

KNUTH-MORRIS-PRATT.

Leçons: 907, 927.

Ref:

Théorème

L'algorithme de recherche (K)MP est correct et termine en $O(m+n)$, où l'on cherche un mot de longueur m dans un mot de longueur n .

Résumé

On cherche $u \in \Sigma^*$ dans $t \in \Sigma^*$ sur un alphabet Σ .

I - Déroulement de l'algorithme sur un exemple significatif.

II - Définition de l'automate des occurrences, reconnaissant $\Sigma^* u$.

III - Étude de la fonction Bord pour un calcul efficace de l'automate.

IV - Algorithme de Morris-Pratt + l'amélioration de Knuth.

I • Prénoms $t = A B A D A B A B A C$

$u = A B A C$.

• L'algorithme de Morris-Pratt donne:

$t = A B A D A B A \text{ } \underline{C} A C \dots$

1] $\underline{A}^{\vee} B^{\vee} \underline{A}^{\vee} C^{\vee}$

2] $\underline{A}^{\vee} B^{\vee} x$

3] $A^{\vee} x$

4] $A^{\vee} B^{\vee} A^{\vee} C^{\vee}$

• L'idée, c'est que quand il y a un problème, on peut décaler la fenêtre de plus d'un caractère.

• En se rend compte que le motif déjà matché avant décalage, c'est à la fois un préfixe de u et un suffixe de $t[k+1]$.

II. Le déroulement de KMP peut être vu comme le parcours de l'automate suivant : $A = (\mathbb{Q}, \varepsilon, u, \delta)$ où :

- * \mathbb{Q} est l'ensemble $\text{Pref}(u)$ des préfixes de u . "ce qu'on a déjà matché".
- * ε (le mot vide) est l'état initial "on n'a encore rien matché".
- * la fonction de transition δ est :

$$\delta(v, a) = \begin{cases} va & \text{si } va \in \text{Pref}(u) \\ \text{Bnd}(va) & \text{sinon. où } \text{Bnd}(w) \text{ est le plus long bord strict de } w. \end{cases}$$

- * u l'état final : "on a reconnu le mot entier".

Ce qu'il se passe, c'est que l'état indique la partie du mot qu'on a déjà matché, puis :

- soit on lit la lettre a qui suit v dans u , et on continue vers va .
 - soit on lit une autre lettre a , et là, on veut éviter de recommencer depuis ε , donc on recommence après un préfixe de va , qu'on avait déjà reconnu à la fin de ce que l'on vient de lire, c'est-à-dire à la fin de va . (sans que ça soit va !)
- C'est donc $\text{Bnd}(va) = \text{Le plus long bord strict de } va$.

III. Mais comment calculer $\text{Bnd}(va)$ efficacement ?

• Grâce au lemme :

$$\text{Bnd}(va) = \begin{cases} \text{Bnd}(v)a & \text{si } \text{Bnd}(v)a \in \text{Pref}(u) \\ \text{Bnd}(\text{Bnd}(v)a) & \text{sinon.} \end{cases}$$

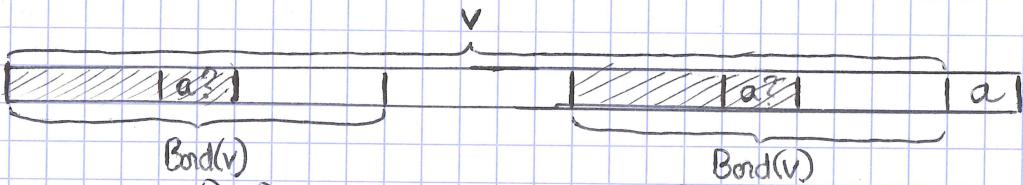
• Intuitivement, on a : → si $va \in \text{Pref}(u)$ on m'a qu'à continuer.

→ si $\text{Bnd}(v)a \in \text{Pref}(u)$, on a :

$\text{Bnd}(v)$	$ $	a	$ $	$\text{Bnd}(v)$	$ $	a
()						

il y a un a ici car $\text{Bnd}(v)a \in \text{Pref}(u)$ et $v \in \text{Pref}(u)$.

→ si non, $\text{Bord}(v)a \notin \text{Pref}(u)$, et on cherche un préfixe de u qui finisse par a . C'est-à-dire un mot au début de v , suivi d'un a , qui apparaît aussi à la fin de uva et donc de $\text{Bord}(v)a$.[⊕]



(les deux zones hachurées sont identiques)

[⊕] si ce n'était pas dans $\text{Bord}(v)a$, on serait dans le premier cas car on aurait un mot, bord de va , finissant par a , et donc un bord de v plus long que $\text{Bord}(v)$, c'est-à-dire ... $\text{Bord}(v)$!

• On peut maintenant prouver le lemme :

- soit $w = \text{Bord}(va)$, il s'écrit $w'a = w$ car c'est un suffixe de va .

→ $w \in \text{Pref}(va)$ donc $w' \in \text{Pref}(v)$

→ $w \in \text{Suff}(va)$ donc $w' \in \text{Suff}(v)$

⇒ de plus, w est maximal donc w' est le plus grand bord de v tel que $w'a$ est un préfixe de u .

- cela laisse deux possibilités :

1 - $\text{Bord}(v)a \in \text{Pref}(u)$, auquel cas $w' = \text{Bord}(v)$ par maximilité.

2 - w' est un préfixe strict de $\text{Bord}(v)$, et comme $w'a = \text{Bord}(va)$, w' est aussi à la fin de v et donc de $\text{Bord}(v)$.

De là, w' est un préfixe de $\text{Bord}(\text{Bord}(v))$, d'où $w = \text{Bord}(\text{Bord}(v)a)$.

• En itérant ce lemme, on obtient

$$\text{Bord}(v)a = \begin{cases} \text{Bord}(v)a & \text{si } \text{Bord}(v)a \in \text{Pref}(u) \\ \text{Bord}(\text{Bord}(v))a & \text{sinon si } \text{Bord}(\text{Bord}(v))a \in \text{Pref}(u) \\ \dots & \\ a & \text{sinon si } a \in \text{Pref}(u) \\ \epsilon & \text{sinon tout court.} \end{cases}$$

III]. On peut traduire tout ça en algorithme directement, mais pour cela, on va seulement garder en mémoire la taille du bord des prefixes de la forme $u[1 \dots k]$.

On définit aussi $\text{JT}(k) = |\text{Bord}(u[1 \dots k])|$.

CalculT(u):

$$\text{JT}(1) \leftarrow 0 \quad k \leftarrow 0.$$

Pour $i = 2 \dots |u|$:

Tant que $k > 0$ et $u[k+1] \neq u[i]$:

$$k \leftarrow \text{JT}(k).$$

Si $u[k+1] = u[i]$:

$$k \leftarrow k + 1$$

$$\text{JT}(i) \leftarrow k.$$

• Terminaison: on a $\text{JT}(k) < k$ donc le while termine et CalculT aussi.

• Correction: c'est le lemme 1, on calcule bien les bords !

• Complexité: on remarque que k est toujours positif, le while ne peut pas désincrémenter k de plus que sa valeur.

De plus, k augmente de 1 au plus une fois par passage dans

le for, dans lequel on passe n fois. On passe donc moins de n fois dans le while, d'où une complexité en $O(n)$, avec $n = |u|$.

- On peut (enfin) écrire l'algorithme de recherche : il suffit de chercher un bord de $u\#\bar{t}$, pour $\#\notin\Sigma$, de taille $|u|$.

Or en effet :

$\forall k \geq 0$, $J_t(k) \leq |u|$ et $J_t(k) = u$ si et seulement si u apparaît dans t en position $k - 2|u| - 1$.

- Or, d'une part, $\#\notin\Sigma$, donc $J_t(k) \leq |u|$ si et seulement si u apparaît une occurrence de $\#$ dans t .

- D'autre part, si $J_t(k) = |u|$, alors u est un préfixe et... un suffixe de $u\#\bar{t}[1..k]$. C'est donc un suffixe ~~de t~~ de $t[1..k-|u|-1]$ en position $k-2|u|-1$.

→ on trouve donc une occurrence de u dans t en :

$$\mathcal{O}(|u| + |\bar{t}| + 1) = \mathcal{O}(|u| + |t|).$$