

Knuth-Morris-Pratt algorithm

The algorithm we present now is an improvement due to Knuth of the Morris-Pratt algorithm, based on avoiding situations which lead to certain letter comparison failures.

	$j-i$		$j-i+p$		$j-1$	j
$t :$...	BORDER($x[1 \dots i-1]$)	...	BORDER($x[1 \dots i-1]$)	b	
	1		p		$i-1$	i
$x :$	BORDER($x[1 \dots i-1]$)	...	BORDER($x[1 \dots i-1]$)	a		
			1	$\beta_x(i-1)$	$s_x(i)$	
			slided $x :$	BORDER($x[1 \dots i-1]$)	a	

The observation is that if $a = a$ then the sliding would lead to a comparison failure. So let us enforce that a sliding to compare $x[s_x(i)]$ to $t[j]$ is done if and only if $x[s_x(i)] \neq x[i]$.

Knuth-Morris-Pratt/Maximum disjoint borders

If $x[s_x(i)] = x[i]$, what should we do?

Let us note $y = x[1 \dots i - 1]$. We should consider successively $\text{BORDER}^2(y) \cdot a$, $\text{BORDER}^3(y) \cdot a$ etc., until we find a k such that $\text{BORDER}^k(y) \cdot a \not\leq x$ or $\text{BORDER}^k(y) = \varepsilon$. Such a border $\text{BORDER}^k(y)$ is called a **maximum disjoint border** of y in x .

If it is non-empty, we slide the word so that we compare $x[s_x^k(i)]$ and $t[j]$.

This disjointedness constraint can be precomputed on the pattern x *alone* before comparing it to the text t : all we have to do is to change the failure function s_x of the algorithm of Morris and Pratt and provide a new one. *The search algorithm itself does not change.*

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

The disjointedness implicitly entails that it is a relative concept: we consider the disjoint border of a *proper prefix* of another word. Thus, the letter which follows the prefix is used as a constraint, i.e., it must not follow the disjoint border.

Let $ya \leq x$. Let us note $\overline{\text{BORDER}}_x(y)$ the maximum disjoint border of y in x . In this case, the letter a is used to constrain the definition of the disjoint border, as we must have $\overline{\text{BORDER}}_x(y) \cdot a \not\leq y$.

What if $y = x$, then? In this case, there is no right-context, like the letter a above, to constrain the maximum disjoint border. In this case, we still can take the maximum border, i.e., $\overline{\text{BORDER}}_y(y) = \text{BORDER}(y)$.

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

More precisely, if the maximum border of y is not followed by a , then it is also the maximum disjoint border we are looking for. In other words:

$$\overline{\text{BORDER}}_x(y) = \text{BORDER}(y) \quad \text{if } ya \leq x \text{ and } \text{BORDER}(y) \cdot a \not\leq y$$

If the maximum border of y is followed by a we must take the maximum disjoint border of the maximum border:

$$\overline{\text{BORDER}}_x(y) = \overline{\text{BORDER}}_x(\text{BORDER}(y)) \quad \text{if } ya \leq x \text{ and } \text{BORDER}(y) \cdot a \leq y$$

By extension, if $x = y$, we take the maximum border:

$$\overline{\text{BORDER}}_y(y) = \text{BORDER}(y)$$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

In summary, assuming $ya \leq x$ **and** $y \neq \varepsilon$

$$\overline{\text{BORDER}}_y(y) = \text{BORDER}(y)$$

$$\overline{\text{BORDER}}_x(y) = \begin{cases} \overline{\text{BORDER}}_x(\text{BORDER}(y)) & \text{if } \text{BORDER}(y) \cdot a \leq y \\ \text{BORDER}(y) & \text{otherwise} \end{cases}$$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

By taking the length of each side of the equations,

$$|\overline{\text{BORDER}}_x(y)| = \begin{cases} |\text{BORDER}(y)| & \text{if } x = y \text{ or } \text{BORDER}(y) \cdot a \not\leq y \\ |\overline{\text{BORDER}}_x(\text{BORDER}(y))| & \text{otherwise} \end{cases}$$

Let $\gamma_x(i)$ be the length of the disjoint maximum border of $x[1 \dots i]$:

$$\gamma_x(i) = |\overline{\text{BORDER}}_x(x[1 \dots i])| \quad 1 \leq i \leq |x|$$

or, equivalently,

$$\gamma_x(|y|) = |\overline{\text{BORDER}}_x(y)| \quad y \leq x$$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

Therefore

$$\gamma_x(|y|) = \begin{cases} \beta_x(|y|) & \text{if } x = y \text{ or } \text{BORDER}(y) \cdot a \not\leq y \\ \gamma_x(|\text{BORDER}(y)|) & \text{otherwise} \end{cases}$$

that is to say

$$\gamma_x(|y|) = \begin{cases} \beta_x(|y|) & \text{if } x = y \text{ or } \text{BORDER}(y) \cdot a \not\leq y \\ \gamma_x(\beta_x(|y|)) & \text{otherwise} \end{cases}$$

If $|y| = i$, we can write instead

$$\gamma_x(i) = \begin{cases} \beta_x(i) & \text{if } x = y \text{ or } \text{BORDER}(y) \cdot a \not\leq y \\ \gamma_x(\beta_x(i)) & \text{otherwise} \end{cases}$$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

The condition can also be rewritten in terms of β_x and i , just as we did with the Morris-Pratt algorithm, pages 35 and 19. Let $|y| = i$, then

$$\begin{aligned}\text{BORDER}(y) \cdot a \not\leq y &\iff x[\beta_x(i) + 1] \neq x[i + 1] \\ x = y &\iff |x| = i\end{aligned}$$

So, finally, for $1 \leq i \leq |x|$,

$$\gamma_x(i) = \begin{cases} \beta_x(i) & \text{if } i = |x| \text{ or } x[\beta_x(i) + 1] \neq x[i + 1] \\ \gamma_x(\beta_x(i)) & \text{otherwise} \end{cases}$$

which allows us to naturally extend γ_x on 0: $\gamma_x(0) = \beta_x(0) = -1$.

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

Before going further, let us check by hand the following values of $\gamma_x(i)$ and compare them to $\beta_x(i)$: we must always have $\gamma_x(i) \leq \beta_x(i)$, since we plan an optimisation.

x		a	b	c	a	b	a	b	c	a	c
i	0	1	2	3	4	5	6	7	8	9	10
$\beta_x(i)$	-1	0	0	0	1	2	1	2	3	4	0
$\gamma_x(i)$	-1	0	0	-1	0	2	0	0	-1	4	0

One difference between β_x and γ_x is that, in the worst case, there is always an empty maximum border ε , i.e., $\beta_x(i) = 0$, whereas there may be no maximum disjoint border at all, i.e., $\gamma_x(i) = -1$. For instance, the prefix abc of x has an empty maximum border, i.e., $\beta_x(3) = 0$, but has no maximum disjoint border, i.e., $\gamma_x(3) = -1$, since $x[1] = x[4]$.

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

This definition of γ_x relies on β_x , more precisely, the computation of $\gamma_x(i)$ requires $\beta_x^p(i)$, i.e., values $\beta_x(j)$ with $j < i$, since

$$\beta_x(0) = -1 \quad \beta_x(i) = 1 + \beta_x^k(i-1) \quad 1 \leq i$$

where k is the smallest non-zero integer such that

- either $1 + \beta_x^k(i-1) = 0$ or $x[1 + \beta_x^k(i-1)] = x[i]$

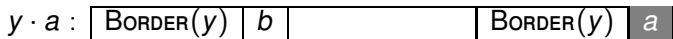
or, equivalently,

$$\text{BORDER}(ya) = \begin{cases} \text{BORDER}^k(y) \cdot a & \text{if } \text{BORDER}^k(y) \cdot a \leq y \\ \varepsilon & \text{if } \text{BORDER}^{k-1}(y) = \varepsilon \end{cases}$$

where k is the smallest non-zero integer such that $\text{BORDER}^k(y) \cdot a \leq y$ or $\text{BORDER}^{k-1}(y) = \varepsilon$.

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

Therefore, let us find another definition of $\text{BORDER}(ya)$ which relies on $\text{BORDER}(y)$ but **not** on $\text{BORDER}^q(y)$ with $2 \leq q$. This can be achieved by considering again the figure



Indeed, if $a = b$ then $\text{BORDER}(ya) = \text{BORDER}(y)$. Else, $a \neq b$, and thus the maximum border can be found among the *disjoint* borders of $\text{BORDER}(y)$, otherwise $\text{BORDER}(ya) = \varepsilon$.

$$\text{BORDER}(ya) = \begin{cases} \overline{\text{BORDER}_x^q}(\text{BORDER}(y)) \cdot a & \text{if } \overline{\text{BORDER}_x^q}(\text{BORDER}(y)) \cdot a \leq y \\ \varepsilon & \text{if } y = \varepsilon \text{ or } \overline{\text{BORDER}_x^{q-1}}(\text{BORDER}(y)) = \varepsilon \end{cases}$$

where $ya \leq x$ and q is the smallest integer such that

$$\overline{\text{BORDER}_x^q}(\text{BORDER}(y)) \cdot a \leq y \text{ or } \overline{\text{BORDER}_x^{q-1}}(\text{BORDER}(y)) = \varepsilon.$$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

By taking the lengths and letting $0 \leq i \leq |x| - 1$ and $|y| = i$,

$$\begin{aligned}\beta_x(0) &= -1 \\ \beta_x(i+1) &= \begin{cases} \gamma_x^q(\beta_x(i)) + 1 & \text{if } x[\gamma_x^q(\beta_x(i)) + 1] = x[i+1] \\ 0 & \text{if } i = 0 \text{ or } \gamma_x^{q-1}(\beta_x(i)) = 0 \end{cases}\end{aligned}$$

where q is the smallest integer such that

- either $x[\gamma_x^q(\beta_x(i)) + 1] = x[i+1]$
- or $\gamma_x^{q-1}(\beta_x(i)) = 0$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

Since $\gamma_x(i) = -1$ if and only if $i = 0$, we have

$$\begin{aligned}\gamma_x^{q-1}(\beta_x(i)) = 0 &\iff \gamma_x(\gamma_x^{q-1}(\beta_x(i))) = \gamma_x(0) \iff \gamma_x^q(\beta_x(i)) = -1 \\ &\iff \gamma_x^q(\beta_x(i)) + 1 = 0\end{aligned}$$

Therefore we can simplify the new definition of β_x as

$$\begin{aligned}\beta_x(0) &= -1 \\ \beta_x(1) &= 0 \\ \beta_x(i+1) &= \gamma_x^q(\beta_x(i)) + 1 \quad 1 \leq i \leq |x| - 1\end{aligned}$$

where q is the smallest integer such that

- either $x[\gamma_x^q(\beta_x(i)) + 1] = x[i + 1]$
- or $\gamma_x^q(\beta_x(i)) + 1 = 0$

Knuth-Morris-Pratt/Maximum disjoint borders/Imperative

```
GAMMA( $x$ )
1   $\gamma_x[0] \leftarrow -1$ ;  $offset \leftarrow -1$   $\triangleright offset = \beta_x(0)$ 
2  for  $i \leftarrow 1$  to  $|x| - 1$ 
3      do repeat
4           $offset \leftarrow \gamma_x[offset]$ 
5          until  $offset = -1$  or  $x[offset + 1] = x[i]$ 
6           $offset \leftarrow offset + 1$   $\triangleright offset = \gamma_x^q(\beta_x(i - 1)) + 1 = \beta_x(i)$ 
7          if  $i = |x|$  or  $x[offset + 1] = x[i + 1]$ 
8              then  $\gamma_x[i] \leftarrow offset$   $\triangleright \gamma_x(i) = \beta_x(i)$ 
9              else  $\gamma_x[i] \leftarrow \gamma_x[offset]$   $\triangleright \gamma_x(i) = \gamma_x(\beta_x(i))$ 
10  $GAMMA \leftarrow \gamma_x$ 
```

Knuth-Morris-Pratt/Better failure function

With the algorithm of Morris and Pratt, it was convenient to use a failure function $s(i) = 1 + \beta(i - 1)$. Similarly, we need here another failure function, noted r , such that $r(i) = 1 + \gamma(i - 1)$, for $1 \leq i$.

In order to finish, we need to give the algorithm for the failure function r .

We could then keep `BETA` and create a separate function `KMP-FAIL` simply by sticking to the above definition:

$\text{KMP-FAIL}(\gamma_x, i)$
1 $\text{KMP-FAIL} \leftarrow 1 + \gamma_x[i - 1] \qquad \triangleright r_x[i] \leftarrow 1 + \gamma_x[i - 1]$

But this is not efficient: it would be better to precompute and store all the values of `KMP-FAIL` in an array.

Knuth-Morris-Pratt/Better failure function (cont)

Let us recall the definition of γ_x . Let $1 \leq i \leq |x|$ and

$$\gamma_x(0) = -1$$

$$\gamma_x(i) = \begin{cases} \beta_x(i) & \text{if } i = |x| \text{ or } x[\beta_x(i) + 1] \neq x[i + 1] \\ \gamma_x(\beta_x(i)) & \text{otherwise} \end{cases}$$

Then, for $1 \leq i \leq |x| - 1$

$$r_x(1) = 1 + \gamma_x(0) = 0$$

$$r_x(1 + i) = 1 + \gamma_x(i) = \begin{cases} 1 + \beta_x(i) & \text{if } x[\beta_x(i) + 1] \neq x[i + 1] \\ 1 + \gamma_x(\beta_x(i)) & \text{otherwise} \end{cases}$$

Knuth-Morris-Pratt/Better failure function (cont)

This can be slightly simplified into

$$r_x(1) = 0$$
$$r_x(1 + i) = \begin{cases} r_x(1 + \beta_x(i)) & \text{if } x[1 + \beta_x(i)] = x[1 + i] \\ 1 + \beta_x(i) & \text{otherwise} \end{cases}$$

where for $1 \leq i \leq |x| - 1$.

Now, we need to express β_x in terms of r_x instead of γ_x .

Knuth-Morris-Pratt/Better failure function (cont)

First, let us prove by induction on $0 \leq p$ that

$$r_x^p(1 + i) = 1 + \gamma_x^p(i) \quad 0 \leq i \leq |x| - 1$$

Because

$$r_x^0(1 + i) = 1 + \gamma_x^0(i) \iff 1 + i = 1 + i$$

and, assuming it is true up to p , we have

$$r_x^{p+1}(1 + i) = r_x(r_x^p(1 + i)) = r_x(1 + \gamma_x^p(i)) = 1 + \gamma_x(\gamma_x^p(i)) = 1 + \gamma_x^{p+1}(i)$$

which is the property at rank $p + 1$.

Knuth-Morris-Pratt/Better failure function (cont)

So, the definition of β_x

$$\beta_x(0) = -1 \quad \beta_x(1) = 0 \quad \beta_x(1+i) = 1 + \gamma_x^q(\beta_x(i)) \quad 1 \leq i \leq |x| - 1$$

where q is the smallest integer such that either $x[\gamma_x^q(1 + \beta_x(i))] = x[1 + i]$ or $1 + \gamma_x^q(\beta_x(i)) = 0$, is equivalent to

$$\beta_x(0) = -1 \quad \beta_x(1) = 0 \quad \beta_x(1+i) = r_x^q(1 + \beta_x(i)) \quad 1 \leq i \leq |x| - 1$$

where q is the smallest integer such that

- either $x[r_x^q(1 + \beta_x(i))] = x[1 + i]$
- or $r_x^q(1 + \beta_x(i)) = 0$

Knuth-Morris-Pratt/Better failure function (cont)

Or, by changing $i + 1$ into i ,

$$\beta_x(0) = -1 \quad \beta_x(1) = 0 \quad \beta_x(i) = r_x^q(1 + \beta_x(i-1)) \quad 2 \leq i \leq |x|$$

where q is the smallest integer such that

- either $x[r_x^q(1 + \beta_x(i-1))] = x[i]$
- or $r_x^q(1 + \beta_x(i-1)) = 0$

Knuth-Morris-Pratt/Better failure function/Imperative

```
KMP-FAIL( $x$ )
1   $r_x[1] \leftarrow 0$ ;  $offset \leftarrow 0$                                  $\triangleright offset = 1 + \beta_x(1 - 1)$ 
2  for  $i \leftarrow 1$  to  $|x| - 1$ 
3      do repeat
4           $offset \leftarrow r_x[offset]$ 
5          until  $offset = 0$  or  $x[offset] = x[i]$ 
6           $offset \leftarrow offset + 1$                                  $\triangleright offset = 1 + \beta_x(i)$ 
7                                   $\triangleright offset = 1 + r_x^q(1 + \beta_x(i - 1))$ 
8          if  $x[offset] = x[1 + i]$ 
9              then  $r_x[i] \leftarrow r_x[offset]$                      $\triangleright r_x(i) = r_x(1 + \beta_x(i))$ 
10             else  $r_x[i] \leftarrow offset$                          $\triangleright r_x(i) = 1 + \beta_x(i)$ 
11  KMP-FAIL  $\leftarrow r_x$ 
```

Knuth-Morris-Pratt/The code

The Knuth-Morris-Pratt algorithm is exactly the Morris-Pratt algorithm, except that we use a better failure function r instead of s :

```
KMP( $x, t$ )
1   $r \leftarrow \text{KMP-FAIL}(x)$ 
2   $i \leftarrow 1; j \leftarrow 1$ 
3  while  $i \leq |x|$  and  $j \leq |t|$ 
4      do if  $i = 0$  or  $x[i] = t[j]$ 
5          then  $i \leftarrow i + 1; j \leftarrow j + 1$ 
6          else  $i \leftarrow r[i]$ 
7  if  $|x| < i$ 
8      then ...            $\triangleright$  Occurrence of  $x$  in  $t$  at position  $j - |x|$ .
9      else ...            $\triangleright$  No occurrences.
```

Knuth-Morris-Pratt/The code (cont)

If we allow arguments to be functions, we can elegantly factorise the two text search algorithms into one:

```
SEARCH( $x, t, f$ )  
1   $i \leftarrow 1; j \leftarrow 1$   
2  while  $i \leq |x|$  and  $j \leq |t|$   
3      do if  $i = 0$  or  $x[i] = t[j]$   
4          then  $i \leftarrow i + 1; j \leftarrow j + 1$   
5          else  $i \leftarrow f[i]$   
6  if  $|x| < i$   
7      then ... ▷ Occurrence of  $x$  in  $t$  at position  $j - m$ .  
8      else ... ▷ No occurrences.
```

Knuth-Morris-Pratt/The code (cont)

Then

$$\text{MP}(x, t) = \text{SEARCH}(x, t, \text{MP-FAIL})$$

$$\text{KMP}(x, t) = \text{SEARCH}(x, t, \text{KMP-FAIL})$$

Knuth-Morris-Pratt/Maximum disjoint borders (cont)

Example. Consider the following table for the search of the word $x = \text{abacabac}$ in the text $t = \text{babacacabacaab}$.

j	1	2	3	4	5	6	7	8	9	10	11	12	13	14
	b	a	b	a	c	a	c	a	b	a	c	a	a	b
i	1	1	2	3	4	5	6	1	2	3	4	5	6	2
	0						1						1	
							0							

Knuth-Morris-Pratt/Example

1 2 3 4 5 6 7 8 9 10 11 12 13 14
t :

b	a	b	a	c	a	c	a	b	a	c	a	a	b
---	---	---	---	---	---	---	---	---	---	---	---	---	---

x :

a	b	a	c	a	b	a	c				
	a	b	a	c	a	b	a	c			
		a	b	a	c	a	b	a	c		
			a	b	a	c	a	b	a	c	
				a	b	a	c	a	b	a	c

Comparison of Morris-Pratt and Knuth-Morris-Pratts

For example, here is the table comparing the values of the two failure functions for the same pattern:

x	a	b	a	c	a	b	a	c
i	1	2	3	4	5	6	7	8
$s_x(i)$	0	1	1	2	1	2	3	4
$r_x(i)$	0	1	0	2	0	1	0	2