

Survey and Comparison of String Matching Algorithms

Chayapathi A R^a, G Sunil Kumar^b, Manjunath Swamy BE^c, Thriveni J^d, Venugopal K.R^e

^a Information Science Department, Visvesvaraya Technological University, Acharya Institute of Technology Bengaluru, Karnataka, India, archayapathi@gmail.com

^b Computer Science Department, Visvesvaraya Technological University, Bangalore University, UVCE Bengaluru, Karnataka, India, gsuneel.k@gmail.com

^c Computer Science Department, Don Bosco Institute of Technology, Bengaluru, Karnataka, India,

^{d,e} Computer Science Department, Bangalore University, UVCE Bengaluru, Karnataka, India

Article History: Received: 11 January 2021; Revised: 12 February 2021; Accepted: 27 March 2021; Published online: 23 May 2021

Abstract: There are many applications which makes use of pattern matching algorithm. Most of current websites implements pattern matching algorithm in order to display the results faster. There exist different kind of data such as image, text, video, audio. In order to deal with such kind of data different pattern matching algorithms are used. One algorithm performs well in particular type of data, while it degrades in other kind of data. Our aim is to find best pattern matching algorithm. One of the key aspects of any string-matching algorithm is how fast the string matching is done along with the degree of search performance. This paper offers a survey on various String-matching algorithms along with the comparative analysis to provide a brief idea regarding the better algorithm for improving the search performance..

Keywords: Brute Force, Rabin-Karp, Boyer-Moore, Knuth-Morris, Aho-Corasick, Commentz-walter, Smith-Waterman, Needleman-Wunsch, Hamming Distance, Levenshtein Distance

1. Introduction

In the current world any websites with or without internet connected will implement search options in their web applications. This is implemented to get the results with less time without searching whole website. Pattern matching algorithms had made its roots in many domains such as medical, information technology, data mining, machine learning, forensics, network, defence, space. Pattern matching algorithm is a technique which accepts two parameters such as the pattern and the large set of data or document which may or may not contain given pattern, then the pattern is matched against the document to find whether it exist in that document or not. Required actions are taken based on the results.

String matching algorithms are identified in various methods. Such as Approximate and Exact string matching algorithms. Exact string matching is searching for the same pattern in the text and approximate string matching is searching for the most similar pattern in the text. And, the search can be made on the basis of the pattern occurrence in the given text. They are Single pattern search and Multiple pattern search. Single pattern search is searching for the single and first incidence of the pattern in the text and the process of identifying the many existence of the same given pattern in the text is Multiple pattern search.

The main job of pattern matching algorithm is to find whether given pattern exist in the large set of data. Based on the match one can take required decisions. Algorithms implemented must be in such a way that it should meet the requirements such as time complexity, space complexity, memory and fetch the results faster.

There are several Pattern matching algorithms namely Boyer Moore algorithm, Rabin-Karp algorithm, naïve string search algorithm, Needleman-Wunsch algorithm, Hamming distance and Levenshtein algorithm, Commentz-Walter algorithm etc. that can be applied for exact or approximate searches to be made accordingly. All of these string matching algorithms play a vital role in implementing the above-mentioned applications in the real-world scenarios.

2. Survey On Pattern Matching Algorithms

Brute Force Algorithm

Brute force algorithm popular as Naïve algorithm. It is very direct approach to search any text string. It keeps iterating through the text, and the pattern is compared with the first few characters of text for the length of pattern. If mismatch occurs shift the pattern one step right and with the first character of pattern compare next character of text and if match occurs proceed comparison with next characters of both text and pattern. Continue the above process, if match occurred for the entire length of pattern that means pattern occur in the text string hence return the position where the match occurred. Time complexity is $O(m*n)$ as both worst case and best case, where (m) is the length of text string and n is the length of the pattern.

Input: Patt [1...m] is a pattern string; Txt[1...n] is a text string.

Output: Position of the sub-string of text matching Patt or -1 if not matched then its returned

for j ← 0 to n-m do

i ← 0

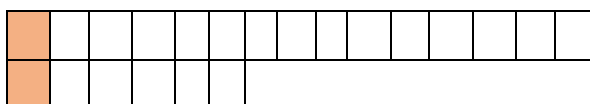
While i < m && Patt[i] == Txt[j+i] do

i ← i+1

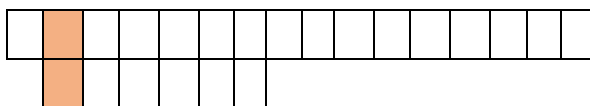
If i == m return j //match successful

Return -1 // match unsuccessful [2]

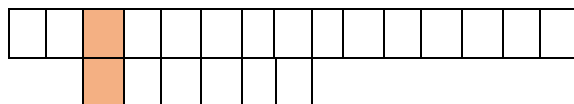
Consider an example where text be “CAT IS A MAMMAL” and pattern be “MAMMAL”.



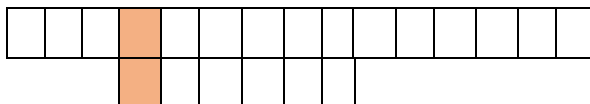
C! = M hence shift pattern by 1 and compare pattern from next character of text.



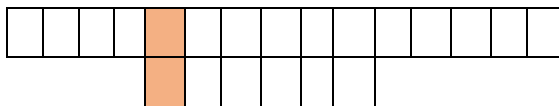
A! = M and hence shift pattern by 1 and compare pattern from next character of text.



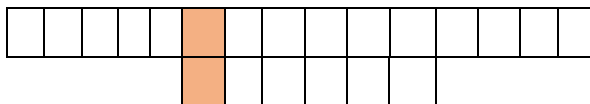
T! = M and hence shift pattern by 1 and compare pattern from next character of text.



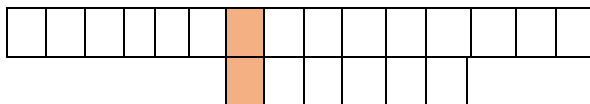
‘!’ = M and hence shift pattern by 1 and compare pattern from next character of text.



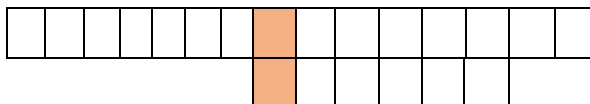
I! = T and hence shift pattern by 1 and compare pattern from next character of text.



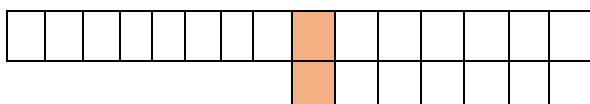
S! = M and hence shift pattern by 1 and compare pattern from next character of text.



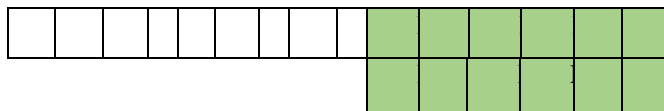
‘!’ = M and hence shift pattern by 1 and compare pattern from next character of text.



A! = M and hence shift pattern by 1 and compare pattern from next character of text.



‘!’ = M and hence shift pattern by 1 and compare pattern from next character of text.



Now all the characters of pattern, match with the txt character, hence algorithm returns the position where match is successful.

Applications

The brute-force algorithm is used to determine the matches between the decimal RGB frames and the secret text in video steganography. [3]

Advantages

- Brute force algorithm is a basic and simple algorithm mainly used when search happens in small amount of data.
- It does not require pre-processing.

Disadvantages

- It is not efficient algorithm hence not possible to implement where data is in huge quantity.
- It fails solving the problem which contains hierarchical structured data and the data contains logical operations.
- It is not efficient when there are lots of matching prefixes ex: if pattern is “ddde” and text is “ddddddddddde”.

B. Rabin-Karp Algorithm

Rabin-Karp Algorithm works based on the hashing technique. It is similar to brute force comparison except it improves the speed of comparison. First step is to calculate the hash value of the given pattern. It makes window of size length of pattern, and this window is made movement right to the text each time when hash values become unequal. Second step is to calculate the hash value of characters inside the window. Then the algorithm iterates through the text string. If hash values of pattern and window become equal then only it starts comparison of each character in the window with each character of pattern and if all the characters of window matches with the characters of pattern then return the position of pattern in the text. If characters mismatch then it stops comparison and moves to the right by one character and continue the above process. $O(m*n)$ is the worst case Time-Complexity and $O(m+n)$ as average case.

Algorithm

```

rabinKarpSearch(txt, patt, prm)
Begin
patternLen := pattern-Length
patternHash := 0 and stringHash := 0, h := 1
stringLen := string Length
mxChar := total no of characters in the character set
for index k of all character in patt, do
hsh := (h*mxChar) mod prm
done
for all character index k of patt, do
patternHash := (mxChar*patternHash + patt[k]) mod prm
stringHash := (mxChar*stringHash + txt[k]) mod prm
done
for k := 0 to (stringLen - patternLen), do
if patternHash = stringHash, then
for chrIndex := 0 to patternLen - 1, do

```

```

    if txt[k+chrIndex] ≠ patt[chrIndex], then
        breaktheloop
    done
    if chrIndex = patternLen, then
        print the location as pattern found at k position.
        if k < (stringLen - patternLen), then
            stringHash := (mxChar*(stringHash - txt[k]*hsh)+txt[k+patternLen]) mod prm, then
                if stringHash < 0, then
                    stringHash := stringHash + prm
            done
        End [4]

```

For example, consider text = “acbfabcgef” and pattern = “abc”.

First calculate the hash value of the pattern. Let prime number be 3. Let the values for alphabets be 1 to 26 for a to z respectively.

$$\text{Hash value} = x_1 * \text{prime}^0 + x_2 * \text{prime}^1 + \dots + x_n * \text{prime}^n.$$

Where, {x1, x2, ..., xn} are the characters of the txt string, n is the length of pattern.

1. hash (abc) = 1*3^0+2*3^1+2*3^2 = 34.

hash of first three characters of text is hash (acb) = 1*3^0+3*3^1+2*3^2 = 28

28! = 34 hence calculate hash value of next three characters of text.

2. In order to make efficient algorithm calculate the hash value using rolling hash function

Let x = oldhash value – previous character value

$$x = x / \text{prime}$$

$$\text{newhash value} = x + \text{value of last character in the window} * \text{prime}^{\text{length}(\text{pattern})} - 1.$$

Therefore, hash(cbf) is

$$x = 28 - 1 = 27$$

$$x = 27 / 3 = 9$$

$$\text{hash}(cbf) = 9 + 6 * 3^2 = 63$$

63! = 34 hence calculate hash value of next three characters of text.

3. hash(bfa) is

$$x = 63 - 3 = 60$$

$$x = 60 / 3 = 20$$

$$\text{hash}(bfa) = 20 + 1 * 3^2 = 29$$

29! = 34 hence calculate hash value of next three characters of text.

4. hash(fab) is

$$x = 29 - 2 = 27$$

$$x = 27 / 3 = 9$$

$$\text{hash}(fab) = 9 + 2 * 3^2 = 27$$

27! = 34 hence calculate hash value of next three characters of text.

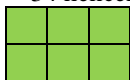
5. hash(abc) is

$$x = 27 - 6 = 21$$

$$x = 21 / 3 = 7$$

$$\text{hash}(abc) = 7 + 2 * 3^2 = 34.$$

34 == 34 hence now compare each character of pattern with the chosen text characters



All the characters match with the pattern hence stop iteration and return the position of pattern in the text that is

5.

Applications

- Detecting plagiarism.[5]

- Text processing
- Bioinformatics
- Compression [6]

Advantages

- It increases the speed when compared to brute force algorithm.
- Since it compares hash value first it skips character comparison against the pattern character and calculation of hash value takes less time.
- It can deal with multiple pattern matching hence good for plagiarism.

Disadvantages

- It performs inefficient when compared to brute force algorithm if hash values become equal and the characters are not same as pattern.
- It requires additional space.

C.Aho-CorasickAlgorithm

The Aho-Corasick algorithm is a popular dictionary matching algorithm. Here matching of all the dictionary words in a single iteration of text input is accomplished. Given all the dictionary words as the input, the algorithm firstly pre-processes them to build an automaton once and save for later data stream to match.

Aho-Corasick algorithm works by building a state machine using a string for comparison. The state machine will begin with a null empty root node which is by non-attendance unmatched state. Each pattern to be compared appends states to the machine, initially from the root node till pattern end is reached. By the traversal of state machine failure pointers are detected and inserted from each node to the highest prefix of the node.

First step is to build tree which is a tree like structure, tree ends with leaf and each leaf gives the various dictionary words. Next step is to construct failure function. Failure function is built in such a way that if the proper suffix of the current node is also a proper prefix then add a link from current node to the node which is also a proper prefix. If there is no suffix or if there is no proper prefix for the current node's proper suffix then add link to the starting node or the root node. It has three important functions success transaction, failure transaction and finally output matching. **Words for each tree node will be set up using bread first search traversal on the tree. The success transactions follow the edge in the tree to find the children of current tree node. The failure transaction set up links between failed string matches and the node on other branches which share the longest common suffix. The output list stores all the words ending at current node and its failure node.**

While running the algorithm it traverses the graph starting by success transaction to child node. If the pattern does not exist then follow failure transaction to its proper suffix node. If the algorithm reaches the node where output keyword is not empty, then algorithm will return all the matched characters that ends at the current character position of the input text string. It has time complexity of $O(m+n)$.

Algorithm

buildTree (patList, s)

Input: The list of all patterns, and the size of the list

Output: Transition map is generated to find the patterns

Begin

initialize elements to output-array to 0

initialize elements to fail-array to -1

initialize elements to goto matrix to -1

s := 1 //at first there is only one state(s)

for every pattern 'i' in the patList, do

word := patList[i]

present := 0

for every character 'ch' of word, do

```

    if goto[present, chr] = -1 then
goto[present, chr] := state
s := s + 1
present:=goto[present, chr]
    done
    out[present] := out[present] OR (shift left 1 for i times)
done
for every characters chr, do
    if goto[0, chr] ≠ 0 then
        fail[goto[0,chr]] := 0
        insert goto[0, chr] into a Queue q
    done
while q is not empty, do
newState := first element of q
    delete from q
    for every character chr, do
        if goto[newState, chr] ≠ -1 then
failure := fail[newState]
        while goto[failure, chr] = -1, do
failure := goto[failure, chr]
        done
fail[goto[newState, chr]] = failure
out[goto[newState, chr]] :=out[goto[newState,ch]] OR out[failure]
        insert goto[newState, chr] into q.
    done
done
return s
End

```

getNextState(presState, nextChar)

Input:the present state character and the next character to findthe next state

Output: the next state

Begin

answer := presState

ch := nextChar

while goto[answer, chr] = -41, do

answer := fail[answer]

done

return goto[answer, chr]

End

patternSearch(patList, s, text)

Input: List of patterns, size of the list and the main text

Output: The indexes of the text where patterns are found

Begin

 call buildTree(patList, s)

presState := 0

for every indexes of the text, do

 if out[presState] = 0

 ignore the next portion and go for next iteration

for every patterns in the patList, do

 if the pattern is found using output array, then

 print the location where pattern resides

 done

done

End [7]

Consider an example where finite set of patterns be {HONEY, MOON, MONEY and NET}

Automata for the above patterns is shown in the fig1

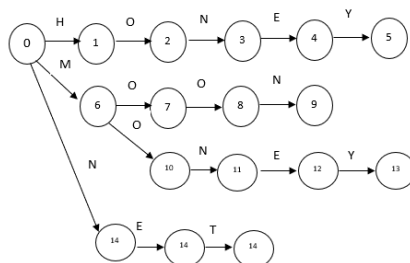


Fig. 1. Automata

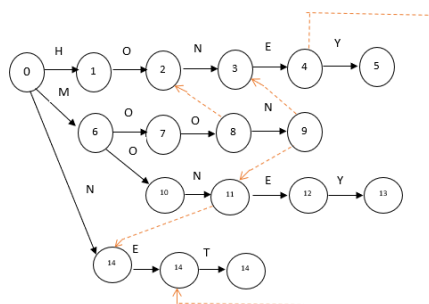


Fig. 2. Failure function for the automata

Then failure function is constructed as shown below fig2.

Output function transition is shown in the fig3 and output function table is shown in the fig4.

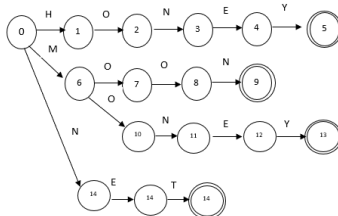


Fig.03.0Output0function

FINAL STATE	OUTPUT
NODE 5	HONEY
NODE 9	MOON
NODE 12	MONEY
NODE 15	NET

Fig. 4. Output function table

Finally, pattern is searched in the constructed automata. the searching phase of ahcorasick is simple while scanning the text it walks through automata if any transition found, it getstransition, else check the failure function.

If text is HONEYPOTNET then search is done as shown in the fig5. [8].

STATE	CHARACTER	TRANSITION	FAILURE	COMMENT
0	H	0 -> 1	-	TRANSITION FOUND
1	O	1 -> 2	-	TRANSITION FOUND
2	N	2 -> 3	-	TRANSITION FOUND
3	E	3 -> 4	-	TRANSITION FOUND
4	Y	4 -> 5	-	TRANSITION FOUND
5	P	-	0	TRANSITION NOT FOUND
6	O	-	0	TRANSITION NOT FOUND
7	T	-	0	TRANSITION NOT FOUND
8	N	0 -> 13	-	TRANSITION FOUND
9	E	13 -> 14	-	TRANSITION FOUND
10	T	14 -> 15	-	TRANSITION FOUND

Fig. 5. Searching transition table of automata

From the fig 5 there exist two meaningful words from the given text, hence this algorithm can be used to identify any bad packets entering into the network.

Various Applications are

- Intrusiondetection mechanism
- Detection of Plagiarism
- Deploy Bioinformatics tools
- Applications of Digitalforensic
- Textmining arena

Advantages

- Everycharacteroftextisanalyzedonlyonetime.
- Despite of the input symbols. Thedeterministictransitionstepisachievedbetweenstates[9]

Disadvantage

- Algorithm makes use of more storage to store transition rules of the deterministic finite state machine. [10]

D. Boyer-Moore Algorithm

Boyer-Moore algorithm compares the characters starting from right to the left of the pattern against the text in the same direction as like pattern, starting with the index equal to the length of pattern-1. It matches the tail of the pattern rather than head. This algorithm makes use of bad match table which is the main cause to reduce the time complexity.

Construction of bad match table

1. This table must not have value less than 1.
2. Keep comparing the pattern to the text starting with the right most character in the pattern.
3. Make a table rows representing value and columns representing characters of the pattern.
4. The table must not contain repetitive character, if the pattern contains repeated character update the value corresponding to that character.
5. Value for last character will be length of pattern if that character was not existing before otherwise leave the same value.
6. Other character which is not present in the pattern is represented by * in the table and value assigned will be the length of the pattern.

This algorithm has time complexity of $O(m/n)$ as best case, $O(m*n)$ as worst case and $O(m/|\Sigma|)$ as average case.

Algorithm

```
fullSuffixMatch(shiftArr, borderArr, pattern)
```

```
Begin
```

```
n := pattern length
```

```
i := n
```

```
k := n+1
```

```
borderArr[i] := k
```

```
while i > 0, do
```

```
    while k <= n AND pattern[i-1] ≠ pattern[k-1], do
```

```
        if shiftArr[k] = 0, then
```

```
            shiftArr[k] := k-i;
```

```
            k := borderArr[k];
```

```
        done
```

```
    decrease i and k by 1
```

```
borderArr[i] := k
```

```
done
```

```
End
```

```
partialSuffixMatch(shiftArr, borderArr, pattern)
```

```
Begin
```

```
n := pattern length
```

```
j := borderArr[0]
```

```
for index of all characters „i“ of pattern, do
```

```
    if shiftArr[i] = 0, then
```

```
shiftArr[i] := j
  if i = j then
j := borderArr[j]
done
End
searchPattern(txt, patt)
Begin
patternLen := patt length
stringLen := txt size
for all entries of shiftArr, do
  set all entries to 0
done
call fullSuffixMatch(shiftArr, borderArr, patt)
call partialSuffixMatch(shiftArr, borderArr, patt)
shift := 0
while shift <= (stringLen - patternLen), do
j := patternLen - 1
  while j >= 0 and patt[j] = txt[shift+j], do
    decrease j by 1
  done
if j < 0, then
  print the shift as, there is a match
shift := shift + 0shiftArr[0]
  else
    shift := shift + shiftArr[j+1]
done
End0[11]
```

Consider an example, let text be "THIS IS A BOOK" and pattern be "BOOK"

Construct a bad match table as shown in the fig6:

Length of pattern = 4

Pattern	B	O	O	K
Index	0	1	2	3

Value for B = 4 - 0 - 1 = 3

Pattern	B	O	O	K	*
Value	3				

Value for O = 4 - 1 - 1 = 2

Pattern	B	O	O	K	*
Value	3	2			

Value for O = 4 - 2 - 1 = 1

Pattern	B	O	O	K	*
Value	3	2	1		

K is the last character and k is not available earlier in the pattern hence it will have the value equal to length of the pattern that is 4

Pattern	B	O	O	K	*
Value	3	2	1	4	

Any other character is represented by * and value will be length of the pattern that is 4

Fig. 6. Searching transition table of automata

Next compare with the text string considering bad match table.

T	H	I	S		I	S		A		B	O	O	K
B	O	O	K										

S! = K and S is not the character of pattern that means it is other character hence see the value for the other character in bad match table, it is 4 hence shift the pattern by 4.

T	H	I	S		I	S		A		B	O	O	K
				B	O	O	K						

Space! = K and space is other character and hence value for the other character in the bad match table is 4 hence shift the pattern by 4.

T	H	I	S		I	S		A		B	O	O	K
										B	O	O	K

B! = K and B exist in the pattern and bad match table and value for B is 3, Hence shift the pattern by 3.

T	H	I	S		I	S		A		B	O	O	K
										B	O	O	K

K == K hence compare previous characters of text with the previous characters of pattern until the length of pattern.

T	H	I	S		I	S		A		B	O	O	K
										B	O	O	K

Now all the characters of pattern are matched with text characters hence return the position.

Applications:

- Text editors
- Commands substitutions [12]
- Intrusion Detection System.

Advantages:

- Boyer-Moore algorithm pre-process only the pattern not the text.
- Algorithm runs faster as length of pattern increases.
- It skips many characters at the same instance instead of searching of every character hence it is efficient algorithm.

Disadvantage:

- Mismatch character will give small shift in some condition, if match not occurs after many matches [13].
- Unable to process small size patterns properly. [14]

E.Knuth-Morris-prattAlgorithm

Knuth-Morris algorithm contrast the characters of pattern and text from left to right. It works based on prefix and suffix match within the given pattern. Compare each character of text with each character of pattern, if all symbols of pattern matched with the text substring of length pattern, then return starting position of text string where pattern exist. If there is no match of particular character then find substring in the pattern which must be suffix as well as prefix in that substring. If no found then compare next character of text with starting character of pattern and continue the process. If suffix and prefix found then compare next character of text with next character immediately after the prefix substring and continue the process. This method avoids backward movement for comparison and also reduces time complexity. It has time complexity of (m) where m is the length of text string.

The algorithm can be made more efficient if temporary array is built. This array contains from which position comparison need to takes place. Time complexity to build array is $O(n)$ where n is length of pattern. Hence over all it has time complexity of $O(m+n)$.

Algorithm:

```
findprefix(patt, m, prefixArr)
```

```
Begin
```

```
len := 0
```

```
prefixArray[0] := 0
```

```
for all character index k of pattern, do
```

```
  if patt[k] = patt[len], then
```

```
    increase len by 1
```

```
prefixArray[k] := len
```

```
  else
```

```
    if len ≠ 0 then
```

```
len := prefixArr[len - 1]
```

```
  decrease k by 1
```

```
  else
```

```
prefixArr[k] := 0
```

```
  done
```

```
End
```

```
Kmp_Algorithm(txt, patt)
```

```
Begin
```

```
N1 := size of text
```

```
M1 := size of pattern
```

```
call findprefix(patt, M1, prefixArr)
```

```
while k < N1, do
```

```
  if txt[k] = patt[j], then
```

```
    increase k and j by 1
```

```
  if j=M1, then
```

```
    printthelocation
```

```

(k-j)asthepatternisthere
  j:=prefixArr[j-1]
elseifk<N1ANDpatt[j]≠txt[k]then
  ifj≠then
    j:=prefixArr[j-1]
  else
    increasekby1
done
End[15]

```

For example, consider text be “abgabfabfabx” and pattern be “abfabx”.

Temporary array for pattern must be created before comparison as shown in the figure 7. Initially the values for first pattern will be zero.

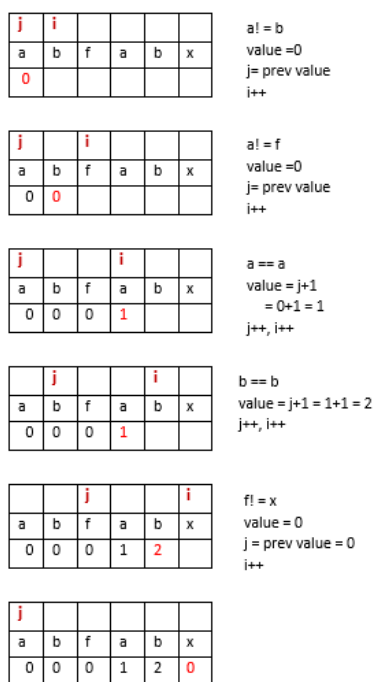


Fig .7. Temporary table

After construction of temporary table pattern is matched with text as shown in the figure 8.

a	b	g	a	b	f	a	b	f	a	b	x
a	b	f	a	b	x						
0	0	0	1	2	0						

$g \neq f$ and there is no suffix which is also a prefix in both pattern and text, check for the previous value of 'f' that is 0, hence start next comparison of pattern from index zero against the character of text where match was unsuccessful that is 'g'.

a	b	g	a	b	f	a	b	f	a	b	x
		a	b	f	a	b	x				
		0	0	0	1	2	0				

$g \neq a$ and there is no suffix which is also prefix in both pattern and text, there is no previous value for 'a' hence shift by one and compare the pattern with next character of text.

a	b	g	a	b	f	a	b	f	a	b	x
			a	b	f	a	b	f	a	b	x
			0	0	0	1	2	0			

$f \neq x$ and 'ab' is a suffix and also a prefix in text as well as pattern hence shift the pattern to the position where it is also a suffix and start comparison from index position 2.

a	b	g	a	b	f	a	b	f	a	b	x
						a	b	f	a	b	x
						0	0	0	1	2	0

Fig .8. Text pattern comparison according to Knuth-Morris algorithm

Applications

- Parallel Knuth-Morris is to be used in parallel image processing applications [16]
- DNA sequence analysis.

Advantages

- It is more efficient than rabin karp and naïve algorithm.
- The execution time of KMP algorithm is $O(m+n)$ which is very fast.
- Algorithm not required moving in backwards direction of the text string. [17]
- This algorithm works better if text length increases hence this algorithm is implemented where search need to be done in large documents.

Disadvantages

- It won't work so well as the alphabet size enhances. Due to which the odds of disparity is more. [18]

F. Commentz-Walter Algorithm

The string probing Commentz-Walter algorithm is proposed by Beate Commentz-Walter. It is a combined with several notes from Aho–Corasick with the fast matching of the Boyer Moore string search algorithm [19]. As in the Aho–Corasick string matching algorithm, at once it can investigate for multiple patterns. It suits best for the applications that possess pattern that are shorter than the text or where it carries on through several probes. The Boyer–Moore algorithm uses information gathering during the pre-process step to skip sections of the text, resultant in a lower steady factors than many other string based search algorithms. From a generic perspective the execution of the algorithm speeds up with increase in the length of the patterns.

The important step in this string matching algorithm is when the string matching process finds a mismatch in the end of the pattern then it skips the text instead of probing every symbol in the given text. If the characters are not matching with any of the characters in the text no need arises to continue backward searching along the text. If the symbols in the probing text do not match with the pattern text, then the next character in the text to verify is found n characters farther along the text, where n is the length of the pattern. The length of the pattern can be formulated through a bad character table. A partial shift is initiated based on the presence of a character in the

text. Then set up along with the matching character and the process is iterated. This method of jumping along the text for comparisons instead of verifying every symbol in the text results in decrease in the number of comparisons. This enhances the competence of the algorithm. The Commentz-Walter algorithm has a time complexity $O(N+M+Z)+O(MN)$ for execution.

Algorithm:[25]

Compute function last

$a \leftarrow k-1$

$b \leftarrow k-1$

Repeat

If $P[b]=T[a]$ then

if $b=0$ then

return a // we have a match

else

$a \leftarrow a-1$

$0 \leftarrow b-1$

else

$a \leftarrow a+k - \text{Min}(b, 1+\text{last}[T[a]])$

$b \leftarrow k-1$

until $a > n-1$

Return "no match"

Example:

Input: MainString: "ABAAABCDBBABCDDDEBCABC", Pattern: "ABC"

Outputs/Results:

Search Pattern occurs in location: 4

Search Pattern occurs in location: 1

Search Pattern occurs in location: 18

Applications

- Text editors
- command substitutions

Advantage

- This algorithm is the fastest when pattern is moderately sized.

Disadvantage

- But the pre-processing time that is taken in this algorithm is considered to be a disadvantage as it requires more time.

G. Waterman Algorithm

The Smith Waterman algorithm is based on the principle of dynamic programming. It computes the optimal local alignment of two sequences [2]. The Smith Waterman algorithm is for detecting local alignments of sequence. Also it ensures detection of identical regions prevailing between two nucleotide or protein sequences. The algorithm is used to compare segments of all possible lengths to arrive at optimal similarity. On comparing with the Needleman Wunsch Algorithm, the algorithm ensures that the negative scoring matrix cells are set to zero. Thus for backtracking only positive scores are visible. The algorithm functions by starting with maximum scoring matrix cell and progress until zero-recorded cell is obtained. Finally it produces the local alignment with highest score.

The steps of operation for two sequences A and B are illustrated below:

1. Before and after alignment the symbols in a sequence should be in the identical order.
2. Establishing Alignment a symbol from a sequence with another is always possible.
3. Alignments are denoted by a blank ('-')
4. Alignment of two blanks is not permitted

Smith Waterman Algorithm relies on Gapped alignments to find the optimal distance between sequences by aligning with the gaps. Smith Waterman algorithm has a time complexity of $O(MN)$ for execution.

Algorithm

1. Determine the substitution matrix and also the gap in penalty scheme. $s(a,b)$ is the similarity score for the elements having 2 sequences. Here k is the penalty of a gap with length- k

2. Create a matrix H of scores and assign it to the first row and first column. The scoring matrix size is given by the term $(n+1)*(m+1)$. Also the matrix employs a based indexing.

$$H_k = H_l \text{ for } k < n \text{ and } l < m$$

3. Enter the scoring matrix using the equation below

$$H_{ij} = \max(\text{Choice 1} \leftarrow H(i-1, j-1) + S(B(i), A(j)) \text{ \{score of aligning } a_i \text{ and } b_j\}$$

$$\text{Choice 2} \leftarrow H(i-1, j) + d \text{ \{score of } a_i \text{ along with gap\}}$$

$$\text{Choice 3} \leftarrow H(i, j-1) + d \text{ \{score of } b_j \text{ along with gap\}}$$

\{no similarity upto a_i and b_j \})

4. Traceback starts with uppermost scores in the H - the score matrix and culminates at a matrix cell possessing a score of traceback which is relied on the origin of every score to produce recursive best local alignment [25].

Applications

- Biometrics

Advantage

- As it implies to the local alignment problems Optimal local alignment can be achieved.

Disadvantage

- But the time complexity and the space complexity for this algorithm is comparatively high.

H.Needleman-Wunsch Algorithm

The Needleman-Wunsch algorithm works on the principle of optimal matching results. This is a basic algorithm employed for solving the problems of sequence alignment [21]. The Needleman-Wunsch algorithm operates by performing global alignment of two sequences. Moreover it is employed in the arena of bioinformatics for aligning protein and nucleotide sequences. This algorithm referred as optimal matching Algorithm and also is an example of dynamic programming. The aligned character scores are procured by using similarity matrix. Also the Linear Gap d is found

by using similarity matrix. The Needleman-Wunsch Algorithm comprises three stages:

1. Score Matrix Initialization
2. Score calculation and completing the trace back matrix.
3. Draw Inference using alignment of the trace back matrix [25].

The two types of matrices employed in Needleman-Wunsch Algorithm are: the score and the trace back matrices.

Traceback matrix algorithm:

1. Traceback employs a method of drawing inference of the paramount alignment through traceback matrices.
2. Traceback process compulsorily starts at the last cell and it is positioned as bottom right cell.
3. Its movement is based on the traceback value provided in the cell.

4. Three potential traversal occurring are: diagonal, left or up.
5. The traceback process is completed when the top-left cell is indicated by- “done”.

Best Alignment:

1. The traceback path based values are employed to infer the Alignments. Also, the values of the traceback matrix are taken into account.

2. The letters from two sequences are aligned in traceback matrix. Further Gap is created based on the sequence orientation. “Left” creates a Gap in the left sequence and a gap is created in the top sequence if it is “Up”. And, thus procured sequences have a backward alignment [25].

The Needleman-Wunsch algorithm has proven to produce best alignment for two sequences. It starts the traceback is accomplished from the right-lower corner position in the traceback matrix and further culminates at the left-top most cell position of the matrix. This is irrespective of the length or complexity of sequences. The algorithm has proven to function identically and guarantees best alignment for different sequences. The Needleman Wunsch algorithm has a time complexity of $O(MN)$ for execution.

Algorithm

```

fork=0 to length(B)-1
  F(k,0) ← d*k
end for
for l=length(A)-1
  F(0,l) ← d*l
end for
for k=1 to length(B)
  for l=1 to length(A)
    Choice1 ← F(k-1,l-1)+S(B(k),A(l))
    Choice2 ← F(k-1,l)+d
    Choice3 ← F(k,l-1)+d
    F(k,l) ← max(Choice1,Choice2,Choice3)
  end for
end for

```

To compute alignment, start from right bottom cell from the matrix and choose the possible choices

if Choice1, then A(l) and B(k) are aligned

if Choice2, then A(l) is aligned with a gap

if Choice3, then B(k) is aligned with a gap

Applications

- Bioinformatics to align nucleotide sequence

Advantage

- This search algorithm considers order of sequence of characters while comparing which makes it more efficient.

Disadvantage

- But requires same length of string ie. The pattern and the text for comparing.

H. HammingDistanceAlgorithm

Hamming Distance Algorithm is an approximate matching algorithm which allows defined distinction in the sample and the text during string matching. Estimated match is allowed for a limited number of errors or edit operations required for the search pattern to match with the text [24]. The mismatches can occur due to any difference in the character called 'mismatch/substitution' or an extra character called 'insertion' or a missing character called 'deletion'. Considering two strings of the same length, hamming distance between the two strings can be defined as the minimum number of replacements one should make to turn one of the strings as another. Hamming distance is measured by tracking the number of positions where corresponding symbols differ from each other. For alphabetical strings and DNA sequences the distance also works.

. Hamming distance model has the time complexity $O(N^2)$ for execution.

Algorithm:[25]

//Initialization

i=0 count=0

while str1[i]!=str2[i]

count++

i++

endwhile

return count

Example

In this example two DNA sequences considered are: AACTCCA and AGCTAAC, the Hamming distance occurring is 4, since symbol mismatch occurs at positions 2, 5, 6 and 7.

Applications

- Systematics as a measure of genetic distance

Advantage

- Suitable for exact string matching and allows Single-bit error detection and correction.

Disadvantage

- But it requires more execution time.

I. Levenshtein Distance Algorithm

Levenshtein Distance is an approximate matching algorithm which allows certain differences in the pattern and the text while string matching. String resemblance comprises of wide array of applications, prominent ones are: web search, text comparison, plagiarism detection. Also the different computation methods exist, salient ones are: the longest common, edit distance, and, substring algorithms [22]. Based on approximate matching, a restricted number of faults or correction operations are identified for the pattern searched in the matching process. The mismatches can occur due to any difference in the character called 'mismatch/substitution' or an extra character called 'insertion' or a missing character called 'deletion'. Considering two strings of the same length, Levenshtein edit distance between the two strings can be defined as the minimal number of replacements that should be made to turn one of the strings to the other which includes substitution, insertion as well as deletion. The difference between Hamming distance and edit distance is that, here we are not considering distance and the strings no longer need to be of the same length as they go through in sections and deletions as well. This algorithm is possessing $O(N+M)$ time complexity for execution.

Algorithm:[25]

//initialization

for q ← to m do

E(q,0) ← q

endfor

```

form←0tondo
E(0,m)←0
endfor
//editdistanceE(q,m)
forq←0tomdo
  form←0tondo
    if(T(m)=P(q))then
      E(q,m)←(q-1,m-1)
    else
      min←MIN[E(q-1,m),E(q,m-1)]
      E(q,m)←min+1
    end if
  endfor
endfor
returnE

```

Example

Levenshtein distance between *barking* and *dark*, these transformations are accomplished:

1. The word *barking* → (indicated as) *barkin* (with deletion of *g*)
2. The word *barkin* → *barki* (with deletion of *n*)
3. The word *barki* → *bark* (with deletion of *i*)
4. The word *bark* → *dark* (with substitution of *b*)

Thus it can be concluded that Levenshtein distance obtained for the two word strings is 4.

Application

- Spell checkers

Advantage

- As this algorithm uses the Edit distance which allows insertion and deletion along with substitution like the Hamming distance algorithm makes it much more efficient.

Disadvantages

- But this algorithm does not consider order of sequence of characters.

3. Comparative Analysis

Algorithm	Comparison	Pre-processing	Time Complexity
Brute Force	Right side to Left side	None	$O(n*m)$
Rabin-Karp	Right side to Left side	$O(n)$	avg $O(m+n)$ worst $O(m*n)$
Aho-Corasick	Not applicable	$O(m+n)$	$O(N + L + Z)$
Boyer-Moore	Right side to	$O(n+ \Sigma)$	$O(n), \Omega(m/n)$

	Left side		
Knuth-Morris	Right side to Left side	O(n)	O(m)
Commentz-walter	Right to Left	none	O(N+M+Z)+O(MN)
Smith-Waterman	Right side to Left side	-	O(MN)
Needleman-Wunsch	Right side to Left side	-	O(MN)
Hamming Distance	Right side to Left side	-	O(N ²)
Levenshtein Distance	Right side to Left side	-	O(N+M)

Table .1. Comparative analysis of the algorithms

Table 1.1 shows pre-processing time, comparison order, and time complexity for all the ten algorithms. Time complexity is different for each algorithm. When compared to all algorithms Knuth Morris algorithm has less time complexity. Hence Knuth Morris algorithm is an efficient algorithm.

4. Conclusion

From the survey the conclusion is that Boyer Moore and Knuth Morris algorithms have less time complexity. Both the algorithms have similar time complexity. Boyer-Moore algorithm works better if the pattern length is large. Whereas the Knuth-Morris algorithm is efficient when length of text string is larger and pattern has repeated patterns. Boyer-Moore algorithm is better to use if the pattern length is large and Knuth-Morris algorithm is better to use if length of text string is larger.

References

1. Jiji. N ,Dr. T Mahalakshmi ,Survey of Exact String Matching Algorithm for Detecting Patterns in Protein Sequence, Advances in Computational Sciences and Technology ISSN 0973-6107 Volume 10, Number 8 (2017) pp. 2707-2720
2. Vibha Gupta, Maninder Singh, Vinod K. Bhalla ,Pattern Matching Algorithms for Intrusion
3. Detection and Prevention System: A Comparative AnalysisInternationalConferenceon Advances in Computing,Communicationsand Informatics (ICACCI),2014
4. KhuloodAbuMaria,Mohammad A. Alia, MaherA. AlsarayrehandEman Abu Maria UN-Substituted Video Steganography, KSII TRANSACTIONS ON INTERNET AND INFORMATION SYSTEMS VOL. 14, NO. 1, January 2020.
5. Sheshasayee, A., &Thailambal, G. A comparitive analysis of single pattern matching algorithms in text mining. 2015 International Conference on Green Computing and Internet of Things (ICGCIoT), (2015).
6. en.wikipedia.org/wiki/Rabin%E2%80%93Karp_algorithm.
7. www.cs.rit.edu/~lr/courses/alg/student/1/Rabin_Karp.pdf.
8. Alfred V. Aho and Margaret J. Corasick Efficient String Matching: An Aid to Bibliographic Search.
9. Saima Hasib, Mahak Motwani, Amit Saxena, International Journal of Computer Science and Information Technologies, Vol. 4 (3) (2013)
10. HyunJin Kim, A Memory-Efficient Deterministic Finite Automaton-Based Bit-Split String Matching Scheme Using Pattern Uniqueness in Deep Packet Inspection.
11. Zeeshan Ahmed Khan, R.K Pateriya,Multiple Pattern String Matching Methodologies: A Comparative Analysis, International Journal of Scientific and Research Publications, Volume 2, Issue 7, July 2012 3 ISSN
12. Sheshasayee, A., &Thailambal, G. (2015). A comparitive analysis of single pattern matching algorithms in text mining. International Conference on Green Computing and Internet of Things (ICGCIoT) 2015.
13. www.personal.kent.edu/~rmuhamma/Algorithms/MyAlgorithms/StringMatch/ boyerMoore.htm

14. Vivek Srivastava, B K Trapathi, V K Pathak, A Novel Hybrid Intelligent Model for Classification and Pattern Recognition Problems, (IJCSIS) International Journal of Computer Science and Information Security, Vol. 10, No. 2, February 2012
15. S. Antonatos, K. G. Anagnostakis, M. Polychronakis, and E. P. Markatos, "Performance analysis of content matching intrusion detection systems," in Proc 4th IEEE/IPSJ Symposium on Applications and the Internet, 2004, pp. 208-215.
16. SS. Swapna, Yashdeep Jha, Syed Zaheed, Keertik Dewangan, Sayyed Mujahid Pasha, A Survey on Different Pattern Matching Algorithms of Various Search Engines, International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)
17. SercanAygün , EceOlçayGüneş, LidaKouhalvandi,Python Based Parallel Application of Knuth–Morris–Pratt Algorithm, IEEE 4th Workshop on Advances in Information, Electronic and Electrical Engineering (AIEEE) 2016
18. www.slideshare.net/sabiyasabiya/knuth-morris-pratt-string-matching-algo
19. Kranthi Kumar Mandumula, Knuth-Morris-Pratt,Indiana State University Terre Haute IN, USA
December 16, 2011
20. Beate Commentz- Waite:"A String Matching Algorithm fast on the Average", TR 79.09.007 Heidelberg ScientificCenter, IBM, Germany Sept. 1979.
21. Hsien-Yu Liao, Meng-Lai Yin, Yi Cheng,"A Parallel Implementation of the Smith-Waterman Algorithm for Massive Sequences Searching",in Proceedings of the 26th Annual International Conference of the IEEE EMBS San Francisco, CA, USA • September 1-5, 2004.
22. Bailong FENG and Jing GAO, "Distributed Parallel Needleman-Wunsch Algorithm on Heterogeneous Cluster System", in Proceedings of the 2015 International Conference on Network and Information Systems for Computers, 2015.
23. ShengnanZhang , Yan Hu , GuangrongBian "Research on String Similarity Algorithm based on Levenshtein Distance", School of Computer Science and Technology, Wuhan University of Technology, Hubei Province, Wuhan, China, Department of Aviation Ammunition, Air Force College of Service, Jiangsu Province, Suzhou, China, 2017.
24. Prince Mahmud, Md. Sohel Rana, Kamrul Hasan Talukder, "An Efficient Hybrid Exact String Matching Algorithm to Minimize the Number
of Attempts and Character Comparisons", 21st International Conference of Computer and Information Technology, 2018.
26. Solon P. Pissis and Ahmad Retha, "Generalised Implementation for Fixed-Length Approximate String Matching underHamming Distance & Applications", IEEE International Parallel and Distributed Processing Symposium Workshops, 2015.
27. www.wikipedia.org