

I Recherche de motifs

1) Algorithme naïf

On cherche à trouver les occurrences d'une chaîne de caractères $M[1..m]$ dans un texte $T[1..n]$ ($m \leq n$).

L'algorithme naïf associé est la méthode dite des fenêtres glissantes:

$T: a b c a b a a b c a b a c$
 $M \rightarrow \rightarrow \rightarrow a b a a$

Recherche naïve (T, M)

pour s de 0 à $m-m$

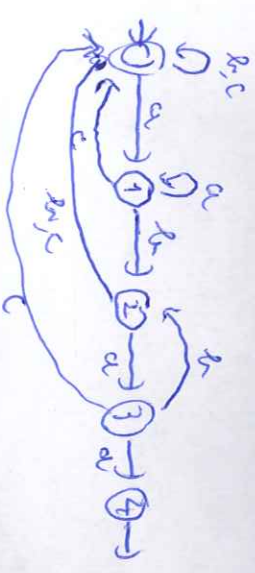
$n_i P[1..m] = T[s+1..s+m]$

afficher: "M apparaît à la position" s

Complexité temporelle: $O((n-m+1)m)$

2) Automates des occurrences, Knuth-Morris-Pratt

Idee: Un constructeur reconnaît récursivement le motif M et on prend T en mot d'entrée



La construction de la fonction de transition se fait à l'aide de la fonction suivante $O(nq) = m \times \{a, c, \dots\}$ appelée $M[q, a]$ en fait:

$$f(q, a) = \sigma(M[1..q]a)$$

Complexité temporelle:

• Construction de l'automate $O(m^2 | \Sigma |)$

• Recherche: $O(m)$

Amélioration: algorithme de Knuth-Morris-Pratt

Principe: On évite le calcul de la fonction de transition en introduisant une nouvelle fonction nécessaire en temps linéaire:

$T[q] >$ la longueur du plus long préfixe de M qui est un suffixe propre de $M[1..q]$.

Une série d'algorithmes:

1 $q = 0$
 $KMP(T, M)$

2 pour i de 1 à m

3 tant que $q > 0$ et $M[q+1] \neq T[i]$

4 $q = T[q]$

5 si $M[q+1] = T[i]$

6 $q = q+1$

7 si $q = m$

8 afficher "le motif apparaît en position" $i-m$

9 $q = T[q]$

3) Algorithme de Boyer-Moore

Principe: On applique une variante de la méthode des fenêtres glissantes sur le caractère par caractère de la fin du motif au texte à chaque tentative.



Notons Z le plus long suffixe commun à M et $T[i, \dots, j]$.
Si $Z \neq M$, on doit effectuer un décalage:

- si $k-Z$ apparaît dans M , on aligne son occurrence dans T avec son occurrence la plus à droite dans M
- sinon on aligne le plus grand préfixe de M qui est un suffixe de Z .

Pq: le décalage du décalage à effectuer est indépendant du texte, et peut être précalculé à partir de M .

Complexité temporelle

Décalage: $O(m | \Sigma |)$
 Recherche: $O(m)$

4) Table des suffixes

Principe: On construit une table luee en ordre lexicographique contenant les suffixes de T et on recherche M par dichotomie.

Complexité de la recherche: $O(m \log m)$.

Le prétraitement sur texte ne peut être calculé à l'avance
méthode de la fin en fin en $O(m \log m)$

5) Automate d'applications

- Boyer-Moore est l'algorithme le plus communément utilisé (diff + F1 graph)
- Sur table des suffixes ne prête à des recherches multiples sur un même texte.

II Comparaison de chaînes de caractères

1) Plus longue sous séquence commune

- soit $X[1..m]$ et $Y[1..n]$ deux chaînes de caractères. Une plus longue séquence commune (PLSC) de X et Y est une chaîne $Z[1..k]$ telle qu'il existe p, q croissants tels $X[p..q] = Y[p..q]$.

Propriété: toute structure d'une PLSC

- si $X[m]$ et $Y[n]$, $Z[k] = X[m]$ et $Z[1..k-1]$ est une PLSC de $X[1..m-1]$ et $Y[1..n-1]$
- si $X[m] \neq Y[n]$, Z est une PLSC de X et Y
- si $X[m] = Y[n]$, Z est une PLSC de X et Y

De cette propriété, on déduit un algorithme de programmation dynamique permettant de calculer une PLSC, pour une complexité temporelle en $O(m \cdot n)$

DVP T

2) Distance d'édit

- Un caractère x pour une chaîne X à une chaîne Y en un nombre minimum d'opérations de type :
- remplacer un caractère par un autre
 - insérer un caractère
 - supprimer un caractère

On étend un algorithme de programmation dynamique calculant la distance pour tous les préfixes de X et Y en ne servant du calcul de la distance pour les préfixes plus petits.

Exemple : pour tous de POMME à POIAGE

	P	O	M	M	E
P	0	1	2	3	4
O	1	0	1	2	3
M	2	1	0	1	2
M	3	2	1	0	1
E	4	3	2	1	0

↳ substitution

NO : On remplit le tableau de gauche à droite et de bas en haut

Application : - correction orthographique

- commande "diff"

3) Application à la recherche d'indices : recherche avec R. différences

Un caractère ou mot d'une chaîne de la recherche d'un motif M dans un texte T , mais on ne permet de numérotage des occurrences qui diffèrent de M par au plus k opérations.

Pour cela, on applique une variante de l'algorithme de distance d'édit :

- la ligne d'initialisation de T vaut 0 partout, pour permettre à la recherche de prendre en compte les occurrences de M partant de T
- à la dernière ligne, on a "oui" M dans T , ainsi toutes les cases de cette ligne indiquent une distance inférieure à k uniquement par référence à occurrences approchées.

III Compression : Codage de Huffman

Problème : Un caractère ou ensemble de caractères doit être fréquemment d'opération. On cherche un codage préfixe optimal de cet ensemble, c'est à dire une fonction associant à chaque caractère une suite binaire telle que :

- aucune suite n'est préfixe d'une autre
- le nombre moyen de caractères transmis est minimal pour un codage préfixe.

Un algorithme pour cela est l'algorithme planteur de Huffman :

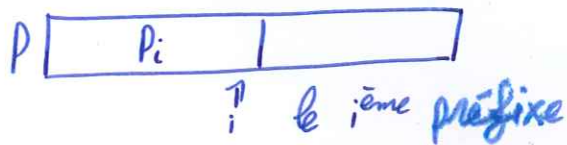
- On construit un ensemble de feuilles situées par les caractères et leur fréquence d'apparition
- Tant que il y en a au moins 2, on choisit les deux plus petits, que k on combine pour former un arbre dont les racines est étiquetée par la somme des fréquences.



Arbre final :

Automate des occurrences

P motif de longueur m , T texte de longueur n



\mathcal{A} $Q = \{0, \dots, m\}$, $q_0 = 0$, $F = \{m\}$, $\delta(q, a) = \sigma(P_q a)$

déterministe

où $\sigma(x) = \max \{k \in \{0, \dots, m\} \mid P_k \sqsupseteq x\}$

↳ bien défini car $P_0 = \epsilon \sqsupseteq x \quad \forall x \in \Sigma^*$

Thm: \mathcal{A} reconnaît $\Sigma^* P$ et \mathcal{A} est minimal.

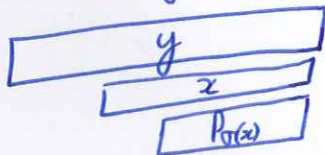
Montrons plus précisément que : $\forall i \in \{0, \dots, m\}$, $\delta(q_0, T_i) = \sigma(T_i)$.

Ce qui donnera le résultat car :

il y a une occurrence de P qui fini à la position i dans T

ssi $P \sqsupseteq T_i$ ssi $\sigma(T_i) = m$.

Lm 1 $x \sqsupseteq y \Rightarrow \sigma(x) \leq \sigma(y)$



$P_{\sigma(x)} \sqsupseteq y$ donc $\sigma(x) \leq \sigma(y)$
par maximalité de $\sigma(y)$.

Lm 2 $\sigma(xa) \leq \sigma(x) + 1$

Si $\sigma(xa) = 0$ ok.



$P_{\sigma(xa)-1} \sqsupseteq xa \Rightarrow \sigma(xa) - 1 \leq \sigma(x)$
par maximalité de $\sigma(x)$.

Lm 3 $\sigma(xa) = \sigma(P_{\sigma(x)}a)$

*
$$\begin{array}{l} \boxed{x} \quad \boxed{a} \\ \boxed{P_{\sigma(x)}} \quad \boxed{a} \\ \boxed{P_{\sigma(xa)}} \end{array}$$

$$\frac{P_{\sigma(x)}a}{\alpha} \supset xa \Rightarrow \sigma(\alpha) \leq \sigma(xa)$$

(lm 1)

*
$$\underbrace{P_{\sigma(xa)}}_{\beta} \supset xa$$

Mais $|\alpha| = \sigma(x) + 1 \geq |\beta| = \sigma(xa)$
(lm 2)

donc $\beta \supset \alpha$

d'où $\sigma(xa) \leq \sigma(\alpha)$ par maximalité de $\sigma(\alpha)$.

Conclusion réc sur i

$i=0$ $d^L(q_0, \epsilon) = q_0$ ok

$i \rightarrow i+1$
$$\begin{aligned} d^L(q_0, T_{i+1}) &= d^L(\underbrace{d^L(q_0, T_i)}_{\sigma(T_i)}, a_{i+1}) = \sigma(P_{\sigma(T_i)} a_{i+1}) \\ &= \sigma(T_i a_{i+1}) \quad (\text{lm 3}) \\ &= \underline{\sigma(T_{i+1})}. \end{aligned}$$

Minimalité Pour $k \in \{0..m\}$ on pose $P = P_k Q_k$,
on a $d^L(q_0, P_k) = k$.

Soient $i < j \in \{0..m\}$. MQ $i \neq j$, cād $L_i \neq L_j$.

*
$$\underbrace{d^L(q_0, P)}_m = d^L(\underbrace{d^L(q_0, P_j)}_j, Q_j)$$
 d'où $Q_j \in L_j$.

* $d^L(q_0, P_i Q_j) = d^L(i, Q_j) \neq m$ car $|P_i Q_j| < m$
d'où d minimal.

Tri des Suffixes

Naïf: $\mathcal{O}(n \times n \log n)$
durée ↑ nb comparaisons
1 comparaison

$w \in \Sigma^*$, $|w| = n$, $S = \{w_i \mid i \in \{1..n\}\}$, $\Sigma = \text{alph}(w)$.



donc $|S| = \mathcal{O}(n)$

On cherche à trier S par ordre lexicographique : on veut $\pi \in \mathcal{O}_n$ tq $w_{\pi(i)}$ soit le i^{e} me élément de S trié.

k premières lettres

On note $w[i..k] = w[i..n]$ si $k \geq n$.

$S_k = \{w[i..i+k-1] \mid i \in \{1..n\}\}$ ($S_n = S$, $|S_k| \leq |S| = n$).

$R_k[i]$ = le rang de $w[i..i+k-1]$ dans S_k trié, $R_k[i] = 0$ si $i > n$

$\pi_k \in \mathcal{O}_n$ tq $w[\pi_k(i).. \pi_k(i)+k-1]$ soit le i^{e} me élément de S_k

\cup à rang égal : stable vis à vis de π_{k-1}

Rq: $R_k[i] \in \{1..|S_k|\}$ alors que $\pi_k(i) \in \{1..n\}$

Idee: si $|S_k| = n = |S|$, $R_k = R_n$ a n elt^s \neq donc $\pi = R_n^{-1}$

k=1 $R_1[i]$ = rang de $w[i]$ dans Σ

k → 2k Lm de dédoublement: $R_{2k}[i]$ est le rang de $(R_k[i], R_k[i+k])$ dans la liste couples $(i \in \{1..n\})$ triée.

Preuve $w[i..i+2k-1] = \underbrace{w[i..i+k-1]}_{u(i)} \underbrace{w[i+k..i+2k-1]}_{v(i)}$

- $w(i) < w(j)$ si $u(i)v(i) < u(j)v(j)$
- si $(u(i), v(i)) < (u(j), v(j))$
- si $(R_k[i], R_k[i+k]) < (R_k[j], R_k[j+k])$

d'où le résultat.

Rq importante : Poser $R[\varepsilon_i] = 0$ pour $i > n$ revient à considérer le mot wa^∞ où a est une lettre inférieure à toutes les autres \rightarrow compatible avec l'ordre lexicographique.

Ex $w = ababca$, $n = 6$, $S = \{w, babca, abca, bca, ca, a\}$

* $S_1 = \{a, b, c\}$

$R_1 = \boxed{1|2|1|2|3|1} \rightarrow (12)(21)(12)(23)(31)(10)$

$\pi_1 = \boxed{1|3|6|2|4|5}$

* $S_2 = \{ab, ba, bc, ca, c\}$

$R_2 = \boxed{2|3|2|4|5|1} \rightarrow (22)(34)(25)(41)(50)(10)$

$\pi_2 = \boxed{6|1|3|2|4|5}$

* $R_4 = \boxed{2|4|3|5|6|1}$

$\pi = \pi_4 = \boxed{6|1|3|2|4|5}$

Algorithme

Utilise le tri par dénombrement :

$TRID(\pi_1, n, R)$ renvoie $\pi_2 \in \mathcal{O}_n$ tel que $R[\varepsilon_{\pi_2(i)}]$ est le i ème élément de $\{R[\varepsilon_i] \mid i \in \{1, \dots, n\}\}$ et stable vis à vis de π_1 (si $R[\varepsilon_{\pi_1(i)}] = R[\varepsilon_{\pi_1(j)}]$ avec $\pi_1(i) < \pi_1(j)$ alors $\pi_2(i) < \pi_2(j)$).

↳ s'exécute en $\mathcal{O}(n^2)$ (en supprimant $\{R[\varepsilon_i]\} \subseteq \{1, \dots, n\}$).

$TRID(\pi_1, n, A) =$

pour $i = 1$ à n faire $C[i] \leftarrow 0$ // $n = \max_i R[\varepsilon_i]$

pour $j = 1$ à n faire $C[A[j]] \leftarrow C[A[j]] + 1$ // $n = |w|$

pour $i = 1$ à n faire $C[i] \leftarrow C[i] + C[i-1]$

pour $j = n$ à 1 faire $\pi_2(j) \leftarrow C[A[\pi_1(j)]]$

$C[A[\pi_1(j)]] \leftarrow C[A[\pi_1(j)]] - 1$

$\boxed{\mathcal{O}(n)}$

TRI-SUFFIXES (w) =

- $k \leftarrow 1$
- $\pi \leftarrow \text{TRID}(\text{id}, n, w[1]) \quad \mathcal{O}(n)$
- $\alpha \leftarrow 1, R_1[\pi(1)] \leftarrow 1$
- pour $i=2$ à n faire
 - si $w[\pi(i)] \neq w[\pi(i-1)]$ alors $\alpha \leftarrow \alpha + 1$) $\mathcal{O}(n)$
 - $R_1[\pi(i)] \leftarrow \alpha$
- tant que $\alpha < n$ faire // $\alpha = |S_k|$
 - $\pi \leftarrow \text{TRID}(\pi, n, R_k[+k])$) $\mathcal{O}(n)$
 - $\pi \leftarrow \text{TRID}(\pi, n, R_k)$) $\mathcal{O}(n)$
 - $\alpha \leftarrow 1, R_{2k}[\pi(1)] \leftarrow 1$
 - pour $i=2$ à n faire $\mathcal{O}(n)$
 - si $(R_k[\pi(i)], R_k[\pi(i)+k]) \neq (R_k[\pi(i-1)], R_k[\pi(i-1)+k])$ alors $\alpha \leftarrow \alpha + 1$
 - $R_{2k}[\pi(i)] \leftarrow \alpha$
 - $k \leftarrow 2k$
- renvoyer π \square

Il y a au pire $\log_2 n$ itérations de la boucle tant que

$$\Rightarrow \boxed{\mathcal{O}(n \log n)}$$