Given a text tat[0,1,..., N-1] and a pattern pat[0,1,..., M-1], write an algorithm search that prints all occurrences of pat[] ~~and~~ in tat[]. You may assume text N>M.

Example:   tat[] = "AABAACAADAABAABA"

pat[] = "AABA"

— Pattern found at ~~index~~ indices 0, 9 and 12

$i = 0 + \frac{\alpha}{2} \cdot 3$

$j = 0 + M-2?$

$N=4, M=2$

$4-2=2$

• Naive Algorithm (exhaustive search)

Algorithm-Naive (tat[], pat[])

1. for i = 0 to N-M
2. if pat[0,1,...M-1] == tat[i, i+1, ... i+M-1] ~~then~~
        then print "pattern found"
3.
4. endif
5. end for

complexity: $O(m(n-m))$

• KMP (Knuth, Morris, Patterson) Algorithm.

Matching overview
    tat = "AAAAABAAABA"
    pat = "AAAA"
we compare first window of tat and pat
tat = AAAAABAAABA
pat = AAAA
We find a match (same as naive approach)
In the next step, we compare next window
                        of tat and pat
tat = AAAAABAAABA
pat = AAAA (Pattern shifted one position)

This is where KMP does optimization over naive. In this 2nd window, we only compare fourth A of pattern ~~with~~ with fourth character of current window of text to decide whether current window matches or not. Since we know, first 3 characters will anyway match, we skipped matching first 3 characters

# Preprocessing :

$lps[i] = \sqrt{}$ <u>size of</u> the longest proper prefix of $pat[0,1,...,i]$ which is also a suffix of $pat[0,1,...,i]$.

- "ABC" → <u>proper prefix</u> are "A", "AB"
  → suffix are "C", "BC", "ABC"

$pat[] = \{A, A, B, A\}$

$pat = "AAAA"$ , $lps[] = \{0, 1, 2, 3\}$

$pat = "ABCDEF"$ , $lps[] = \{0, 0, 0, 0, 0\}$

$pat = "AABAACAABAA"$ , $lps[] = \{0, 1, 0, 1, 2, 0, 1, 2, 3, 4, 5\}$

$pat = "AAACAAAAAC"$ , $lps[] = \{0, 1, 2, 0, 1, 2, 3, 3, 3, 4\}$

# Algorithm: <u>lps calculation</u>

1. $len \leftarrow 0$ , $lps[0] \leftarrow 0$

   while($i <$ size of pat string)

2. If $pat[len]$ and $pat[i]$ match, then ~~also increment~~ ~~len by 1 and also then~~ $len = len+1$

3. $len = len + 1$
   $lps[i] = len$ , $i = i+1$

4. 

5. If $pat[len]$ and $pat[i]$ do not match and $len > 0$, then
   $len = lps[len-1]$
   If $len = 0$, then $i = i+1$

6.

7.

for $i = 1$ to $M-1$ do

# Example:

$pat = "AAACAAAA"$   $lps[] = \{0, 1, 2, 0, 1, 2, 3, 3\}$

1. We start the comparison of pat[j] with j=0
with characters of the current window of that

2. We keep matching txt[i] and

2.

1.  $i \leftarrow 0, \; j \leftarrow 0$

2.  If txt[i] and pat[j] match, term  $\parallel$ for current
3.       $i = i + 1, \; j = j + 1$                            window

4.  If txt[i] and pat[j] do not match

   - we know that pat[0,1,...,j-1] match with
                                          txt[i-j,...i-1]

   - We also know lps[j-1] is the count of
     characters of pat[0,1,...j-1] that are
     both proper prefix and suffix.

   - From the above two points, we can
   conclude that we do not need to match
   these lps[j-1] character with txt[i-j,...i-1]

Time: $O(N)$.