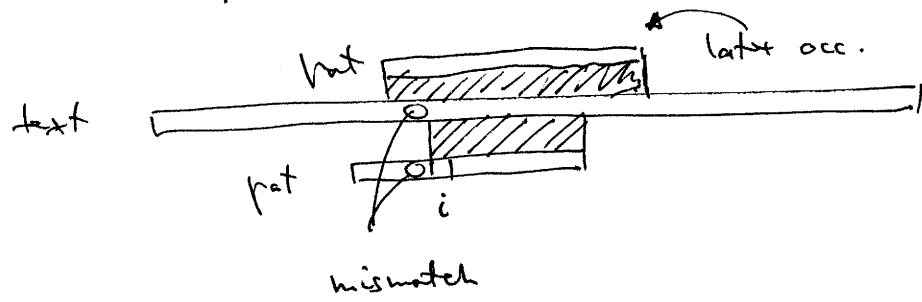
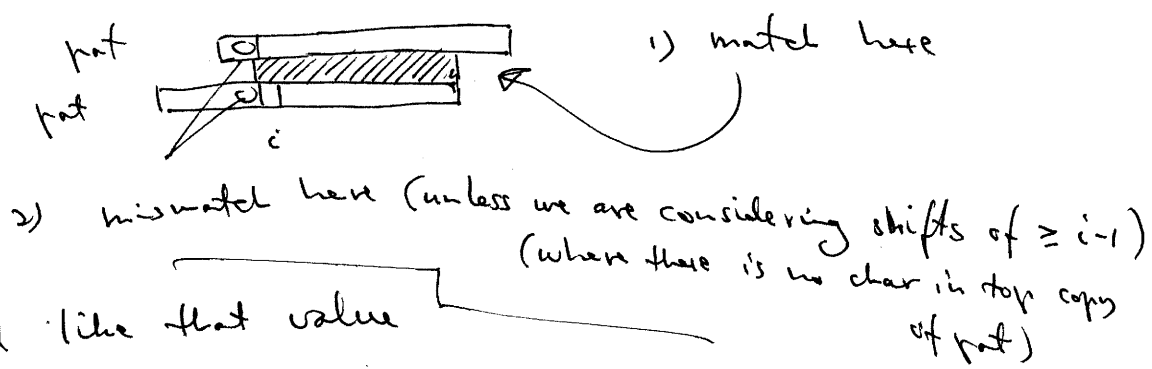


# Preprocessing for BM

Recall that after a mismatch, we for the def. of BM-shift consider a next possible occurrence:



This leads us to consider least possible shift s.t.:



We would like that value

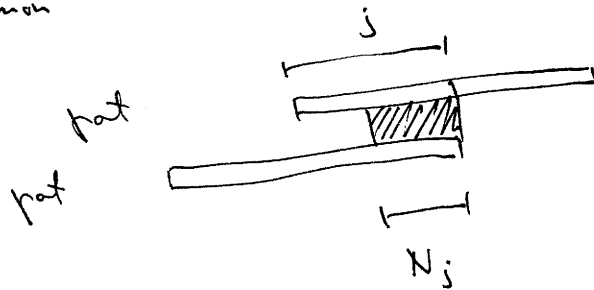
for all  $i \geq 2$ , [BM-shift(i) is then that value.]  
(and  $i \leq m+1$ )

For  $i = 1$ , we have an occ., and consider ~~larger~~ least possible shift  $\geq 1$  such that 1) holds [2) will be in "unless" case].  
[BM-shift(1) is then that value]

Note:  $i$  here is one larger than there corresponding value (called  $j$ ) in textbook BM alg. on page 28 (as it is there position of ~~the~~ mismatch, not last match).  
Our exposition here corresponds to handout from Gusfield book

(2)

First, define for  $j = 1, 2, \dots, |pat|$   $N_j$  as largest <sub>common</sub> suffix of  $pat[1..j]$  and  $pat$ .

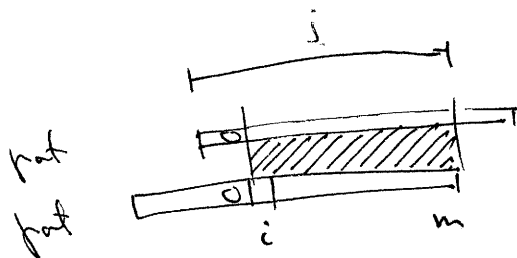


Observe that (by left-right mirror-flip of figure) the  $N_j$  values can (be found from) the prefix table values  $Z_i$  for the reverse string (and for  $i = \overset{m-j+1}{\cancel{m-j+1}}$ ).

Hence, the  $N_j$  values can be found in  $O(|pat|)$  time.

for  $i \geq 2$  and  $i \leq m+1$

Let  $L'(i)$  be largest  $j \leq m$  s.t.  $pat[1..j]$  has  $pat[i..m]$  as suffix and  $j - (m - i + 1) \geq 1$  and  $pat[i-1] \neq pat[j - (m - i + 1)]$

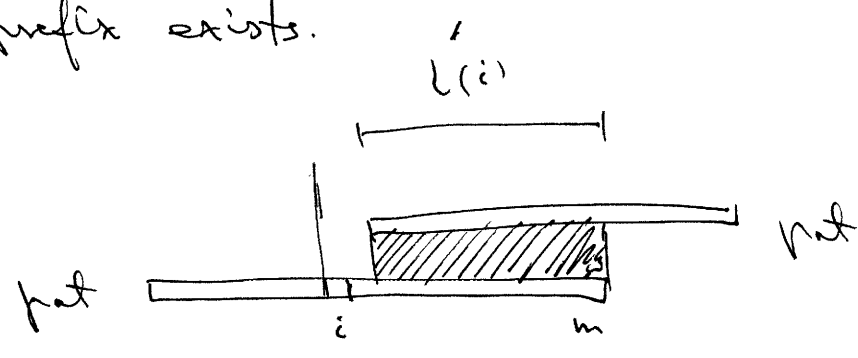


( $m = |pat|$ )

Let  $L'(i) = 0$  if no such  $j$  exists. Let  $L(1) = 0$ .

I.e., if  $L'(i) > 0$ , this gives  $BM\_Shift(i)$  for  $i \geq 2$ .  
 (if 2) is not in unless part) (as  $m - L'(i)$ ).

Let  $L'(i)$  = length of largest suffix of  $pat[i..m]$  being a prefix of  $pat$ . Let  $L'(i) = 0$  if no such prefix exists.

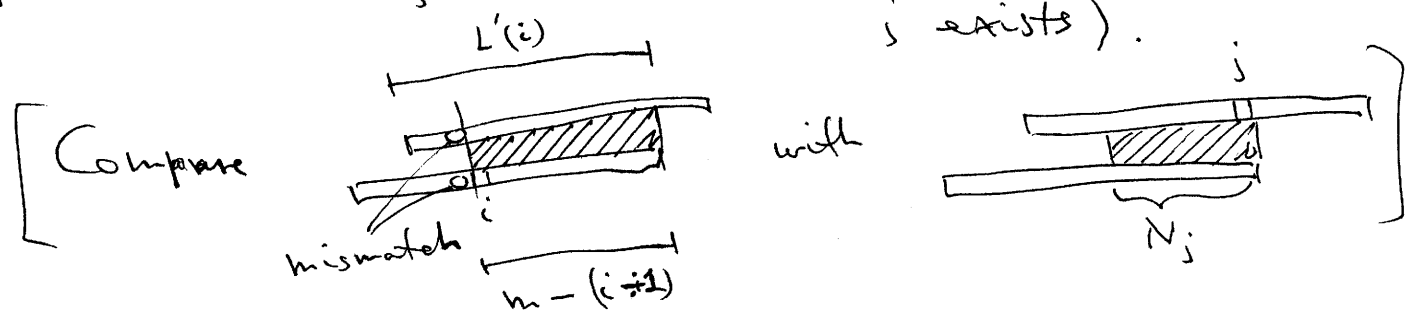


I.e., this gives  $BM\_Shift(i)$  [as  $m - L'(i)$ ] if  $L'(i) = 0$  [if our "shift search" enters the "unless" part of 2)].

Note:  $L'(1) = m$  always. <sup>by above def.</sup> This corresponds to occ., and a shift of  $\geq 1$  is needed. This shift is  $L'(2)$ . So we should set  $L'(1) = L'(2)$ .

Observe:  $L'(i) =$  largest  $j < m$  where

$j \geq m - i + 2$  AND  $N_j = m - i + 1$  (or  $= 0$  if no such  $j$  exists).



$$N_j = m - i + 1 \quad \left( \text{AND } j \geq m - i + 2 \right)$$

$$i = m - N_j + 1 \quad (*)$$

Note: We can leave out the "AND" (and hence the "if" in code) if we allow  $L'$  to cover 1st "unless" case of  $\geq$  [I think]

Given the  $N_j$  table (prefix table), the non-zero  $L'(i)$  values [those where a  $j$  exists] will be those hit by  $*$ ) when  $j$  goes from 1 to  $m-1$ . For each of these, we should keep the largest  $j$  as  $L'(i)$ .

So the alg.

```

for i = 1 to m+1:
  L(i) = 0
for j = 1 to m-1:
  i = m - N_j + 1
  if j >= m - i + 2
    L(i) = j

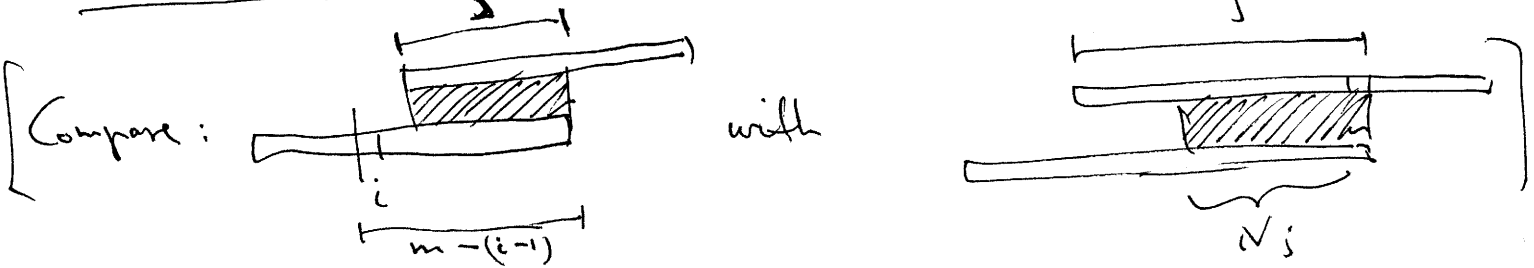
```

Note: For these  $j$  ( $\leq m-1$ ), the possible values of  $N_j$  are  $0, 1, \dots, m-1$ , so possible values for  $i$  are  $2, 3, \dots, m+1$ . Which is what we wanted.

is correct. Clearly  $O(m)$ .

(or zero if no such  $j$  exists)

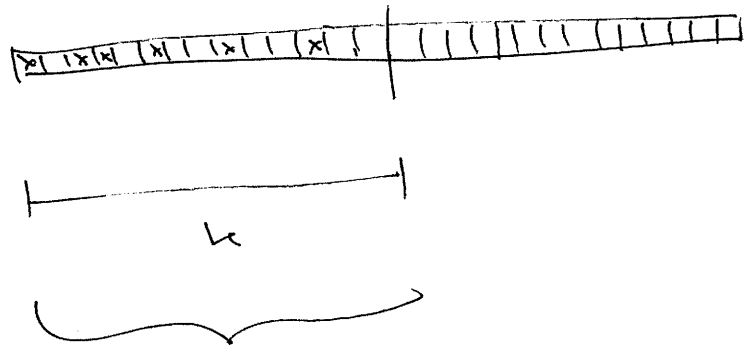
Observe:  $L'(i) =$  "largest  $j \leq m - i + 1$  s.t.  $N_j = j$ "



Let  $k = m - i + 1 \iff i = m - k + 1$

Let  $A(j) = "N_j == j"$

Truth-table for  $A(j)$ :



$$L'(i) = \max_{\substack{j \leq k \\ \uparrow \\ m-i+1}} \{N_j = j\}$$

So alg. for  $L'(i)$ :

max = 0  
 for k = 1 to m-1  
   if  $N_k = k$   
     max = k  
    $L'(m-k+1) = \text{max}$   
     i

$L'(1) = L'(2)$   
 $L(m+1) = 0$

Recall:  
 $L'(1) = m$  by above def. of  $L'$  [see page 3],  
 so we should set  
 $L'(1) = L'(2)$ , as  
 we need a shift  $\geq 1$ .

If  $i = m+1$ , there are no possible unless-case matches so set  $L(m+1) = 0$

Clearly  $O(m)$ .

6

Recall: If  $L'(i) > 0$ , we should use  $L'(i)$ .  
Else, we should use  $L(i)$ .

(By our defs.  $L'(1) = 0$  and  $L'(1) = L'(2)$ ,  
this also covers the case  $i = 1$  [i.e., an  
occurrence].)

In other words, we should use

$$\max \{ L'(i), L(i) \}$$

(if  $L'(i) > 0$ , then  $L'(i) \geq L(i)$ ,  
by the way we define our shifts (for a  
given  $i$ , the  $j$ 's tried for  $L'$  are  
bigger than those for  $L$  (the 'unless' case  
in  $\Rightarrow$ )))

So finally, create BM-shift by:

For  $i = 1$  to  $m+1$ .

$$\text{BM-shift}(i) = m - \max \{ L'(i), L(i) \}$$

Total time is  $O(m)$ .