



---

# Alignement de séquences

---

LU3IN003

Léa Cohen-Solal

Lysa Rudy

*Dirigé par Olivier Spanjaard*

# Table des matières

<b>1</b>	<b>Le problème d'alignement des séquences</b>	<b>2</b>
1.1	Alignement de deux mots . . . . .	2
<b>2</b>	<b>Algorithmes pour l'alignement de séquences</b>	<b>4</b>
2.1	Méthode naïve par énumération . . . . .	4
<b>3</b>	<b>Programmation dynamique</b>	<b>7</b>
3.1	Calcul de la distance d'édition par programmation dynamique . . . . .	7
3.2	Calcul d'un alignement optimal par programmation dynamique . . . . .	9
<b>4</b>	<b>Amélioration de la complexité spatiale du calcul de la distance</b>	<b>14</b>
<b>5</b>	<b>Amélioration de la complexité spatiale du calcul d'un alignement optimal par la méthode "diviser pour régner"</b>	<b>17</b>

Le code des taches A, B, C et D se trouve dans les fichiers *tache\_A.py*, *tache\_B.py*, *tache\_C.py* et *tache\_D.py* respectivement . Les tests des fonctions sont dans les fichiers *essai\_A.py*, *essai\_B.py*, *essai\_C.py* et *essai\_D.py*

# 1 Le problème d'alignement des séquences

## 1.1 Alignement de deux mots

### Question 1

Montrons que si  $(\bar{x}, \bar{y})$  et  $(\bar{u}, \bar{v})$  sont des alignements de  $(x, y)$  et  $(u, v)$ , alors  $(\bar{x}.\bar{u}, \bar{y}.\bar{v})$  est un alignement de  $(x.u, y.v)$ .

$(\bar{x}, \bar{y}) = (x_1 \dots x_k, y_1 \dots y_k)$  un alignement de  $(x, y) = (x'_1 \dots x'_{k'}, y'_1 \dots y'_{k''})$

$(\bar{u}, \bar{v}) = (u_1 \dots u_n, v_1 \dots v_n)$  un alignement de  $(u, v) = (u'_1 \dots u'_{n'}, v'_1 \dots v'_{n''})$  avec  $n', n'' \leq n$  et  $k', k'' \leq k$ .

$(x.u) = x'_1 \dots x'_{k'}.u'_1 \dots u'_{n'}$  avec  $|x.u| = n' + k'$

On obtient  $x$  en supprimant tous les gaps de  $\bar{x}$  et  $u$  en supprimant tous les gaps de  $\bar{u}$ . Comme on le remarque  $x.u$  s'obtient de la même manière.

Soient  $\bar{X} = (\bar{x}.\bar{u})$  et  $\bar{Y} = (\bar{y}.\bar{v})$

— Montrons que  $\pi(\bar{x}.\bar{u}) = x.u$  et  $\pi(\bar{y}.\bar{v}) = y.v$

On sait que  $\pi(\bar{x}) = x$ ,  $\pi(\bar{y}) = y$ ,  $\pi(\bar{u}) = u$  et  $\pi(\bar{v}) = v$

Par définition  $\pi(a.b) = \pi(a).\pi(b)$

$\implies \pi(\bar{x}.\bar{u}) = \pi(\bar{x}).\pi(\bar{u}) = x.u$  et  $\pi(\bar{y}.\bar{v}) = \pi(\bar{y}).\pi(\bar{v}) = y.v$

— Montrons que  $|\bar{X}| = |\bar{Y}|$  c'est-à-dire que  $|\bar{x}.\bar{u}| = |\bar{y}.\bar{v}|$

$(\bar{x}.\bar{u}) = x_1 \dots x_k u_1 \dots u_n$  et  $(\bar{y}.\bar{v}) = y_1 \dots y_k v_1 \dots v_n$

Comme  $(\bar{x}, \bar{y})$  et  $(\bar{u}, \bar{v})$  sont respectivement des alignements de  $(x, y)$  et  $(u, v)$ ,  $|\bar{x}| = |\bar{y}|$  et  $|\bar{u}| = |\bar{v}|$

Donc  $|\bar{x}.\bar{u}| = |\bar{x}| + |\bar{u}| = k + n = |\bar{y}| + |\bar{v}| = |\bar{y}.\bar{v}|$

Donc  $|\bar{X}| = |\bar{Y}|$

— Montrons que  $\forall i \in \{1, \dots, n+k\} \bar{X}_i \neq \_ \text{ ou } \bar{Y}_i \neq \_$

On sait que

$\forall i \in \{1, \dots, k\} \bar{x}_i \neq \_ \text{ ou } \bar{y}_i \neq \_$

$\forall i \in \{1, \dots, n\} \bar{u}_i \neq \_ \text{ ou } \bar{v}_i \neq \_$

Supposons qu'il y ait un  $j \in \{1, \dots, n+k\}$  tel que  $\bar{X}_j = \_ \text{ et } \bar{Y}_j = \_$

- Si  $j \leq k$   $\bar{X}_j = \bar{x}_j = \_ \text{ et } \bar{Y}_j = \bar{y}_j = \_ \implies \bar{y}_j = \_ \text{ et } \bar{x}_j = \_ \text{ donc } (\bar{x}, \bar{y})$  n'est pas un alignement de  $(x, y)$

- Si  $j > k$   $\bar{X}_j = \bar{u}_{j-k} = \_ \text{ et } \bar{Y}_j = \bar{v}_{j-k} = \_ \implies \bar{u}_{j-k} = \_ \text{ et } \bar{v}_{j-k} = \_ \text{ donc } (\bar{u}, \bar{v})$  n'est pas un alignement de  $(u, v)$

$\implies$  Contradiction.

Donc  $\forall i \in \{1, \dots, n+k\} \overline{X}_i \neq \_ \text{ ou } \overline{Y}_i \neq \_$

On en conclut donc que  $(\overline{x}.\overline{u}, \overline{y}.\overline{v})$  est un alignement de  $(x.u, y.v)$

## Question 2

Soient  $|x| = n$  et  $|y| = m$  et  $(\overline{x}, \overline{y})$  un alignement de  $(x, y)$ .

Par définition de l'alignement  $|\overline{x}| = |x| + k$  et  $|\overline{y}| = |y| + k'$  avec  $k$  et  $k' \geq 0$  où  $k$  est le nombre de gaps ajouté à  $x$  et  $k'$  le nombre de gaps ajoutés à  $y$ .

donc  $|\overline{y}| \geq m$  et  $|\overline{x}| \geq n$  Or  $|\overline{y}| = |\overline{x}|$  Donc  $|\overline{x}| \geq \max(n, m)$

Montrons que le nombre de gaps ne peut pas dépasser  $\min(n, m)$ .

*Supposons ici que  $n \leq m$  ( Le raisonnement est le même si  $m \leq n$  ) :*

$\max(n, m) \leq |\overline{x}| \leq n + \min(n, m) = 2n$

$\max(n, m) \leq |\overline{y}| \leq m + \min(n, m) = n + m \max(2n, n+m) = n+m.$

Donc  $|\overline{x}| = |\overline{y}| \leq n + m$ .

En effet supposons que  $|\overline{x}| = n + m + 1$ .

$\overline{x}$  compte donc  $m+1$  gaps  $\{g_1, \dots, g_{m+1}\}$  où  $g_i$  est la position du  $i$ -ième gap dans  $\overline{x}$ .

$\overline{y}$  compte donc  $n+1$  gaps  $\{g_1, \dots, g_{n+1}\}$  où  $g_i$  est la position du  $i$ -ième gap dans  $\overline{y}$

$\overline{y}$  a plus de gaps que  $\overline{x}$  n'a de lettre  $\implies \exists i \in \{1, \dots, n+m+1\}$  tel que  $\overline{y}_i = \_$  et  $\overline{x}_i = \_ \implies (\overline{x}, \overline{y})$  n'est pas un alignement de  $(x, y)$ .

Donc  $|\overline{x}| = |\overline{y}| \leq n + m$ .

Si  $|\overline{x}| = n + m$

$\overline{x}$  compte donc  $m$  gaps  $\{g_1, \dots, g_m\}$

$\overline{y}$  compte donc  $n$  gaps  $\{g_1, \dots, g_n\}$

Pour chaque lettre  $l$  de  $x$   $\{l_{x1}, \dots, l_{xn}\}$  il est possible de placer la lettre  $l_{xi}$  dans  $\overline{x}$  et le gap  $g_i$  dans  $\overline{y}$  au même indice.

Par ce processus les lettres de  $\overline{x}$  seront placées aux mêmes indices que les gaps de  $\overline{y}$  et les lettres de  $\overline{y}$  seront placées aux mêmes indices que les gaps de  $\overline{x}$ . La longueur maximale d'un alignement de  $(x, y)$  est donc  $n+m$ .

## 2 Algorithmes pour l'alignement de séquences

### 2.1 Méthode naïve par énumération

#### Question 3

Nous cherchons le nombre de mots obtenus en ajoutant  $k$  gaps à un mot de  $n$  lettres. On obtient alors un mot de  $n+k$  lettres. Pour approcher le problème d'une autre manière on pourrait se demander : Pour  $k \in \{1, \dots, n+k\}$ , combien de combinaisons différentes de  $k$  chiffres pouvons nous créer ? En prenant en compte le fait que la combinaison  $\{1, 3, 7\}$  est la même que  $\{3, 1, 7\}$  nous cherchons donc le nombre de combinaisons de  $k$  éléments parmi  $n+k$  éléments, soit :

$$\binom{n+k}{k}$$

#### Question 4

Une fois que l'on a ajouté  $k$  gaps à  $x$  pour obtenir  $\bar{x}$ , on obtient l'égalité :

$n + k = m + k'$  où  $k'$  est le nombre de gaps que l'on rajoute à  $y$ . Les gaps de  $\bar{x}$  ne doivent pas être au même endroit que les gaps de  $\bar{y}$ . Ainsi les  $k'$  gaps peuvent se situer seulement sur les  $n$  emplacements restants. Alors, il y a  $\binom{n}{k'}$  façons d'insérer les  $k'$  gaps dans  $y$  pour satisfaire cette contrainte. On a :  $\binom{n}{k'} = \binom{n}{n+k-m}$ .

On observe que pour  $k$  gaps, on obtient  $\binom{n+k}{k} \times \binom{n}{k'}$  alignements possibles pour  $(x, y)$ . Enfin, d'après la question 2, au maximum,  $\bar{x}$  compte  $m$  gaps ( $k=m$ ) et  $\bar{y}$  compte  $n$  gaps ( $k'=n$ ). On obtient donc le nombre d'alignements possibles de  $(x, y)$  :  $\sum_{k=0}^m \left( \binom{n+k}{k} * \binom{n}{n+k-m} \right)$ .

Si  $n=15$  et  $m=10$ , on obtient que le nombre d'alignements possibles de  $(x, y)$  est : 298 199 265

#### Question 5

En supposant que le nombre d'alignements possibles d'un couple de mots  $(x, y)$  de longueurs respectives  $n$  et  $m$  est  $\sum_{k=0}^m \left( \binom{n+k}{k} * \binom{n}{n+k-m} \right)$ , la complexité temporelle pour énumérer tous les alignements de  $(x, y)$  est donc de  $\sum_{k=0}^m \left( \binom{n+k}{k} * \binom{n}{n+k-m} \right)$

$$\binom{n+k}{k} * \binom{n}{n+k-m} = \frac{(n+k)!}{k!n!(n+k-m)!(m-k)!} = \frac{(n+k)!}{k!(n+k-m)!(m-k)!} \leq (n+k)! \leq (n+m)!$$

Donc  $\sum_{k=0}^m \left( \binom{n+k}{k} * \binom{n}{n+k-m} \right) \leq \sum_{k=0}^m (n+m)! = m * (n+m)!$

La complexité est donc de  $O(m \times (n+m)!)$ .

## Question 6

Comme nous avons besoin de stocker toutes les distances d'édition avant de trouver la plus petite. Pour chaque alignement trouvé et donc pour chaque distance trouvée on doit stocker  $x$  ( $O(n)$ ) et  $y$  ( $O(m)$ ). Comme vu à la question 5, le nombre maximal d'alignements possible est en  $O(m * (n+m)!)$ . La complexité spatiale totale est donc :

$$\text{nombre d'alignements possible de } (x,y) * \text{taille de } x \text{ et } y = O(m * (m + m)! * (n + m))$$

Les complexités spatiales et temporelles de cette méthode se rapprochent de la complexité exponentielle, elle n'est donc pas optimale.

### **Tache A**

#### TEST DE DIST\_NAIF

Les fonctions *DIST\_NAIF* et *DIST\_NAIF\_REC* se trouvent dans le fichier *Tache\_A.py* Après vérification, on obtient bien les valeurs de retour 10, 8 et 2 quand on applique *dist\_naif* sur *Inst\_0000010\_44.adn*, *Inst\_0000010\_7.adn* et *Inst\_0000010\_8.adn*

Pour *Inst\_0000010\_7.adn* Dist = 8

Pour *Inst\_0000010\_44.adn* Dist = 10

Pour *Inst\_0000010\_8.adn* Dist = 2

#### TEMPS MAX DE DIST\_NAIF

En effectuant des tests sur plusieurs instances, on s'aperçoit qu'après  $l_x = 12$ , on dépasse la minute. On remarque tout de même que lorsque  $l_y$  se rapproche de  $l_x$ , les temps d'exécutions s'allongent considérablement comme on peut remarquer entre les fichiers 12\_13 (et 12\_32 ) et 12\_56 .

```
Instances_genome/Inst_0000010_44.adn 0.03739595413208008
Instances_genome/Inst_0000010_7.adn 7.484240770339966
Instances_genome/Inst_0000010_8.adn 2.916300058364868
Instances_genome/Inst_0000012_13.adn 13.385396957397461
Instances_genome/Inst_0000012_32.adn 13.467100143432617
Instances_genome/Inst_0000012_56.adn 96.79302024841309
```

### CONSOMMATION MEMOIRE DE DIST\_NAIF

Nous avons mesuré la quantité de mémoire utilisée par la fonction DIST\_NAIF pour les 5 premières instances ( les quantités sont en octets ou bytes).

Pour cela nous avons utiliser le module *tracemalloc*.

Inst\_0000010\_44.adn Temps = 0.0508267879486084 consommation mémoire 19329

Inst\_0000010\_7.adn Temps = 10.213174819946289 consommation mémoire 27069

Inst\_0000010\_8.adn Temps = 4.136791229248047 consommation mémoire 29832

Inst\_0000012\_13.adn Temps = 20.95844578742981 consommation mémoire 29967

Inst\_0000012\_32.adn Temps = 21.820443868637085 consommation mémoire 30623

Inst\_0000012\_56.adn Temps = 172.80903887748718 consommation mémoire 30757

### 3 Programmation dynamique

#### 3.1 Calcul de la distance d'édition par programmation dynamique

##### Question 7

- Si  $\bar{u}_l = \_$ , par définition de l'alignement on ne peut pas avoir  $\bar{v}_l = \_$  donc  $\bar{v}_l$  est le dernier caractère de  $y$  soit  $y_j$ . Donc si  $\bar{u}_l = \_$ ,  $\bar{v}_l = y_j$ .
- Si  $\bar{v}_l = \_$ , de la même manière on aura  $\bar{u}_l = x_i$ .
- Si  $\bar{u}_l \neq \_$  et  $\bar{v}_l \neq \_$ , il n'y a qu'une seule possibilité,  $\bar{u}_l = x_i$  et  $\bar{v}_l = y_j$

##### Question 8

$$C(\bar{u}, \bar{v}) = \begin{cases} c_{ins} + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) & \text{si } \bar{u}_l = \_ \\ c_{del} + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) & \text{si } \bar{v}_l = \_ \\ c_{sub}(\bar{u}_l, \bar{v}_l) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) & \text{si } \bar{v}_l \neq \_ \text{ et } \bar{u}_l \neq \_ \end{cases}$$

##### Question 9

On a  $D(i,j) = d(x_{[1....i]}, y_{[1....j]}) = \min\{C(\bar{x}_{[1....i]}, \bar{y}_{[1....j]}) \mid (\bar{x}_{[1....i]}, \bar{y}_{[1....j]}) \text{ est un alignement de } (x_{[1....i]}, y_{[1....j]})\}$ . Nous avons déterminé dans la question précédente la formule récursive permettant de calculer la distance d'édition de  $(\bar{x}, \bar{y})$ . Sachant que pour  $(x,y)$  il y a plusieurs manières d'obtenir  $(\bar{x}, \bar{y})$ , on souhaite trouver l'alignement ayant la distance d'édition la plus faible. Prenons un exemple : Soit  $x_{[1,2,3,4]} = ATGA$  et  $y_{[1,2,3]} = GGC$  Pour trouver  $d(x_{[1,2,3,4]}, y_{[1,2,3]})$  on prend la dernière lettre de chaque mot (ici A et C) et on les confronte à un choix :

- Soit on place un gap à la fin de  $\bar{x}$  donc  $\bar{x}_l = \_$  et dans ce cas la,  $\bar{y}_l = y_j = y_3 = C$ . D'après q8 on a  $C(\bar{u}, \bar{v}) = c(\bar{x}_l, \bar{y}_l) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) = c(\_, C) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) = c_{ins} + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]})$ . Donc la lettre  $y_3$  a été utilisé, on doit continuer le processus avec les mots  $x_{[1,2,3,4]} = ATGA$  et  $y_{[1,2]} = GG$ . Plus formellement si on place un gap à la fin de  $\bar{x}$  on aura distance d'édition  $(x_{[1....i]}, y_{[1....j-1]}) = c_{ins} + \text{distance d'édition}(x_{[1....i]}, y_{[1....j-1]})$
- Soit on place un gap à la fin de  $\bar{y}$  donc  $\bar{y}_l = \_$  et  $\bar{x}_l = x_i = A$  dans ce cas la, d'après q8 on a  $C(\bar{u}, \bar{v}) = c(\bar{x}_l, \bar{y}_l) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) = c(A, \_) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) = c_{del} + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]})$  De la même manière on aura distance d'édition  $(x_{[1....i]}, y_{[1....j]}) = c_{del} + \text{distance d'édition}(x_{[1....i-1]}, y_{[1....j]})$
- Soit on ne place aucun gap à la fin de  $\bar{x}$  et  $\bar{y}$  et dans ce cas la on aura  $C(\bar{u}, \bar{v}) = c(\bar{x}_l, \bar{y}_l) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) = c(A, C) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) = c_{sub}(A, C) + c(\bar{u}_{[1....l-1]}, \bar{v}_{[1....l-1]}) \implies \text{distance d'édition}(x_{[1....i]}, y_{[1....j]}) = c_{sub}(x_i, y_j) + \text{distance d'édition}(x_{[1....i-1]}, y_{[1....j-1]})$

Nous pouvons donc en déduire une formule de  $D(i,j)$  :



$$D(i,j) = \min ( c_{ins} + D(i, j-1) , c_{del} + D(i-1, j) , c_{sub}(x_i, y_j) + D(i-1, j-1) )$$

En prenant  $i' = i-1$  et  $j' = j-1$  on retrouve :

$$D(i,j) = \min ( c_{ins} + D(i, j') , c_{del} + D(i', j) , c_{sub}(x_i, y_j) + D(i', j') )$$

### Question 10

$$D(0,0) = d(x_{[1...0]}, y_{[1...0]}) = d(\emptyset, \emptyset) = 0.$$

En effet l'alignement de 2 mots vides est unique car, si on a x et y vide, la seule chose que l'on peut faire pour obtenir  $\bar{x}$  et  $\bar{y}$  est d'ajouter des gaps. Cependant si l'on ajoute un gap à x et y on obtiendra  $\bar{x}_1 = \bar{y}_1 = \_$  ce qui ne peut pas correspondre à un alignement de (x,y). Le seul alignement possible de (x,y) est donc  $\bar{x} = \bar{y} = \emptyset$  La distance d'édition est donc nulle.

### Question 11

$D(0,j) = d(\emptyset, y_{[1...j]})$ . Il n'y a qu'une seule façon d'obtenir un alignement de (x,y) avec x le mot vide. Il faudra rajouter j gap dans x et 0 gap dans y. En effet pour que  $(\bar{x}, \bar{y})$  puisse être un alignement de (x,y) il faut que  $|\bar{x}| = |\bar{y}| \geq |y|$ . Comme  $|y| = j$ ,  $\bar{x}$  doit au moins avoir j gaps. Si on rajoute j+1 gap à x, on devra rajouter un gap à y et on aura donc forcément un  $i \in \{1, \dots, j+1\}$  pour lequel  $\bar{x}_i = \bar{y}_i$ . Donc on a  $\bar{x}$  composé de j gaps et  $\bar{y} = y_{[1...j]}$

On a donc :

$$D(0,j) = C(\bar{x}, \bar{y}_{[1,...,j]}) = \sum_{k=1}^j c(\bar{x}_k, \bar{y}_k) = \sum_{k=1}^j c(\_, y_k) = \sum_{k=1}^j c_{ins} = \sum_{k=1}^j 2 = 2j$$

## Question 12

```
DIST_1
Entrée : x, y et une matrice T vide de taille (n+1)*(m+1)
Sortie : distance d'édition de (x,y)
n = |x|
m = |y|
T = [n+1][m+1] On crée un tableau de n+1 lignes et m+1 colonnes qu'on remplit de 0
For i from 0 to n :
    T[i][0] = c_del * i
For j from 0 to m :
    T[0][j] = c_ins * j

For i from 1 to n :
    For j from 1 to m :
        T[i,j] = min(T[i][j-1] + c_ins , T[i-1][j] + c_del, T[i-1][j-1] + c_sub(x[i-1],y[j-1]))

retourner T[n][m]
```

## Question 13

Dans *DIST\_1*, on alloue une matrice de  $n + 1$  lignes et  $m + 1$  colonnes. La taille totale est donc de  $n+1 * m+1$  cases. La complexité spatiale de *DIST\_1* est  $O(n * m)$

## Question 14

L'algorithme se compose de :

- une boucle de  $n$  itérations permettant de donner la valeur  $D(i,0)$  pour  $i$  allant de 0 à  $n$
- une boucle de  $m$  itérations permettant de donner la valeur  $D(0,j)$  pour  $j$  allant de 0 à  $m$
- 2 boucles imbriquées ( une à  $n$  itérations et l'autre à  $m$  itérations) contenant chacune des opérations élémentaires en  $O(1)$ .

Donc la complexité totale est de  $O(n + m + n * m) = O(n * m)$

## 3.2 Calcul d'un alignement optimal par programmation dynamique

### Question 15

Montrons que :

Si  $D(i,j) = D(i-1,j-1) + c_{sub}(x_i, y_j)$  alors  $\forall (\bar{s}, \bar{t}) \in Al(i-1,j-1)$ ,  $(\bar{s}.x_i, \bar{t}.y_j) \in Al(i,j)$

On suppose que  $D(i,j) = D(i-1,j-1) + c_{sub}(x_i, y_j)$  et soit  $(\bar{s}, \bar{t}) \in Al(i-1,j-1)$ .

Montrons que  $(\bar{s}.x_i, \bar{t}.y_j) \in Al(i,j)$  :

$(\bar{s}.x_i, \bar{t}.y_j) \in Al(i,j)$  si il respecte les 2 conditions suivantes :

- $(\bar{s}.x_i, \bar{t}.y_j)$  est un alignement de  $(x_{[1,...,i]}, y_{[1,...,j]})$
- $C(\bar{s}.x_i, \bar{t}.y_j) = d(x_{[1,...,i]}, y_{[1,...,j]})$

Montrons que  $(\bar{s}.x_i, \bar{t}.y_j)$  est un alignement de  $(x_{[1,...,i]}, y_{[1,...,j]})$  :

- $\Pi(\bar{s}.x_i) = \Pi(\bar{s}).\Pi(x_i)$  et  $\Pi(\bar{t}.y_j) = \Pi(\bar{t}).\Pi(y_j)$   
Comme  $x_i$  et  $y_j \in \Sigma$  alors  $\Pi(x_i) = x_i$  et  $\Pi(y_j) = y_j$ .  
On a  $(\bar{s}, \bar{t})$  est un alignement de  $(x_{[1,...,i-1]}, y_{[1,...,j-1]})$  alors  $\Pi(\bar{s}) = x_{[1,...,i-1]}$  et  $\Pi(\bar{t}) = y_{[1,...,j-1]}$   
On a donc bien  $\Pi(\bar{s}.x_i) = x_{[1,...,i-1]}.x_i = x_{[1,...,i]}$  et  $\Pi(\bar{t}.y_j) = y_{[1,...,j-1]}.y_j = y_{[1,...,j]}$
- $|\Pi(\bar{s}.x_i)| = |\bar{s}| + |x_i| = 1 + |\bar{s}|$   
 $|\Pi(\bar{t}.y_j)| = |\bar{t}| + |y_j| = 1 + |\bar{t}|$   
Comme  $(\bar{s}, \bar{t})$  est un alignement de  $(x_{[1,...,i-1]}, y_{[1,...,j-1]})$   $|\bar{s}| = |\bar{t}| = k$   
On a donc  $|\Pi(\bar{s}.x_i)| = |\bar{s}| + 1 = |\bar{t}| + 1 = |\Pi(\bar{t}.y_j)|$ .
- Pour faciliter la notation on écrira  $(\bar{s}.x_i, \bar{t}.y_j) = (\bar{u}, \bar{v})$ .  
Comme  $(\bar{s}, \bar{t})$  est un alignement de  $(x_{[1,...,i-1]}, y_{[1,...,j-1]})$  on a  
 $\forall i \in [1, \dots, |\bar{t}|] \bar{x}_i = \_ \text{ ou } \bar{y}_i = \_$   
Soit  $k \in [1, \dots, |\bar{t}| + 1]$ .  
- Si  $k \leq |\bar{t}|$  alors  $\bar{u}_k = \bar{s}_k$  et  $\bar{v}_k = \bar{t}_k$ . Donc  $k \leq |\bar{t}| \implies \bar{u}_k = \_ \text{ ou } \bar{v}_k = \_$   
- Si  $k = |\bar{t}| + 1$  alors  $\bar{u}_k = x_i \neq \_ \text{ et } \bar{v}_k = y_j \neq \_$   
On a donc bien  $\forall i \in [1, \dots, |\bar{t}.y_j|] \bar{u}_i = \_ \text{ ou } \bar{v}_i = \_$   
 $(\bar{s}.x_i, \bar{t}.y_j)$  est donc bien un alignement de  $(x_{[1,...,i]}, y_{[1,...,j]})$

Montrons que  $C(\bar{s}.x_i, \bar{t}.y_j) = d(x_{[1,...,i]}, y_{[1,...,j]})$

On rappelle qu'ici on notera  $(\bar{s}.x_i, \bar{t}.y_j) = (\bar{u}, \bar{v})$ . Soit  $|\bar{u}| = n$

$$C(\bar{s}.x_i, \bar{t}.y_j) = C(\bar{u}, \bar{v}) = \sum_{k=1}^n c(\bar{u}_k, \bar{v}_k) = c(\bar{u}_n, \bar{v}_n) + C(\bar{u}_{[1,...,n-1]}, \bar{v}_{[1,...,n-1]})$$

Ici on a  $(\bar{u}_n, \bar{v}_n) = (x_i, y_j)$ , avec  $x_i \neq \_ \text{ et } y_j \neq \_$ , donc  $c(\bar{u}_n, \bar{v}_n) = c_{sub}(\bar{u}_n, \bar{v}_n)$

$$C(\bar{u}, \bar{v}) = c_{sub}(x_i, y_j) + C(\bar{u}_{[1,...,n-1]}, \bar{v}_{[1,...,n-1]}) = c_{sub}(x_i, y_j) + C(\bar{s}, \bar{t}).$$

$$(\bar{s}, \bar{t}) \in Al(i-1, j-1) \text{ donc } C(\bar{s}, \bar{t}) = d(x_{[1,...,i-1]}, y_{[1,...,j-1]}) = D(i-1, j-1)$$

$$C(\bar{u}, \bar{v}) = c_{sub}(x_i, y_j) + D(i-1, j-1).$$

$$\text{Or } D(i-1, j-1) + c_{sub}(x_i, y_j) = D(i, j), \text{ donc } C(\bar{u}, \bar{v}) = D(i, j) = d(x_{[1,...,i]}, y_{[1,...,j]})$$

On a donc bien montré que :

$$\text{Si } D(i, j) = D(i-1, j-1) + c_{sub}(x_i, y_j) \text{ alors } \forall (\bar{s}, \bar{t}) \in Al(i-1, j-1), (\bar{s}.x_i, \bar{t}.y_j) \in Al(i, j)$$

## Question 16

SOL\_1

```

Entrée x, y et T la matrice contenant les valeurs de D
Sortie : alignement optimal de (x,y)
i ← |x|
j ← |y|
s ← ""
t ← ""

while i>0 and j>0
    if T[i,j] == T[i-1,j-1] + c_sub(x[i-1],y[j-1])
        s ← x[i-1]+s
        t ← y[j-1]+t
        i ← i-1
        j ← j-1
    elif T[i,j] == T[i,j-1] + c_ins
        s ← "_" + s
        t ← y[j-1]+t
        j ← j-1
    elif T[i,j] == T[i-1,j] + c_del
        s ← x[i-1]+s
        t ← "_" + t
        i ← i-1

while i>0
    if T[i,j] == T[i-1,j] + c_del
        s ← x[i-1]+s
        t ← "_" + t
        i ← i-1

while j>0
    if T[i,j] == T[i,j-1] + c_ins :
        s ← "_" + s
        t ← y[j-1]+t
        j ← j-1
retourner (s,t)

```

## Question 17

Pour résoudre le problème ALI avec les ressources que nous avons à notre disposition, nous devons d'abord construire la matrice des valeurs de D grâce à la fonction *DIST\_1* puis construire l'alignement associé avec la fonction *SOL\_1*

Soit x et y de taille respective n et m.

— **Complexité de DIST\_1 =  $O(n.m)$**

— **Complexité de SOL\_1 :**

La fonction *SOL\_1* est composée de 3 boucles :

- La première boucle fait au plus  $c1 = \min(n,m)$  itérations.

- Si  $|x| > |y|$  :

On passe seulement par la 2e boucle de complexité  $c2 = n - c1$

- Si  $|x| < |y|$  :

On passe seulement par la 3e boucle de complexité  $c2 = n - c1$

La complexité totale est de  $c_1 + c_2 + \text{Complexité de } DIST\_1 = O(n.m + n) = O(n^2)$  car  $m \leq n$ .

### Question 18

Comme vu dans la question 13, la complexité spatiale de  $DIST\_1$  est  $O(n * m)$  et dans  $SOL\_1$  on utilise la mémoire du tableau T de taille  $(n + 1) * (m + 1)$  et les mots x et y donc la complexité spatiale de  $SOL\_1$  est  $O(n * m + n + m)$  La complexité totale est  $O(2 * n * m + n + m) = O(n * m)$

### Tache B

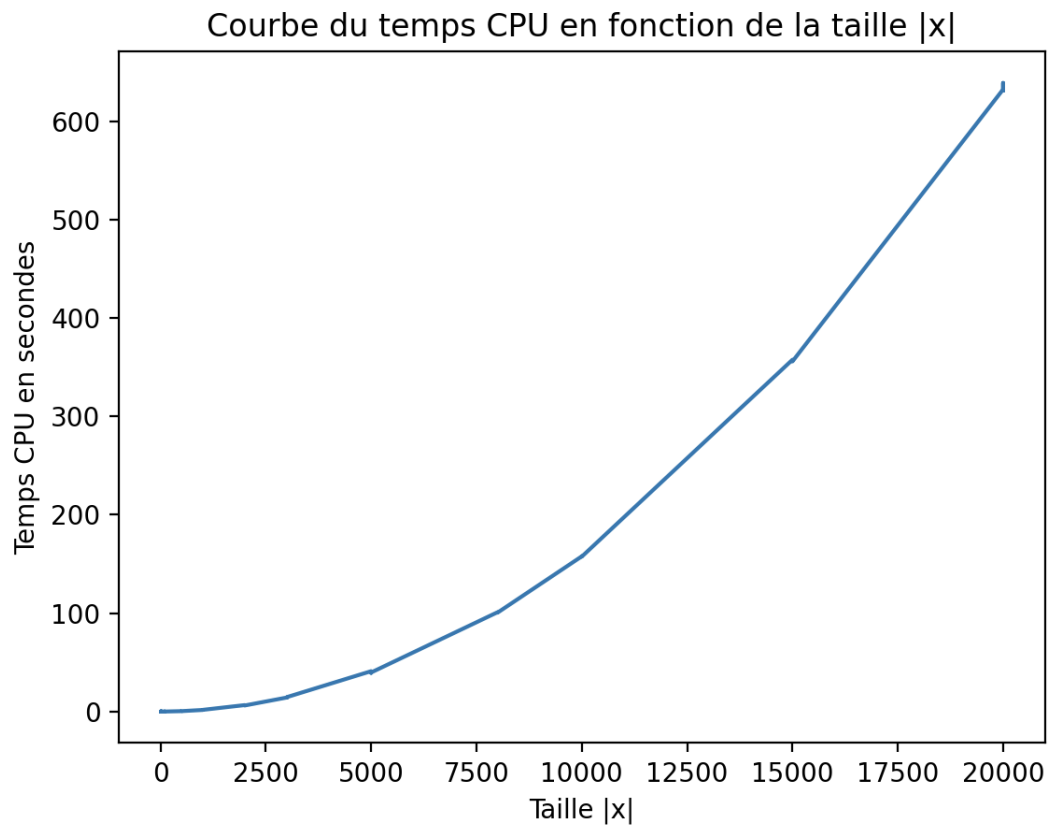
#### TEST DE PROG\_DYN

```
PROG_DYN de Inst_0000010_44.adn = DIST = 10 alignement = ( TATATGAGTC , TAT_T___T_
PROG_DYN de Inst_0000010_7.adn = DIST = 8 alignement = ( TGGGTGCTAT , GGGGTTCTAT
PROG_DYN de Inst_0000010_8.adn = DIST = 2 alignement = ( AACTGTCTTT , AACTGT_TTT
PROG_DYN de Inst_0000012_13.adn = DIST = 9 alignement = ( CTGGAAGTGCG , CT_G_AACTG_G
PROG_DYN de Inst_0000012_32.adn = DIST = 6 alignement = ( CCATTGATTTTC , _CATTGA_TTT_
PROG_DYN de Inst_0000012_56.adn = DIST = 9 alignement = ( GCTTAACTAACG , GCTAACT_ACT
PROG_DYN de Inst_0000013_45.adn = DIST = 2 alignement = ( CGGGGGGTAGCAA , CGGGGGGT_GCAA
PROG_DYN de Inst_0000013_56.adn = DIST = 9 alignement = ( GCTTAACTAAC_GA , GCATAACT__CAGA
PROG_DYN de Inst_0000013_89.adn = DIST = 8 alignement = ( TCCGCAAGCGTGT , TCGGCAA_CGTCT
PROG_DYN de Inst_0000014_23.adn = DIST = 15 alignement = ( CCCCACAGGGCTAG , __AGACA_GGC_AG
```

#### CALCUL DU TEMPS CPU DE PROG\_DYN

Après calcul, nous nous limitons aux mots x de longueur au plus 20 000 car après ceux-ci le calcul de prog\_dyn(x,y) prend plus de 10 minutes.

À l'aide des instances fournies, on obtient la courbe suivante :



On observe que la courbe correspond bien à la complexité théorique trouvée à la question 17 :  $O(n^2)$

#### CONSOMMATION MEMOIRE DE PROG\_DYN

On a mesuré la quantité de mémoire utilisée pour l'instance *Inst\_0015000\_3.adn* et on obtient 244077 octets. Pour des petites instance on obtient environ 13000 octets.

Consommation memoire pour Instances\_genome/Inst\_0000013\_45.adn : 19329

## 4 Amélioration de la complexité spatiale du calcul de la distance

### Question 19

On a  $T[i,j] = \min(T[i,j-1] + c_{ins}, T[i-1,j] + c_{del}, T[i-1,j-1] + c_{sub}(x_i, y_j))$ .

Pour avoir  $T[i,j]$  il faut avoir accès à  $T[i-1,j]$ ,  $T[i-1,j-1]$  et  $T[i,j-1]$  c'est-à-dire aux lignes  $i-1$  et  $i$  et aux colonnes  $j-1$  et  $j$ .

### Question 20

```
DIST_2
Entrée : x et y
Sortie : distance d'édition de (x,y)
n = |x|
m = |y|
T = [2][m+1]      On crée un tableau de 2 lignes et de m+1 colonnes
For i from 0 to 1 :
    T[i][0] = c_del*i
For i from 0 to m :
    T[0][j] = c_ins * j

For i from 1 to n :
    For j from 1 to m :
        T[1][j] = min(T[0][j] + c_del , T[1][j-1] + c_ins, T[0][j-1] + c_sub(x[i-1],y[j-1]))
    For k from 0 to m :
        T[0][k] = T[1][k]      La deuxième ligne du tableau devient la premiere ligne

retourner T
```

### Tache C

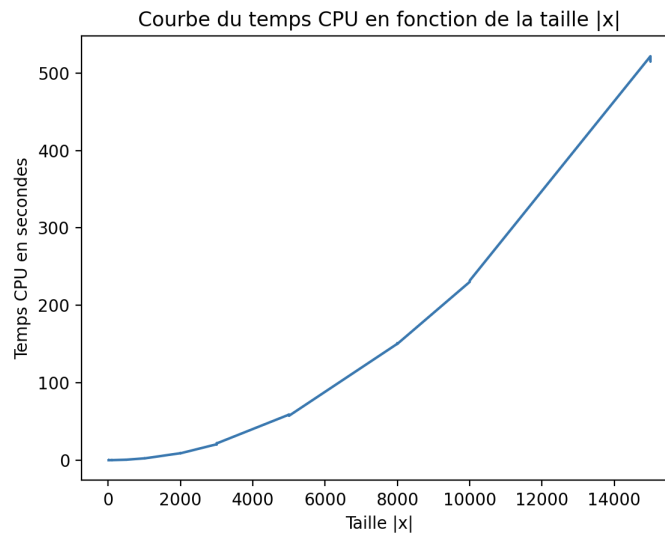
#### TEST de DIST\_2

```
DIST_2 pour l'instance Inst_0000010_44.adn : 10
DIST_2 pour l'instance Inst_0000010_7.adn : 8
DIST_2 pour l'instance Inst_0000010_8.adn : 2
DIST_2 pour l'instance Inst_0000012_13.adn : 9
DIST_2 pour l'instance Inst_0000012_32.adn : 6
DIST_2 pour l'instance Inst_0000012_56.adn : 9
DIST_2 pour l'instance Inst_0000013_45.adn : 2
DIST_2 pour l'instance Inst_0000013_56.adn : 9
DIST_2 pour l'instance Inst_0000013_89.adn : 8
DIST_2 pour l'instance Inst_0000014_23.adn : 15
DIST_2 pour l'instance Inst_0000014_7.adn : 8
DIST_2 pour l'instance Inst_0000014_83.adn : 8
DIST_2 pour l'instance Inst_0000015_2.adn : 8
DIST_2 pour l'instance Inst_0000015_4.adn : 23
DIST_2 pour l'instance Inst_0000015_76.adn : 11
DIST_2 pour l'instance Inst_0000020_17.adn : 14
DIST_2 pour l'instance Inst_0000020_32.adn : 8
```

### COURBE DE TEMPS CPU DE DIST\_2

Après calcul, nous nous limitons aux mots  $x$  de longueur au plus 15 000 car après ceux-ci le calcul de  $DIST_2(x,y)$  prend plus de 10 minutes.

À l'aide des instances fournies on obtient la courbe de temps CPU de  $DIST_2$  suivante :

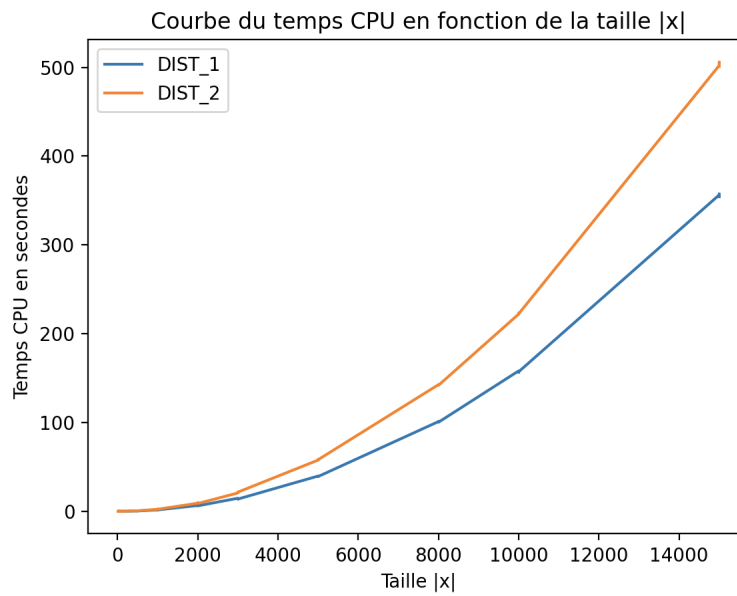


La complexité théorique nous donne :  $O(2 + m + 1 + n * (m + m)) = O(n * m)$ .  
Donc  $DIST_2$  est en  $O(n^2)$  car  $m \leq n$ , ce que l'on obtient bien sur la courbe.

### COMPARAISON $DIST_1$ et $DIST_2$

Si l'on compare les performances de  $DIST_1$  et  $DIST_2$ , on remarque que  $DIST_1$  est plus rapide  $DIST_2$  comme représenté sur le graphique ci-dessous :





### CONSOMMATION MEMOIRE DE DIST\_2

On a mesuré la quantité de mémoire utilisée pour l'instance *Inst\_0015000\_3.adn* et on obtient 244077 octets. Pour des petites instance on obtient environ 13000 octets.

Consommation memoire pour Instances\_genome/Inst\_0000013\_45.adn : 19329

Consommation mémoire de Inst\_0015000\_3.adn 244077

## 5 Amélioration de la complexité spatiale du calcul d'un alignement optimal par la méthode "diviser pour régner"

### Question 21

```
Mot_gaps
    Entrée : k un entier
    Sortie : mot composé de k gaps
    s = mot vide
    For i from 1 to k
        s = s + "_"
    retourner s
```

### Question 22

Nous avons implémenté 2 fonctions auxiliaires :

— **teste\_lettre(x,l)**

Cette fonction teste si la lettre l existe dans x. Si c'est le cas elle renvoie l'indice de la première occurrence de l dans x, sinon elle renvoie -1

— **teste\_complémentaire(x,l)**

Cette fonction teste si il existe une lettre complémentaire à l dans x. Si c'est le cas elle renvoie la première occurrence de celle ci, sinon elle renvoie -1

```
teste_lettre
    Entrée x un mot et l une lettre
    Sortie : retourne l'indice de la lettre l dans x, retourne -1 si l n'est pas dans x
    for i from 0 to |x|
        if x[i] == l
            retourner i
    retourner -1

teste_complémentaire
    Entrée : un mot x et une lettre l
    Sortie : retourne l'indice de la lettre l dans x, retourne -1 si le complémentaire de l n'est pas dans x
    for i from 0 to n
        if (x[i] in {'A','T'} and l in {'A','T'}) or (x[i] in {'C','G'} and l in {'C','G'})
            retourner i
    retourner -1
```

```

align_lettre_mot
  Entré : x et y
  Sortie : alignement optimal de (x,y)

  lettre ← teste_lettre(y,x[0])
  longueur ← |y|
  res ← mot_gaps(longueur)
  complémentaire ← teste_complémentaire(y,x[0])

  if lettre >= 0          il y a une lettre similaire entre x et y
    res[lettre] ← x[0]
    retourner (res,y)

  if complémentaire >= 0
    res[complémentaire] ← x[0]
    return (res,y)

  res[0] ← x[0]          On place la lettre x[0] à la première place dans res
  return (res,y)

```

### Question 23

Soient  $(x^1, y^1) = (\text{BAL}, \text{RO})$  et  $(x^2, y^2) = (\text{LON}, \text{ND})$ .

On remarque facilement que :

1. Un alignement optimal de  $(x^1, y^1)$  est  $(\text{BAL}, \text{RO}_-)$  avec une distance d'édition de 13. On le notera  $(\bar{s}, \bar{t})$ .
2. Un alignement optimal de  $(x^2, y^2)$  est  $(\text{LON}_-, \text{ND})$  avec une distance d'édition de 9. On le notera  $(\bar{u}, \bar{v})$ .

Montrons que  $(\bar{s}, \bar{u}, \bar{t}, \bar{v})$  n'est pas un alignement optimal de  $(x, y)$

$(\bar{s}, \bar{u}, \bar{t}, \bar{v}) = (\text{BALLON}_-, \text{RO}_- \text{ND})$  avec une distance d'édition de 22.

Or si on prend  $\bar{n} = \text{BALLON}_-$  et  $\bar{m} = \text{R}_- \text{ND}$ , on a bien  $(\bar{n}, \bar{m})$  un alignement de  $(x, y)$  avec une distance d'édition de  $17 < 22$ .

Donc  $(\bar{s}, \bar{u}, \bar{t}, \bar{v})$  n'est pas un alignement minimal.

## Question 24

```

sol_2(x,y)
    Entrée : x et y
    Sortie : alignement optimal de (x,y)
    n ← |x|
    m ← |y|
    if m==0
        s ← mot_gaps(n)
        retourner (x,s)
    if n==1 and m>0
        s ← align_lettre_mot(x,y)
        retourner s

    i ← partie entière inférieur (n/2)
    j ← coupure(x,y)

    (x1,y1) ← sol_2(x[0:i],y[0:j])
    (x2,y2) ← sol_2(x[i:n+1],y[j:m+1])

    retourner tuple( map( lambda i,j : i+j, (x1,y1) , (x2,y2) ))

```

## Question 25

Coupure  
 Entrée : x et y  
 Sortie : indice j associé à i.

```

n ← |x|
m ← |y|
i ← partie entière inférieure de (n/2)
T ← np.zeros((2,m+1),int)
I ← np.zeros((2,m+1),int)
for i2 from 0 to 1
    T[i2][0] ← c_del*i2
    I[i2][0] ← 0
for j2 from 0 to m :
    T[0][j2] ← c_ins * j2
    I[0][j2] ← j2

for k from 1 to n :
    for j from 0 to m :
        if j==0 :
            T[1][j] ← T[0][j] + c_del
        else:
            T[1][j] ← min(T[0][j] + c_del , T[1][j-1] + c_ins, T[0][j-1] + c_sub(x[k-1],y[j-1]))

        if k>i :
            if j==0 :
                I[1][j] ← I[0][j]
            else :
                if T[1][j] == T[0][j] + c_del :
                    I[1][j] ← I[0][j]
                elif T[1][j] == T[1][j-1] + c_ins :
                    I[1][j] ← I[1][j-1]
                elif T[1][j] == T[0][j-1] + c_sub(x[k-1],y[j-1]) :
                    I[1][j] ← I[0][j-1]

    for k2 from 0 to m :
        if k> i :
            I[0][k2] ← I[1][k2]
        T[0][k2] ← T[1][k2]
retourner I[1][m]

```

### Question 26

Dans *coupure*(*x*, *y*), on alloue seulement deux tableaux de deux lignes et *m* + 1 colonnes, un tableau pour *D*(*i*, *j*) et un tableau pour calculer l'indice *j*\* associé à *i*\*. Donc la complexité spatiale de *coupure*(*x*, *y*) est  $O(2 * (2 * m + 1)) = O(m)$

### Question 27

Dans l'exécution de *SOL\_2* :

- Si *y* est le mot vide alors on retourne *mot\_gaps*(*n*) donc  $O(n)$
  - Si *x* est un mot de longueur 1, on retourne *align\_mot\_lettre*(*x*, *y*) en  $O(m)$
  - Sinon on fait appel à *coupure*, donc avec une complexité spatiale en  $O(m)$ .
- On a l'inégalité suivante :

$$\frac{n}{2^k} \geq 1$$

En effet pour chaque appel récursif on divise le mot *x* par 2 jusqu'à ce que  $|x| = 1$ . On a donc  $k \leq \log_2(n)$ . Il y a au plus  $2 * \log_2(n)$  appels récursifs de *SOL\_2* ( car on appelle 2 fois *SOL\_2* à chaque fois). Chaque appel récursif appelle *coupure* en  $O(m)$  et alloue 2 tableaux de tailles  $2*(m+1)$  donc chaque appel est en  $O(m + 2*(m+1)) = O(m)$ .

La complexité spatiale de *SOL\_2* est  $C(n, m) = \begin{cases} O(n) & \text{si } |y| = 0 \\ O(m) & \text{si } |x| = 1 \\ O(m * \log_2(n)) & \text{sinon} \end{cases}$

### Question 28

- Les instructions des lignes 1 à 5 sont en  $O(1)$
- La 1e boucle est en  $O(1)$ , la 2e en  $O(m)$  avec des instructions en  $O(1)$
- La 3e boucle est en  $O(n*m)$  avec des instructions en  $O(1)$
- La dernière boucle est en  $O(m)$  avec des instructions en  $O(1)$

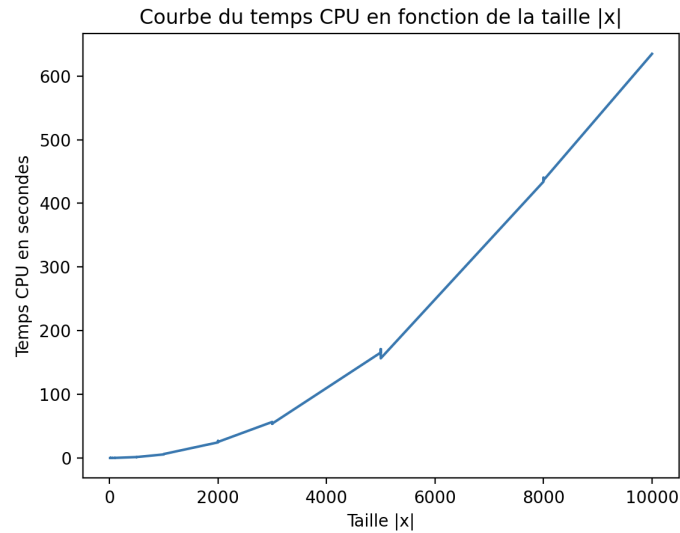
La complexité totale est de  $O(m + n*m) = O(n * m)$

### Tache D

#### TEMPS CPU SOL\_2

Après calcul, nous nous limitons aux mots *x* de longueur au plus 10 000 car après ceux-ci le calcul de *SOL\_2* prend plus de 10 minutes.

À l'aide des instances fournies, on obtient la courbe suivante :



### CONSOMMATION MÉMOIRE SOL\_2

Nous avons mesuré la consommation mémoire de la fonction *SOL\_2* sur :

- Une instance de petite taille ( 12 caractères ) : Environ 19 000 octets

Consommation memoire pour Instances\_genome/Inst\_0000013\_45.adn : 19329

- Des instances de grande taille ( 500 caractères ) : Environ 45 000 octets

Consommation memoire pour Instances\_genome/Inst\_0000500\_8.adn : 44884

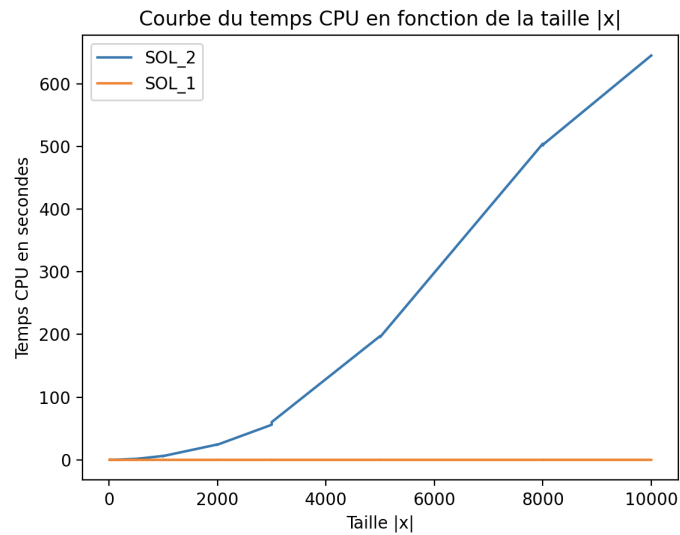
- Des instances de grande taille ( 500 caractères )

Consommation memoire pour Instances\_genome/Inst\_0002000\_3.adn : 139821

### Question 29

On compare l'efficacité temporelle de *SOL\_1* avec *SOL\_2* en effectuant des mesures de temps et en traçant une courbe du temps CPU pour chaque fonction.

Pour des instances identiques, on obtient :



On observe que, bien que la complexité spatiale soit améliorée dans  $SOL_2$  par rapport à  $SOL_1$  car on ne calcule pas une matrice entière, la complexité temporelle est énormément dégradée. On voit sur les courbes ci-dessus que  $SOL_1$  est bien plus rapide que  $SOL_2$ . Il serait plus intéressant de comparer  $PROG\_DYN$  et  $SOL_2$  car les deux fonctions construisent la matrice des valeurs  $D(i,j)$ . On obtient les courbes suivantes :

