

## TD5 : Algorithmes sur les textes

Version imprimable

### Exercice 1 : Chercher un mot

Soit une chaîne de caractères  $d$  de longueur  $n$  et un mot  $t$  de longueur  $m < n$ .

1. Donner l'algorithme naïf retournant la position de la première occurrence du mot  $t$  dans la chaîne  $d$  (ou  $-1$  s'il est absent). Quelle est sa complexité?
2. Donner de même l'algorithme retournant la position de toutes les occurrences du mot  $t$  dans la chaîne  $d$  (ou une liste vide s'il est absent). Quelle est sa complexité?
3. (\*) On suppose qu'on peut tester si un caractère  $c$  appartient au motif  $t$  en temps constant. Proposez un algorithme plus efficace que l'algorithme naïf.

### Exercice 2 : Compter les mots

1. 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).

2. (\*) Dessiner l'automate fini correspondant.

### Exercice 3 : Palindrome

1. 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-la!").

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

Quelle est sa complexité?

2. (\*) Pouvez-vous définir un automate fini capable de reconnaître les palindromes?

### 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 permettent de transformer d1 en d2, avec les opérations suivantes :

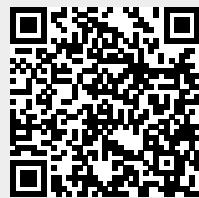
- **ins**(a) → insertion du caractère a
- **perm** (a, b) → remplacement de a par b
- **del** (a) → suppression du caractère a

Il existe différentes manières de transformer la chaîne d1 en d2. On peut par exemple supprimer tous les caractères de d1 et insérer tous les caractères de d2, mais c'est rarement le nombre d'opérations optimal ( $|d1|+|d2|$ ).

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. La résolution de ce problème repose sur les principes de la programmation dynamique.  
Résoudre sur papier l'algorithme de calcul pour "cloche" et "hochet"
3. (\*) Écrire l'algorithme général permettant de calculer la distance d'édition entre deux chaînes quelconques d1 et d2. Quelle est sa complexité?
4. (\*\*) Est-il possible de réduire cette complexité? Si oui, dans quels cas?

From:

<https://wiki.centrale-med.fr/informatique/> - WiKi informatique



Permanent link:

[https://wiki.centrale-med.fr/informatique/tc\\_info:td3](https://wiki.centrale-med.fr/informatique/tc_info:td3)

Last update: **2019/11/14 11:02**