

TD3 : Algorithmes sur les textes

Exercice 1 : Chercher un mot

Soit une chaîne de caractères t et un mot w de longueur $m \leq n$.

1. Écrire un algorithme retournant la position de la première occurrence du mot w dans la chaîne t (ou -1 s'il est absent). Quelle est sa complexité?
2. Écrire un algorithme retournant la position de toutes les occurrences du mot w dans la chaîne t (ou une liste vide s'il est absent). Quelle est sa complexité?

Exercice 2 : Compter les mots

Écrire un algorithme qui compte le nombre de mots dans un texte.

Remarque : On considère comme caractère d'espacement tout caractère qui n'est pas alphanumérique (alphabétique accentué ou non et chiffres).

Exercice 3 : Palindrome

Écrire un algorithme récursif permettant de savoir si un tableau de caractères est un palindrome (un palindrome se lit "à l'endroit" et "à l'envers" de la même façon, comme par exemple "à l'étape, épate-lai").

Remarque : on ne considère ni la ponctuation, ni les espaces, ni les accents.

Quelle est sa complexité?

Exercice 4 : Distances entre chaînes de caractères

On cherche à exprimer une distance entre deux chaînes de caractères.

Une distance entre 2 textes d_1 et d_2 est telle que :

- $\text{dist}(d_1, d_2) = \text{dist}(d_2, d_1)$
- $\text{dist}(d_1, d_2) \geq 0$
- $\text{dist}(d_1, d_2) + \text{dist}(d_2, d_3) \geq \text{dist}(d_1, d_3)$

Distance de Hamming

La distance de Hamming entre deux chaînes de même taille est définie comme le nombre de caractères non appariés. Ainsi la distance de Hamming entre "passoire" et "pastèque" est égale à 4. Peut-on généraliser cette distance à des chaînes de taille différente?

Distance d'édition

La distance d'édition est définie, pour deux chaînes de longueur quelconque, comme le nombre minimal d'opérations permettant de transformer d_1 en d_2 , avec les opérations suivantes :

- **ins**(a) \rightarrow insertion du caractère a
- **perm**(a, b) \rightarrow remplacement de a par b
- **del**(a) \rightarrow suppression du caractère a

auxquelles on ajoute l'identité **id**(a, a) de coût nul.

Il existe différentes manières de transformer la chaîne d_1 en d_2 . On peut par exemple supprimer tous les caractères de d_1 et insérer tous les caractères de d_2 , mais c'est rarement le nombre d'opérations optimal ($|d_1| + |d_2|$).

La résolution de ce problème repose sur les principes de la programmation dynamique.

1. Essayez de calculer "à la main" la distance d'édition
 - entre "robe", "arbre", et "porte".
 - entre "cloche", "hochet" et "louche".
 - Vérifiez sur ces exemples que l'inégalité triangulaire est bien respectée.
2. Résoudre au tableau l'algorithme de calcul pour "cloche" et "hochet"
3. Écrire l'algorithme général permettant de calculer la distance d'édition entre deux chaînes quelconques d_1 et d_2 . Quelle est sa complexité?
4. Est-il possible de réduire cette complexité? Si oui, dans quels cas?

From:

<https://wiki.centrale-marseille.fr/informatique/> - **WIKI informatique**

Permanent link:

https://wiki.centrale-marseille.fr/informatique/tc_info:td3

Last update: **2018/10/03 11:58**

