

Tri des suffixes

Ref: Crochemore: Algorithmique du texte

L'objectif du classement est de calculer une permutation p des indices sur y qui satisfait la condition $y[p[0]..n-1] <_{\text{lex}} y[p[1]..n-1] <_{\text{lex}} \dots <_{\text{lex}} y[p[n-1]..n-1]$

Prop: Une telle permutation peut être calculée en temps $O(n \log n)$.

Démo • On définit $\text{pre}_k(u) = u[0..k-1]$ et on va comparer les pre_k des suffixes, pour k à croissance exponentielle. Plus précisément, on définit $R_k(i)$ le rang de $\text{pre}_k(y[i..n-1])$ dans la liste ordonnée des mots de l'ensemble $\{\text{pre}_k(u) : u \text{ suffix } y, u \neq \epsilon\}$

• On peut voir R_k comme une relation d'équivalence sur les positions via $i \equiv_k j$ si $R_k(i) = R_k(j)$. La propriété clé étant que \equiv_{k+1} est un raffinement de \equiv_k . On va donc partir de \equiv_1 , qui revient à identifier les lettres de y , pour arriver à \equiv_n , où chaque suffixe sera seul dans sa classe.

• Lemme: (Dedoublement). Pour k, i entiers, $k > 0$ et $0 \leq i < n$, $R_{2k}(i)$ est le rang du couple $(R_k(i), R_k(i+k))$ dans la liste classée en ordre lexico de ces couples (où on prendra $R_k(i) = -1$ pour $i \geq n$).

Démo: * Poser $R_k(i) = -1$ pour $i \geq n$ revient à considérer le mot infini ya^∞ , où a est une lettre strictement inférieure à toutes celles qui apparaissent dans y . Ainsi, pour $i \geq n$, le facteur de taille k correspondant est a^k , qui est bien inférieur à tous les autres, d'où la convention

* Par définition, $R_{2k}(i)$ est le rang de $\text{pre}_{2k}(y[i..n-1])$ dans la liste classée des facteurs de longueur $2k$ du mot ya^∞ . On note

$$u(i) = \text{pre}_k(y[i..n-1])$$

$$v(i) = \text{pre}_k(y[i+k..n-1])$$

On a $\text{prez}_k(y[i..n]) = u(i) \cdot v(i)$, d'où pour $0 \leq i \neq j < n$,
 l'inégalité $\text{prez}_k(y[i..n]) < \text{prez}_k(y[j..n])$ est équiv
 à $u(i) \cdot v(i) < u(j) \cdot v(j)$
 $\Leftrightarrow (u(i), v(i)) < (u(j), v(j))$
 $\Leftrightarrow (R_k(i), R_k(i+k)) < (R_k(j), R_k(j+k))$

D'où le résultat □

• On obtient alors l'algo suivant:

TRI-SUFFIXES (y, n)

```

    POUR  $n \leftarrow 0$  à  $n-1$  FAIRE  $p[n] \leftarrow n$ 
    k ← 1
    POUR  $i \leftarrow 0$  à  $n-1$  FAIRE  $R_1(i) \leftarrow$  rang de  $y[i]$  dans la liste classée des lettres de  $y$ 
    p ← TRI( $p, n, R_1, 0$ )
    i ← max  $R_1$ 
    TANT QUE  $i < n-1$  FAIRE
        p ← TRI( $p, n, R_k, k$ )
        p ← TRI( $p, n, R_k, 0$ )
         $R_k \leftarrow$  RANG( $p$ )
        i ← max  $R_k$ 
        k ← 2k
    Renvoyer p
  
```

Dédoublement
 $k \leftarrow 2k$

• $\text{Tri}(p, n, R, k)$ trie la suite d'entiers $p[0]+k, p[1]+k, \dots, p[n-1]+k$ suivant leur clé R . De plus Tri doit vérifier une condition de stabilité qui est : si $R[p[i]+k] = R[p[j]+k]$, alors $p[i] < p[j]$ ssi $p'[i] < p'[j]$.
 \hookrightarrow on peut prendre le tri par bac qui est linéaire

• Pour prouver la correction, on montre l'invariant de base suivant:
 $(\text{prez}_k(y[p[n]..n-1])) : n = 0..n-1$ est croissante.

• Pour la complexité : on a au plus $\log_2 n$ étapes, chacune en $O(n)$
 d'où le $O(n \log n)$ □