

NOM : PEGATTE

Prénom : Timothée

Jury :

Algèbre ← Entourez l'épreuve → Analyse

Sujet choisi : 907. Algorithmique du texte: exemples et applications.

Autre sujet :

0. Introduction des définitions basiques :
 Σ un alphabet fini

Def 1: (words) Σ^* . On écrit $w \in \Sigma^*$ comme $w = w_1 \dots w_n$ avec $w_i \in \Sigma$

Def 2: Soit $w \in \Sigma^*$. On dit que $x \in \Sigma^*$ est :
 • un préfixe de w ssi $\exists z \in \Sigma^*$ tq $w = xz$
 • un suffixe de w ssi $\exists z \in \Sigma^*$ tq $w = zc$
 • un facteur de w ssi $\exists y, z \in \Sigma^*$ tq $w = yxz$

Ex 3: $w = abab$. Pref(w) = $\{ \epsilon, a, ab, aba, abab \}$
 Suf(w) = $\{ \epsilon, b, ab, abab \}$
 Fact(w) = $\{ \epsilon, a, b, ab, abab, b, ab, abab \}$

Def 3: (factor) $|w| = n$ ssi $w = w_1 \dots w_n$
 $\hookrightarrow w$ se code en $\Theta(n \log |\Sigma|)$, soit $\Theta(n)$ si Σ est supposé fixe.

I. Recherche d'un motif dans un texte :

I.1. Problème :
 • SEARCH(m, t) : • Entrées : un motif m et un texte $t \in \Sigma^*$, $|m| = m, |t| = n$
 • Sortie : réponse à "il est factif ?"

Rmq 5: Variante en ligne

I.2. Algorithme naïf :
 Algorithme POUR : ALGORITHME NAÏF
 POUR $j \leftarrow 1$ FAIRE
 TRAITER t (j, j+m et t[j]=m) FAIRE
 SI $t[j \dots j+m-1] = m$ FAIRE
 L renvoie t[j] est factif

Rmq 7: Cet algorithme résout SEARCH(m, t) en temps $\mathcal{O}(n \cdot m)$

Rmq 8: Dans un modèle de Random-Access, cet algorithme tourne en temps $\mathcal{O}(n)$.

I.3. Algorithme parallèle: SHIFT-AND
 • On considère qu'une opération élémentaire est le "ET" bit à bit sur des registres machines de taille p telle que si $m \leq p$

Def 5: (Motif) Étant donné un mot m et un préfixe pre de m , pour $c \in \Sigma$ on définit le bit de pre et c ssi $pre \cdot c = m$

Ex 10: $m = aaaaa$. $B_0(m) = 1$, $B_1(m) = 0$, $B_2(m) = 0$, $B_3(m) = 0$, $B_4(m) = 0$

Algorithme SHIFT-AND(m, t)
 PRE-FACTUL
 POUR $c \in \Sigma$ FAIRE $B[c] \leftarrow 0$
 POUR $j = 1 \dots m$ FAIRE $B[j] \leftarrow B[j-1] \oplus 2^{m-j}$
 RECHERCHE
 POUR $i = 1 \dots n$ FAIRE
 POUR $k = 0 \dots m-1$ FAIRE
 SI $(t[i-k] \oplus B[k]) = 0$ ALORS renvoie t[i-k] est factif

Rmq 12: Shift-AND résout SEARCH(m, t) en temps $\mathcal{O}(n)$.

I.4. Algorithme KMP :

Def 13: (Bordure) Étant donné un mot $u \in \Sigma^*$, on dit que w est une bordure de u ssi w est un préfixe de u et w est un suffixe de u , et $|w| \leq |u|$ (ie $w \in u$).

Ex 14: Bordures d'un mot

Def 14: (Bordure) Étant donné $u = a_1 \dots a_n$, on construit l'ensemble des facteurs f de u qui sont des préfixes de u et des suffixes de u . On se pose la question de savoir si f est un facteur de u .

Ex 16: On a des facteurs de $ababab$ qui sont des préfixes et des suffixes de $ababab$.

Rmq 17: Soit f l'ensemble des facteurs de u . Alors Fact(u) \subseteq Fact(f)

Rmq 18: Si f est l'ensemble des préfixes et des suffixes de u , on a Fact(u) \subseteq Fact(f)

Algo 19: BOM(u, t)

PROCHAUN
 L ← fonction de transition de l'occure des lettres de u
 REVERSER

```

    pos ← 0
    TANT QUE pos ≤ n-m FAIRE
        q ← état initial de l'automate
        j ← m
        TANT QUE j > 0 ET q ≠ q FAIRE
            q ← δ(q, tpos+j)
            j ← j-1
        SI j = 0 ALORS retourner occurrence à la position pos+1
        pos ← pos + j-1
    
```

Prop 20: BOM(t) retour SERCH(u, t) en temps O(mn) en pire cas, et en temps O(m log(m)) dans un modèle de données R.

II.6. Algorithme de Karp-Rabin:

Algo 21: KR(u, t)

```

    a ← h(u)
    POUR i = 1..n-m+1 FAIRE
        SI h(t[i..i+m-1]) = a ALORS
            [ SI h(u[i..i+m-1]) = a ALORS retourner à la pos i
            ]
    
```

Prop 22: KR(u, t) retour SERCH(u, t) en temps O(mn) en pire cas

II. Indirection:

Soit t un texte fixe dans l'alphabet.

II.1. Problèmes:

- Pr 23: (Appartenance) Soit $a \in \Sigma^*$ trouver le plus long préfixe de $a \in \text{Fact}(t)$
- Pr 24: (Position) Soit $a \in \text{Fact}(t)$, déterminer la position de sa première occurrence dans t
- Pr 25: (Occurrences) Soit $a \in \text{Fact}(t)$, déterminer le nombre d'occurrences de a dans t
- Pr 26: (Parties) Soit $a \in \text{Fact}(t)$, déterminer la séquence partielle des occurrences de a dans t
- Ex 27: $t = \text{ababababc}$, $a = \text{ab}$
- Pr 27: $t = \{1, 3, 6\}$, # occurrences = 3, FirstPos = 1, LastPos = 6

II.2. Table des suffixes:

Def 28: La table des suffixes de t est un tableau contenant tous les suffixes de t par ordre décroissant.

Prop 29: Cette table peut être construite en temps O(n log n)

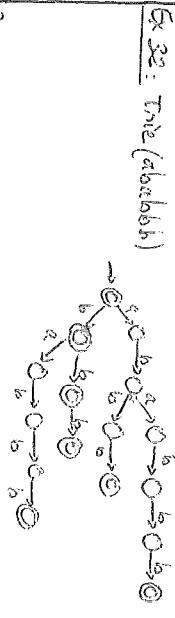
Prop 30: Etant donné la table des suffixes:

Appartenance	se retour en temps	$O(m + \log n)$
Position	"	$O(m + \log n)$
# occurrences	"	$O(m + \log n)$
Parties	"	$O(m + \log n + \# \text{ occurrences})$

II.3. Arbre des suffixes:

Def 31: (Trie) Le trie d'un texte est un automate acroreccit définie par:

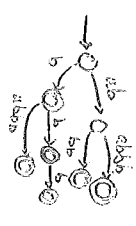
- Les feuilles de t sont les états
- La fonction de transition δ est définie par $\delta(i, a) = j$ si ia est un préfixe de t.
- Les états branches sont les suffixes de t.



Prop 33: Le Trie des suffixes peut être construit en temps et en espace $O(n^2)$

Def 34: L'arbre des suffixes est obtenu en supprimant les nœuds de degré ≤ 1 du Trie des suffixes.

Ex 35: Tc (ababbb):



Prop 36: Tc peut être stocké en $O(n \log n)$.

Def 37: (L'arbre des suffixes) $\forall a \in \Sigma$, $\text{ls}(a) = \Sigma$ pour chaque état z de l'automate

DEF 38:

Algorithme d'Ukkonen pour construire Tc en temps linéaire en $\log n$

Def 38: $\text{max}(a) = \max \{ |z| : z \in \Sigma^* \text{ et } \#(z) \neq 0 \}$
 $\text{min}(a) = \min \{ |z| : z \in \Sigma^* \text{ et } \#(z) \neq 0 \}$

Prop 39: $\text{first-pos}(a) = |t| - \text{max}(a) - |a|$

$\text{last-pos}(a) = |t| - \text{min}(a) - |a|$

Prop 40: $T(n)$ pouvant être constant $\in O(n)$, on peut répondre à Appelbeurre, Ration, # courses, pour un mot z , au temps $O(|z| \times \log |H|)$ et en mémoire $O(n)$.

Prop 41: Si les états terminaux de $T(n)$ sont atteints par une suite de suffixes consécutifs, la liste des positions des courses de z peut se retourner en temps $O(|z| \log |H| + \text{courses})$.

III. Compression : codes de Huffman

III.1. Codes préfixes

Def 42: Un langage Σ^* est dit préfixe si $\forall x, y \in \Sigma^*$, x n'est pas préfixe de y .

Ex 43: $\Sigma = \{a, b, c, d, e\}$
 $\Sigma^* = \{ \emptyset, a, b, c, d, e, aa, ab, ba, bb, \dots \}$

Def 43: (arbre d'un langage) L'arbre d'un langage est défini comme la Trie des codes $fc(a)$, $a \in \Sigma^*$.

Ex 44:



Prop 45: L'arbre d'un langage préfixe optimal est complet.

Prop 46: (codé) Pour $e \in \Sigma$, soit $f(e)$ la fréquence de e dans le texte, alors $|f(e)| \leq |f(a)| \leq |f(c)|$.

III.2. Construction algorithmique d'un langage préfixe optimal:

Algo 47: HUFFMAN(Σ) (avec $n = |\Sigma|$)

$Q \leftarrow$ file de priorité vide

POUR $i := 1 \dots n$ FAIRE INSERTER($Q, \Sigma_i, P(\Sigma_i)$)

POUR $i := 2 \dots n$ FAIRE

$\Sigma_i \leftarrow$ concat nouveau nœud

$P(\Sigma_i) \leftarrow$ EXTRAIRE-MIN(Q)

distance $P(\Sigma_i) \leftarrow$ EXTRAIRE-MIN(Q)

INSERTER($Q, \Sigma_i, P(\Sigma_i) + P(\Sigma_{i-1})$)

retourner EXTRAIRE-MIN(Q)

Lemme 47: Soit un ensemble de caractères $d_1, d_2, \dots, d_n \in \Sigma$ les deux caractères du mot $d_{i-1}d_i$.

Il existe un langage préfixe optimal des lang. $|L_1| = |L_2|$ et $(d_1, d_2) \in L_1$ n'a pas de préfixe commun.

Lemme 48: Les caractères d'un langage préfixe optimal vérifient la propriété de sous-structure optimale.

Prop 50: L'algorithme HUFFMAN(Σ) produit un langage préfixe optimal.

IV. La notion de similitude : un aperçu

IV.1. FBS : large sous-séquence commune:

Def 51: (large-séquence) Pour $x_1, y_1 \in \Sigma^*$, on dit que x est une sous-séquence commune de x_1, y_1 si x existe $(i, j) \in \mathbb{N}^2$ tq $x_i = y_j$, $\forall i \in \mathbb{N}$.

Ex 52: BCBAB est une sous-mot de ABCBABC

Prop 53: (LCS) Étant donné $x_1, y_1 \in \Sigma^*$, trouver une plus large sous-séquence commune z de x_1, y_1 .

Prop 54: LCS peut être calculé par programmation dynamique en temps $O(|x_1| \cdot |y_1|)$.

IV.2. Distance d'édition:

Prop 55: Étant donné $x_1, y_1 \in \Sigma^*$, la distance d'édition est le nombre minimum d'opérations pour passer de x_1 à y_1 .

Prop 56: Venant de Σ^* avec opérations permises

Ex 57: Opérations : insertion d'une lettre, suppression d'une lettre, changer une lettre. On peut passer de x_1 à y_1 en 6 opérations.

Prop 58: L'exemple 57 peut être résolu par prog. dyn. en temps $O(|x_1| \cdot |y_1|)$.

Rem 59: L'algorithme de séquence se réduit au $Prop 57$.

IV.3. Recherche de motifs approchés:

Prop 60: Étant donné un motif x et un texte t , on cherche toutes les occurrences de x dans t avec une distance d'édition $\leq k$ pour k fixe.

Prop 61: Le $Prop 60$ peut être résolu par prog. dyn. en temps $O(|x| \cdot |t| \cdot k)$.

Refs: • Cormen : Huffman et k -DIP et prog. dyn.

• Navarro : Editable Pattern Matching in Strings. Algo de recherche comparative, recherche approchée

• Crochemore : Text Algorithms. Algo d'indices : Huffman

• Crochemore : Algorithmes de séquence et similitude. Le mot

* Éléments d'algèbre linéaire : chp 10

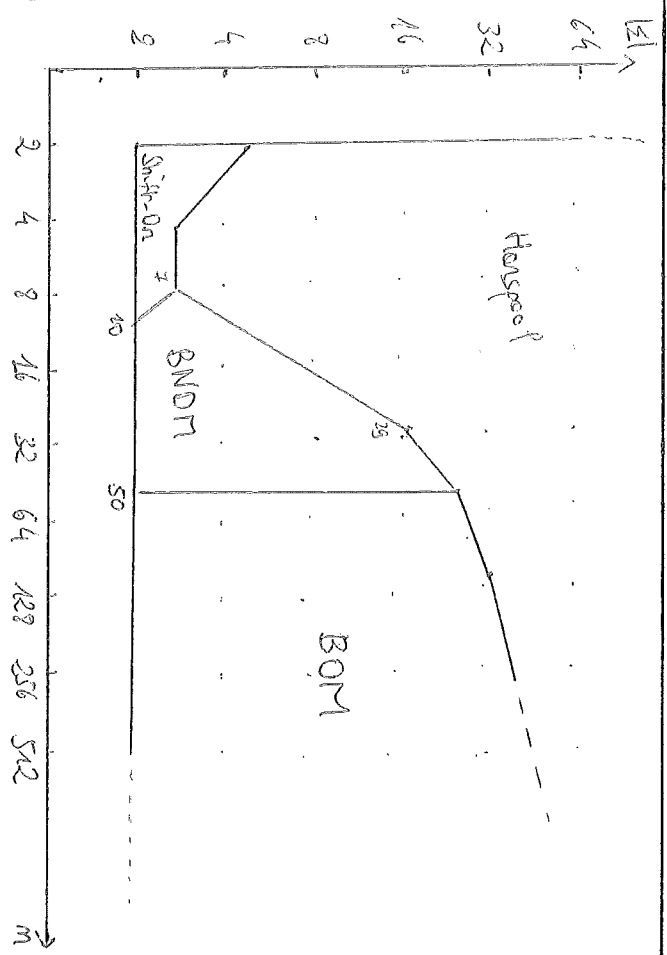


Figure 1: Comparatif des algorithmes de recherche

	(B)	D	(C)	A	(B)	(E)
A	0	01	01	1	← 1	1/2
(B)	0 ₁	← 1				
(C)			0 ₂	← 2		
(D)					0 ₃	
(E)						1 ₃
D						3 ₁
(A)						3 ₁
B						1 ₄

Figure 2: Construction de la plus longue sous-séquence commune.