

COÛT D'ALIGNEMENT ET DISTANCE DE LEVENSHTEIN

Pour parler de coût d'alignement et pour la suite de distance entre les mots, on se ramène à donner un coût aux opérations.

SUB On considère que la substitution présente un coût qui ne dépend que des lettres que l'on échange, et non de la position concernée.

De plus on considère assez naturellement ce coût nul si on substitue $a \rightarrow a$, et symétrique.

On le note $\text{sub}(a, b)$

DEL On considère le coût de cette opération constante, on le note del
INS _____ ins

Rq Dans certains modèles il peut être intéressant de considérer qu'en début ou en fin, les suppressions et les insertions sont moins coûteuses, notamment lorsque les deux mots que l'on cherche à aligner sont de longueurs très différentes. Cf. alignem't ac mèches.

ex

$x = \text{ABRACADABRA}$

$y = \text{BRAS}$

Comme $|x| > |y|$ on "fait une prime" sur la suppression en début ou en fin de mot.

aurait faire $\begin{array}{c} \text{A BRA - CADABRA} \\ \text{--- BRAS -----} \end{array}$ ou $\begin{array}{c} \text{ABRACADABA -} \\ \text{----- BRAS} \end{array}$

seront moins coûteuses que $\begin{array}{c} \text{ABRACADABRA -} \\ \text{B ----- RAS} \end{array}$.

Def

Le coût d'un alignement est la somme des coûts des opérations d'édition auxquelles il correspond, on le note $C(\tilde{x}, \tilde{y})$.

Pté

$$C(\tilde{x}, \tilde{y}) = \sum_{i=1}^{|\tilde{x}|} c(x_i, \tilde{y}_i) \text{ où } c(a, b) = \begin{cases} \text{del si } b = - \\ \text{ins si } a = - \\ \text{sub}(a, b) \text{ sinon} \end{cases}$$

ex

$$\begin{aligned} \text{del} &= \text{ins} = 1 \\ \text{sub}(a, b) &= 1 \text{ si } a \neq b \\ &= 0 \text{ si } a = b \end{aligned}$$

$$\begin{aligned} \tilde{x} &= \text{SALLE} \\ \tilde{y} &= \text{SE-L-S} \end{aligned}$$

$$C(\tilde{x}, \tilde{y}) = 0 + 1 + 1 + 0 + 1 + 1 = 4$$

Def

On introduit la "distance" de Levenshtein comme étant pour deux mots x et y .

$$\text{Lev}(x, y) = \min \{ C(\tilde{x}, \tilde{y}) \mid (\tilde{x}, \tilde{y}) \text{ est un alignement de } (x, y) \}$$

C'est le coût minimal d'un alignement entre x et y .

C'est aussi le coût minimal d'une succession d'opérations transformant x en y .

ex

$$\begin{aligned} \text{pour } x &= \text{SALLE} \\ y &= \text{SELS} \end{aligned}$$

$$C\left(\begin{array}{c} \tilde{x} = \text{SA L L E} \\ \tilde{y} = \text{SE - L S} \end{array}\right) = 0 + 1 + 1 + 0 + 1 = 3 < 4.$$

On ne peut pas faire moins donc $\text{Lev}(\text{SALLE}, \text{SELS}) = 3$

Pti

Si $\text{del} = ins > 0$ et si sub est une distance sur Σ
alors Lev est effectivement une distance sur Σ^*

Preuve • sub en tant que distance est symétrique.

Parce que $\text{del} = ins$ on est assuré que c est symétrique.

De suite C et Lev sont symétriques.

- Si $\text{Lev}(x, y) = 0$, il existe (\tilde{x}, \tilde{y}) un alignement de (x, y) tel que $C(\tilde{x}, \tilde{y}) = 0$, alors $\forall i \in [1..|\tilde{x}|] \quad c(\tilde{x}_i, \tilde{y}_i) = 0$
Comme del et ins sont > 0 , cela implique $\tilde{x}_i \neq -$ et $\tilde{y}_i \neq -$.
Donc nous $x = \tilde{x}$ et $y = \tilde{y}$.
De plus comme sub est une distance $\text{sub}(\tilde{x}_i, \tilde{y}_i) = 0 \iff \tilde{x}_i = \tilde{y}_i$.
D'où $\tilde{x} = \tilde{y}$. Finalement $x = y$.
D'où le caractère défini de Lev .

- Si on considère x, y, z trois mots, en enchainant les opérations justifiant $\text{Lev}(x, y)$, et celles justifiant $\text{Lev}(y, z)$, on passe de x à z pour un coût $\text{Lev}(x, y) + \text{Lev}(y, z)$. Or $\text{Lev}(x, z)$ est moindre (définition).

D'où l'inégalité triangulaire.

Ainsi LEV est bien une distance.