

# CALCUL DE LA DISTANCE DE LEVENSHTEIN ET D'UN ALIGNEMENT

Considérons les coûts *del* et *ins* finies.

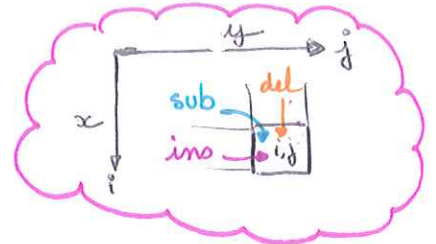
Considérons aussi sur une distance sur  $\Sigma$  finie.

Soient  $x = x_1 \dots x_n$  et  $y = y_1 \dots y_m$  deux mots sur  $\Sigma$ .

On introduit  $l_{i,j} = \text{Lev}(x_1 \dots x_i, y_1 \dots y_j)$  pour  $(i,j) \in [0..n] \times [0..m]$ .

Pté

•  $l_{m,m} = \text{Lev}(x,y)$



•  $\forall i \in [0..m] \quad l_{i,0} = i \times \text{del}$   
 •  $\forall j \in [0..m] \quad l_{0,j} = j \times \text{ins}$

•  $\forall (i,j) \in [1..m] \times [1..m] \quad l_{i,j} = \min \begin{cases} l_{i-1,j-1} + \text{sub}(x_i, y_j) \\ l_{i-1,j} + \text{del} \\ l_{i,j-1} + \text{ins} \end{cases}$

Preuve

a) Définition seule

b) passer de  $x_1 \dots x_i$  à  $\epsilon$  requiert au minimum de supprimer  $i$  lettres, et cela suffit, d'où  $\text{Lev}(x_1 \dots x_i, \epsilon) = i \times \text{del}$

c) à l'inverse passer de  $\epsilon$  à  $y_1 \dots y_j$  requiert au minimum d'insérer  $j$  lettres, et comme cela suffit si on insère les bonnes lettres dans l'ordre c-à-d  $y_1$  en pos 1,  $y_2$  en pos 2 ... jusqu'à  $y_j$  en position  $j$ , on a  $\text{Lev}(\epsilon, y_1 \dots y_j) = j \times \text{ins}$ .

d) Considérons  $(\tilde{x}, \tilde{y})$  un alignement de  $x_1 \dots x_i$  et  $y_1 \dots y_j$  de longueur  $a$ .  
 $(\hat{x}, \hat{y}) = (\tilde{x}, \tilde{x}_a, \tilde{y}, \tilde{y}_a)$  reste un alignement, mais pas  $x$  et  $y$ .

• Si  $\tilde{y}_a = -$  alors  $(\hat{x}, \hat{y})$  est un alignement de  $(x_1 \dots x_i, y_1 \dots y_j)$  qui au minimum sera de coup  $l_{i,j}$ , auquel il faudra ajouter le coût de la dernière opération :  $c(\hat{x}_a, \hat{y}_a) = c(\hat{x}_a, -) = \text{del}$   
 donc au min  $l_{i,j} = l_{i-1,j} + \text{del}$

• Si  $\tilde{x}_a = -$  alors  $(\hat{x}, \hat{y})$  est un alignement de  $(x_1 \dots x_i, y_1 \dots y_{j-1})$  dont le coût est au minimum  $l_{i, j-1}$ , donc au min  $l_{i,j} = l_{i, j-1} + \text{ins}$

• Si  $\tilde{x}_a \neq -$  et  $\tilde{y}_a \neq -$   $(\hat{x}, \hat{y})$  est un alignement de  $(x_1 \dots x_{i-1}, y_1 \dots y_{j-1})$ , et la dernière opéra est une substitution de coût  $\text{sub}(x_i, y_j)$ , donc  $l_{i,j}$  vaut au minimum  $l_{i-1, j-1} + \text{sub}(x_i, y_j)$

d) (suite). On conclut en se rappelant que  $l_{i,j}$  est le coût d'un alignement minimal, il faut donc considérer le min entre ces trois options.

## Algo

La propriété précédente indique et prouve l'algorithme de programmation dynamique suivant pour calculer la distance de Levenshtein entre deux mots.

LEV(x, y)

$n \leftarrow$  longueur de  $x$

$m \leftarrow$  longueur de  $y$ .

$L \leftarrow$  tableau indexé par  $[0..m] \times [0..m]$

Pour  $i$  allant de 0 à  $n$

$L[i, 0] \leftarrow i \times \text{del}$

Pour  $j$  allant de 0 à  $m$

$L[0, j] \leftarrow j \times \text{ins}$

Pour  $i$  allant de 1 à  $n$

Pour  $j$  allant de 1 à  $m$

$L[i, j] \leftarrow \min(L[i-1, j] + \text{del}; L[i, j-1] + \text{ins}; L[i-1, j-1] + \text{sub}(x_i, y_j))$

Retourner  $L[n, m]$

## Complexité

\* Cet algorithme a une complexité temporelle en  $O(m \times m)$ .

\* tel qu'il est écrit sa complexité temporelle est en  $O(m \times m)$ , mais on peut facilement l'améliorer, en ne gardant en mémoire que la ligne courante et la précédente.

Donc la complexité spatiale est plutôt en  $O(m)$



# Amélioration

Quitte à calculer à chaque étape le minimum entre les 3 opérations possibles, autant retenu laquelle (une de celle en fait) on a choisi, afin de pouvoir finalement retourner un alignement réalisant la distance annoncée

## Algo bis

LEV\_bis(x, y)

$n \leftarrow |x|$   
 $m \leftarrow |y|$

$L \leftarrow$  tableau indexé par  $[0..n] \times [0..m]$

$M \leftarrow$  \_\_\_\_\_  $[1..n] \times [1..m]$

Pour  $i$  allant de 0 à  $n$

$L_{i,0} \leftarrow i \times \text{del}$

Pour  $j$  allant de 0 à  $m$

$L_{0,j} \leftarrow j \times \text{ins}$

Pour  $i$  allant de 1 à  $n$

Pour  $j$  allant de 1 à  $m$

$D \leftarrow L[i-1, j] + \text{del}$

$I \leftarrow L[i, j-1] + \text{ins}$

$S \leftarrow L[i-1, j-1] + \text{sub}(x_i, y_j)$

Si  $\min(I, D, S) = I$

alors  $\begin{cases} L[i, j] \leftarrow I \\ M[i, j] \leftarrow (i, j-1) \end{cases}$

Si  $\min(I, D, S) = D$

alors  $\begin{cases} L[i, j] \leftarrow D \\ M[i, j] \leftarrow (i-1, j) \end{cases}$

Si  $\min(I, D, S) = S$

alors  $\begin{cases} L[i, j] \leftarrow S \\ M[i, j] \leftarrow (i-1, j-1) \end{cases}$

$\tilde{x} \leftarrow \varepsilon ; i \leftarrow n$   
 $\tilde{y} \leftarrow \varepsilon ; j \leftarrow m$

tant que  $i > 1$  ou  $j > 1$

match  $M[i, j]$  with

$I(i, j-1) \rightarrow \begin{cases} \tilde{x} \leftarrow \tilde{x} \\ \tilde{y} \leftarrow y_j \tilde{y} \end{cases}$

$I(i-1, j) \rightarrow \begin{cases} \tilde{x} \leftarrow x_i \tilde{x} \\ \tilde{y} \leftarrow \tilde{y} \end{cases}$

$I(i-1, j-1) \rightarrow \begin{cases} \tilde{x} \leftarrow x_i \tilde{x} \\ \tilde{y} \leftarrow y_j \tilde{y} \end{cases}$

$(i, j) \leftarrow M[i, j]$

retourner  $(\tilde{x}, \tilde{y}, L[n, m])$

Cplx complexité temporelle  $O(n \times m)$   
matrice  $O(n \times m)$

esc

$\Sigma$  l'alphabet romain.

$\text{ins} = \text{del} = 1$

$$\text{sub}(\alpha, \beta) = \begin{cases} 0 & \text{si } \alpha = \beta \\ 1 & \text{si } \alpha \text{ et } \beta \text{ sont 2 voyelles} \\ 1 & \text{consonnes} \\ 2 & \text{sinon} \end{cases}$$

NB c'est une distance!

$x = \text{RAPE} \quad m = 4$   
 $y = \text{LAPIN} \quad m = 5$

	$j \rightarrow$	0	1	2	3	4	5
		$\epsilon$	L	A	P	i	N
$i \downarrow$	0	$\epsilon$	1	2	3	4	5
1	R	1	1	2	3	4	5
2	A	2	2	1	2	3	4
3	P	3	3	2	1	2	3
4	$\epsilon$	4	4	3	2	2	3

ex  $1 = \min(1 + \text{ins}, 1 + \text{del}, 0 + \text{sub})$   
 $\text{sub}(L, R)$

$\rightarrow \text{LEV}(x, y) = 3$

M		1	2	3	4	5
1		0,0	(1,1)	(1,2)	(1,3)	(1,4)
2		(1,1)	(1,1)	(2,2)	(2,3)	(2,4)
3		(2,1)	(2,2)	(2,2)	(3,3)	(3,4)
4		(3,1)	(3,2)	(3,3)	(3,3)	(4,4)

$\tilde{x} = \epsilon, \tilde{y} = \epsilon, i = 4, j = 5$

comme  $M[4,5] = (4,4) = (4, 5-1)$   $\begin{cases} \tilde{x} = - \\ \tilde{y} = y_5 = N \end{cases}$

comme  $M[4,4] = (3,3) = (4-1, 4-1)$   $\begin{cases} \tilde{x} = x_4 = \epsilon \\ \tilde{y} = y_4 = N = IN \end{cases}$

puis de  $\tilde{m} \begin{cases} \tilde{x} = PE_- \\ \tilde{y} = PiN \end{cases}$

puis  $\begin{cases} \tilde{x} = APE_- \\ \tilde{y} = APiN \end{cases}$

puis  $\begin{cases} \tilde{x} = RAPE_- \\ \tilde{y} = LAPiN \end{cases}$

D'où

RAPE_	3
LAPIN	
1 0 0 1 1	