

TP Informatique 11

Exercice 1

On définit la *distance de Levenshtein* ou *distance d'édition* notée lev entre deux chaînes de caractères X et Y par

$$\text{lev}(X, Y) = \begin{cases} \max(|X|, |Y|) & \text{si } \min(|X|, |Y|) = 0 \\ \text{lev}(X[1:], Y[1:]) & \text{si } X[0] = Y[0] \\ 1 + \min(\text{lev}(X[1:], Y), \text{lev}(X, Y[1:]), \text{lev}(X[1:], Y[1:])) & \text{sinon} \end{cases}$$

On peut montrer que la distance de Levenshtein est le nombre minimal d'opérations élémentaires (substitution, insertion ou suppression d'un caractère) pour transformer une chaîne de caractères en une autre.

Pour deux chaînes de caractères X, Y de tailles respectives n et m , on pose

$$\forall (i, j) \in \llbracket 0; n \rrbracket \times \llbracket 0; m \rrbracket \quad d(i, j) = \text{lev}(X[:i], Y[:j])$$

Les fonctions seront testées au fur et à mesure de leur saisie.

1. Écrire une fonction $\text{lev1}(X, Y)$ d'arguments X, Y des chaînes de caractères et qui renvoie la distance de Levenshtein entre X et Y en réalisant le calcul de manière descendante sans mémorisation.
2. Écrire une fonction $\text{lev2}(X, Y)$ d'arguments X, Y des chaînes de caractères et qui renvoie la distance de Levenshtein entre X et Y en réalisant le calcul de manière descendante avec mémorisation.
3. Préciser $d(i, 0)$ pour $i \in \llbracket 0; n \rrbracket$ et $d(0, j)$ pour $j \in \llbracket 0; m \rrbracket$.
4. Pour $(i, j) \in \llbracket 1; n \rrbracket \times \llbracket 1; m \rrbracket$, déterminer une expression de $d(i, j)$ fonction de $d(i-1, j)$, $d(i, j-1)$, $d(i-1, j-1)$ et de $X[i-1]$, $Y[j-1]$.
5. Compléter fonction $\text{lev3}(X, Y)$ d'arguments X, Y des chaînes de caractères pour qu'elle calcule la matrice $(d(i, j))_{(i,j) \in \llbracket 0; n \rrbracket \times \llbracket 0; m \rrbracket}$ de manière ascendante et renvoie $d(n, m)$.

Exercice 2

Un *alignement* de deux chaînes de caractères X et Y est formé d'un couple de chaînes \hat{X} et \hat{Y} de même taille, obtenues respectivement à partir des chaînes X et Y avec éventuellement des insertions d'espaces (caractère ' ') mais pas sur des positions identiques.

Un *désaccord* dans un alignement est une différence entre deux caractères sur un indice donné.

Un *alignement optimal* de deux chaînes de caractères a et b est un alignement de ces chaînes dont le nombre de désaccords est minimal.

On peut montrer que la distance de Levenshtein entre deux chaînes de caractères est le nombre minimal de désaccords dans un alignement optimal.

Les fonctions seront testées au fur et à mesure de leur saisie.

1. Écrire une fonction `lev_tab(X,Y)` d'arguments `X`, `Y` des chaînes de caractères et qui renvoie la matrice $(d(i,j))_{(i,j) \in [0;n] \times [0;m]}$ définie dans l'exercice précédent.
2. Écrire une fonction `lev_align(X,Y)` d'arguments `X`, `Y` des chaînes de caractères et qui renvoie un alignement optimal de ces chaînes en utilisant le résultat de `lev_tab(X,Y)`.
3. Compléter la fonction `lev_align_all(X,Y)` d'arguments `X`, `Y` des chaînes de caractères pour qu'elle renvoie la liste de tous les alignements optimaux de ces chaînes.