

# LIEN ENTRE PLUS LONGUE SOUS-CHAÎNE COMMUNE ET DISTANCE DE LEVENSHTEIN

en plus  
Cochemare  
7.3 p 242.

Pb Le problème de la plus longue sous chaîne commune (PLSCC) est de trouver, étant données en entrée deux chaînes = mots  $x$  et  $y$ , un mot  $z$  qui soit à la fois sous-chaîne / sous-mot de  $x$  et sous-mot de  $y$ , et de longueur maximale.

Dans un premier temps on pourra se chercher que la longueur maximale d'une sous-chaîne commune.

⚠ sous mot  $\neq$  facteur      sous-mot  $\approx$  suite extraite.

Déf Soit  $x = x_1 \dots x_n$  un mot.  
 $z = z_1 \dots z_l$  est un sous mot de  $x$   
 $\Leftrightarrow$  il existe une fonction strictem<sup>t</sup> croissante  $\varphi$  de  $[1..l]$  vers  $[1..n]$  telle que  $\forall i \in [1..l] z_i = x_{\varphi(i)}$ .  
aut dit il existe une "extra"  $\ell$  telle que  $z = x_{\varphi(1)} \dots x_{\varphi(\ell)}$

NB : facteur  $\Rightarrow$  sous-mot

ex LABO est un sous mot de LABRADOR

ici  $m=8, l=4$  et  $\varphi = \begin{cases} 1 \mapsto 1 \\ 2 \mapsto 2 \\ 3 \mapsto 3 \\ 4 \mapsto 8 \end{cases}$

mais LABO n'est pas un facteur de LABRADOR.

En revanche LAB, OR, ABRA sont des facteurs de LABRADOR et donc des sous-mots.

# Alignement → SSC

Un alignement fournit naturellement une sous-chaine commune (SSC) il suffit de ne considérer que les lettres qui sont "support" d'une substitution triviale (changer une lettre en elle-même).

ex  $x = \text{BLABLABLA} \quad m=9$   
 $y = \text{ABRACAPABRA} \quad m=12$

$\tilde{x} : \ominus \text{B} \text{L} \ominus \text{A} \text{B} \text{L} \text{A} \ominus \ominus \text{B} \text{L} \text{A}$   
 $\tilde{y} : \text{A} \text{B} \ominus \text{R} \text{A} \text{C} \ominus \text{A} \text{D} \text{A} \text{B} \text{R} \text{A}$

↳ BAABA est une SSC à  $x$  et  $y$ .

Cependant si un alignement ne contient pas de substitution non triviale, cela s'exprime encore plus facilement: il suffit de considérer les lettres qui ne sont pas-face à un trait.

Considérons  $\text{ins} = \text{del} = 1$  et  $\text{sub}(a,b) = \begin{cases} 3 & \text{si } a \neq b \\ 0 & \text{si } a = b \end{cases}$  jusqu'à la fin de cette fiche.

Puisque la substitution\* est strictement plus chère qu'une insertion et une suppression, un alignement minimal sans substitution\*  
 \*= non triviale. Sous-entendu dans la suite.

ex  $x = \text{BLABLABLA} \quad m=9$   
 $y = \text{ABRACADABRA} \quad m=12$

$\tilde{x} \ominus \text{B} \text{L} \ominus \text{A} \text{B} \text{L} \text{A} \ominus \ominus \text{B} \text{L} \text{A}$   
 $\tilde{y} \text{A} \text{B} \ominus \text{R} \text{A} \text{C} \ominus \text{A} \text{D} \text{A} \text{B} \ominus \text{R} \text{A}$   
 $\cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot$

→ BAABA  
 $L=5$

↳ coût de l'alignement  $\tilde{C}(x,y) = \underbrace{4}_{\text{Nb del}} + \underbrace{6}_{\text{Nb ins}} = 10$

# SCC -> alignement

Et l'inverse une SCC  $z$  de  $x$  et  $y$  fournit un alignement de  $x$  et  $y$ .

On écrit  $z = z_1 \dots z_L$ . Il existe, puisque  $z$  est sous-mot de  $x$  et de  $y$ , deux fonctions  $\varphi$  et  $\psi$ , st  $\rightarrow$  telles que  $z = (x_{\varphi(i)})_{i \in \{1, \dots, L\}} = (y_{\psi(i)})_{i \in \{1, \dots, L\}}$ .

Alors on peut décomposer  $\begin{cases} x = \mu_0 x_{\varphi(1)} \mu_1 \dots x_{\varphi(L)} \mu_L & \text{où} \\ y = \nu_0 y_{\psi(1)} \dots y_{\psi(L)} \nu_L \end{cases}$

les  $(\mu_i)$  et les  $(\nu_i)$  sont des mots, éventuellement vides.

Alors on a l'alignement  $\begin{cases} \tilde{x} = \mu_0 \text{---} x_{\varphi(1)} \mu_1 \text{---} \dots \text{---} x_{\varphi(L)} \mu_L \text{---} \\ \tilde{y} = \text{---} \nu_0 y_{\psi(1)} \text{---} \nu_1 \dots \dots y_{\psi(L)} \text{---} \nu_L \end{cases}$

NB: L'alignement n'est ni unique ni canonique.

ESC

$x = \text{BLA BLA BLA}$

$\varphi$

$z = \text{BAABA}$

$\mu_0 = \epsilon$   
 $\mu_1 = L$   
 $\mu_2 = BL$   
 $\mu_3 = \epsilon$   
 $\mu_4 = L$   
 $\mu_5 = \epsilon$

$y = \text{ABRACADABRA}$

$\psi$

$z = \text{BAABA}$

$\nu_0 = A$   
 $\nu_1 = R$   
 $\nu_2 = C$   
 $\nu_3 = DA$   
 $\nu_4 = R$   
 $\nu_5 = \epsilon$

d'où l'alignement

$\tilde{x}$	-	B	L	-	A	B	L	-	A	-	-	B	L	-	A
$\tilde{y}$	A	B	-	R	A	-	-	C	A	D	A	B	-	R	A

## Leve -> L

Si  $z$  est une SCC de  $x$  et  $y$  "correspondant" à l'alignement  $(\tilde{x}, \tilde{y})$   
 alors  $2|z| = |x| + |y| - C(\tilde{x}, \tilde{y})$

Preuve Considérons  $N =$  le nombre total de symbole dans l'alignement  $(\tilde{x}, \tilde{y})$  (ie sur les 2 lignes)  
 $N = |x| + |y| + N_{ins} + N_{del}$  // en comptant les lettres puis les traits  
 $= 2(|z| + N_{ins} + N_{del})$  // en regardant comment chaque "colonne" est colorié  
 Donc  $2|z| = |x| + |y| + N_{ins} - 2N_{ins} + N_{del} - 2N_{del} = |x| + |y| - (1 \times N_{del} + 1 \times N_{ins}) = C(\tilde{x}, \tilde{y})$

Cr

$$L(x,y) = \frac{|x| + |y| - \text{Lev}(x,y)}{2}$$

où  $L$  désigne la longueur d'une PLSCC

Preuve 1)  $L(x,y) \geq \frac{|x| + |y| - \text{Lev}(x,y)}{2}$   $n = |x|, m = |y|$

Pour chaque alignement de coût  $C$  il existe une PLSCC de longueur  $L$  où  $2L = n + m - C$ . Comme  $L \leq L(x,y)$  par def. on a  $n + m - C \leq 2L(x,y)$ . En particulier pour un alignement de coût min on a  $n + m - \text{Lev}(x,y) \leq 2L(x,y)$  soit  $L(x,y) \geq \frac{n + m - \text{Lev}(x,y)}{2}$

2)  $\frac{|x| + |y| - \text{Lev}(x,y)}{2} \geq L(x,y)$ .

Etant donné une SCC de taille  $L$  on peut construire un alignement de coût  $C$  où  $2L = n + m - C$ . Comme  $C \geq \text{Lev}(x,y)$  on a  $2L \leq n + m - \text{Lev}(x,y)$ . En particulier pour une PLSCC on a  $2L(x,y) \leq \frac{n + m - \text{Lev}(x,y)}{2}$

On conclut par double inégalité.

CCP

Calculer la distance de Levenshtein revient à calculer la longueur d'une plus longue sous-suite commune lorsque qu'on ait fixé

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

$$\text{sub}(a,b) = \begin{cases} 3 & \text{si } a \neq b \\ 0 & \text{si } a = b \end{cases}$$

De plus étant donné un alignement réalisant la distance de Levenshtein on a une PLSCC, et réciproquement étant donné une SCC de longueur  $L(x,y)$ , on a un alignement de coût minimal.

minimiser  $C(\hat{x}, \hat{y}) \Leftrightarrow$  maximiser  $|z|$