

mettre des " (Admix) "

Notations Σ alphabet = en. fin.
 $n \in \mathbb{Z}^+$ n lettres i-ème lettre
 $n \in [i, j]$ factor de i à j
 n préfixe de longueur n

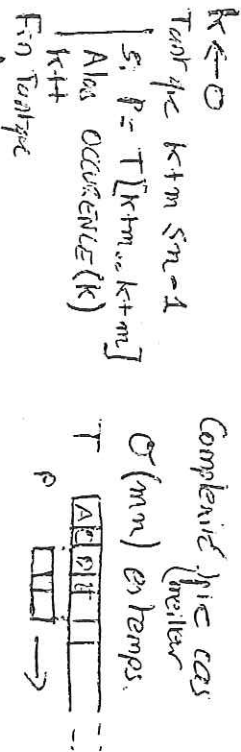
Motivations : - recherche et remplacement : traitement de texte

- traitement automatique de journaux / scripts.
- détection information : audio, bio-informatique
- transmission de l'information
- compilation

I - RECHERCHE DE MOTIFS

1 - Motif exact

Texte $T \in \Sigma^+$ $|T| = n$ Motif $P \in \Sigma^+$ $|P| = m \leq n$
Algorithme naïf

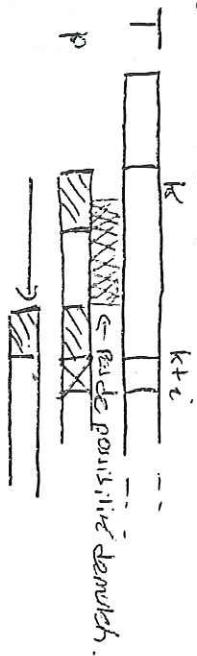


Utilisation des bordes : algorithme de
Norani-Raitt.

DF-Bord (u) = plus long préfixe strict de u aussi suffixe de u .

Exemple Bord(abcaba) = abca

Idee: être plus astucieux en cas d'échec on obtient la structure de P.



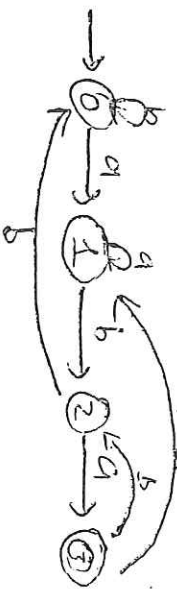
nécessaire de calculer la table des bordes de P en temps $\Theta(m|Z|)$.

Remarque: cette méthode est encore améliorée par l'algorithme de Knuth-Morris-Pratt.

c - Automate des occurrences

Idee: Trouver dans T char haves - tous les T_i ayant P pour suffixe ie. $i \in \mathbb{Z}^+$ $P \in \Sigma^+$

On construit donc un automate par exemple P = abca



On note $\Phi(T_i)$ l'état après lecture de $T[0..i]$ et $\sigma(T_i)$ le plus long préfixe de P suffixe de T_i .

Théorème $\forall i \in \{0, 1, \dots, n\}$ $\Phi(T_i) = \sigma(T_i)$
 La table des occurrences est minimal pour reconnaître Σ^* .

Extension plusieurs mots = Aho-Corasick

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99

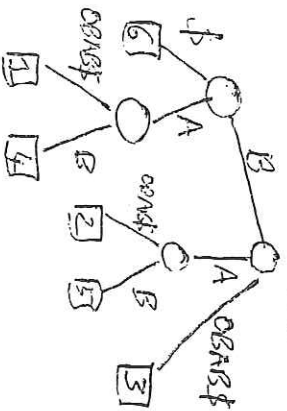
DEV 4

1- Arbre des suffixes

Idee : Texte préfixé dans lequel on va faire un grand nombre de recherches incomplètes à l'avance = grand retranchement.

Définition L'arbre des suffixes du texte T est un arbre engendré à n feuilles (par chaque préfixe) équilibré par des facteurs de T tels que la concaténation des étiquettes de la racine à la feuille k est le suffixe d'indice k.

Exemple T = BAABABAB\$
1 2 3 4 5 6



Théorème [Ukkonen 85]

L'arbre des suffixes pour s[1..n] construit en temps linéaire $\Theta(n)$

La recherche de P (padding spécifique) dans s[1..n] est faite en parcourant le sous-arbre.

e) Bilen

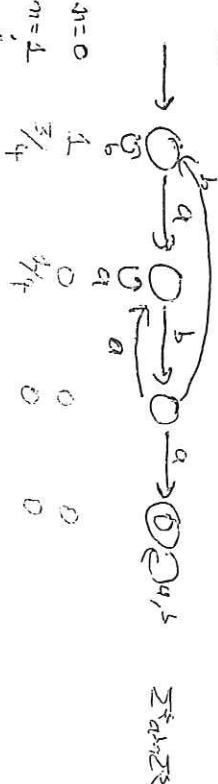
	Rechercher	Exécution	Espace
Naïf	-	$\Theta(mn)$	$\Theta(1)$
MP	$\Theta(m Z)$	$\Theta(n)$	$\Theta(m)$
kmf	$\Theta(m)$	$\Theta(m)$	$\Theta(m)$
Algorithme	$\Theta(m Z)$	$\Theta(m)$	$\Theta(m Z)$
Aho	$\Theta(m)$	$\Theta(m)$	$\Theta(m \log n)$

2) Probabilité d'occurrence d'un motif

Idee Source aléatoire (ex: A 20%, B 17%, ...) et T_n issu de cette source. On s'intéresse toujours, est-ce significatif?

On peut modifier l'algorithme des occurrences et le hauteur m et en hauteur de Markov

Ex: a b a a 1/4 1/4



Rq cela fonctionne car l'algorithme est déterministe.

3) Cas général : opérations régulières.

En informatique, on cherche souvent P sous forme déterministe régulière.

Idee: construire l'algorithme et le déterminer / minimiser $P_m = \text{cost étendu}$ exponentiel.

logiciel : GREP

II COMPARAISON

→ T3 cible en biologie: homologies → orthographe
→ traitement de signal → recherche de virus.

1) Distances orthogonales

- Mots de même taille = distance de Hamming (nbr de différences)
- Distance d'édition ou de Levenshtein : nbr minimum de substitutions, insertions, suppressions.
- Matrices de substitutions : raison en bioinfo.

DEUX Calcul de la distance d'édition par programmation dynamique.

myqt.

$$T[i, j] = \min(x \in \{0, 1\}, y \in \{0, 1\})$$

$$T[-1, -1] = 0$$

$$T[i, -1] = T[i-1, -1] + 1$$

$$T[-1, j] = T[-1, j-1] + 1$$

$$T[i, j] = \min \begin{cases} T[i-1, j-1] + 1 & \text{(substitution)} \\ T[i-1, j] + 1 & \text{(suppression)} \\ T[i, j-1] + 1 & \text{(insertion)} \end{cases}$$

Utilisation : correction orthographe

3) Détection de répétitions

Ex: on cherche les facteurs qui se répètent 2x dans BAAGAB.

- ① construire l'arbre des suffixes
 - ② parcourir l'arbre et compter le nombre de feuilles de chaque sous-arbre (récurrent)
 - ③ Parcours à profondeur, calculer les racines ≥ 2
- Ex = B et BA et A (et C).

III COMPRESSION : CODE DE HUFFMAN

On a soit un texte écrit avec Σ ou alors une source aléatoire émettant émetteur du lettres de Σ . On note f_i la fréquence / probabilité d'occurrence de la lettre i .

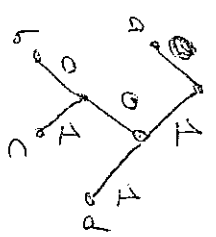
On veut encoder en binaire Σ de façon à minimiser le nombre de bits: $\sum_{a \in \Sigma} f_a \times \text{longueur } a$.

Huffman: on construit un arbre d'encodage de manière glabonne.

- ① $|\Sigma|$ feuilles de poids f_i .
- ② Trouver les deux sous arbres de poids minimum et fusionner, somme les poids des sous-arbres.

Ex

Σ	f_i
a	70
b	20
c	3
d	37



Théorème de l'information.

- Information d'une lettre i : $-\log_2 p_i$ (a.s.s)
- Entropie de la source: $H = -\sum_{x \in \Sigma} p_x \log_2 p_x$
- = information moyenne.
- Théorème de Shannon: on ne peut pas faire mieux que H.
- Huffman: codage optimal au sens de la H et si de l'information.

Notations

$P \in \Sigma^*$ $T \in \Sigma^*$
 $|P| = m \leq |T| = n$
 \exists "est suffixe de"
 p_i : ième préfixe
 $\sigma(x) = \max\{|P_n| \exists x\}$

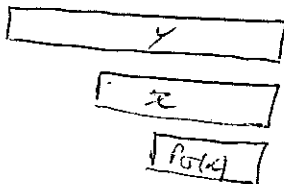
Automate \mathcal{A}

$Q = \llbracket 0, m \rrbracket$
 $I = 0$ $F = \{m\}$
 $\delta(q, a) = \sigma(P_q a)$
 $\phi: \begin{cases} \phi(\epsilon) = 0 \\ \phi(xa) = \delta(\phi(x), a) \end{cases}$

Théorème $\forall i \in \llbracket 0, n \rrbracket, \phi(T_i) = \sigma(T_i)$. En particulier $\mathcal{L}(\mathcal{A}) = \Sigma^* P$.
 De plus \mathcal{A} est minimal!

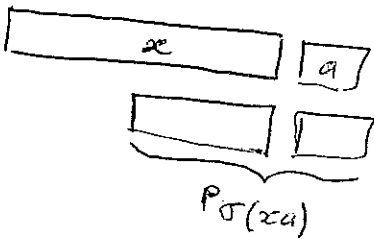
Démonstration Correction

Lemme 1 $\forall x \in \Sigma^* \forall y \in \Sigma^* x \supset y \Rightarrow \sigma(x) \leq \sigma(y)$



$P_{\sigma(x)} \supset x$
 $x \supset y$ } $P_{\sigma(x)} \supset y$ donc $\sigma(x) \leq \sigma(y)$

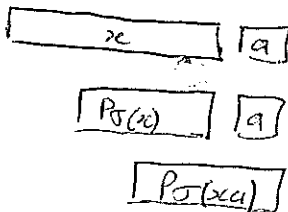
Lemme 2 $\forall x \in \Sigma^* \forall a \in \Sigma, \sigma(xa) \leq \sigma(x) + 1$



$P_{\sigma(xa)} \supset xa$
 donc $P_{\sigma(xa)-1} \supset x$
 donc $\sigma(xa)-1 \leq \sigma(x)$

Lemme 3 $\forall x \in \Sigma^*, \forall a \in \Sigma, \sigma(P_{\sigma(x)} a) = \sigma(xa)$

$P_{\sigma(x)} \supset x \Rightarrow P_{\sigma(x)} a \supset xa \xRightarrow{\text{(lemme 1)}} \sigma(P_{\sigma(x)} a) \leq \sigma(xa)$



On a $P_{\sigma(xa)} \supset xa$
 et $P_{\sigma(x)} a \supset xa$

donc $|P_{\sigma(xa)}| = \sigma(xa) \leq \sigma(x) + 1 \stackrel{\text{(lemme 2)}}{\leq} |P_{\sigma(x)}| + 1 = |P_{\sigma(x)} a|$

d'où $|P_{\sigma(xa)}| \leq |P_{\sigma(x)} a|$ et $P_{\sigma(xa)} \supset P_{\sigma(x)} a$

d'où $\sigma(xa) \leq \sigma(P_{\sigma(x)} a)$

ve du théorème Réécriture n° 1

$$i=0 \quad T_0 = \varepsilon \quad \phi(\varepsilon) = 0 \quad \text{et } \sigma(\varepsilon) = \max\{i \mid p_i \in E\} = 0 \quad \phi(T_0) = \sigma(T_0)$$

$$\begin{aligned} \text{si } i \rightarrow i+1 \\ a = T_{i+1} \\ \phi(T_{i+1}) = \phi(T_i a) = \delta(\phi(T_i), a) = \delta(\sigma(T_i), a) = \sigma(p_{\sigma(T_i)} a) \stackrel{\text{lemme 3}}{=} \sigma(\cancel{T_i} a) \\ \sigma(T_i a) = \sigma(T_{i+1}). \end{aligned}$$

Minimalité

Soient i, j q.c.q vérifiant $0 \leq i < j \leq m$

On pose $Q_k \in \Sigma^*$ défini par $p_k Q_k = P$

On note L_k le langage reconnu à partir de l'état k .

$$\phi(p_j Q_j) = \phi(P) = m \quad \text{donc } Q_j \in L_j$$

$$\phi(p_i Q_j) < m \quad \text{car } |p_i Q_j| < m \quad \text{donc } Q_j \notin L_j$$

Ainsi $i \neq j$ dans l'équivalence de Nerode. Tous les états sont distingués
2 à 2 donc A est minimal.