

Ch: 8 String matching

↳ String matching algorithm are normally used in text processing. (done in compilation of program).

↳ String matching means finding one or more generally all the occurrences of a string in the text.

↳ These occurrence are called as pattern. Hence, sometime string matching algorithm are also called as pattern matching algorithms.

→ In this section we will discuss various string matching algorithms such as -

- ① The Naive method
- ② Rabin-Karp method
- ③ Finite Automata for string matching
- ④ KMP-Algorithm (Knuth Morris Pratt).

(1) The Naive method

↳ This is the simplest method which works using Brute force approach. (Straightforward approach).

Example:-

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
text:	a	c	m	a	n			i	k	e	s		m	a	n	g	o

	0	1	2	3	4
Pattern:	m	a	n	g	o

↳ Hence seen index 12, because it match with the pattern is found from that location in text.

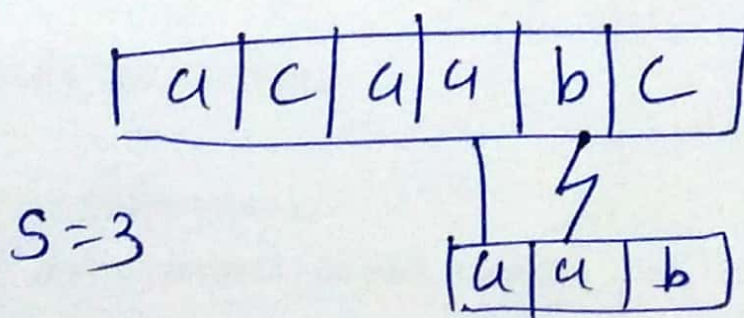
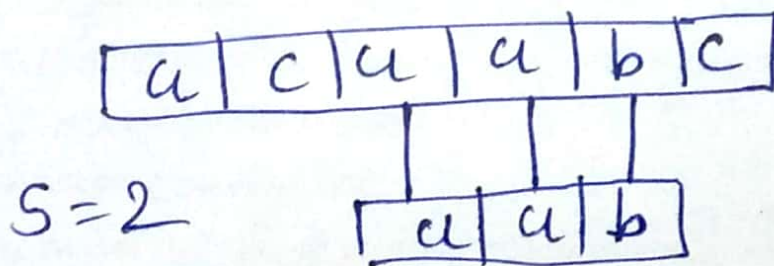
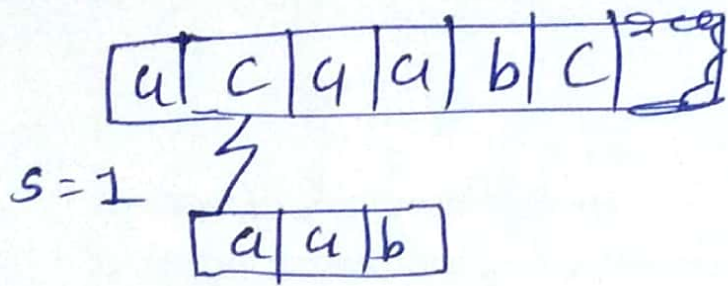
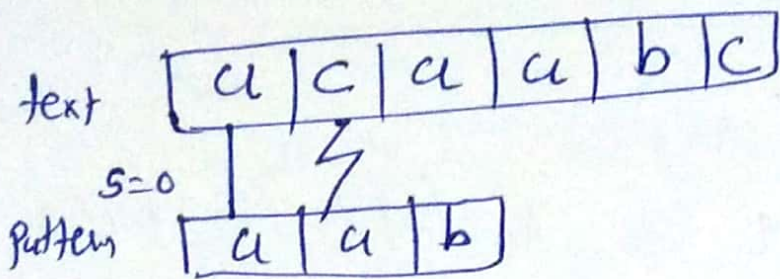
→ $O(mn)$

① The Naive String Matching algorithm

Naive-string-matcher(T, P)

1. $n = T.length$
2. $m = P.length$
3. for $s = 0$ to $n - m$
4. if $P[1 \dots m] == T[s+1 \dots s+m]$
5. Print "Pattern occurs with shift" s .

Example of Naive method



→ one occurrence of the pattern at shift $s=2$.

② The Rabin Karp Algorithm

↳ This algorithm assumes each character to be a digit in radix-d notation.

↳ $O((n-m+1)m)$

Example Using the Rabin-Karp algorithm for text $T = 3141592653589793$ find for pattern $P = 26$ with modulo $q = 11$.

↳ Text =

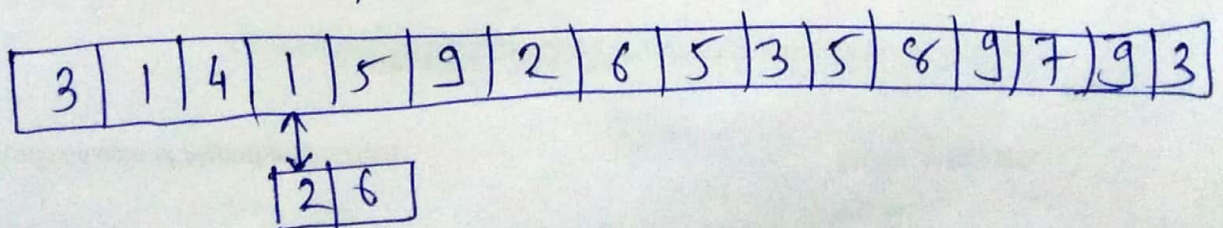
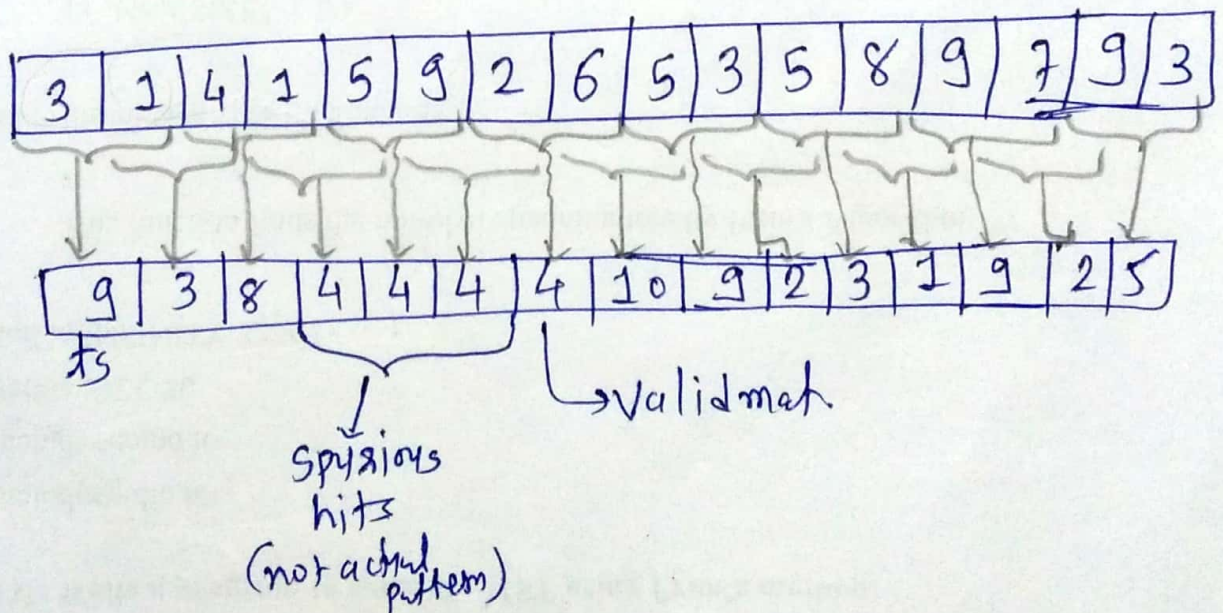
3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Pattern =

2	6
---	---

Here $26 \bmod 11 = \underline{4} \leftarrow P$

↳ We will now obtain mod 11 value for each text.



3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

↓

2	6
---	---

3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

↓

2	6
---	---

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3

↓ ↓

2	6
---	---

Match found

↳ At 7th position we find pattern matching.

Ex: Text = 2359023141526739921

Pattern = 31415

$q = 13$

② The Rabin Karp Algorithm

Rabin-Karp-Matcher (T, P, d, q)
text pattern
base of num.

1. $n = T.length$

2. $m = P.length$

3. $h = d^{m-1} \bmod q$ // hash value

4. $P = 0$

5. $t_0 = 0$

6. for $i = 1$ to m

7. $P = (dP + P[i]) \bmod q$ } modulo stored
8. $t_0 = (dt_0 + T[i]) \bmod q$ } in P & t_0

9. for $s = 0$ to $n - m$

10. if $P == t_s$ // modulus value match!

11. if $P[1..m] == T[s+1..s+m]$

12. Print "Pattern occurs with shift" s

13. if $s < n - m$

14. $t_{s+1} = (d(t_s - T[s+1]h) + T[s+m+1]) \bmod q$

↳ Running time of Algorithm $((n - m + 1)m)$

* String Matching with Finite Automata

↳ Finite Automata is a very effective tool which is used in string matching algorithms.

↳ Definition of finite Automata:-

A finite automata is a five tuples.

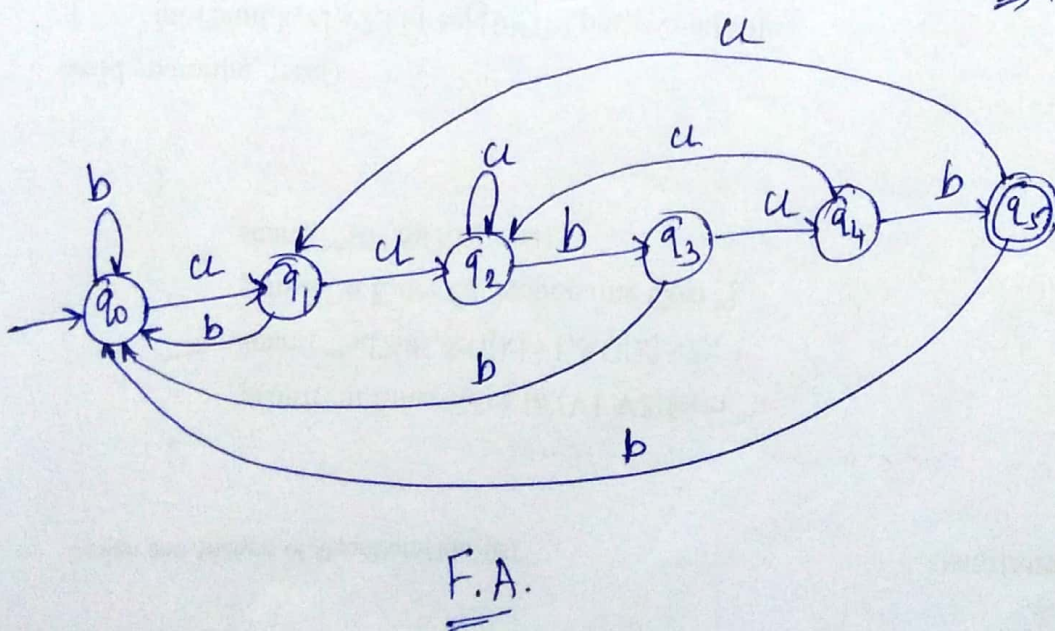
$$M = (Q, \Sigma, q_0, A, \delta)$$

Where

- Q is a set of states.
- q_0 is a start state
- A is a final or Accepting state
- Σ is a finite Input set
- δ is a function from $Q \times \Sigma \rightarrow Q$.

Example:- Text = aababababababab.
Pattern = aabab.

↳ We can construct a finite Automata for Pattern = aabab.



⇒ Transition table

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

Text=

a	a	a	b	a	b	a	a	b	a	a	b	a	b	a	a	b	
a	a	b	a	b													
	a	a	b	a	b												
		a	a	b	a	b											
			a	a	b	a	b										
				a	a	b	a	b									
					a	a	b	a	b								
						a	a	b	a	b							
							a	a	b	a	b						
								a	a	b	a	b					
									a	a	b	a	b				
										a	a	b	a	b			
											a	a	b	a	b		
												a	a	b	a	b	
													a	a	b	a	b

match
nd →

match
nd →

Finite Automaton - matcher (T, δ, m)

1. $n = T.length$

2. $q = 0$

3. for $i = 1$ to n

$q = \delta(q, T[i])$

4. $q = m$

5.

if $q == m$

print "Pattern occurs with shift " $i-m$

6.

text(1) -- text(2) --

last position

of Pattern (Accepting state)

(*) KMP (Kuth Mossis Pratt Algo)

↳ The Basic idea behind this algorithm is to build a prefix array.

↳ Some times this array is also called π array.

↳ This prefix array is built using the prefix and suffix ~~are~~ information of pattern.

↳ The overlapping prefix and suffix is used in KMP algorithm.

Example : text = budbubybabudab
 Pattern = ababada

~~Prefix~~

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
text :-	b	a	d	b	a	b	a	b	a	b	a	d	a	a	b
Pattern :-	a	b	a	b	a	d	a								

↳ first find out prefix table.

→ Consider string
"a"

prefix = ϵ
suffix = ϵ

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0						

↳ Consider string
"ab"

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0	0					

Prefix = Λ, a

Suffix = Λ, b

↳ Consider string
"aba"

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0	0	1				

Prefix = Λ, a, ab

Suffix = Λ, a, ba

↳ Consider string
"abab"

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0	0	1	2			

Prefix = $\Lambda, a, ab, abab$

Suffix = Λ, b, ab, bab

↳ Consider string
"ababa"

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0	0	1	2	3		

Prefix = $\Lambda, a, ab, aba, abab$

Suffix = $\Lambda, a, ba, aba, babab$

↳ Consider string:
"ababad"

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0	0	1	2	3	0	

Prefix = $\Lambda, a, ab, aba, abab, ababa$

Suffix = $\Lambda, d, ad, bad, abad, babad$

↳ Consider string "ababad"

0	1	2	3	4	5	6
a	b	a	b	a	d	a
0	0	1	2	3	0	1

Prefix = $\Lambda, a, ab, aba, abab, ababa, ababad$

Suffix = $\Lambda, a, da, bada, abada, babada$

We will Compute the Prefix table for the Pattern "ababada".

	1	2	3	4	5	6	7
$P[i]$	0	1	2	3	4	5	6
	a	b	a	b	a	d	a
$\pi[i]$	0	0	1	2	3	0	1

→ Compute $text[0]$ with $Pattern[0]$

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Text	b	a	d	b	a	b	a	b	a	b	a	d	a	a	b

Pattern	a	b	a	b	a	d	a
	0	1	2	3	4	5	6

→ Here b & a, as it is not matching.

→ Compute $text[1]$ with $Pattern[0]$ {Prefix-table[0] = 0}

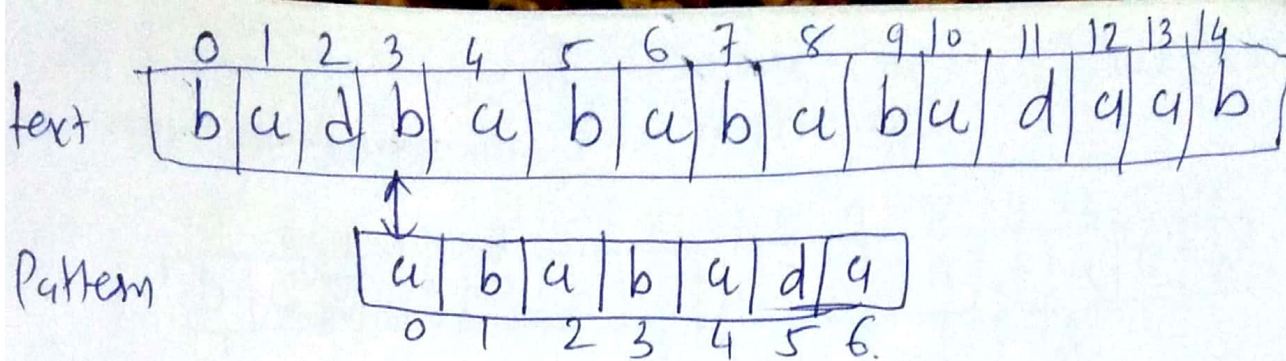
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Text	b	a	d	b	a	b	a	b	a	b	a	d	a	a	b

Pattern	a	b	a	b	a	d	a
	0	1	2	3	4	5	6

→ $text[1]$ is matching with $Pattern[0]$,
We will now compare $text[2]$ with $Pattern[1]$

→ But $Text[2]$ is not match with $Pattern[1]$.

So, Back track on Pattern and Compare $Pattern[0]$ with $text[2]$ because we consult $Prefix-table[0] = 0$, So, Consider $Pattern[0]$ is Compared with $Text[2]$.

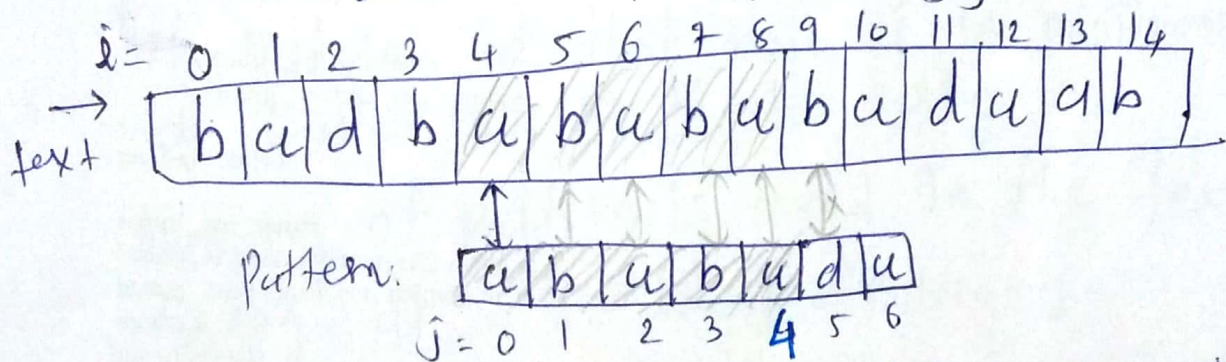


↳ Here text[3] is not matching with Pattern[0].

↳ We will then ask Prefix table[0] for the location of Pattern.

Prefix-table[0] = 0.

So, we will compare pattern[0] with text[4].



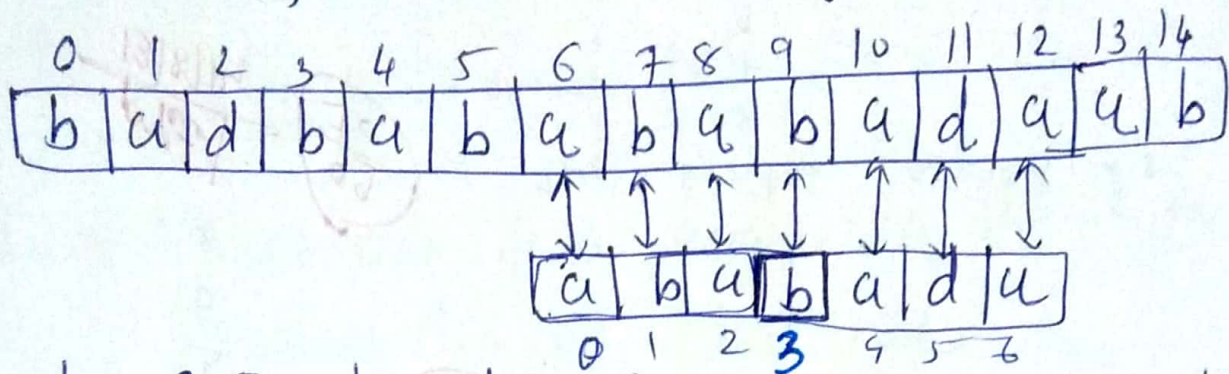
- text[4] matches with Pattern[0], increment i & j
- text[5] matches with Pattern[1], increment i & j
- text[6] matches with Pattern[2], increment i & j
- text[7] matches with Pattern[3], increment i & j
- text[8] matches with Pattern[4], increment i & j

But text[9] is not matching with Pattern[5].

↳ Hence we must backtrack on Pattern, that means j will be positioning on location 4. (Pattern 5 to decrease)

→ Consider Prefix-table[4] = 3, first three (K=3) characters matched that indicates, compare pattern[3] with current i position, means text[9]

→ Hence we will compare text[9] with Pattern[3], which is matching.



- text[10] with pattern[4], matching \therefore increment i and j
- text[11] with pattern[5], matching \therefore increment i and j
- text[12] with pattern[6], matching \therefore increment i and j.

→ Thus, we have reached on the last character of Pattern, at that time i is positioned at location 12 in the text array.

matching of pattern is found in the text array at = $i - \text{length of pattern} + 1$

$$= 12 - 7 + 1$$

$$= 6$$

→ This required pattern matches at location 6.

Ex. find prefix table for "abudab"

a	b	a	d	a	b
0	0	1	0	1	2

⊛ The Knuth-Morris-Pratt Algorithm

KMP_Matcher(T, P)

1. $n = T.length$
2. $m = P.length$
3. $\pi = \text{Compute-Prefix-Function}(P)$
4. $q = 0$ // number of characters matched
5. for $i = 1$ to n // scan the text from left to right
6. while $q > 0$ and $P[q+1] \neq T[i]$
7. $q = \pi[q]$ // next character does not match
8. if $P[q+1] == T[i]$
9. $q = q + 1$ // next character matches
10. if $q == m$ // is all of P matched?
11. Print "Pattern occurs with shift" $i - m$
12. $q = \pi[q]$ // look for the next match

Compute - prefix - function (P)

1. $m = P.length$
2. let $\pi[1 \dots m]$ be a new array
3. $\pi[1] = 0$
4. $k = 0$
5. for $q = 2$ to m
6. while $k > 0$ and $P[k+1] \neq P[q]$
7. $k = \pi[k]$
8. if $P[k+1] == P[q]$
9. $k = k + 1$
10. $\pi[q] = k$
11. return π