

Notation: On notera l'alphabet sur lequel on travaille Σ , supposé fini. Une lettre sera notée $T = t_1 t_2 \dots t_{n-1}$; où $\forall i \in \{0; n-1\}$, $t_i \in \Sigma$ et n est la longueur du texte (nombre de caractère). On note de même un motif $M = m_0 \dots m_{k-1}$, de longueur k . k et n sont supposés toujours finis.

I Algorithmes de recherche de motif dans un texte

Problème: Soit un texte T et un motif M . On cherche le plus petit indice $j \in \{0; n-1\}$ où n est la longueur du texte tel que $t_j t_{j+1} \dots t_{j+k-1} = M$, k la longueur du motif.

Algorithme 1: Algorithme naïf

Entrée: T texte, M motif

Sortie: indice j minimal tel que $t_j t_{j+1} \dots t_{j+k-1} = M$

Recherche naïve ($T; M$)

$n, k = \text{longueur}(T); \text{longueur}(M);$

Pour $j = 0$ à $n-k$:

 Si $M = t_j t_{j+1} \dots t_{j+k-1}$:
 renvoyer j

renvoyer "erreur: T ne contient pas M "

La complexité temporelle en pire cas est en $O(nk)$

→ Exemple 2:

Si on cherche $aaab$ dans $aaaaaa \dots aab = T$
grand nombre de a

On voit apparaître la complexité

Remarque 3: Dans le cas de texte en français, on peut estimer que le texte ne contient pas beaucoup de préfixes du motif. Dans le cas de la recherche de séquence dans un génome, on peut atteindre la complexité en pire cas.

Algorithme 4: Rabin-Karp.

On suppose $\Sigma = \{0; n-1\} \subset \mathbb{N}$.

Principe: On va faire un pré-traitement sur T . Si on cherche un motif de taille k , on va créer un tableau de taille $n-k$ où la case i contient $t_i \dots t_{i+k-1}$. Soit $m = \sum_{i=0}^{k-1} m_i$. Il suffira alors de regarder pour les i tel que $T[i] = m$ si $M = t_i t_{i+1} \dots t_{i+k-1}$

Complexité: Ce paradigme est en pratique plus efficace que la recherche naïve mais, en pire cas, la complexité est en $O(nk)$. De plus, il faut un tableau de taille $n-k$ en mémoire.

Exemple 3: Soit $\Sigma = \{0, 1\}$. On cherche $M = 011$.
 Si $T = (002) \dots (002) 011$,

on voit apparaître la complexité en pire cas: $O(nk)$.

Definition 6: Soit T un texte; M un motif sur un alphabet Σ .

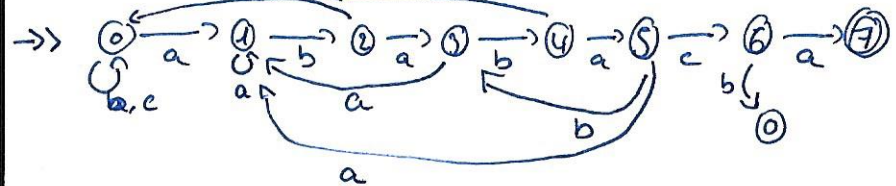
On définit l'automate des occurrences par:

- l'alphabet Σ
- $Q = \{0, \dots, k\}$ les états, où k est la longueur du motif.
- $\{0\}$ l'état initial
- $\{k\}$ l'état final
- δ la fonction de transition, tel que $\forall q \in Q, \forall a \in \Sigma$, $\delta(q, a)$ renvoie l'indice j maximal tel que $m_0 \dots m_{j-1}$ soit préfixe de M et suffixe de $m_0 m_1 \dots m_{q-1} a$.

Remarque 7: Un tel objet a besoin d'un espace mémoire en $O(|\Sigma|k)$

Exemple 8: Soit $P = ababaca$; $\Sigma = \{a; b; c\}$.

L'automate des occurrences associées est



Si en 1, 2, 3, 4, 6; on lit un c , on est renvoyé à l'état 0. Une fois en 7, l'automate s'arrête.

Algorithme 9: Morris-Pratt

Morris-Pratt ($T; M; \Sigma$)

$\delta = \text{créer_delta}(M; \Sigma)$

$n = \text{longueur}(T)$

$q = 0$

pour $i = 0$ à $n-1$:

$q = \delta(q, T[i])$

 si $q = m$

 renvoyer $i - k$

renvoyer "T ne contient pas M"

La complexité hors pré-traitement est linéaire en la taille du texte. Le pré-traitement (contenu dans créer_delta, qui renvoie la fonction δ de transition de l'automate des occurrences associée à M) nécessite une place mémoire en $O(|\Sigma|M)$ et s'effectue en $O(|\Sigma|k^2)$ opérations.

creer_delta(M; Σ):

k = longueur(M)

Pour q = 0 à m

pour tout a ∈ Σ

k = min(m+1, q+1)

répéter

k = k-1

jusqu'à M_k J M_q a

δ(q, a) = k

renvoyer δ

On note "M_k J M_q a" le fait que m_0...m_{k-1} soit suffixe de m_0...m_{q-1} a.

Les 2 premières boucles comptent pour O(m|Σ|) opérations. La boucle répéter peut se faire au plus q ≤ k fois, la vérification M_k J M_q a se fait en au plus k opérations → complexité en O(k^3|Σ|).

Remarque 10 (admis): IP est possible de faire tomber la complexité de creer_delta à O(k|Σ|) opérations.

Dev 1: Algorithme de Knuth-Morris-Pratt, validité. L'algorithme a une complexité spatiale en O(k), temporelle en O(n).

Remarque 11: Les algorithmes Morris-Pratt et Knuth-Morris-Pratt peuvent se généraliser au cas d'ensemble de mots finis.

II Distance d'édition et recherche approchée de motif

Exemple 12: Si la recherche n'est pas sensible à la casse, on peut retrouver "Abc" au lieu de "abc". On va généraliser la recherche approchée de motif.

Définition 13: On pose les fonctions suivantes:

Sub: Σ × Σ → ℝ⁺; où sub(a;b) est le coût de substitution de la lettre a par b. sub(a;a) = 0, (a;b) ∈ Σ²

Del: Σ → ℝ⁺; où del(a) est le coût de la suppression de a, a ∈ Σ

Ins: Σ → ℝ⁺; où Ins(a) est le coût d'insertion de la lettre a.

Définition 14: On note ∀ m; y ∈ Σ*, mots sur Σ, Σ_m y l'ensemble des suites d'opérations menant de m au mot y. C'est non vide, car on peut supprimer toutes les lettres de m puis ajouter toutes celles de y.

Définition 15: Lev: Σ* × Σ* → ℝ⁺

(m; y) ↦ min(cout de T, T ∈ Σ_m y)

est appelée la distance d'édition.

Théorème 16: Lev est une distance sur Σ* si et seulement si sub est une distance sur Σ

et
• ∀ a ∈ Σ, del(a) = Ins(a) > 0

Définition 17: Un alignement entre x et y ∈ Σ* est un mot z sur l'alphabet (Σ ∪ {ε})² / {ε; ε} dont la projection sur la première composante est x; y sur la deuxième composante. On note z = (x_0 ... x_n / y_0 ... y_n)

Exemple 18: Si x = ACGA; y = ATGCTA, un alignement est z = (ACGεεA / ATGCTA). (A) représente une substitution, (ε) représente la suppression.

Définition 19: Le coût d'un alignement est la somme des coûts de chaque opération représentée.

Définition 20: Un alignement z entre m; y ∈ Σ* est dit optimal si son coût est Lev(m; y).

Dev 2 (Bashien) Calcul de l'alignement optimal entre 2 mots sur Σ*

Problème: Soit un texte T, un motif M, de longueur respective n et k. On souhaite trouver les séquences M', de longueur k' ≥ k, avec un nombre maximal de différences, ie où Lev(M, M') ≤ d, d est finie, d ≤ k. On suppose que toute substitution, suppression et ajout coûte 1

Dev 3 (Rami) Recherche des facteurs de T à une distance au plus k de M par la programmation dynamique. Amélioration avec les cases l-spéciales. Complexité.

Exemple 21: Ce problème convient dans le cas de recherche de séquences ADN dans un génome en tenant compte d'éventuelles mutations.

III Compression

Remarque 22: On ne peut pas compresser indéfiniment sans perte d'information. On ne peut pas forcément stocker l'information de n bits en $n-1$ bits.

Définition 23: L'algorithme de Lempel Ziv (LZ) est un algorithme de compression par substitution de facteurs. Il remplace des séquences par l'indice des séquences dans un dictionnaire. Il fonctionne en une seule lecture de l'entrée.

Exemple 24: "Un bon gars est un gars bon." donnera alors "Un bon gars est [1] [3] [2]"

Remarque 25: Cette base d'algorithme est aujourd'hui majoritairement utilisée

Algorithme 26: LZ 77

Basé sur une fenêtre qui glisse de gauche à droite sur le texte, divisée en 2 parties: Une qui contient le dictionnaire, l'autre qui renvoie le texte en premier (tampon de lecture)

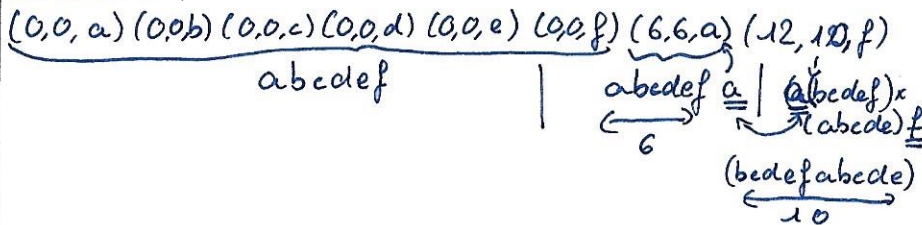
Initialement, la première partie est vide, la dernière contient T. A chaque itération, on cherche la plus longue séquence préfixe de T contenue dans le dictionnaire, codé par (i, j, c) où i est la distance entre le début du hupon et la position de la répétition dans le dictionnaire

→ j est la longueur de la répétition

→ c est le premier caractère du hupon différent du caractère correspondant dans le dictionnaire.

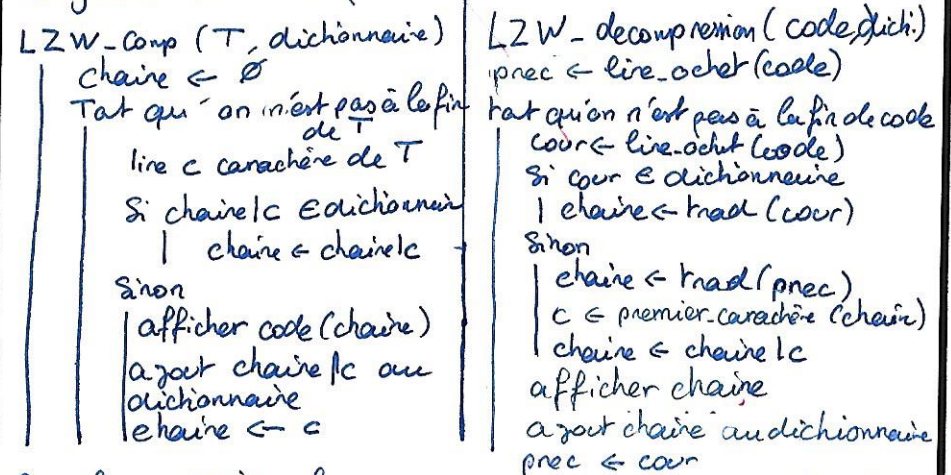
La fenêtre coulisse de $j+1$ caractères.

Exemple 27: On applique LZ 77 à abcdefabcdefabcdefabcdef. On obtient:



Algorithme 28 Lempel-Ziv-Welch (LZW)

Il consiste à ne coder que l'indice n dans le dictionnaire. Il est nécessaire d'avoir un dictionnaire initial (Par exemple la table ASCII). On remplace chaque groupe de caractères déjà connu par son code et on ajoute au dictionnaire un nouveau groupe formé par ce groupe suivi du prochain caractère.



Pour la compression, le si-sinon gère le seul cas problématique: la répétition de chaîne (cas où cour ∉ dictionnaire). On peut avec les mêmes dictionnaires initiaux.

IV Codes correcteurs

Problème: On veut vérifier algorithmiquement si un message reçu par un ordinateur est correct.

Exemple 29: Code de parité

On envoie un message de n bits. On va envoyer un $n+1$ ième bit de sorte à ce que la somme des éléments soit pair.

010 → 010 1 || Permet de détecter une erreur
 011 → 011 0 || mais ne permet pas de corriger.

Exemple 30: Code de répétition.

Soit m le message envoyé. On va répéter 3 fois m .

1001 → 1001 1001 1001

Si on envoie 1000 1001 1001, on peut détecter et corriger l'erreur.