

## Unit-6:-

### String Matching:-

#### 1. Naive string matching

Ex: Text  $T = \underline{AABA}CAAD\underline{AABA}BAA$

pattern  $p = AABA$  The pattern is found at index using

Naive string matching.

Sol: Text  $T = \boxed{A} \boxed{A} \boxed{B} \boxed{A} AC AAD AABA BAA$   
pattern  $p = \boxed{A} \boxed{A} \boxed{B} \boxed{A}$  pattern is found at index 1.

step 2:-

$i$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A} AC AAD AABA BAA$   
 $j$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  pattern is not matching.

step 3:-

$\boxed{A} \boxed{A} \boxed{BAAC} AAD AABA BAA$   
 $\boxed{A} \boxed{B} \boxed{A} \boxed{A}$  Not matched.

then shift one step to right side.

step 4:-

$\boxed{A} \boxed{A} \boxed{B} \boxed{A} \boxed{A} \boxed{C} \boxed{A} AAD AABA BAA$   
 $j$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  not matching, then shift into one step right.

step 5:-

$\boxed{A} \boxed{A} \boxed{B} \boxed{A} \boxed{A} \boxed{C} \boxed{A} \boxed{A} DAA BAA BAA$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  not matching, then shift into one step right

step 6:-

$\boxed{A} \boxed{A} \boxed{B} \boxed{A} \boxed{A} \boxed{C} \boxed{A} \boxed{A} \boxed{D} AABA BAA$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  not matching.

step 7:-

$\boxed{A} \boxed{A} \boxed{B} \boxed{A} \boxed{A} \boxed{C} \boxed{A} \boxed{A} \boxed{D} \boxed{A} ABA BAA$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  not matching, then shift to right.

step 8:-

$\boxed{A} \boxed{A} \boxed{B} \boxed{A} \boxed{A} \boxed{C} \boxed{A} \boxed{D} \boxed{A} \boxed{A} BAA BAA$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  not matching.

step 9:-

$\boxed{A} \boxed{A} \boxed{B} \boxed{A} \boxed{A} \boxed{C} \boxed{A} \boxed{D} \boxed{A} \boxed{A} \boxed{B} AA BAA$   
 $\boxed{A} \boxed{A} \boxed{B} \boxed{A}$  not matching.



Step 10:-

A A B A A C A A D A A B A A

A | A | B | A

A | A | B | A

pattern is found at index

Step 11:-

A A B A A C A A D A A B A A

A | B | A | A

A | A | B | A

Not matching.

Step 12:-

A A B A A C A A D A A B A A

B | A | A | B

A | A | B | A

Not matching.

Step 13:-

A A B A A C A A D A A B A A

A | A | B | A

A | A | B | A

pattern is found at

Index 8.

Step 14:-

A A B A A C A A D A A B A A

A | B | A | A

A | A | B | A

Not matching.

The window is shifting to one step to right there is no text.

→ pattern is found at Index 1, 10, 13.

Algorithm:-

Algorithm Naive ( ) — 0

{

n ← length (text) — 1

m ← length (pattern) — 1

for s ← 1 to n - m + 1 do — 0 (n - m + 1)

do P[1...m] = T[s+1...s+m] — 1

Then pattern "pattern occurs" with shift — 1

} — 0

Best case:-

Text: A A B A A C A A

pattern: F | A | A — 0(n)

worst case

Text: A A A A A A A A

patterns: A A A A — 0(m \* (n - m + 1))

length of pattern window.



Text: a c a a b c

patterns: aab

Find the pattern in a given text using Naive algorithm.

\* Tries: Also called as Radix Tree, Digital Tree, prefix Tree.

→ There are three types of Tries. → Information Retrieval data structure.

1. Standard Trie:-

string  $s = \{ball, box, bomb, basket, stock, stop\}$

→ standard trie is a tree which stores the strings.

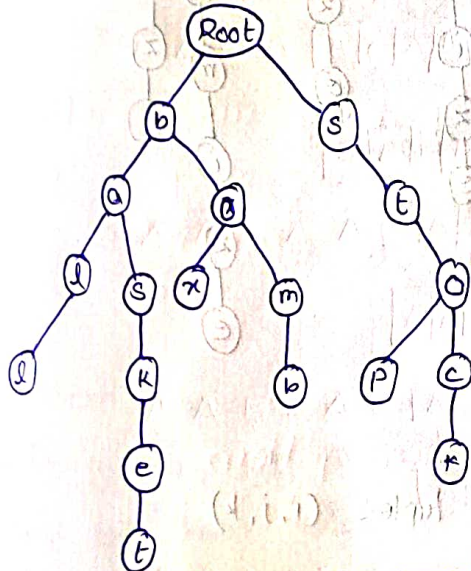
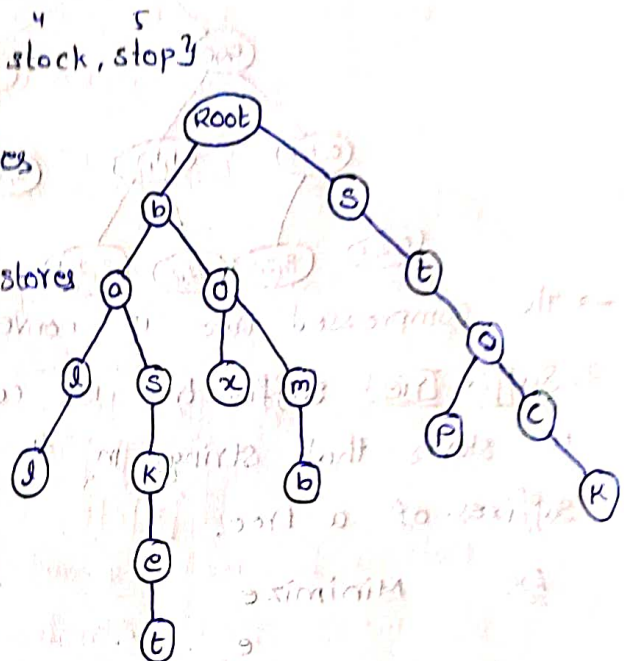
→ Each and every node in a tree stores the letter from the given string.

2. Compressed Trie:-

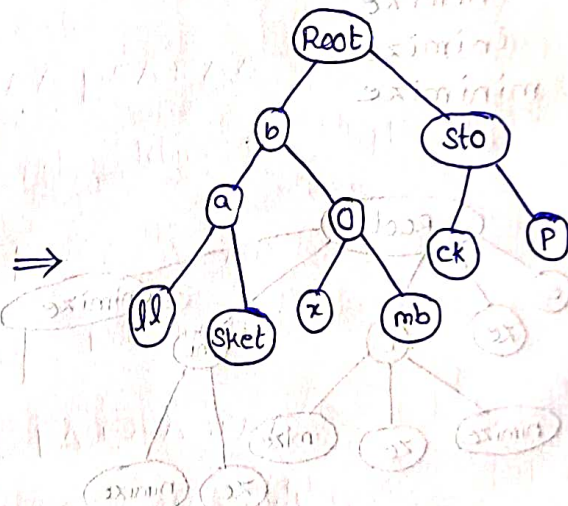
→ Advanced version of standard trie which stores the strings in the form of tree structure.

Each and every node stores the letters from the given string.

→ so we are merging the leaf nodes which contains one node or group of nodes.



⇒



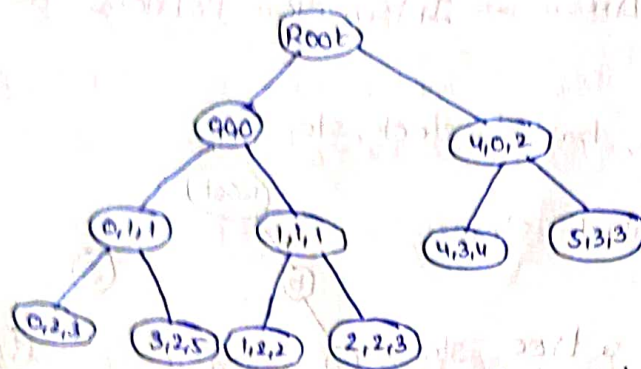
0(0) = b a  
 8(1) = b' o  
 8(2) = b o m b  
 3(3) = b a s k e t  
 s(4) = s t o c k  
 s(5) = s t o p

Tuple  $(i, j, k)$

$i$  is the string present in the index.

$j$  is the prefix which is starting at index.  $[chr(i+c)]$

$k$  is the prefix which is ending at index.



→ The compressed tree is converted into the form of tuple  $(i, j, k)$

3. Suffix Tree - suffix tree is compressed tree which can be used to store that strings in the form of tree structures by using suffixes of a tree.

Ex: Minimize

e

ze

ize

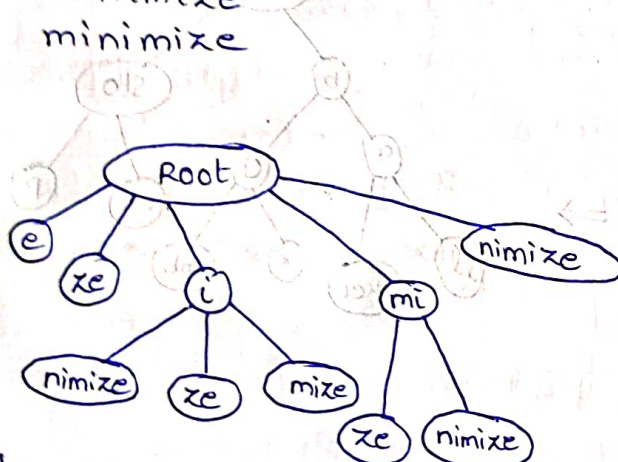
mize

imize

nimize

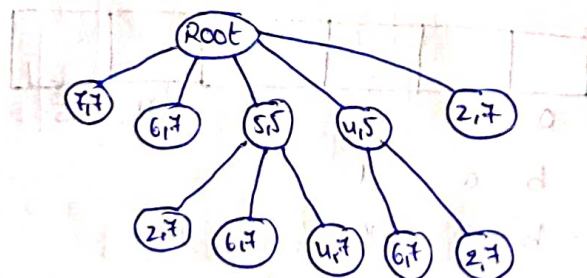
inimize

minimize



→ The tree is converting in the form of tuples  $(i, j, k)$

$[ctrl+v]$





# \* Robin-Karp Algorithm [1987] :-

→ Robin-Karp algorithm is finding or matching in the given text by using hash value.

→ pattern characters are converted into hash values that is matching with text window hash value, both the hash values are same then only we can check the characters.

Text  $T = A B C C D D A E F G$

pattern  $P = C C D$

$n$  is the size of pattern

$m$  is the size of text

$n=3$   $m=10$ ,  $d=10$

$d$  no. of character in a given text

→ Find the hash value of pattern

$$h(P) = \sum (v \times d^{n-1}) \mod 13$$

$$= (3 \times 10^{3-1} + 3 \times 10^{3-2} + 4 \times 10^{3-3}) \mod 13$$

$$= (3 \times 10^2 + 3 \times 10 + 4) \mod 13$$

$$= 334 \mod 13$$

→ Find the hash value of Text window

hash (Text) = A B C

$$\text{hash}(T) = (1 \times 10^{3-1} + 2 \times 10^{3-2} + 3 \times 10^{3-3}) \mod 13$$

$$= (10^2 + 2 \times 10 + 3) \mod 13$$

$$= (100 + 20 + 3) \mod 13$$

$$= (129) \mod 13$$

$$= 6$$

pattern hash value and text window hash value both are not matching so shift the text one step right.

→ Find the hash value of text window

hash (text) = B C C

$$= (2 \times 10^{3-1} + 3 \times 10^{3-2} + 3 \times 10^{3-3}) \mod 13$$

$$= (2 \times 10^2 + 3 \times 10 + 3) \mod 13$$

$$= (200 + 30 + 3) \mod 13$$

$$= (233) \mod 13$$

$$= (233) \mod 13 \quad 232 \mod 13 = 12$$

A - 1

B - 2

C - 3

D - 4

E - 5

F - 6

G - 7

H - 8

I - 9

J - 10

$$\begin{array}{r} 13 \overline{) 232} \quad (17) \\ 13 \\ \hline 102 \\ 91 \\ \hline 11 \end{array}$$



pattern hash value and text window value not matching then text window is moving one step to right side.

→ Text  $T = \boxed{CCD}$

pattern  $p = CCD$

Find the hash value of text window "CCD"

$$\begin{aligned} \text{hash}(T) &= (3 \times 10^{3-1} + 3 \times 10^{3-2} + 4 \times 10^{3-3}) \bmod 13 \\ &= (3 \times 10^2 + 3 \times 10 + 4) \bmod 13 \\ &= (334) \bmod 13 \end{aligned}$$

∴ pattern hash value = 9

pattern hash value and text window hash value are same then compare characters.

Text  $T = A, B, \boxed{C, C, D}, D, A, E, F, G$

$\downarrow \quad \downarrow \quad \downarrow$   
 $\boxed{C, C, D}$

Pattern is matching at index 3.

\* Knauth - Morris Pratt Algorithm [KMP]

prefix table, called as pi table ( $\pi$ )

→ prefix is same as suffix.

Text: a b a b c a b a

prefix: a, ab, aba, abab, ...

suffix: a, ba, aba, abac, ...

1	2	3	4	5	6	7	8
a	b	a	b	c	a	b	a
0	0	1	2	0	1	2	1

→ Matching the pattern in given text by using prefix table.

→ Main idea of KMP algorithm is reducing the no. of comparisons.

is there any prefix values is same as suffix values.

Ex: Text  $T = a b a b a b d$

pattern  $p = a b a b d$

sol: first prepare prefix table for the given pattern

0	1	2	3	4	5
	a	b	a	b	d
	0	0	1	2	0

→ prefix table

$$\begin{array}{r} \pi[0] = 0 \\ \pi[1] = 0 \\ \pi[2] = 0 \\ \pi[3] = 1 \\ \pi[4] = 2 \\ \pi[5] = 0 \end{array}$$

$$\pi[0] = 0$$

$$\pi[1] = 0$$

$$\pi[2] = 0$$

$$\pi[3] = 1$$

$$\pi[4] = 2$$

$$\pi[5] = 0$$





## \* Boyer Moore Algorithm:-

- Boyer Moore Algorithm is the fastest way to finding or matching the pattern in a string.
- If the characters are not matching then we can check the bad match table.
- The bad match table is deciding how many times the pattern is moving from current position to the next matching position depends on values of the bad match table.

Ex:- Text T = THIS IS A TEST  
 pattern P = TEST Find the pattern using boy moore Algorithm.

Sol:- step 1:- First create bad match table for the pattern "TEST"

letter	T	E	S	*
value	3	2	1	4

b	a	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Find the values of each and every letter in the pattern

Value = length of the pattern - index of letter - 1

$$\text{Value (T)} = 4 - 0 - 1 = 3$$

$$\text{Value (E)} = 4 - 1 - 1 = 2$$

$$\text{Value (S)} = 4 - 2 - 1 = 1$$

Step 2:- Text: T H I S I S A T E S T

pattern = T E S T Not matching

The letter s is available in bad match table

→ the value of s is = 1 means the pattern is jumping one position to the right.

Step 3:- Text: T H I S I S A T E S T

pattern = T E S T Not matching

The empty letter is not available in the bad match table then consider "\*" symbol. The value of "\*" is 4.

Step 4:- Text: T H I S ! S . A T E S T

pattern = T E S T Not matching.

→ the letter not available in the bad match table, then consider "\*" value = 4. pattern jumping 4 positions to right.



Step 5:-

Text:

T	H	I	S		I	S		A		T	E	S	T
---	---	---	---	--	---	---	--	---	--	---	---	---	---

x not matching.

T	E	S	T
---	---	---	---

→ s is available in bad match table. The value of s = 1  
text window letter

Step 6:-

Text :-

T	H	I	S		I	S		A		T	E	S	T
---	---	---	---	--	---	---	--	---	--	---	---	---	---

T	E	S	T
---	---	---	---

Ex:-

Text T:- WELCOME TO MY WORLD

pattern P:- WORLD

W	O	R	L	D	*
4	3	2	1	0	5

Step 1:-

W E L C O M E T O M Y W O R L D

W O R L D x not matching

Step 2:-

W E L C O M E T O M Y W O R L D

W O R L D x not matching

Step 3:-

W E L C O M E T O M Y W O R L D

W O R L D x not matching.

Step 4:-

W	E	L	C	O	M	E	T	O	M	Y	W	O	R	L	D
											↓	↓	↓	↓	↓
											W	O	R	L	D