

Recherche des facteurs à distance d'édition au plus k

Réf: Crochemore: Algorithms on strings

Th: Étant donné un texte t et un motif x , on peut trouver toutes les occurrences à distance au plus k de x dans t : en temps $O(kn)$, où $n = |t|, m = |x|$.

Démo • Première couche de programmation dynamique: on relâche le problème en posant $D[i, j] = \min_p (\text{dist}(x[1..i], t[p..j])) : p = 1, \dots, j+1$
Ainsi, chaque j tq $D[i, j] \leq k$ correspond à une occurrence de x à distance $\leq k$ à la position $j-m$.

Et le tableau D vérifie la relation de récurrence suivante:

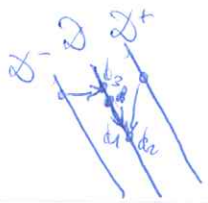
$$\begin{aligned} * D[0, j] &= 0 \quad \forall j \leq n & * D[i, 0] &= i \quad \forall i \leq m \\ * D[i, j] &= \min \begin{cases} D[i-1, j] + 1 & \leftarrow \text{suppression} \\ D[i, j-1] + 1 & \leftarrow \text{insertion} \\ D[i-1, j-1] + 1_{x[i] \neq t[j]} & \leftarrow \text{évaluation substitution} \end{cases} \end{aligned}$$

On obtient un algorithme dynamique qui résout le problème en temps $O(mn)$.

• Première idée d'amélioration: dans la relation de récurrence, on voit que si $x[i] = t[j]$, alors $D[i, j] = D[i-1, j-1]$ dans certains cas.
Pour pouvoir améliorer l'algorithme, il faut d'abord le modifier un peu: on voudrait calculer de manière itérative toutes les cases contenant j , pour $j = 0 \dots k$, ainsi on évite de remplir les cases dont la distance est $> k$.

• On définit ainsi une case l -spéciale: une case est l -spéciale si elle contient la valeur l et que aucune case plus basse sur sa diagonale ne contient la valeur l .

Voyons maintenant comment calculer les cases $(l+1)$ -spéciales étant donné les cases l -spéciales. Étant donné une diagonale et sa case l -spéciale, la case $(l+1)$ -spéciale provient soit de la diagonale au-dessus, soit de celle d'en



dessous, soit de la même diagonale. On calcule donc les trois cas candidates d_1, d_2, d_3 et on garde la plus basse. Puis on descend le long de la diagonale tant que on voit $l+1$, i.e tant que $x[j] = t[j]$.

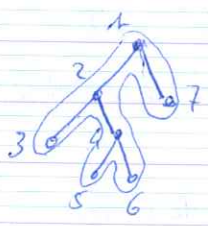
- Comme pour tout l , on a au plus $n+m$ cases l -spéciales (une par diag), si on cherche à faire la descente de diagonale en $O(1)$, notre algo tournera en $O(k \times (n+m)) = O(kn)$
nombre d'étapes cases l-spec

- On cherche donc à résoudre le problème suivant : étant donné i et j on veut calculer le plus ^{long} préfixe entre $x[i+1..m]$ et $t[j+1..n]$, qui sont respectivement suffixe de x et de t . On va donc construire l'arbre des suffixes de la chaîne $z = x\#t\$$ et on veut l'intersection commune entre les nœuds correspondant à $x[i+1..m]$ et $t[j+1..n]$, ce qui peut se faire en temps constant après un pré calcul linéaire (cf calcul de l'intersection commune via tableaux)

- Ainsi, étant donné x et t , on calcule l'arbre des suffixes en temps linéaire, puis le pré calcul pour l'intersection commune aussi linéaire. Enfin on lance l'algo dynamique qui fait des requêtes en $O(1)$, d'où une complexité en $O(k \cdot n)$.

□

Annexe: Arbre commun : on étiquette l'arbre selon un parcours séquentiel de gauche à droite, ce qui nous donne un tableau représentant l'arbre.



→ [1, 2, 3, 2, 4, 5, 4, 6, 4, 2, 1, 1]