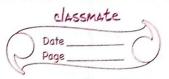
classmate

-	A = A + A + A + A + A + A + A + A + A +														
ANI	Given String: abadab														
	posses to company but sages whigher														
da	(i) Initially put o o 1 2 3 4 5 in oth location a b a d a b														
0 39	in oth Jocation abab														
	august who Plot I be to														
	(ii) String: ab														
	(ii) String: ab Prefix: E, a Suffix: E, b P 0 0														
	5affix: E, b P 10101 1														
	2112														
	did d di camitali														
(Prefix: E, a, ba P, 0 0 1 Suffix: E, a, ba P, 0 0 1														
adot n	Prefix: E, a, ab a d a b Suffix: E a ba P 0 0 1														
	Suffix: E a ba P 0 0 1 1 1														
	arthog by the detent to pattern														
	(iv) String: abad														
	Prefix: E, a, ab, aba abadab														
	Suffix: E, d, ad, bad F 0101110														
	The still to the state of the s														
1000	the man beauty of the control of the														
	Prefix: E, a) da, ada, bada 00101"														
	Prefix: Elajab, aba, abad ab ad ab														
Z I X	uffix: E, a) da, ada, bada 1010111011														
	A STATURE OF THE STATE OF THE S														
	refix: E, a, ab, aba, abad, abada ob a da b uffix: E, b, ab, dab, adab, badab o o 1 0 1 2														
1	efix: E, a, ab, aba, abad, abada 00 1012														
2	uffix: E, b, ab, dab, adab, badab [0]0]1]2]														
1-0168	anife patriote to sail - compat to sail to sail to														



	It is defined as the size of the largest
	It is defined as the size of the largest prefix of the P[0j-1] that is also
	a suffix of P[1j]
0.1	1011111111110111223121111111111111
0	LPS[j] = Longest Proper Profix of Pattern
1 2	Pro j] which also a Suffix of
1	Pattern P[oj]
	V [01721107] =
01 _	It also indicates how much of last companison
7 4	can be reused, if it fails.
-	It enables avoiding bart-tracking on Text Strings
	Algorithm Calculate LPS (Pat[], m)
	\$ 11901111111 (A 11 6/10)
	1. K=0 // Length of previous longest pattern
	2. p[0] = 0
	3. i = 1 // Pointer on Pattern
	4. while (i <m)< td=""></m)<>
	the day a mine to a day
	5. if (Pat[i] == Pat[k]) // If match found
	case ()
	6. K++
	7. P[i]=k
	8. i°++
	1, } , , , , , , , , , , , , , , , , , ,
	g. else
1	10. $f(k)=0$ Mismatch & $k = 0$ case 2
	k = P(k-1)
	12. else // Mismatch & K=0
	13. 1 p[i]=0
	19. 3++
	3 1

and the second s	Example?
and a supplied	Example:
	Pattern = aaacaaaaaaa
	0 1 2 3 4 5 6 7 8 9 10
	Pat a a a a a a a a a
	0 0 1 2 0 1 2 3 3 3 4 5
	Ti NIC NICT
=	Proressing
	1:123334567788910
	K: 0, 1 2 1 0 0 1 2 3 2 3 2 3 4 5
	Total dead adjusted to the state of the stat
	nse 1 1 2 2 3 1 1 1 2 1 2 1 1 1
	/ A CITE A NICE Wheling on although (A I C=
Etica i	ernal averyand to attend to a previous lenner
	0 = [0] 0 = 0
	Example 2:
	- Construction to the
	abadbae abf
P	0 0 0 0 0 1 0 1 2 0
(1) 201	
-	8 + + 5
	Example 3:
100.10	+ + 1 6
and section of the section	aaabaacd
	8-012301200
	A STREET OF STREET OF STREET OF STREET
Contract of the second section	



Ē	Search Alanxidhan andriada
	Search Algorithm - In this LPS is used to decide next
	character to be matched
	- Logic is - Do not match character
	that are definitely going to match.
->	How to find number of characters to be
cient.	SKIDDER
	- Initially set pointer i at initial location
	- Initially set pointer i at initial location of Text and set pointer j at initial location
asteld.	of pattern.
	- treep incrementing i and j until Text[i]
	and Pattern[j] matches
	- when there is mismatch at index j of
	Pattern it implies
mod c	@ Pattern [0 j-1] = Text[i-ji-1]
	6 LPS [j-1) is count of characters
	of Pattern [0j-1] which are
	both Prefix and Suffix.
	1
	> we do not need to match LPS[j-1]
	characters with Text[i-ji-1] as
- V	they will anyway match.
4017	THE COLUMN TO THE PARTY OF THE



```
A)gorithm
  KMP-Matcher (Pat[], Text[])
      m = pat. length
     n = Text. length
  3. Calculate_LPS (Pat[],m)
  4. i=0 // Pointer on Text
  5. j=0 // Pointer on Pattern
    while (i<n)
        if (Pat[j]==Text[i]) // Case Do Match
  8 .
  9.
       if (j==m) // Case 2: Pattern Found
  of.
          print ("Pattern Found at", i-j)
 ))-
         j=P[j-1]
       if (i<n and Pat[j]!=Text[i])
  12
            if (j!=0) // Case3: Mismatch d
j=p[j-1]
 14
  15
                  1 Case A. Mismatch & j=0
 16
 17
```

6	u e s	noite	1:		F	aft	ern	: (ao	a							
			_		rex (:	0	0	2 (h on	0	01 (3 (2	
3	0)0	fion	•			and the second second second					Name of the latest						
					ngth												
				Ler	gth	0~	-	Tex	t		ים ב	- 11		-			
	Co	m 191	Jd P	L	PS		for	r	pat	t er	0						
		,,,	710			,											
					0	*	, 2										
		Po	rot	(2 0	a	a										
			P		0)	†	2										
,															100		
(0	1	1	2	3	4	1	5	6	7	8	1	9	10 a	
1			<u>a</u>	a		a]	C	<u>a</u>) 0	۲]	a	a	C		α	aj	
)		1					-								
											-						
	5	0	1	2	3	2	2	4	5	6	7	2	2	Q	a	10	
	9																
	1				2						2	2	1			-	
Co	se	1	1	1	3	3	4	1	1	1	1	3	3 4	4	1	2	
				2						= 2	2						
				Pattern)													
Pattern Found:0								Found:4									
1	3-320								7-3-4								
	Pattern										2						
											und						
			-							8	-3=	5					
	No.																