

CALCUL DE LA PLUS LONGUE SOUS-CHAÎNE COMMUNE

Goedemaer
p242
(7.3).

Pour intro : consultez LIEN ENTRE LA PLUS LONGUE SOUS-CHAÎNE COMMUNE ET LA DISTANCE DE LEVENSHTAIN.

En s'inspirant de ce qu'on a fait pour calculer la distance de Levenshtein on fait un algo de programmation dynamique.

Cependant notez quelques changements :

- Les trois options **INS**, **DEL** et **SUB** n'ont plus le même statut. On ne réalise une substitution que si elle est triviale, en fait on la réalise soit $x_i = y_j$. Si $x_i \neq y_j$ on départage **INS** et **DEL** par un max.
- Un max plutôt qu'un min? Oui parce que maintenant on stocke la longueur de la PLSC construite jusqu'à là. Se rappeler que maximiser la longueur d'une SSC c'est minimiser le coût d'alignement.
- L'amélioration consistant à ne retenir que la colonne courante et la précédente pour que l'espace utilisé soit réduit de $O(n \times m)$ à $O(n)$ est enfin réalisée.

Algo 1

PLSCC(x, y)

$m \leftarrow |x|$

$m \leftarrow |y|$

$C_1 \leftarrow$ tableau indexé de 0 à m init à 0

// colonne "mineure"

$C_2 \leftarrow$

// colonne "courante"

Pour j allant de 1 à m

Pour i allant de 1 à m

Si $x_i = y_j$

alors $C_2[i] \leftarrow C_1[i-1] + 1$ // une lettre de plus dans notre SSC

sinon $C_2[i] \leftarrow \max(C_1[i], C_2[i-1])$

$C_1 \leftarrow C_2$ // la colonne courante est "mise en mémoire".

Retourner C_2 . // on pourrait a priori retourner $C_1[m]$ mais on aura besoin de C_1 encore pour la suite.

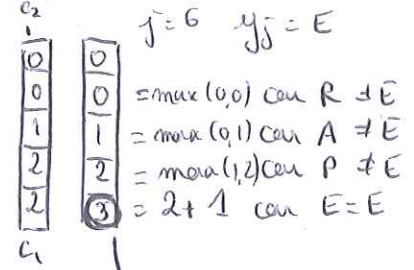
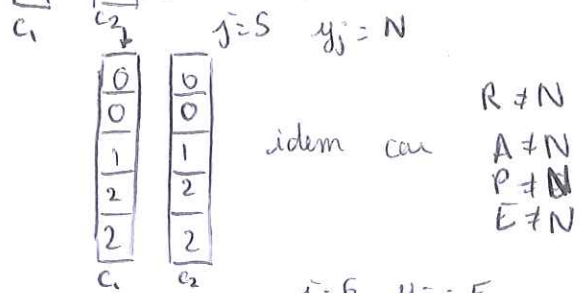
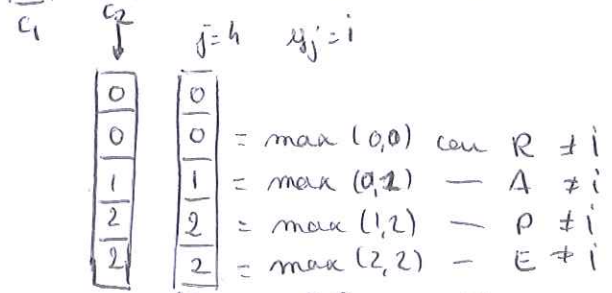
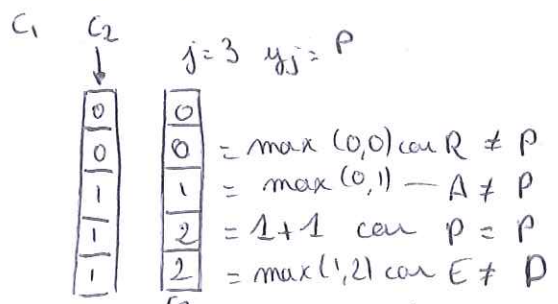
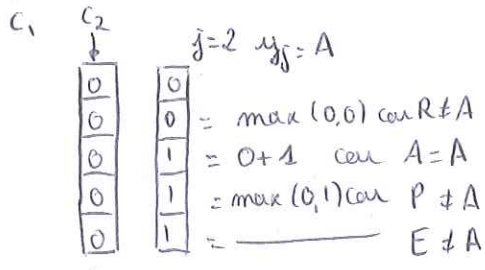
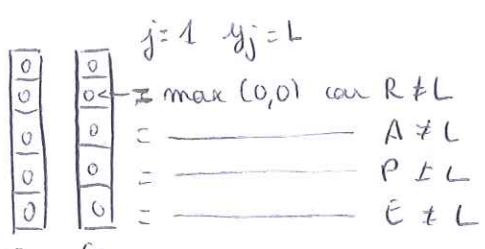
Complexité

Complexité temporelle en $O(m \times m)$

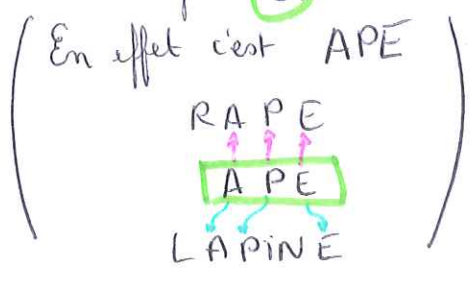
Complexité spatiale en $O(m)$

ex

x = RAPE
y = LAPINE
m = 4
n = 6



→ la plus longue sous chaîne est de longueur 3.



Comme précédemment on va au delà de la longueur de la PLSCC et on demande une PLSCC.

L'idée est de stocker d'où vient notre max ou si on a fait une substitution triviale dans un tableau qui nous permettra de remonter la solution, l'alignement qui justifie cette longueur de SCC.

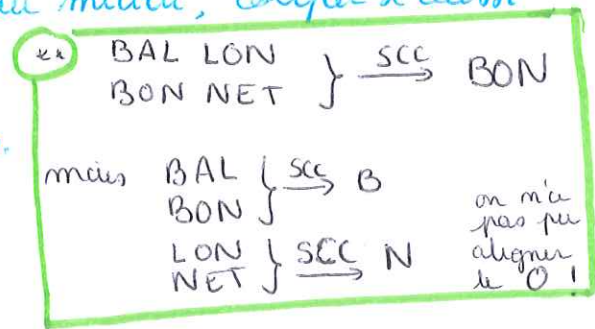
C'est la même idée que lorsqu'on veut fournir un alignement minimal, justifiant la distance si diversifiée annoncée.

Mais cet algorithme nécessiterait un espace mémoire en $O(n \times m)$, puisque l'astuce des 2 colonnes n'est plus possible, parce qu'on veut garder tout le tableau en mémoire pour remonter.

Pour palier à cela on utilise le paradigme pour séparer, on va chercher à aligner des morceaux plus petits.

Mais Δ si on veut couper y en deux au milieu, couper x aussi au milieu n'est pas une bonne idée.

On ^{doit} alors calculer là où on doit couper x , c'est à quoi sert P dans l'algo ci-dessous.



Algo 2

PLSCC-2(x, y)

$n \leftarrow$ longueur(x)
 $m \leftarrow$ longueur(y)
 $d \leftarrow \lfloor m/2 \rfloor$

Si $m=0$ ou $n=0$
 alors retourner E

Si $m=1$
 alors pour j allant de 1 à n
 si $y_j = x_1$
 alors retourner x_1
 retourner E

Si $n=1$
 alors pour i allant de 1 à m
 si $x_i = y_1$
 alors retourner y_1
 retourner E

$C_1, C_2, P_1, P_2 \leftarrow$ tableaux indexés de 0 à m initialisés à 0

$C_1 \leftarrow$ PLSCC-1(x, y, -y_d)

Pour i allant de d à m
 $P_1[i] \leftarrow i$

Pour j allant de $d+1$ à m

Pour i allant de 1 à n

Si $x_i = y_j$
 alors $C_2[i] \leftarrow C_1[i-1]$
 $P_2[i] \leftarrow P_1[i-1]$

sinon si $C_2[i-1] > C_1[i]$
 alors $C_2[i] \leftarrow C_2[i-1]$
 $P_2[i] \leftarrow P_2[i-1]$

sinon $C_2[i] \leftarrow C_1[i]$
 $P_2[i] \leftarrow P_1[i]$

$C_1 < C_2, C_2 < 0$
 $P_1 < P_2, P_2 < 0$

$k \leftarrow P_2[m]$ // oublier C_1, C_2, P_1, P_2

$u \leftarrow$ PLSCC-2($x_1 - x_k, y_1 - y_d$)

$v \leftarrow$ PLSCC-2($x_{k+1} - x_n, y_{d+1} - y_m$)
 retourner uv.

sub triviale

DEL

INS

exc

$x = \text{GIRAFE } m=6$
 $y = \text{GRAFITI } m=7$

↘ = substitution triviale
 ↓ = suppression (doxc)
 → = insertion (doxc)



	/	G	R	A	F	i	T	i
/	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1
i	0	1	1	1	1	2	2	2
R	0	1	2	2	2	2	2	2
A	0	1	2	3	3	3	3	3
F	0	1	2	3	4	4	4	4
E	0	1	2	3	4	4	4	4

R_9 Si on coupait aussi x au milieu on ne pourrait pas aligner le A:

$\left. \begin{matrix} \text{GIR} \\ \text{GRA} \end{matrix} \right\} \rightarrow \begin{matrix} \text{G i R} \\ \text{G - R A} \end{matrix}$
 $\left. \begin{matrix} \text{AFE} \\ \text{FITi} \end{matrix} \right\} \rightarrow \begin{matrix} \text{A F E} \\ \text{- P - i T i} \end{matrix}$

PLSCC_2(x,y)

$m=6; m=7$
 $d = [3,5] = 3$

$C_1 = \text{PLSCC}_1(\text{GIRAFE}, \text{GRA}) = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 2 \\ 3 \\ 3 \\ 3 \end{bmatrix}; P_1 = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix} \dots C_i = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 2 \\ 3 \\ 4 \\ 4 \end{bmatrix}; P_i = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 4 \\ 4 \\ 4 \end{bmatrix}$

$k=4$

$u_1 = \text{PLSCC}_2(\text{GIRA}, \text{GRA})$

$m=4; m=3$

$d=1$

$C_1 = \text{PLSCC}_1(\text{GIRA}, \text{G}) = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}; P_1 = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \dots C_i = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 2 \\ 3 \end{bmatrix}; P_i = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 2 \\ 2 \end{bmatrix}$

$k=2$

$u_{1,1} = \text{PLSCC}(\text{Gi}, \text{G}) = \text{G}$ car $m=1$ et G apparaît ds Gi

$u_{1,2} = \text{PLSCC}(\text{RA}, \text{RA}) \dots = \text{RA}$

= GRA

$u_2 = \text{PLSCC}_2(\text{FE}, \text{FITi})$

$m=2; m=3; d=1; C_1 = \text{PLSCC}_1(\text{FE}, \text{F}) = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}; P_1 = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \dots C_i = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}; P_i = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$

$k=1$

$u_{2,1} = \text{PLSCC}(\text{F}, \text{FITi}) = \text{F}$ car $m=1$ et F apparaît ds FITi

$u_{2,2} = \text{PLSCC}(\text{E}, \text{FITi}) = \text{E}$ ————— E n' ————— pas —————

= F

= GRAF

Manque complexité Δ