

Exercice : Détermination de la plus longue sous-chaine commune à deux chaînes de caractères

Eléments de correction

En génétique par exemple, il peut s'avérer très important de savoir quels sont les points communs de deux génomes. Le problème de la plus longue sous-chaine permet d'apporter une réponse à cette question.

On appelle **sous-chaine** d'une chaîne de caractères $a = a_1 \dots a_n$ toute chaîne de caractères s extraite de a telle que $s = a_i \dots a_k$ où $i \leq \dots \leq k < n$ et $\forall p \in \{i, \dots, k\}, a_p \in a$. Les caractères de s n'apparaissent pas nécessairement de manière consécutive dans la chaîne a .

Soient $a = a_1 \dots a_n$ et $B = b_1 \dots b_p$ 2 chaînes de caractères non vides. On appelle **plus longue sous-chaine commune** à a et b toute sous-chaine commune à a et b de longueur maximale.

Si l'une des chaînes a ou b est vide ou si a et b n'ont aucune sous-chaine commune, la chaîne vide est alors l'unique plus longue sous-chaine commune à a et b .

Exemple : les chaînes de caractères "AAA" et "TAA" sont les plus longues sous-chaînes communes aux chaînes de caractères "ATAGA" et "TAACA".

1. Pour les chaînes $a="AATGCG"$ et $b="TATTAGC"$? Donner les solutions des plus longues suites communes.

Solution : ATGC et AAGC.

2. Écrire une **fonction de signature** `est_sousChaine(ch : str, sch : str) -> bool` dont les `ch1` et `sch` paramètres sont deux chaînes de caractères. Cette fonction renvoie True si `sch` est une sous-chaine de `ch`, False sinon.

Solution :

```
def est_sousChaine(ch, sch):
```

```
    j = 0
    for i in range(len(ch)):
        if ch[i] == sch[j]:
            j += 1
        if j == len(sch):
            return True
    return False
```

3. Écrire une **fonction de signature** `est_commune(ch1 : str, ch2 : str, sch : str) -> bool` où les paramètres sont des chaînes de caractères. Cette fonction renvoie True si `sch` est une sous-chaine commune à a et b .

Solution :

```
def est_commune(ch1, ch2, sch):
    return est_sousChaine(ch1, sch) and est_sousChaine(ch2, sch)
```

4. Le problème de la plus longue sous-chaine commune entre a et b est noté $L(a, b)$, son résultat est la longueur maximale d'une sous-chaine commune à a et b .

Formuler le **problème $L(a, b)$ récursivement** afin de pouvoir justifier de sa sous-structure optimale.

Solution :

Expression récursive - Sous-structure optimale

On se donne une chaîne a , de longueur n et on note a_1, \dots, a_n ses caractères.

De même on se donne une seconde chaîne b de longueur m et b_1, \dots, b_m ses caractères.

On a alors 2 situations :

- le **dernier caractère est commun** : $c = a_n = b_m$. Il est alors le dernier caractère de la plus longue sous-séquence commune et cette sous-séquence commune est obtenue en concaténant la plus longue sous-séquence commune de a_1, \dots, a_{n-1} et b_1, \dots, b_{m-1} et c .
- le **dernier caractère est différent** : l'un au moins de x_n et y_m n'est pas dans cette plus longue sous-séquence commune.

			b_m
		$(n-1, m-1)$	$(n-1, m)$
a_n		$(n, m-1)$	(n, m)

On peut le retirer. La plus longue sous-séquence commune est la plus grande entre :

→ la plus longue sous-séquence commune entre a_1, \dots, a_{n-1} et b_1, \dots, b_m ,

→ la plus longue sous-séquence commune entre a_1, \dots, a_n et b_1, \dots, b_{m-1} ,

On est donc ramené à chercher des sous-séquences maximales entre des sous-chaînes de a et b : on a bien l'apparition de sous-structure optimale dans la solution.

Une approche récursive directe va de nouveau amener à un problème de complexité exponentielle.

Chevauchement des sous-problèmes

Le calcul de la plus longue sous-séquence commune d'une portions de chacune des deux chaînes est utilisé par beaucoup d'autres calculs de portions plus grande. On a donc bien un chevauchement très important des sous-problèmes, ce qui invite de nouveau à un algorithme de programmation dynamique.

Soit $L(i, j)$ la longueur de la plus longue sous-chaîne des chaînes extraites $a_1 a_2 \dots a_i$ et $b_1 b_2 \dots b_j$ des chaînes a et b.

On regarde si l'une des 2 chaînes est vide.

Si non, on regarde si le dernier élément de chaque chaîne est le même ou pas

On obtient la récurrence :

Pour tout $i \in [0, n]$ et pour tout $j \in [0, m]$,

- Si $i = 0$ ou $j = 0$, $L(i, j) = 0$
- Si $a_i \neq b_j$ $L(i, j) = \max(L(i-1, j), L(i, j-1))$
- Si $a_i = b_j$ $L(i, j) = 1 + L(i-1, j-1)$