

A **finite automaton** is a quintuplet $M = (Q, q_0, A, \Sigma, \delta)$

- Q : finite set of states
- q_0 : initial state
- $A \subset Q$: set of final states
- Σ : finite alphabet
- δ is a function from $Q \times \Sigma$ in Q called the transition function.

The **suffix function** associated with a pattern $P[1..m]$:

$$\sigma : \Sigma^* \rightarrow \{0, 1, \dots, m\} \quad (1)$$

$$t \rightarrow \sigma(x) = \max\{k \mid P[1..k] \text{ suffix of } t\}$$

The number $\sigma(t)$ is the size of the largest prefix of the pattern being searched for, which is the suffix of the text t .

Example : For $P = ab$, one have $\sigma(\epsilon) = 0$, $\sigma(ccaca) = 1$, $\sigma(ccab) = 2$.

If x is suffix of y , $\sigma(x) \leq \sigma(y)$.

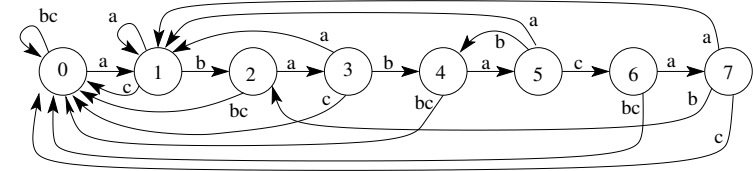
b is suffix of ab , $\sigma(b) \leq \sigma(ab)$

a is suffix of aa , $\sigma(a) \leq \sigma(aa)$.

Automaton associated with a pattern. This is the automaton for which we are in state q if and only if the largest prefix we have just read is $P[1..q]$.

- $Q = \{0, 1, \dots, m\}$
- $q_0 = \{0\}$
- $A = \{m\}$
- $\delta(q, a) = \sigma(P_q a)$ maximum suffix of the concatenation of P_q with a .

Example : Search for pattern *ababaca*



Once the automaton has been constructed, the text traversal algorithm is as follows :

```

1 def FiniteAutomatonSearch(T, delta, m)
2   n = len(T);
3   q = 0;
4   for (i=1; i<=n; i++) {
5     q = delta(q, T[i]);
6     if q = m then
7       print('The pattern appears with the offset', i-m);
8   }

```

Complexity : $O(n)$

Run the automaton on the string *ababacaba*.

Computing the transition function.

The idea is based on the meaning of the different states of the automaton : state i corresponds to the state where the first i letters of the searched pattern have just been read. To build the automaton, we go through all the states of the automaton (from 0 to n , where n is the length of the word we're looking for) and for each state i , we go through each letter a of the alphabet. We then calculate the longest prefix of the pattern that is a suffix of $P[1..i].a$. The length of this suffix gives the arrival state of the transition starting from i via letter a .

```

1 def Transition_Function_computation (P, Sigma)
2   m = len(P);
3   for q in range(m):
4     for a in Sigma :
5       k = min(m, q+1);
6       while (P[1..k] is not a suffix of P_q.a) :
7         k--;
8       delta(q, a) = k;
9   return(delta);

```

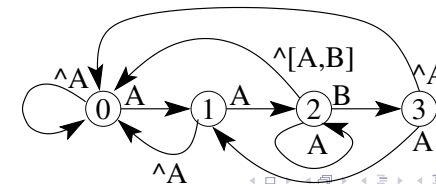
For this function to be correct, the following convention must be used : ϵ is the suffix for all strings.

Complexity analysis :

- lines 6-7 : $O(m^2)$
- line 4 : $O(|\Sigma|)$
- line 3 : $O(m)$
- Global complexity : $O(m^3|\Sigma|)$
- We can do faster.

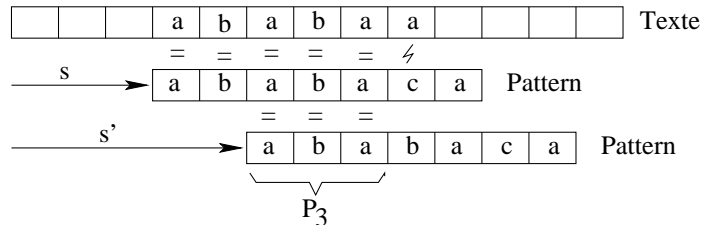
Example : building the automaton for pattern AAB

q	a	k	P_k	$P_q.a$	suffix	δ
q=0	A	$\min(m, q+1) = \min(3, 1) = 1$	A	A	yes	1
	B	1	A	B	no	
		0	ϵ	B	yes	0
q=1	A	$\min(m, q+1) = \min(3, 2) = 2$	AA	AA	yes	2
	B	2	AA	AB	no	
		1	A	AB	no	
q=2	A	$\min(m, q+1) = \min(3, 3) = 3$	AAB	AAA	no	
		2	AA	AAA	yes	2
	B	3	AAB	AAB	yes	3
q=3	A	$\min(m, q+1) = \min(3, 4) = 3$	AAB	AABA	no	
		2	AA	AABA	no	
		1	A	AABA	yes	1
	B	3	AAB	AABB	no	
		2	AA	AABB	no	
		1	A	AABB	no	
	0	ϵ	AABB	yes	0	



This algorithm achieves complexity in $\Theta(n + m)$ by avoiding the transition function δ . It computes an auxiliary function $\pi[1..m]$ precomputed from the pattern in $O(m)$. The array π allows the transition function δ to be computed on the fly if necessary.

Pattern prefix function : Correspondence between the motif and its own shifts.

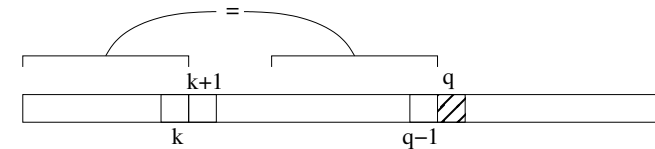


Question : how to calculate s' so that the offset is not invalid?

Answer : Find a suffix P_k of P_q that is a prefix of P .

The prefix function for the P pattern :
 $\Pi : \{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, m-1\}$
 $q \rightarrow \Pi[q] = \max\{k/k < q \text{ and } P_k \text{ is suffix of } P_q\}$
 $\Pi[q]$ is in fact the size of the longest prefix of P which is a proper suffix of P_q .

- We construct the array Π starting from index 0. The initialisation is simple : $\Pi[1] = 0$
- Now let's assume that we have calculated $\Pi[i]$ from $i = 1$ to $q - 1$. To calculate $\Pi[q]$ we have the following situation :



- Since k is the longest prefix that is a suffix of P_{q-1} , the longest prefix that is also a suffix of P_q cannot be longer than P_{k+1} . Furthermore, we have

$$P[k+1] = P[q] \iff \Pi[q] = k+1$$

- If $P[k+1] \neq P[q]$, look for the largest prefix-suffix of P_q . If we don't look at the last letter, the largest prefix-suffix of P_q is also a suffix of P_k . Now we know the largest prefix-suffix of P_k , which is $\Pi[k]$, already constructed. Once we have $P_{\Pi[k]}$, we need to check if it can be extended to the next letter.

```

1 Compute_prefix_Function(P)
2   m = long(P);
3   pi[1] = 0;
4   k = 0;
5   for (q=2; q<=m ; q++) {
6     while (k>0 and P[k+1] != P[q]:
7       {
8         k=pi[k];
9       }
10    if P[k+1] == P[q] then k++;
11    pi[q] = k;
12  }
13  return(pi);

```

1	2	3	4	5	6	7	8	9	10
a	b	a	b	a	b	a	b	c	a

Let's build Π :

Π :	1	2	3	4	5	6	7	8	9	10
	0	0	1	2	3	4	5	6	0	1

- $q = 2, k = 0$ no while ($k=0$)
 $(a=b)?$ no, $k = 0 \implies \Pi[2] = 0$
- $q = 3, k = 0$ no while ($k=0$)
 $(a=a)?$ yes, $k = 1 \implies \Pi[3] = 1$
- $q = 4, k = 1$ no while ($P[k+1]=P[q]$)
 $(b=b)?$ yes, $k = 2 \implies \Pi[4] = 2$
- $q = 5, k = 2$ no while ($P[k+1]=P[q]$)
 $(a=a)?$ yes, $k = 3 \implies \Pi[5] = 3$
- $q = 6, k = 3$ no while ($P[k+1]=P[q]$)
 $(b=b)?$ yes, $k = 4 \implies \Pi[6] = 4$
- $q = 7, k = 4$ no while ($P[k+1]=P[q]$)
 $(a=a)?$ yes, $k = 5 \implies \Pi[7] = 5$
- $q = 8, k = 5$ no while ($P[k+1]=P[q]$)
 $(b=b)?$ yes, $k = 6 \implies \Pi[8] = 6$
- $q = 9, k = 6$ enter into the while
 $while P[7] \neq P[9] (a \neq c) k = \Pi[6] = 4$
 $while P[5] \neq P[9] (a \neq c) k = \Pi[4] = 2$
 $while P[3] \neq P[9] (a \neq c) k = \Pi[2] = 0$
 $(a=c)?$ no, $k = 0 \implies \Pi[9] = 0$
- $q = 10, k = 0$ no while ($k=0 \& P[k+1]=P[q]$)
 $(a=a)?$ yes, $k = 1 \implies \Pi[10] = 1$

Consider the following pattern :

1 2 3 4
S N N S

Let's build the function Π :

Π : 1 2 3 4
0 0 0 1

- $q = 2, k = 0$ no while ($k=0$)
if $P[1] = P[2]$ ($S \neq N$)? no $\implies \Pi[2] = 0$
- $q = 3, k = 0$ no while ($k=0$)
if $P[1] = P[3]$ ($S \neq N$)? no $\implies \Pi[3] = 0$
- $q = 4, k = 0$ no while ($k=0$)
if $P[1] = P[4]$ ($S = S$)? yes, $k = 1 \implies \Pi[4] = 1$

Let's consider $\pi^*[q] = \{q, \pi[q], \pi^2[q], \dots, \pi^t[q]\}$ where t is the first natural number such that $\pi^t[q] = 0$.

Lemma

Let P be a pattern of length m and having the prefix function π . Then, for $q = 1, 2, \dots, m$, one has $\pi^*[q] = \{k \mid P_k \text{ suffix of } P_q\}$

Proof :

1 First inclusion : $i \in \pi^*[q] \implies P_i$ suffix de P_q .

If $i \in \pi^*[q], \exists u \mid \pi^u[q] = i$

- for $u = 0, i = q$ and thus $P_i = P_q$ and P_i is suffix of P_q
- let us suppose $P_{\pi^u[q]}$ suffix of P_q for each $u < u_0$
 $P_{\pi^{u_0}[q]} = P_{\pi[\pi^{u_0-1}[q]]}$ and one have $P_{\pi^{u_0-1}[q]}$ suffix of P_q and $P_{\pi^{u_0}[q]}$ suffix of $P_{\pi^{u_0-1}[q]}$.
Since the suffix relationship is transitive, we have $P_{\pi^{u_0}[q]}$ suffix of P_q .
- Conclusion : $i \in \pi^*[q] \implies P_i$ suffix of P_q .

□

Let's consider $\pi^*[q] = \{q, \pi[q], \pi^2[q], \dots, \pi^t[q]\}$ where t is the first natural number such that $\pi^t[q] = 0$.

Lemma

Let P be a pattern of length m and having the prefix function π . Then, for $q = 1, 2, \dots, m$, one has $\pi^*[q] = \{k \mid P_k \text{ suffix of } P_q\}$

Proof :

2 let us prove that : $\{k \mid P_k \text{ suffix of } P_q\} \subseteq \pi^*[q]$ **Proof by contradiction :**

let us assume that there is an integer in the set

$\{k \mid P_k \text{ suffix of } P_q\} \setminus \pi^*[q]$.

We denote j the largest of the values.

As $q \in \pi^*[q] \cup \{k \mid P_k \text{ suffix of } P_q\} \implies j < q$.

Let j' be the smallest integer in $\pi^*[q]$ that is greater than j .

- P_j is a suffix of P_q since it belongs to $\{k \mid P_k \text{ suffix of } P_q\}$
- $P_{j'}$ is a suffix of P_q since it belongs to $\pi^*[q]$

Thus P_j suffix of $P_{j'}$ (trivial because $j' > j$)

Moreover j is the largest value of $\{k \mid P_k \text{ suffix of } P_q\} \setminus \pi^*[q]$

We should then have $\pi[j'] = j$ and then $j \in \pi^*[q]$.

Contradiction.

□

Lemma

Let P be a pattern of length m and prefix function π . For $q = 1, 2, \dots, m$, if $\pi[q] > 0$ then $\pi[q] - 1 \in \pi^*[q - 1]$.

Proof :

If $k = \pi[q] > 0$ then P_k suffix of P_q .

Thus P_{k-1} suffix of P_{q-1} (by deleting the last character of P_k and P_q)

According to the previous lemma : $k - 1 \in \pi^*[q - 1]$. □

For $q = 2, 3, \dots, m$, we define the subset $E_{q-1} \subseteq \pi^*[q - 1]$ by :

$$E_{q-1} = \{k \mid k \in \pi^*[q - 1] \text{ by } P[k + 1] = P[q]\}$$

Intuitively, E_{q-1} is made up of values $k \in \pi^*[q - 1]$ such that it is possible to extend P_k to P_{k+1} and obtain a suffix of P_q .

corrolary

Let P be a pattern of length m and prefix function pi . For $q = 2, 3, \dots, m$,

$$\pi[q] = 0 \text{ if } E_{q-1} = \{\}$$

$$\pi[q] = 1 + \max\{k \in E_{q-1} \mid P[k+1] = P_q\}$$

Proof :

If $r = \pi[q] > 0$ then P_r suffix of P_q .

And $r \geq 1 \Rightarrow P[r] = P[q]$

According to the previous lemma, if $r \geq 1$, we have :

$$r = 1 + \max\{k \in \underbrace{\pi^*[q-1]}_{E_{q-1}} \mid P[k+1] = P_q\}$$

If $r = 0$, there is no $k \in \pi^*[q-1]$ for which we can extend P_k to P_{k+1} to obtain a suffix of P_q , since we would then have $\pi[q] > 0$.

Thus $E_{q-1} = \{\}$ □

Algorithm validity :

- ① $\pi[1] = 0$ correct because $\pi[q] < q$ for all q
- ② At the start of each loop iteration, we have $k = \pi[q-1]$
 - for the first loop : imposed by $\pi[1] = 0$ and $k = 0$
 - for the others : imposed by $\pi[q] = k$
- ③ while loop : we run through all the values of $\pi^*[q-1]$ until we find one for which $P[k+1] = P[q]$.
At this point, we know that k is the largest value of E_{q-1} ; and from the corollary, we can give to $\pi[q]$ the value $k+1$.
If no such k is found, $k = 0$ □

<pre> 1 KMP1(T,P) 2 n = long[T]; 3 m = long[P]; 4 PI = Calcul_fonct_prefixe(P); 5 q = 0 ; 6 7 Pour i=1 à n faire 8 tant que q>0 & P[q+1]!=T[i] 9 q = PI[q]; 10 si P[q+1]=T[i]: 11 q = q+1 12 si q = m alors 13 print("hit at ", i-m); 14 q = PI[q]; </pre>	<pre> 1 KMP2(T,P) 2 n = long[T]; 3 m = long[P]; 4 PI = Compute_prefix_Function(P); 5 i = q = 0 ; 6 7 while (i<n): 8 if T[i] == P[q]: 9 i++; q++; 10 else 11 if q==0 : 12 i++; 13 else: 14 q = PI[q-1]; 15 if q==m 16 print("hit at ", i-m); 17 q = PI[q-1]; </pre>
--	---

The first version is based on the same idea as the prefix function.

The second version manages two indices in a single loop : one to indicate progress in the text and another to indicate progress in the pattern.

Complexity analysis : requires amortized analysis...