



SNS COLLEGE OF ENGINEERING

Coimbatore-107



COURSE NAME: ANALYSIS OF ALGORITHM

II YEAR/ IV SEMESTER

UNIT – IV

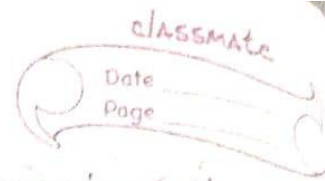
STRING MATCHING ALGORITHM

Topic

Knuth Morris Pratt Algorithm



Unit 4 → Knuth Morris Pratt Algorithm



⇒ The basic idea is to match the text with the pattern.

⇒ It is advantageous than Naïve string.

Disadvantages of Naïve:

⇒ It checks each of pattern and text even if some part has already been matched.

⇒ If mismatch occurs, it shifts the pattern by 1 and starts over.

⇒ Repeated comparisons are done

Step 1:

⇒ Build LPS [Longest prefix suffix Array], before performing match.

⇒ When Mismatch occurs, kmp doesn't start over.

⇒ It uses LPS to skip characters & avoid repeating comparisons.

Example:

Text = "A A A B A"
pattern = "A A B A"
0 1 2 3



Rules:
(1) If $\text{text}[i] = \text{pat}[i]$ (match occurs) +
 \rightarrow increment i ($i++$; $j++$)
(2) If $\text{text}[i] \neq \text{pat}[i]$ (mismatch occurs)
 \rightarrow reset $j = \text{lps}[j-1]$
 \rightarrow then increment i ($i++$).

Example:
Text = A B A A B C A B A D E
Pattern = A B C A B
 0 1 2 3 4

Goal:
find index of pattern in
Text.

Step 1:
Find LPS Array (Longest Prefix
Suffix Array)

Rules: Note (✓)

1. All characters unique
eg: A B C D E

LPS Array = [0 0 0 0 0]

2. All characters same

eg: A A A A

LPS Array = [0 1 2 3]



Step 1:

LPS Array for pattern in Example

PATTERNS - ^{0 1 2 3 4} → Index
A B C A B(i). Split the given character
as A, AB, ABC, ABCA, ABCAB(ii) Calculate the prefix &
suffix for each term

(iii) Update in LPS Array.

(matched character/term length)

(iv) Always don't choose last matched longest Prefix & suffix

LPS Array

Prefix & Suffix for
Each term

1. A → 0 (Initial Always 0) ①. A B
(Left to Right) (Right to Left)

2. B → 0 Prefix Suffix

3. C → 0 1. A ≠ 1. B X

4. A → 1 2. AB X 2. AB X

5. B → 2 Note! A ≠ B; So B → 0 (length of char)

(AB is matched in prefix & suffix of Term AB. But we should not choose last matched longest Prefix & suffix.) Update B → 0 in LPS Array for (mismatch)

②. Prefix & Suffix of "ABC":

Prefix ← Compare → Suffix

1. A ≠ C X

2. AB ≠ BC X

3. ~~ABC~~ Should not consider ~~ABC~~



classmate
Date _____
Page _____

∴ $A \neq C$; Set $C \rightarrow 0$ in LPS Array
(because of mismatch)

③. Prefix & Suffix of "ABCA"
Prefix ← Compare → Suffix

- | | | | |
|--------------------|---|-----|---------|
| 1. A | = | A | ✓ match |
| 2. AB | ≠ | CA | X |
| 3. ABC | ≠ | BCA | X |
| 4. ABCA | | | |
- Don't Consider ~~ABCA~~

∴ $A == A$ (matches) put $A \rightarrow 1$ in LPS Array
(length of char)

④. Prefix & Suffix of ABCAB
Prefix ← Compare → Suffix

- | | | | |
|---------------------|---|------|-----------|
| 1. A | ≠ | B | X |
| 2. AB | = | AB | (match) ✓ |
| 3. ABC | ≠ | CAB | X |
| 4. ABCA | ≠ | BCAB | X |
| 5. ABCAB | | | |
- Don't Consider ~~ABCAB~~

∴ $AB == AB$ (matches) ; put $B \rightarrow 2$ in LPS Array.
(i.e. length of AB (character))

LPS Array					
Index →	0	1	2	3	4
Pattern →	A	B	C	A	B
LPS Array →	0	0	0	1	2



CLASSTIME
Date _____
Page _____

i
 \uparrow
0 1 2 3 4 5 6 7 8 9 10 \rightarrow Index
Text = A B A A B C A B A D E
pattern = A B C A B

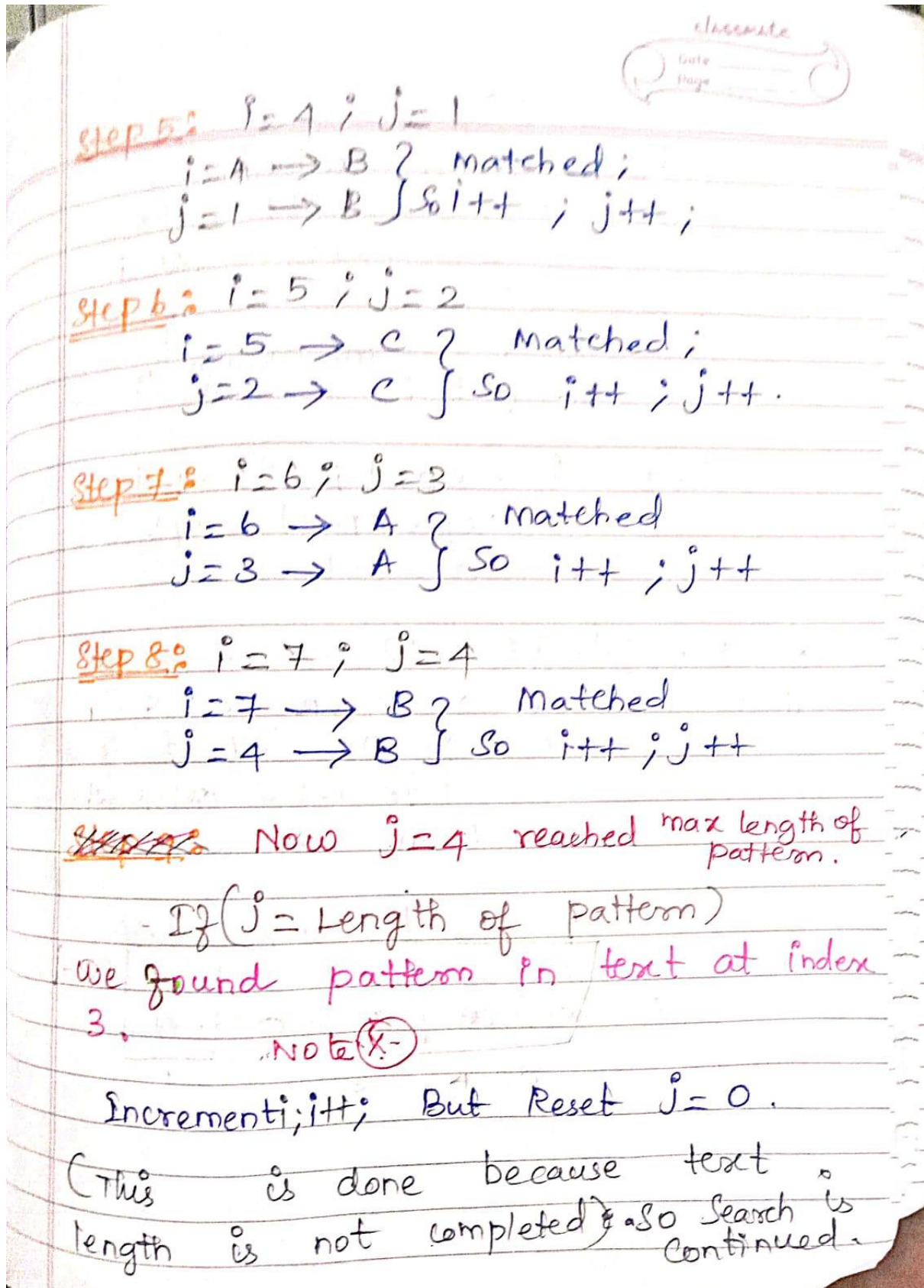
0 1 2 3 4 \rightarrow Index j
 \downarrow
 j

Step 1:
 $i=0 ; j=0 ;$
 $i=0 \rightarrow A$
 $j=0 \rightarrow A$ } Both character matches.
So $i++ ; j++ ;$

Step 2: $i=1 ; j=1 ;$
 $i=1 \rightarrow B$
 $j=1 \rightarrow B$ } matched ;
So $i++ ; j++ ;$

Step 3: $i=2 ; j=2 ;$
 $i=2 \rightarrow A$
 $j=2 \rightarrow C$ } Mismatched ;
So Reset $j = \text{lps}[j-1]$
 $j = \text{lps}[2-1]$
 $j = \text{lps}[1]$
 $j = 0 ; i++ ;$
 $\therefore [\text{lps}[i]=0 \text{ from Lps Array}]$

Step 4: $i=3 ; j=0 ;$
 $i=3 \rightarrow A$
 $j=0 \rightarrow A$ } matched
So $i++ ; j++ ;$





Step 9: $j=0; i=8$

$i=8 \rightarrow A$ } matched
 $j=0 \rightarrow A$ } so $i++$, $j++$.

Step 10: $j=1; i=9$.

$i=9 \rightarrow D$ } mismatched
 $j=1 \rightarrow B$ }

Reset $j = \text{lps}[i-1]$ $\because \text{lps}[0]=0$
 $j = \text{lps}[1-1] = \text{lps}[0]$
 $\therefore j=0; i++$

Step 11: $j=0; i=10$.

Since $i=10$ ^{$\rightarrow B$} End of text.
 $j=0 \rightarrow A$ & mismatch.

← Stop here →

Now we found pattern "ABCAB"
is matched with text at index
3 & 6 in only position.



Algorithm LpsArr(pattern, m, lps[])

```
{
    while (i < m)
    {
        if (pattern[i] == pattern[len])
        {
            len++;
            lps[i] = len;
            i++;
        }
        else {
            if (len != 0)
            {
                len = lps[len-1];
            }
            else {
                lps[i] = 0;
                i++;
            }
        }
    }
}
```

Algorithm kmp(pattern, text)

```
{
    m = length of pattern;
    N = length of text;
    while (i < N)
    {
        if (pattern[i] == text[i])
        {
            i++;
            j++;
        }
    }
}
```



```
if (j == m)
```

```
{
```

```
    // pattern found at index (i-j),
```

```
    j = lps[j-1];
```

```
} else if (i < N && pattern[j] != text[i];
```

```
{
```

```
    if (j != 0)
```

```
        j = lps[j-1];
```

```
    else
```

```
        i++;
```

```
} }
```

Time Complexity:

1) LPS Array Construction = $O(m)$

2) Pattern Search = $O(n)$

Total = $O(n+m)$

* $n \rightarrow$ text length * $m \rightarrow$ pattern length

Efficient:

* No Backtracking in text * Avoids

Rechecking characters * LPS tells

where to resume in pattern during mismatch

* Space Complexity:

* Much faster than brute

force $O(n \times m)$ in worst case

* $O(m)$ for LPS Array