in that style which is actually present in the data, the target ought to be sufficient to contain the patterns where remaining concise should be mined in an acceptable time jump. The source of data is data mining or the data warehouse.

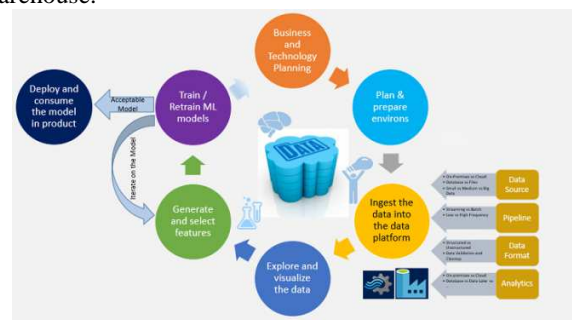


Fig 3 Data mining

### Advantages & disadvantages

DFA's are one of most pragmatic models of calculation, since there is a trifling straight time, steady space, online calculation to reenact a DFA on a surge of information. Additionally, there are effective calculations to discover a DFA perceiving. Complement of language recognized by a given DFA.

1. Union/intersection of languages recognized by two given DFAs.

Because DFAs could be reduced to a canonical, there are also efficient algorithms to determine:

### Research Problem

Syntax error correction is essential part of debugging process. Yet there had been little research investigating how programmers approach syntax error correction & how to help beginner programmers learn to fix errors efficiently.

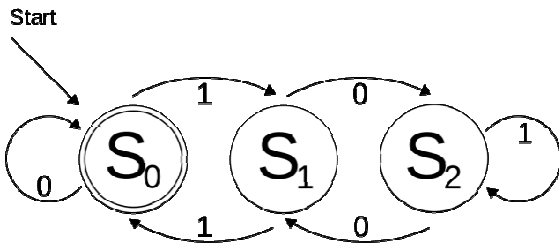


Fig 4 DFA

## II LIRETURE REVIEW

### Ashish Sureka 2011 Detecting Duplicate Bug Report Using

Duplicate reports needs to be identified to avoid a situation where duplicate reports get assigned to multiple developers. Also, duplicate reports could contain complementary information which could be useful for bug fixing.

### V. Neelima1, 2013 Annapurna Bug Detection through Text Data Mining [4]

Content Mining is an interdisciplinary field that draws on data recovery, Data Mining, mach.

ne learning, insights, & computational linguistics. As most information is currently stored as text, Text Mining is believed to have a high commercial potential value.

### Promila D evi1, Rajiv Ranjan 2014 Enhanced Bug Detection by Data Mining Techniques [5]

## III RESEARCH METHODOLOGY

### Pattern matching

Most people use pattern matching in some form. Search engines on the Web use pattern matching to locate information of interest. Patterns can be specific or quite general, using various wildcards that match multiple endings, words, or strings. Many databases have a similar capability, in which a character such as an asterisk is used as a wildcard. . Other systems use methods such as hidden Markov models (HMMs) and neural networks (NNs) for pattern matching.

## IV PROPOSED WORK

In this research we have to enhance performance of existing pattern matching algorithm by modifying them. Objective of our research is to decrease time consumption during pattern matching.

Most people use pattern matching in some form. Search engines on Web use pattern matching to locate information of interest.

### Algorithm of Proposed work

1. Take two string b and a.
2. Here b is pattern and a is actual string.
3. Get the length of a and store in m.
4. Get the length of b and store in n.

5. Set flag variable f equal to 0.
6. Set location variable loc to 0.
7. Set C to 0.
8. Set I to 1 and increment it to m with following step
  - a. Set increment c by 1.
  - b. If string a(i) and b(1) at matches then compare a(i+m-1) with b(n). then set loc=I and flag f to 1.
  - c. Start a nested loop to check inner string of pattern and sub string.
  - d. If data does not match at particular location then again flag is set to 0.
9. At the end if flag f is 1 then string found at loc location.
10. If flag f is 0 then string is not found.

## V EXPERIMENTAL RESULTS

### Implementation of KMP

#### KMP

1. Compares text left-to-right
2. Uses a failure array to shift intelligently
3. takes O(m), where m is length of pattern, to compute failure array
4. takes O(m), space
5. takes O(n), time to search a string

```
> booremoore('abc', 'ghhgghabchjhj')
:elapsed time is 0.000542 seconds.
```

```
uns =
```

```
1
```

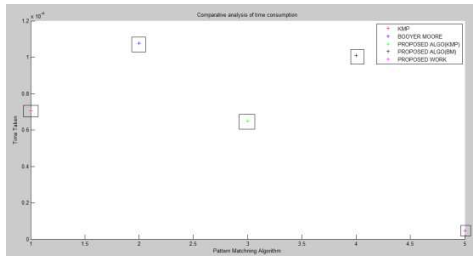
Fig:5 Pattern Is Found In Booremoore Search 1 Would Be Returned

### Proposed implementation

In proposed algorithm we have reduce time consumption secondly there are no of limited character support in Boore moore algorithm here we have increase probability of number of character & reduced time consumption to perform search. Following are result that represent patter matching in case of KMP , BOOYER MOORE And our proposed work.

	Result
<b>KMP BASED</b>	Found
<b>BOOYER MOORE</b>	Found
<b>PROPOSED WORK</b>	Found

FIG:6 Result during pattern matching



**Fig: 7 Mat lab Code To Represent Comparative Analysis**

### CONCLUSION

The Proposed algorithm will be integrated with a new approach which would surely reduce the consumption of time during algorithm tradition and new algorithms .Here Patterns are generally in the form of sequence or structures of trees. The advantages of this pattern matching is consist of out coming locations of the pattern in a token sequence, to get any component with matched pattern, & for substituting the matching pattern with any other token sequence. These Patterns are usually described by using regular expressions & these are matched by using techniques like backtracking. The Knuth-Morris-Pratt algorithm is a very good choice while we desire to search for the same pattern in a repeated way for many different texts. They believe that the answer is true to get as the assumption is every time you run algorithm on the various text preprocessing is only  $O(n)$  whereas for BM this is  $O(n + \text{size of alphabet})$ .

### FUTURE SCOPE

In this type of work finding a particular sequence of the tokens is constituted for the presence of constituents of some type of pattern in to find anomaly. The Proposed algorithm will be integrated with a new approach which would surely reduce the consumption of time during algorithm tradition and new algorithms .Here Patterns are generally in the form of sequence or structures of trees. The advantages of this pattern matching is consist of out coming locations of the pattern in a token sequence, to get any component with matched pattern, & for substituting the matching pattern with any other token sequence.

### REFERENCES

- [1] M. O. Mansur, 2 Mohd. Noor Md. Sap 2005 Outlier Detection Technique in Data Mining: A Research Perspective
- [2] Wahidah Husain1, Pey Ven Low2, Lee Koon Ng3, Zhen Li Ong4, 2011 Application of Data Mining Techniques for Improving Software Engineering
- [3] Ashish Sureka 2011 Detecting Duplicate Bug Report Using Character N-Gram-Based Features we present an approach to identify duplicate bug.
- [4] V. Neelima1, 2013 Annapurna Bug Detection through Text Data Mining Volume 3, Issue 5, May 2013 ISSN: 2277 128X International Journal of

Advanced Research in Computer Science & Software Engineering.

- [5] Promila Devi1, Rajiv Ranjan 2014 Enhanced Bug Detection by Data Mining Techniques.
- [6] Safia Yasmeen 2014 Software Bug Detection Algorithm using Data mining Technique.
- [7] Dhyana Chandra Yadav April 2015 Software Bug Detection using Data Mining International Journal of Computer Applications (0975 – 8887) Volume 115.
- [8] Damini.V.S1, Rabindranath2, Priyashree.K3, .Jayashubhaj.K4 10, May 2016 Optimized Error Detection Analytics within Big data on Cloud International Journal of Innovative Research in Science, Engineering & Technology Vol. 5, Special Issue 10, May 2016.
- [9] Tao Xie & Suresh Thummalapenta, North Carolina State University, David Lo, Singapore Management University, Chao Liu, Microsoft Research —Data Mining in Software Engineering, August, 2009, pp. 55-60.