

Pb Il s'agit encore du problème de la recherche de toutes les occurrences d'un motif $x = x_1 \dots x_m$ dans un texte $t_1 \dots t_n = t$.
On renvoie ces occurrences par leur indice de début.

Idée L'idée est de voir nos lettres comme des chiffres.
Ainsi au lieu de comparer des mots - on compare des nombres.

Codage / Décodage

Soit $\Sigma = [0..(b-1)]$ où $b \in \mathbb{N}^*$

Puisqu'on a coutume de noter $\langle _ \rangle_b$ le codage d'un entier en chiffres selon la base b , on notera ici $\langle _ \rangle_b$ le décodage d'un mot de Σ (une suite de chiffres) selon la base b .

$$\begin{aligned}
 \underbrace{\langle w_1 w_2 \dots w_n \rangle_b}_{\substack{\text{poids} \\ \text{fort}}} &:= w_1 \times b^{n-1} + w_2 \times b^{n-2} \dots + w_{n-1} \times b + w_n \\
 &= \sum_{i=1}^n w_i b^{n-i} \\
 &= P(b) \text{ où } P = \sum_{i=1}^n w_i X^{n-i} \\
 &= \text{EVALUER}(w_n, w_{n-1}, w_{n-2} \dots w_1, b) \\
 &= ((w_1 \times b + w_2) \times b \dots) \times b + w_n
 \end{aligned}
 \quad \left. \vphantom{\begin{aligned} \dots \\ \dots \\ \dots \\ \dots \end{aligned}} \right] \text{cf Horner}$$

décoder($\overbrace{w_1 \dots w_n}^w, b$)

```

m ← |w|
v ← w1
pour i allant de 2 à m
  v ← v × b
  v ← v + wi
retourner v.
  
```

NB: l'algo ci-contre n'est pas utile tel quel mais permet de mieux comprendre les 1^{ères} boucles de l'algo de Rabin-Karp.

ex $\Sigma = [0, 1, 2, 3]$ $b = 4$ $\langle 123 \rangle = 3 + 2 \times 4 + 1 \times 16 = 27$ // calcul à la main
 $= (((1 \times 4) + 2) \times 4) + 3 = (6 \times 4) + 3 = 27$ // Horner.

Dans notre problème on ne voulait non seulement calculer $\succ x \prec_b$ mais aussi $\succ w \prec_b$ pour tout les facteurs w de t de taille m .

Mais du calcul de $\succ w_i \prec_b - w_{i+m-1} \prec_b$ on peut déduire des infos sur $\succ w_{i+1} \prec_b - w_{i+m} \prec_b$ le facteur suivant

ex	$b = 10$	$\succ 1234 \prec = 1234$
	$m = 4$	$\succ 2345 \prec = 234 \times 10 + 5$
		$= (\succ 1234 \prec - 1 \times 10^3) \times 10 + 5$
		$= (\succ 1234 \prec - 1 \times b^{m-1}) \times b + 5$

Pte | Considérons $(m, b) \in \mathbb{N}^*$ fixés, et $\Sigma = [0 \dots (b-1)]$.

Si $w \in \Sigma^m$, $g \in \Sigma$ et $d \in \Sigma$ alors

$$\succ gw d \prec_b = (\succ gw \prec_b - g \times b^{m-1}) + d \quad *$$

Reduction modulo q

L'amélioration par rapport à l'algorithme naïf doit être qu'on puisse comparer d'un côté d'un seul, et contenir d'une fenêtre de taille m et le motif, en comparant les nombres associés. Etat. dit on considère la comparaison de nombre comme une opération constante.

Or le coût de cette opération dépend de la taille des nombres en réalité et ici nos nombres sont bornés par $b^m - 1$, ce qui n'est pas const!!!

En fixant un entier premier q , et en réalisant toutes nos opérations modulo q , nos nombres seront bornés par $q-1$, il est alors justifié de considérer constante la comparaison de nombre, ainsi que les addi. et multipliées.

Mais ! aux faux positifs. En effet

$$\succ x \prec_b \neq \succ t_i - t_{i+m-1} \prec_b [q] \Rightarrow x \neq t_i - t_{i+m-1}$$

$$\text{mais } \succ x \prec_b \equiv \succ t_i - t_{i+m-1} \prec_b [q] \not\Rightarrow x = t_i - t_{i+m-1}$$

D'où une réification (v) dans l'algorithme qui suit.

EN MOYENNE

On fait l'hypothèse que la réduction modulo q du décodage en base b agit comme une transformée aléatoire uniforme entre Σ^* et $\mathbb{Z}/q\mathbb{Z}$.

Supposons le motif x fixe, et T un texte aléatoire sur Σ de n lettres, suivant une loi uniforme.

$$\text{Aut dit } \forall w \in \Sigma^m \quad \mathbb{P}(T=w) = 1/b^m$$

$$\text{et } \forall i \in [1..m], \forall a \in \Sigma \quad \mathbb{P}(T_i=a) = 1/b$$

Puisque les T_i sont unif, $T_i - T_{i+m+1}$ est unif dans Σ^m .

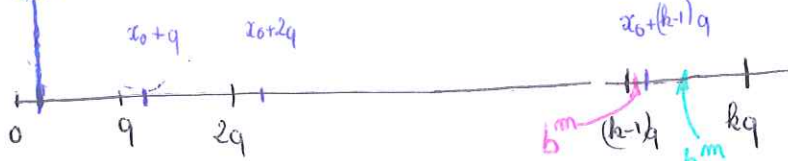
$$\text{Donc } \mathbb{P}(\exists T_i - T_{i+m+1} <_b \equiv x <_b [q]) = \frac{\#\{w \in \Sigma^m \mid \exists w <_b \equiv x < [q]\}}{\#\{w \in \Sigma^m\}}$$

$$= \frac{\#\{s \in [0..b^m-1] \mid s \equiv x <_b [q]\}}{b^m}$$

$$= \frac{k-1}{b^m} \text{ où } k-1 = \lfloor b^m/q \rfloor$$

$$\approx \frac{1}{q}$$

$x_0 \in [0, q[$ tq $x_0 \equiv x <_b [q]$.



Fixons k tq $b^m \in [(k-1)q, kq[$

Selon que b^m se situe "sous" $x_0+(k-1)q$ ou au delà de $x_0+(k-1)q$, il y a $k-1$ ou k s entre 0 et b^m-1 tels que $s \equiv x <_b [q]$ ie tq $s \in x_0\mathbb{Z}$.

On cherche le nombre de fois (nombre moyen) où l'on rentre dans la boucle implicite de (V) car elle est en $O(m)$.

Aut dit on cherche le nombre de tests " $v_t = v_x$ " positifs, soit

$$\mathbb{E} \left[\sum_{i=1}^{n-m+1} \mathbb{1}_{(v_t = v_x)} \right] = \sum_{i=1}^{n-m+1} \mathbb{E} \left[\mathbb{1}_{\exists T_i - T_{i+m-1} <_b \equiv x <_b [q]} \right]$$

$$= \sum_{i=1}^{n-m+1} \mathbb{P}(\text{---})$$

$$\approx (n-m+1) \times \frac{1}{q}$$

D'où une complexité moyenne en pour la boucle B_2 .

$$O \left(\underbrace{n-m+1}_{\substack{\text{test négatif} \\ \text{ou positif} \\ \text{(coût du test)}}} + \underbrace{\frac{(n-m+1) \times m}{q}}_{\substack{\text{test positif} \\ \text{et donc} \\ \text{vérifica (V)} \\ \text{en } O(m)}} \right)$$

Finalement la complexité moyenne est en

$$O \left(m + (n-m+1) \left(1 + \frac{m}{q} \right) \right) = O \left(m + \frac{nm}{q} \right)$$

Algo

Considérons Σ et donc b fixés. Fixons aussi q (év. choisi selon b)

Rabin-Karp (x, t)

$m \leftarrow \text{longueur}(t)$
 $m \leftarrow \text{---}(x)$

$$h \leftarrow b^{m-1} \text{ mod } (q)$$

$$N_x \leftarrow x_1$$

$$N_t \leftarrow t_1$$

Pour i allant de 2 à m

B_1

$$N_x \leftarrow N_x \otimes b \text{ mod } [q]$$

$$N_x \leftarrow N_x \oplus x_i \text{ mod } [q]$$

$$N_t \leftarrow N_t \otimes b \text{ mod } [q]$$

$$N_t \leftarrow N_t \oplus t_i \text{ mod } [q]$$

// $N_x = \Rightarrow x <_b$ et $N_t = \Rightarrow t_1 - t_m <_b$

Pour i allant de 1 à $m-m+1$

B_2

Si $N_t = N_x$

alors $\left(\begin{array}{l} \text{si } t_i \text{ --- } t_{i+m} = x \\ \text{alors afficher } i \end{array} \right) (V)$

Si $i < m-m+1$

alors $N_t \leftarrow (N_t - t_i \times h) + t_{i+m}$ // applica de *

Complexité

AU PIRE

B_1 se fait toujours en $O(m)$

B_2 se fait au pire en $O(\underline{m-m+1} \times \underline{m})$

puisque au pire chaque indice est un début possible de motif et l'on doit faire la vérifica (V) en $O(m)$

ex $x = a^m \quad t = a^m$

La complexité au pire est en $O((m-m+1) \times m)$
soit identique à celle de l'algorithme naïf.

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6\} \quad b = 7.$$

$$x = 116$$

$$t = 425116$$

x	17
0	0
1	17
2	34
3	51
4	68
5	85

Au plus un motif de taille 3 vaut

$$\begin{aligned} >666_7 &= 6 \times 7^2 + 6 \times 7^1 + 6 \times 7^0 \\ &= 6 \times 49 + 42 + 6 \\ &= 300 - 6 + 42 + 6 \\ &= 342 \end{aligned}$$

342 ne tient pas sur un octet (qui vaut au + $2^8 - 1 = 127$)

Si on choisit $-q = 17$ on a $qb = 119 \leq 127$ donc qb tient sur un octet.

$$\begin{aligned} >116_7 &\equiv 1 \times 49 + 1 \times 7 + 6 \times 1 \equiv 62 \equiv 11 [17] \quad \text{à la main} \\ &\equiv \left(\left((1 \times 7) + 1 \right) \times 7 \right) + 6 \quad \text{dans l'algo via Horner.} \end{aligned}$$

$$7; 7+1 \equiv 8 [17]; 8 \times 7 = 56 \equiv 5 [17]; 5+6 \equiv 11 [17]$$

$$\begin{aligned} >425_7 : \quad &4 \times 7 \equiv 28 \equiv 11 [17] && \text{Donc } >425_7 \equiv >116_7 [17] \\ &11+2 \equiv 13 \equiv -4 [17] && \text{il faut tester mais } 4 \neq 1 \\ &7 \times -4 \equiv -28 \equiv -11 \equiv 6 [17] && \text{donc on n'a pas le motif} \\ &6+5 \equiv 11 [17] \end{aligned}$$

$$\begin{aligned} >251_7 : \quad &2 \times 7 \equiv 14 \equiv -3 [17] \\ &-3+5 \equiv +2 [17] \\ &2 \times 7 \equiv 14 \equiv -3 [17] \\ &-3+1 \equiv -2 [17] \neq 11 \quad \text{Donc } >251_7 \neq >116_7 [17] \end{aligned}$$

$$\begin{aligned} >511_7 : \quad &5 \times 7 \equiv 35 \equiv 1 [17] \\ &1+1 \equiv 2 [17] \\ &2 \times 7 \equiv 14 \equiv -3 [17] \\ &-3+1 \equiv -2 [17] \neq 11 \quad \text{Donc } >511_7 \neq >116_7 [17] \end{aligned}$$

$$\begin{aligned} >116_7 : \quad &1 \times 7 \equiv 7 [17] \\ &7+1 \equiv 8 [17] \\ &8 \times 7 \equiv 56 \equiv 5 [17] \\ &5+6 \equiv 11 [17] \end{aligned} \quad \begin{aligned} &\text{On a } >116_7 \equiv >116_7 [17], \\ &\text{On teste et on a bien } t_4 = x_1, t_5 = x_2, t_6 = x_3. \text{ On a trouvé!} \end{aligned}$$