

Soit  $\Sigma$  un alphabet fini.

I) Recherche d'un motif dans un texte

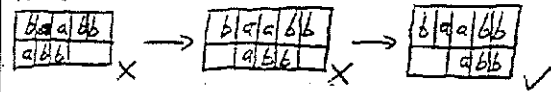
Problème fréquent dans les traitements de texte, la bio-informatique, l'analyse lexicale...  
On note  $P \in \Sigma^*$  le motif de longueur  $m$ , et  $T \in \Sigma^*$  le texte de longueur  $n$ .

1) Algorithme naïf - Fenêtres glissantes [Cormen]

Pour  $k \in \{1, n-m+1\}$ :  

$$\left[ \begin{array}{l} S[T[s..s+m-1]] = P \\ \text{Afficher "occurences de P en k"} \end{array} \right] \text{ RECHERCHE\_NAIVE}(P, T)$$

Exemple: recherche de abb dans baabb



Prop 1 Dans le pire des cas, RECHERCHE\_NAIVE a une complexité temporelle en  $O(m(n-m+1))$  (comptée en nombre de comparaisons)

2) Algorithme de Rabin-Karp [Cormen]

Les prochains algorithmes tentent d'améliorer la complexité par un prétraitement de P ou T. Celui de Rabin-Karp a une complexité au pire de  $O(m(n-m+1))$  aussi, mais il est très utile quand les textes sont courts, et dans la détection de plagiat.

Idee  $\Sigma = \{0, 1, \dots, d\}$ . Soit  $q$  un entier fixe.

- $w \in \Sigma^*$  est un nombre en base  $d$ , noté en minuscule.  
 no calcul par la méthode de Horner: par  $w_1, \dots, w_k$ ,  
 on calcule  $w_k + d(w_{k-1} + d(\dots + d(w_1))) [q]$ .
- On note  $t_s$  le nombre associé à  $T[s+1 \dots s+m]$   

$$t_{s+1} = d \cdot (t_s - T[s+1] \cdot h) + T[s+m-1] \text{ mod } q$$
 où  $h = d^{m-1} [q]$  est calculé à l'avance.
- Dans l'algorithme naïf, on effectue la comparaison  $T[s..s+m-1] = P$  ssi  $p = t_{s+1}$ .

3) Algorithme de Knuth-Morris-Pratt [Cormen]

a) Automate des occurrences

Soit  $\sigma: \Sigma^* \rightarrow \{0, 1, \dots, m\}$  tq  $\sigma(x)$  est la longueur du plus long préfixe de P qui soit suffixe de x.

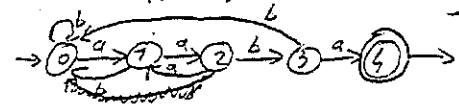
On définit un automate qui reconnaît  $\Sigma^*P$ :

vérifier si P apparaît dans T est alors faisable en  $O(n)$ :

les états sont  $\{0, 1, \dots, m\}$ , où 0 est initial et m final.  
 la fonction de transition  $\delta$  vérifie pour un état  $q$  et une lettre  $a$ :  

$$\delta(q, a) = \sigma(P[1..q] + a)$$

Exemple: Pour aaba



Thm 2  $A_p$  est l'automate minimal reconnaissant  $\Sigma^*P$ .

b) Fonction préfixe

On préfixe utilise une table  $\pi[1..m]$  tq  
 $\pi[i]$  est la longueur du plus long préfixe de P qui est suffixe propre de  $P[1..i]$ .

Prop 3 Dans  $A_p$ , si  $q = m$  ou que  $P[q+1] \neq a$ , on a  

$$\delta(q, a) = \delta(\pi[q], a)$$

On peut donc recalculer  $A_p$  à partir de  $\pi$ .

FUNCTION\_PREFIXE(P): // Renvoie le tableau  $\pi$  associé à P.

```

    pi[1] ← 0
    R ← 0
    Pour q ∈ [2, m]:
        Tant que R > 0 et P[R+1] ≠ P[q]:
            R ← pi[R]
        Si P[R+1] = P[q] alors R ← R+1
        pi[q] ← R
    Renvoyer pi
    
```

Exemple: Pour aabaca:

$\pi = [0, 0, 1, 2, 0, 1]$

Thm 4 FUNCTION\_PREFIXE renvoie  $\pi$  en  $O(m)$ .

Algorithmique du texte: exemples et applications.

⊙ Algorithme

Complexité temporelle en  $O(n^2 + m)$

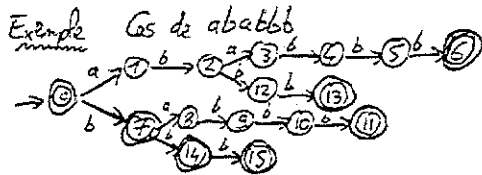
KMP(T, P):

```

π ← FONCTION_PREFIXE(P)
q ← 0
Pour i ∈ [1, n]:
    Tant que q > 0 et P[q+1] ≠ T[i]: q ← π[q]
    Si P[q+1] = T[i]: q ← q+1
    Si q = m:
        Affiche "occurrence de P en i-m"
        q ← π[q]
    
```

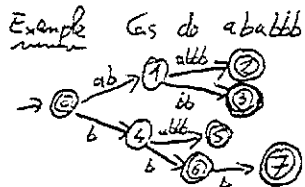
4) Arbre des suffixes [Crochemore]

Def 5 L'arbre des suffixes  $\mathcal{A}$  de T, noté  $\mathcal{A}(\text{Suff}(T))$ , est l'automate déterministe qui reconnaît  $\text{Suff}(T)$  l'ensemble des suffixes de T, et dans lequel deux chemins de même origine ont toujours des fins distinctes.



On insère les suffixes  $T[i..n]$  dans l'ordre croissant des  $i$ . C'est avec cette convention que les sommets ont été numérotés.

Def 6 L'arbre compact des suffixes de T s'obtient à partir de l'arbre des suffixes en supprimant les nœuds de degré 1 non terminaux.



Pour vérifier que P a une occurrence dans T, il suffit de parcourir l'arbre compact: s'il existe un chemin étiqueté par P, c'est bon.

La recherche a donc une complexité  $O(m)$ , c'est le prétraitement qui prend une place immense!

Utilisé quand on a besoin de chercher de nombreux motifs dans T.

Thm 7 (Wiens) On peut calculer l'arbre compact des suffixes de T avec une complexité au pire linéaire en l'espace et en temps.

II Comparaisons de mots

Soit  $x \in \Sigma^m$  et  $y \in \Sigma^n$ .

1) Distance de Levenshtein [Crochemore] (duplt)

Le but est d'étudier la façon la plus rapide de changer x en y. Très utile en bio-informatique et dans les recherches approchées.

On utilise trois opérations d'édition:

- substitution d'une lettre a de x en une lettre b de y (en des positions données) Associée à un coût  $\text{Sub}(a, b)$
- suppression d'une lettre de x, notée a, de coût  $\text{Del}(a)$ .
- insertion d'une lettre de y, notée b, dans x, de coût  $\text{Ins}(b)$ .

Def 8 On note  $\text{Lev}(x, y) = \min \{ \text{coût de } \sigma \mid \sigma \text{ suite d'opérations changeant } x \text{ en } y \}$ , où le coût de  $\sigma$  est la somme des coûts d'opérations utilisées dans  $\sigma$ .

Exemple  $x = \text{ACG-A}$  et  $y = \text{ATGCTA}$ , où  $\text{Ins} = \text{Del} = 1$  et  $\text{Sub}(a, b) = \begin{cases} 0 & \text{si } a=b \\ 1 & \text{sinon} \end{cases}$ .  
 $\text{Lev}(x, y) = 3$ , correspondant à un alignement  $\begin{pmatrix} \text{ACG} & - & - & \text{A} \\ \text{A} & \text{T} & \text{G} & \text{C} & \text{T} & \text{A} \end{pmatrix}$

Prop 9 Lev est une distance sur  $\Sigma^*$

ssi Sub est une distance sur  $\Sigma$  et  $\text{Del}(a) = \text{Ins}(a) > 0$  pour chaque  $a \in \Sigma$ .

Note: un alignement de x et y est une façon de visualiser une suite d'éditions changeant x en y. Formellement,

c'est un mot  $z$  sur l'alphabet  $(A \cup \{\epsilon\}) \times (A \cup \{\epsilon\}) \setminus \{(\epsilon, \epsilon)\}$  dont la projection sur la première composante (resp. seconde composante) est x (resp. y).

Si  $z = (\bar{x}_1, \bar{y}_1) \dots (\bar{x}_p, \bar{y}_p)$ , on écrit plutôt z sous la forme  $\begin{pmatrix} \bar{x}_1 & \dots & \bar{x}_p \\ \bar{y}_1 & \dots & \bar{y}_p \end{pmatrix}$

Le calcul de la distance d'édition se fait par programmation dynamique:

En notant  $T[i, j] = \text{Lev}(x[1..i], y[1..j])$ ,

on a le ~~proposition~~

ne peut pas être de Hoffman

Analyse  
 (L) Pas trop sa place ici ?

On a par  $i, j \geq 2$  :

$$T[i, j] = \min \begin{cases} T[i-1, j-1] + \text{Sub}(x[i], y[j]) \\ T[i-1, j] + \text{Del}(x[i]) \\ T[i, j-1] + \text{Ins}(y[j]) \end{cases}$$

En remarquant que seuls deux lignes ou colonnes de T suffisent au calcul, on a :

Prop 10 On peut calculer Lev en un temps  $O(mn)$  et en espace  $O(\min(m, n))$ .

Note: Pour exhiber la suite d'édition minimale, il faut y ajouter  $O(mn)$  en temps et en espace.

2) Plus longue sous-suite commune [Cochonore]

En notant  $\text{smc}(x, y)$  la longueur de la plus longue sous-suite commune à  $x$  et  $y$ , on voit que ce problème est une spécialisation du précédent, il revient à enlever la substitution des opérations possibles.

Prop 11 Si  $\text{Del} = \text{Ins} = 1$ , et que par  $a, b \in \Sigma$  on a  $\text{Sub}(a, b) \rightarrow \text{Del}(a) + \text{Ins}(b) = 2$ ,

alors  $\text{Lev}(x, y) = n + m - \text{smc}(x, y)$

On peut donc aussi résoudre ce problème par programmation dynamique. Supposons dans la suite que  $m \leq n$ .

SSC\_LONGUEUR(x, y) : // Renvoie la valeur de  $\text{smc}(x[1..i], y)$  par  $i \in [1, m]$ .

```

Pour i ∈ {0, 1, ..., m}: C1[i] ← 0
Pour j ∈ [1, n]:
  C2[0] ← 0
  Pour i ∈ [1, m]:
    Si x[i] = y[j]: C2[i] ← C1[i-1] + 1
    Sinon: C2[i] ← max {C1[i], C2[i-1]}
  C1 ← C2
Retourner C1
  
```

Prop 12 SSC\_LONGUEUR nécessite une complexité  $O(mn)$  en temps et  $O(m)$  en espace.

Thm 13 On peut calculer un plus long sous-mot commun à  $x$  et  $y$  en un temps  $O(mn)$  et avec un espace  $O(\min(m, n))$ .

DEVLPT

Autres algèbres de recherche de motif > Boyer Moore notamment  
 Arbre des suffixes: pour aligner

Compression

Utilisé dans les compressions d'image, de sons ou de vidéos.

On cherche à réduire la taille d'un texte en conservant son information. Présentons ici l'algorithme de Lempel-Ziv-Welch.

On utilise un "dictionnaire"  $d$ , i.e.

une table qui à chaque entier associe une chaîne de caractères.

Le dictionnaire est généralement initialisé avec  $2k$  cases, où  $k = |\Sigma|$ , et les  $k$  premières cases contiennent les lettres de  $\Sigma$ .

LZV(T) : // T est une chaîne de caractères: un mot de  $\Sigma^*$ .

```

W = ε
Pour i ∈ [1, |T|]:
  Si w + T[i] ∈ d: // On a donc une fonction restant si un mot s
    L w = w + T[i] // est dans le dictionnaire d.
  Sinon:
    Ajouter w + T[i] à d
    Imprimer code(w) // code(w) est l'entier i tq d[i] = w.
    W = T[i]
  Imprimer code(w)
  
```

Exemple  $ababcba$  sur un dictionnaire initialisé avec  $\begin{matrix} 0 & 1 & 2 \\ a & b & c \end{matrix}$

On aura le code  $0, 1, 3, 2, 4, 7, 0, 9, 10, 0$ .

Et le dictionnaire finira avec

a	b	c	ab	ba	abc	cb	bab	abab	aa	aaa	aaaa
0	1	2	3	4	5	6	7	8	9	10	11

Prop 14 L'algorithme de décompression n'a pas besoin du dictionnaire rempli en entrée.

On peut en avoir l'intuition avec l'exemple: si on veut décoder  $0, 1, 3, 2$ , en sachant que le dictionnaire comme à  $\begin{matrix} a & b & c \end{matrix}$ , on a  $0 \leftrightarrow a$ ,  $1 \leftrightarrow b$ . On devine donc que  $d[3] = ab$ . Donc  $3 \leftrightarrow ab$ , et  $2 \leftrightarrow c$ . On avait  $ababc$  avant compression.

Références

[Cormen] Algorithmes, Cormen, Leiserson, Rivest, Stein  
 [Cochonore] Algorithmique du texte, Cochonore, Hancart, Decroq  
 [Beauquier] Éléments d'algorithmique, Beauquier, Bostel, Chetanié

## Plus longue sous-suite commune (PLSSC) [Crochemore]

On utilise l'algorithme  $SSC\_LONGUEUR(x[0..m-1], y[0..n-1])$ ,  
 implémentation dynamique de l'idée suivante:

Lemme 1

$$S[i, j] = \begin{cases} 0 & \text{si } i = -1 \text{ ou } j = -1 \\ \text{smc}(x[0..i], y[0..j]) & \text{si } \begin{matrix} i \in \{0, \dots, m-1\} \\ j \in \{0, \dots, n-1\} \end{matrix} \end{cases}$$

alors pour  $i, j \geq 0$ , 
$$S[i, j] = \begin{cases} S[i-1, j-1] + 1 & \text{si } x[i] = y[j] \\ \max(S[i-1, j], S[i, j-1]) & \text{sinon} \end{cases}$$

Preuve Soient  $ua = x[0..i]$  et  $vb = y[0..j]$ , où  $a, b \in \Sigma$ .

• Si  $a = b$ , une PLSSC de  $ua$  et  $vb$  se termine par  $a$ !

Autrement, on pourrait la prolonger par  $a$ , absurde par maximalité de la longueur.

Dans  $S[i, j] = S[i-1, j-1] + 1$ .

• Si  $a \neq b$ , et  $ua$  et  $vb$  possèdent une PLSSC ne se terminant pas par  $a$ ,

on a  $S[i, j] = S[i-1, j]$  (\*)

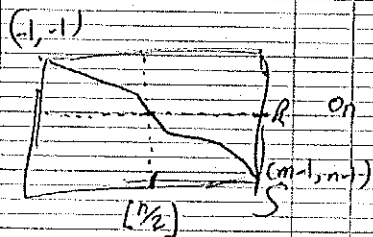
Par symétrie de  $a$  et  $b$ , on a  $S[i, j] = \max(S[i-1, j], S[i, j-1])$

En reliant  $(i, j)$  à  $(i-1, j-1)$  si  $x[i] = y[j]$ ,

à  $(i-1, j)$  si on est dans le cas (\*) du lemme

à  $(i, j-1)$  sinon,

on obtient une ligne reliant  $(-1, -1)$  à  $(m-1, n-1)$  dans  $S$ .



On va chercher la coordonnée  $k$  par laquelle passe la ligne dans la colonne  $\lfloor n/2 \rfloor$ ,  
 pour pouvoir appliquer un algorithme SSC, renvoyant une PLSSC, avec  
 une méthode "divise pour régner".

$k$  est l'indice maximisant 
$$\text{smc}(x[0..k-1], y[0.. \lfloor n/2 \rfloor - 1]) + \text{smc}(x[k..m-1], y[\lfloor n/2 \rfloor .. n-1])$$

Pour le trouver, on calcule d'abord la colonne d'indice  $\lfloor n/2 \rfloor - 1$  de  $S$   
 par  $SSC\_LONGUEUR(x, y[0.. \lfloor n/2 \rfloor - 1])$ .

Puis on termine la construction de  $S$  en mémorisant des pointeurs vers la  
 colonne du milieu dans  $P_1$  et  $P_2$ .

Ces deux tables implémentent la table P vérifiant

$$\left( \begin{array}{l} \forall i \in \{-1, 0, \dots, m-1\}, \forall j \in \{\lfloor \frac{n}{2} \rfloor - 1, \dots, n-1\}, \\ P[i, j] = k \in \{0, \dots, i+1\} \text{ si} \\ \text{smc}(x[0..i], y[0..j]) = \text{smc}(x[0..k-1], y[0..(\frac{n}{2}-1)]) \\ \quad + \text{smc}(x[k..i], y[(\frac{n}{2})..j]) \end{array} \right)$$

Donc  $P[m-1, n-1]$  sera le  $k$  cherché.

lemme 2

$$\left. \begin{array}{l} \bullet \forall i \in \{-1, \dots, m-1\}, \quad P[i, \lfloor \frac{n}{2} \rfloor - 1] = i+1 \\ \bullet \forall j \geq \lfloor \frac{n}{2} \rfloor, \quad P[-1, j] = 0 \\ \bullet \text{ si } i \geq 0 \text{ et } j \geq \lfloor \frac{n}{2} \rfloor, \quad P[i, j] = \begin{cases} P[i-1, j-1] & \text{si } x[i] = y[j] \\ P[i-1, j] & \text{si } x[i] \neq y[j] \text{ et } \\ & S[i-1, j] > S[i, j-1] \\ P[i, j-1] & \text{sinon} \end{cases} \end{array} \right\}$$

Preuve: on se ramène à la définition de P en faisant une récurrence sur  $(i, j)$ , par ordre lexicographique.

- si  $j = \lfloor \frac{n}{2} \rfloor - 1$ ,  $y[\lfloor \frac{n}{2} \rfloor - j] = \epsilon$ , donc  $k = i+1$ .
- si  $i = -1$ , les facteurs de  $x$  sont vides donc  $k = 0$ .
- si  $i \geq 0$  et  $j \geq \lfloor \frac{n}{2} \rfloor$ , supposons  $x[i] = y[j]$  (les autres cas sont similaires)

(lemme 1)  $\text{smc}(x[0..i], y[0..j]) = \text{smc}(x[0..i-1], y[0..j-1]) + 1$

(cas base)  $\text{smc}(x[0..i-1], y[0..j-1]) = \text{smc}(x[0..k-1], y[0..(\frac{n}{2}-1)]) + \text{smc}(x[k..i-1], y[(\frac{n}{2})..j-1])$  où  $k = P[i-1, j-1]$

(lemme 1)  $\text{smc}(x[k..i], y[(\frac{n}{2})..j]) = \text{smc}(x[k..i-1], y[(\frac{n}{2})..j-1]) + 1$

D'où  $P[i, j] = k = P[i-1, j-1]$

Enfin l'algorithme SSC.

La preuve de correction découle essentiellement des deux lemmes 2, et de la justification de SSC\_LONGUEUR.

SSC(x, y):

```

Si m = 1 et x[0] apparait dans y : Renvoie x[0]
Sinon, Si n = 1 et y[0] apparait dans x : Renvoie y[0]
Sinon, Si m = 0 ou m = 1 ou n = 0 ou n = 1 : Renvoie ε
Soyons
C1 ← SSC_LONGUEUR(x, y[0..⌊n/2⌋-1])
Pour i ∈ {1, ..., m-1} : P1[i] ← i+1
Pour j ∈ {⌊n/2⌋, ..., n-1} :
  C2[i-1], P2[i-1] ← 0
  Pour r ∈ {0, ..., m-1} :
    Si x[r] = y[j] :
      C2[r, j] ← C1[r-1] + 1
      P2[r, j] ← P1[r-1]
    Sinon, Si C1[r] > C2[r-1] :
      C2[r, j] ← C1[r]
      P2[r, j] ← P1[r]
    Sinon
      C2[r, j] ← C1[r-1]
      P2[r, j] ← P2[r-1]
  C1 ← C2, P ← P2
R ← P[m-1]
u ← SSC(x[0..R-1], y[0..⌊n/2⌋-1])
v ← SSC(x[R..m-1], y[⌊n/2⌋..n-1])
Retourner u.v
  
```

Par les complexités : SSC s'exécute en temps  $O(mn)$  dans un espace  $O(m)$ .

Preuve : Si  $\mathcal{C}(m, n)$  désigne le nombre la complexité de SSC (pour  $|x|=m, |y|=n$ ), alors  $\mathcal{C}(m, n) = O(mn) + \mathcal{C}(k, \lfloor n/2 \rfloor) + \mathcal{C}(m-k, n - \lfloor n/2 \rfloor)$ .

Ainsi, comme  $O(k \lfloor n/2 \rfloor) + O((m-k)(n - \lfloor n/2 \rfloor)) = O(\frac{mn}{2})$  et que  $\sum_i \frac{mn}{2^i} \leq 2mn$ , on a  $\mathcal{C}(m, n) = O(mn)$ .

• du mémoire nécessaire consiste en la taille de  $C_1, C_2, P_1$  et  $P_2$ , le tout en  $O(m)$ . En effet, comme les appels récursifs n'utilisent pas les informations stockées dans ces tables, on peut réutiliser le même espace dans la suite du calcul.

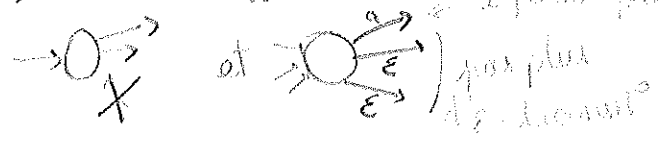
# Questions :

1) Recherche de motif = est-ce qu'on peut détecter l'occurrence d'une expression rationnelle ?

↳ OVP → automate

Pb: pour qu'il soit déterministe  $\forall$  AIS :  $2^n$  états → BoF

↳ Beauchquier = méthode - évite en la taille de l'ape, mais avec des termes ptes. (→ automate normalisé = 2 parts par état)



2) Plus courte sur séquence commune → un mot le plus court possible, qui inclure les deux mots de départ  $(x_1 \dots x_n, y_1 \dots y_m)$

ATGCTA  
ACGA

PLSSC → AGA.  
 $z_1 \dots z_k$   
ATGCTA  
ACGA

→ ATCGCTA  
(ou ACTGCTA)

## Algo:

$w \leftarrow \epsilon$   
tant que  $l < |z|$   
tant que  $x_i \neq z_l$  :  
     $w = w + x_i$   
     $i = i + 1$   
tant que  $y_j \neq z_l$  :  
     $w = w + y_j$   
     $j = j + 1$

$w = w + z_l$   
 $l++$ ;  $i++$ ;  $j++$

$w = w + x[i \dots n] + y[j \dots m]$   
retourne  $w$

← DVPT: PLSSC sans optimisation + cet application

DVPT: automate minimal des occurrences et

DVPT: Rabin Karp

## Distance de Levenshtein sur les mots

Référence : M. CROCHEMORE, W. RYTTER  
*Text Algorithms*

2015-2016

Étant donnés deux mots  $X = x_1 \cdots x_m$  et  $Y = y_1 \cdots y_n$  sur un alphabet  $\Sigma$ , on cherche à savoir quelle est la distance de Levenshtein entre  $X$  et  $Y$ , à savoir le coût minimal pour aller de  $X$  à  $Y$  en effectuant les opérations élémentaires suivantes :

- suppression d'un caractère de  $X$ ,
- ajout d'un caractère de  $X$  dans  $Y$ ,
- substitution d'un caractère de  $X$  dans  $Y$ .

On note  $m$  la longueur de  $X$ ,  $n$  celle de  $Y$  et  $X[i]$  le  $i$ -ème caractère de  $X$ . On suppose  $n \geq m$ .

Considérons  $G$ , un graphe composé de  $m \times n$  nœuds notés  $(i, j)$  (pour  $i \in [0, m]$  et  $j \in [0, n]$ ). Chaque nœud  $(i, j)$  est relié aux nœuds  $(i + 1, j)$  et  $(i, j + 1)$  par une arête de poids 1 et au nœud  $(i + 1, j + 1)$  par une arête de poids  $\delta(X[i + 1], Y[j + 1])$ . La distance de Levenshtein entre  $X$  et  $Y$  est égale au poids du plus court chemin dans  $G$  reliant les sommets  $(0, 0)$  et  $(m, n)$ .

Soit  $EDIT$  un tableau de taille  $m \times n$  défini par :

- $\forall i \in [0, m], EDIT[i, 0] = i$ ,
- $\forall j \in [0, n], EDIT[0, j] = j$ ,
- $\forall i \in [1, m], j \in [1, n], EDIT[i, j] = \min(EDIT[i - 1, j] + 1, EDIT[i, j - 1] + 1, EDIT[i - 1, j - 1] + \delta(X[i], Y[j]))$ .

La distance de Levenshtein entre  $X$  et  $Y$  est alors égale à  $EDIT(m, n)$ , la valeur  $EDIT(i, j)$  représentant la distance minimale entre  $(0, 0)$  et  $(i, j)$  dans  $G$ .

**Theorem 1** *La procédure EDIT a une complexité en temps de  $\mathcal{O}(m \times n)$  et en mémoire de  $\mathcal{O}(m)$ .*

**Proof**

Dans EDIT, on a deux boucles "pour" imbriquées, et dans chacune de ces boucles des opérations en  $\mathcal{O}(1)$ . D'où la complexité en  $\mathcal{O}(m \times n)$ .

On peut toutefois se contenter de garder en mémoire les informations sur la colonne (ou la ligne) actuelle ainsi que la précédente à tout moment du calcul, réduisant ainsi la complexité en mémoire à du  $\mathcal{O}(m)$ . Ceci dit, garder tout le tableau en mémoire nous permet de reconstruire les opérations pour passer de  $X$  à  $Y$ .  $\square$



---

**Algorithm 1** EDIT

---

```
m ← X
n ← Y
for i = 1 to m do
  EDIT[i, 0] ← i
end for
for j = 0 to n do
  c[0, j] ← j
end for
for i = 1 to m do
  for j = 1 to n do
    if  $x_i = y_j$  then
      delta ← 0
    else
      delta ← 1
    end if
    EDIT[i, j] ← min(EDIT[i - 1, j] + 1, EDIT[i, j - 1] + 1, EDIT[i - 1, j - 1] + delta)
  end for
end for
return EDIT[m, n]
```

---

**Example 1** On cherche à comparer les chaînes NICHE et CHIENS. Le tableau construit par EDIT est :

		<i>C</i>	<i>H</i>	<i>I</i>	<i>E</i>	<i>N</i>	<i>S</i>
	<b>0</b>	1	2	3	4	5	6
<i>N</i>	<b>1</b>	1	2	3	4	4	5
<i>I</i>	<b>2</b>	2	2	2	3	4	5
<i>C</i>	<b>3</b>	<b>2</b>	3	3	3	4	5
<i>H</i>	<b>4</b>	3	<b>2</b>	<b>3</b>	4	4	5
<i>E</i>	<b>5</b>	4	3	3	<b>3</b>	4	<b>5</b>

La distance entre NICHE et CHIENS est de 5 et une façon de passer de l'un à l'autre est de supprimer les lettres N et I de NICHE pour obtenir CHE puis de rajouter les lettres I, N et S aux endroits adéquats.