

CALCUL DE LA DISTANCE PAR LES SOUS-MOTS (*)

par Jean-Jacques HEBRARD ⁽¹⁾ et Maxime CROCHEMORE ⁽²⁾

Communiqué par J. BERSTEL

Résumé. – Cet article contient deux méthodes de calcul de la longueur du plus court sous-mot (au sens de sous-suite) permettant de distinguer deux mots différents u et v . L'utilisation d'automates et des structures de données concernant les questions d'« union et recherche » conduit à un algorithme quasi linéaire en la longueur de uv .

Abstract. – This paper gives two methods to compute the shortest subsequence which distinguishes two different words u and v . The use of automata together with data structures for "Union-Find" questions leads to an algorithm almost linear in the length of uv .

Définir des critères de comparaison entre chaînes de caractères (ou mots) est un problème qui apparaît fréquemment, dans des domaines aussi divers que la biochimie, la reconnaissance de la parole, l'informatique (correction de mots mal orthographiés, recherche d'informations non entièrement spécifiées), la reconnaissance de contours, etc. [SK 83].

Parmi les critères les plus couramment utilisés on trouve la distance de Hamming et la distance d'édition [Mo 70, Se 74, WF 74, MP 80]. Ces deux distances sont fondamentalement différentes : la première n'est qu'une stricte comparaison, caractère à caractère, des deux mots, alors que la seconde prend en compte de façon essentielle, les similitudes entre leurs sous-mots (un sous-mot d'un mot est obtenu en supprimant dans ce dernier un certain nombre de caractères).

Il est souvent possible de ramener la comparaison de deux mots u et v à celle de leurs ensembles respectifs de sous-mots. Les deux indices de proximité les plus simples sont alors : la longueur maximale $L(u, v)$ d'un sous-mot

(*) Reçu juillet 1985, révisé décembre 1985.

⁽¹⁾ Laboratoire d'Informatique, Univ. Caen, 14032 Caen Cedex.

⁽²⁾ L.I.T.P.-Rouen et C.S.P., Univ. Paris-Nord, avenue J.-B.-Clément, 93400 Villetaneuse.

commun, et la « distance par les sous-mots » ($d(u, v)$) définie à l'aide de la longueur minimale des mots qui distinguent u et v — c'est-à-dire, des mots qui ne sont pas des sous-mots communs à u et v .

Le calcul de $L(u, v)$ est traité dans de nombreuses études; citons entre autres celle de Hirschberg [Hi 77] et Nakatsu, Kambayashi et Yajima [NKY 82] dont les méthodes sont de complexité maximale quadratique. L'algorithme de Hunt et Szymanski [HS 77], bien que de complexité maximale $O((|u|+|v|)^2 \log(|u|+|v|))$, est d'une grande efficacité dans certaines situations pratiques; il est d'ailleurs utilisé dans la mise en œuvre de la commande DIFF de UNIX.

Les aspects théoriques de l'étude de la « distance par les sous-mots » sont développés par I. Simon dans [Lo 82]. Il en déduit un algorithme linéaire de calcul de $d(u, v)$ qui est malheureusement très complexe et malaisé à mettre en œuvre [Si 84]. Une analyse directe du problème est réalisée dans [He 84]; elle conduit à un algorithme linéaire pour un alphabet à deux lettres.

Nous adoptons une autre approche, fondée essentiellement sur le fait que l'on peut construire en temps linéaire, pour un mot donné u , un automate $\mathcal{S}(u)$ qui reconnaît l'ensemble des sous-mots de u . Le calcul de $L(u, v)$ se réduit alors à la recherche d'un chemin de longueur maximale dans l'automate produit $\mathcal{S}(u) \times \mathcal{S}(v)$; cet algorithme est de complexité quadratique [He 84]. Nous présentons ici les deux algorithmes que l'on obtient en utilisant l'automate des sous-mots pour le calcul de $d(u, v)$. Le premier est l'application à l'automate union $\mathcal{S}(u) \cup \mathcal{S}(v)$ d'une méthode de partitionnement; sa complexité en temps est $O((|u|+|v|) \log(|u|+|v|))$. Le second consiste à ramener le calcul à un problème d'équivalence d'automates; la mise en œuvre est immédiate et très efficace : sa complexité est quasi linéaire en temps et linéaire en espace.

1. AUTOMATES ET SOUS-MOTS

1.1. Définitions générales

Soit (A, \leq) un alphabet fini totalement ordonné. A^* est l'ensemble des mots formés sur A . On note ε le mot vide et $|w|$ la longueur d'un mot w .

On considère l'ordre généalogique sur A^* , noté \leq également, et défini de la façon suivante : pour tout $u, v \in A^*$, $u \leq v$ si et seulement si :

ou bien $|u| < |v|$;

ou bien $|u| = |v|$, $u = ras$, $v = rbt$, avec $r, s, t \in A^*$, $a, b \in A$ et $a \leq b$;

ou bien $u = v$.

Si $u \leq v$, alors soit v est strictement plus long que u , soit u et v sont de même longueur et v est plus grand que u pour l'ordre lexicographique.

Soient $u = u_1 \dots u_m$ ($u_i \in A$) et $v = v_1 \dots v_n$ ($v_j \in A$) deux mots de A^* . On dit que v est un sous-mot de u s'il existe une application $f: \{1, \dots, n\} \rightarrow \{1, \dots, m\}$ strictement croissante telle que $v_1 \dots v_n = u_{f(1)} \dots u_{f(n)}$.

Exemple :

$$A = \{a, b, c\}, \quad u = abcabc, \quad v = bcc.$$

REMARQUE : Un mot peut être sous-mot d'un autre de plusieurs manières. Ainsi abc est un sous-mot de $abcabc$ de quatre façons différentes.

Si v est un sous-mot de u , on dit aussi que v divise u et on note $v \mid u$. Si v ne divise pas u , on note $v \nmid u$.

Un automate fini déterministe complet est un quadruplet $\mathcal{A} = (Q, i, T, t)$, où Q est l'ensemble fini des états, i l'état initial, T l'ensemble des états terminaux et t l'application de $Q \times A$ dans Q qui définit les transitions de \mathcal{A} ; t se prolonge en une application de $Q \times A^*$ dans Q . On emploie la notation abrégée : $\forall q \in Q, \forall w \in A^*, q \cdot w = t(q, w)$.

Le langage reconnu par l'automate \mathcal{A} est l'ensemble $|\mathcal{A}| = \{w \in A^* / i \cdot w \in T\}$.

1.2. Distance par les sous-mots

Étant donné un mot u et un entier l , on note $S(u, l)$ l'ensemble des sous-mots w de u tels que $|w| \leq l$.

$$S(u, l) = \{w \in A^* / w \mid u \text{ et } |w| \leq l\}.$$

On appelle distance par les sous-mots de deux mots u et v , l'élément de $N \cup \{\infty\}$, noté $d(u, v)$, et défini par :

$$d(u, v) = \max \{l \in N / S(u, l) = S(v, l)\} \quad \text{si } u \neq v,$$

$$d(u, v) = \infty \text{ sinon.}$$

Évidemment, $d(u, v)$ n'est pas une distance au sens usuel, mais $\delta(u, v) = 2^{-d(u, v)}$ est une distance ultramétrique. On trouve une étude détaillée de $d(u, v)$ dans [Lo 82].

On dit qu'un mot w distingue deux mots u et v , s'il divise l'un et ne divise pas l'autre. Si w est un mot de longueur minimale qui distingue u et v alors $d(u, v) = |w| - 1$.

Exemple :

$$\begin{aligned} u &= ababa, & v &= aabba, & S(u, 2) &= S(v, 2), \\ & & & & bab \mid u, & bab \mid v, & d(u, v) &= 2. \end{aligned}$$

1.3. Automate des sous-mots

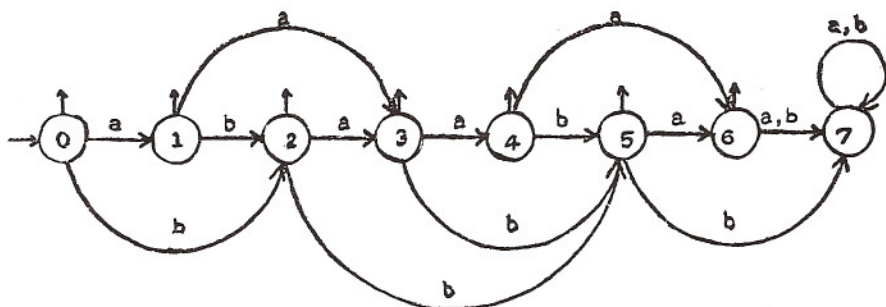
On peut construire, en temps linéaire, l'automate minimal qui reconnaît l'ensemble des sous-mots d'un mot donné.

Étant donné un mot $u = u_1 \dots u_m$ ($u_i \in A$), on note $\mathcal{S}(u)$ l'automate fini déterministe défini par :

$$\begin{aligned} \mathcal{S}(u) &= \{\{0, \dots, m+1\}, 0, \{0, \dots, m\}, t\}, \\ \forall a \in A, \quad t(m, a) &= t(m+1, a) = m+1, \\ \forall a \in A, \quad \forall i \in \{0, \dots, m-1\}, \quad E_a(i) &= \{k \in \{i+1, \dots, m\} / u_k = a\}, \\ t(i, a) &= \begin{cases} \min E_a(i) & \text{si } E_a(i) \neq \emptyset \\ m+1, & \text{sinon.} \end{cases} \end{aligned}$$

Exemple :

$$A = \{a, b\}, \quad u = abaaba, \quad m = 6.$$



On a les propriétés suivantes :

- (1) $\forall i, j \in \{0, \dots, m+1\}, \forall w \in A^*, \quad i \leq j \Rightarrow t(i, w) \leq t(j, w).$
- (2) $\forall i, j \in \{0, \dots, m+1\}, \quad i < j \Rightarrow t(i, u_j) \leq j.$

PROPOSITION 1 : Pour tout $u = u_1 \dots u_m$ ($u_i \in A$), $\mathcal{S}(u)$ reconnaît l'ensemble des sous-mots de u .

Démonstration : Pour tout mot $w = w_1 \dots w_p$ ($w_i \in A$) et tout $i \in \{1, \dots, p\}$, soit $g(i) = t(0, w_1 \dots w_i)$. Il suffit de montrer : $g(p) \leq m \Leftrightarrow w \mid u$. Si $g(p) \leq m$,

alors g est une application strictement croissante de $\{1, \dots, m\}$ qui fait de w un sous-mot de u .

Si $w|u$, alors $w = u_{f(1)} \dots u_{f(p)}$, où f est une application strictement croissante de $\{1, \dots, p\}$ dans $\{1, \dots, m\}$. Par définition de $\mathcal{S}(u)$, $g(1) \leq f(1)$ et pour tout i de $\{2, \dots, p\}$,

$$\begin{aligned} g(i-1) \leq f(i-1) &\Rightarrow g(i-1) < f(i) \\ &\Rightarrow t(g(i-1), u_{f(i)}) \leq f(i) \Rightarrow g(i) \leq f(i) \\ &\text{et donc } g(p) \leq f(p) \leq m. \quad \blacksquare \end{aligned}$$

$\mathcal{S}(u)$ est appelé l'automate des sous-mots de u . Sa construction est réalisée par l'algorithme suivant, en un simple parcours de u de droite à gauche.

ALGORITHME 2 :

Entrée : $u = u_1 \dots u_m (u_i \in A)$

Sortie : t , fonction de transition de $\mathcal{S}(u)$

début

pour chaque lettre a de A faire début $t(m+1, a) \leftarrow m+1$; $t(m, a) \leftarrow m+1$ fin;

pour $i \leftarrow m-1$ à 0 pas -1 faire

début

$t(i, u_{i+1}) \leftarrow i+1$;

pour chaque lettre a de $A \setminus \{u_{i+1}\}$ faire $t(i, a) \leftarrow t(i+1, a)$;

fin

fin.

PROPOSITION 3 : Pour tout mot u , l'algorithme 2 calcule la fonction de transition de $\mathcal{S}(u)$; sa complexité est $O(|A| |u|)$.

Démonstration : Ceci est une conséquence de la définition de $\mathcal{S}(u)$. \blacksquare

2. PARTITIONNEMENT

Le calcul de $d(u, v)$ se ramène à la construction d'une suite d'équivalences $R(k)$ ($k \in \mathbb{N}$) sur l'ensemble des états de l'automate union $\mathcal{S}(u) \cup \mathcal{S}(v)$. Pour tout k , $R(k+1)$ est plus fine que $R(k)$. On montre que $d(u, v) = D$ si et seulement si D est le plus petit entier tel que les états initiaux de $\mathcal{S}(u)$ et $\mathcal{S}(v)$ ne sont pas équivalents pour $R(D+1)$. La construction de la suite R est réalisée en adaptant l'algorithme de partitionnement de Hopcroft [AHU 74], le calcul de $d(u, v)$ est alors effectué en temps $O(|A|(|u| + |v|) \log(|u| + |v|))$.

Considérons les mots $u = u_1 \dots u_m (u_i \in A)$, $v = v_1 \dots v_n (v_j \in A)$ et les automates $\mathcal{S}(u) = (\{0, \dots, m+1\}, 0, \{0, \dots, m\}, t_u)$ et $\mathcal{S}(v) = (\{\bar{0}, \dots, \bar{n}+1\}, \bar{0}$,

$\{\bar{0}, \dots, \bar{n}\}, t_v)$. Soit $\mathcal{S}(u) \cup \mathcal{S}(v) = (Q, I, T, t)$ l'automate union. On convient de confondre $m+1$ et $\bar{n+1}$ en un seul état noté « Puits ».

$$Q = \{0, \dots, m\} \cup \{\bar{0}, \dots, \bar{n}\} \cup \{\text{Puits}\},$$

$$I = \{0, \bar{0}\}, \quad T = Q \setminus \{\text{Puits}\}.$$

Cet automate n'est pas déterministe.

Pour tout entier k , $R(k)$ est la relation d'équivalence définie sur Q par :

$$\forall p, q \in Q, \quad p \equiv q [R(k)]$$

$$\Leftrightarrow (\forall w \in A^*, |w| \leq k \Rightarrow (p.w \in T \Leftrightarrow q.w \in T)).$$

On note $P(k)$ la partition associée à $R(k)$. Ainsi $P(0) = \{\{\text{Puits}\}, Q \setminus \{\text{Puits}\}\}$.

La proposition suivante, dont la preuve est une simple transcription des définitions de d et de la suite d'équivalences R , fournit un moyen de calcul de d . Elle exprime uniquement que la distance $d(u, v)$ entre deux mots u et v est le plus grand entier D pour lequel les états initiaux de $\mathcal{S}(u)$ et $\mathcal{S}(v)$ sont $R(D)$ -équivalents.

PROPOSITION 4 : $\forall D \in N, \quad d(u, v) = D \Leftrightarrow (\forall k \leq D, 0 \equiv \bar{0} [R(k)] \text{ et } 0 \not\equiv \bar{0} [R(D+1)]).$

Démonstration :

$\forall k \in N, 0 \equiv \bar{0} [R(k)] \Leftrightarrow (\forall w \in A^*, |w| \leq k \Rightarrow (w|u \Leftrightarrow w|v)).$ ■

REMARQUES : (1) Pour tout k , $R(k+1)$ est plus fine que $R(k)$ (on note $R(k) \geq R(k+1)$).

(2) Pour tout k , $R(k) \neq R(k+1) \Leftrightarrow (\exists p, q \in Q, \exists a \in A, p \equiv q [R(k)] \text{ et } p.a \not\equiv q.a [R(k+1)]).$

(3) Il existe un entier K tel que $R(0) > \dots > R(K)$ et $\forall k > K, R(k) = R(K)$.

(où $R(k) > R(k') \Leftrightarrow (R(k) \geq R(k') \text{ et } R(k) \neq R(k'))$)

Q étant fini, l'ensemble $E = \{k \in N / R(k) = R(k) = R(k+1)\}$ est non vide, et $K = \min E$.

(4) $R(K)$ est l'équivalence la plus grossière qui soit à la fois plus fine que $R(0)$ et compatible avec t :

$$\forall p, q \in Q, \quad \forall a \in A, \quad p \equiv q [R(K)] \Rightarrow p.a \equiv q.a [R(K)].$$

Le calcul de $R(K)$ est un problème de partitionnement qui peut être résolu en utilisant l'algorithme de Hopcroft [AHU 74]. Cependant celui-ci ne permet

pas de construire effectivement la suite R . Aussi emploierons-nous une méthode dérivée, définie dans [CC 82] pour le partitionnement d'un graphe.

NOTATIONS :

$$- \forall q \in Q, \forall a \in A, qa^{-1} = \{p \in Q / p.a = q\}.$$

$$- \forall B \subset Q, \forall a \in A, Ba^{-1} = \bigcup_{q \in B} qa^{-1}.$$

- $P = \{B_1, \dots, B_h\}$ étant une partition de Q , on note $[B_1, \dots, B_h]$ l'équivalence associée.

- L'intersection de deux équivalences est une opération associative; soient R_1 et R_2 deux équivalences sur Q , on a :

$$\forall p, q \in Q, p \equiv q[R_1 \cap R_2] \Leftrightarrow (p \equiv q[R_1] \text{ et } p \equiv q[R_2]).$$

Pour toute partition P de Q et toute partie T de $\mathcal{P}(Q)$, on note $P \wedge T$ la partition associée à l'équivalence R définie par :

$$\forall p, q \in Q, p \equiv q[R] \Leftrightarrow (\forall B \in P, \forall C \in T \cup \{Q\}, p \in B \cap C \Leftrightarrow q \in B \cap C).$$

Une première construction de la suite $R(k)$ ($k \in \mathbb{N}$) est donnée par la :

$$\text{PROPOSITION 5 : Pour tout } k, R(k+1) = \bigcap_{\substack{B \in P(k) \\ a \in A}} [Ba^{-1}, Q \setminus Ba^{-1}].$$

Démonstration : Soit

$$E = \bigcap_{\substack{B \in P(k) \\ a \in A}} [Ba^{-1}, Q \setminus Ba^{-1}].$$

$$\begin{aligned} \forall p, q \in Q, p \equiv q[E] &\Leftrightarrow \forall a \in A, p.a \equiv q.a[R(k)] \\ &\Leftrightarrow p \equiv q[R(k+1)]. \quad \blacksquare \end{aligned}$$

On améliore cette méthode en distinguant les petites et les grandes classes. Soient

$$k > 0, P(k) = \{B_1, \dots, B_h\} \quad \text{et} \quad P(k-1) = \{B'_1, \dots, B'_{h'}\};$$

toute classe B'_i suivant $R(k-1)$ est union de classes suivant $R(k)$. On définit une application g qui permet d'associer à toute classe B'_i de $P(k-1)$ une classe de $P(k)$ qu'elle contient et qui est de cardinal maximal,

$$g : P(k-1) \rightarrow P(k), \quad \forall i \in \{1, \dots, h'\}, \quad g(B'_i) \subset B_i$$

et

$$\forall B \in P(k), B \subset B'_i \Rightarrow |B| \leq |g(B'_i)|.$$

Les classes qui sont dans l'image de g sont appelées les grandes classes, les autres les petites classes. On note $Pt(k)$ l'ensemble des petites classes de $P(k)$. On a $Pt(0) = \{\{\text{Puits}\}\}$.

La proposition suivante montre qu'il n'est pas nécessaire d'utiliser toutes les classes de $R(k)$ pour calculer $R(k+1)$, les petites classes suffisent.

PROPOSITION 6 : Pour tout k ,

$$R(k+1) = R(k) \cap \left(\bigcap_{\substack{B \in Pt(k) \\ a \in A}} [Ba^{-1}, Q \setminus Ba^{-1}] \right).$$

$$\text{Démonstration : } E' = R(k) \cap \left(\bigcap_{\substack{B \in Pt(k) \\ a \in A}} [Ba^{-1}, Q \setminus Ba^{-1}] \right).$$

Si $k=0$, $Pt(k) = \{\{\text{Puits}\}\}$; $\forall p, q \in Q, p \equiv q[E'] \Leftrightarrow \forall a \in A, p.a \equiv q.a[R(0)]$.

Si $k > 0$, $R(k) = \bigcap_{\substack{C \in P(k-1) \\ a \in A}} [Ca^{-1}, Q \setminus Ca^{-1}]$ (prop. 5) :

$$P(k-1) \wedge Pt(k) = P(k),$$

donc

$$E' = \bigcap_{\substack{B \in P(k) \\ a \in A}} [Ba^{-1}, Q \setminus Ba^{-1}] = R(k+1). \quad \blacksquare$$

REMARQUE : Pour tout k , on a : $Pt(k) = \emptyset \Leftrightarrow R(k) = R(k-1)$.

L'algorithme de calcul de $d(u, v)$ est décrit par le schéma suivant :

ALGORITHME 7. — Fonction $d(u, v)$

```

début
   $k \leftarrow 0$ ;  $R(k) \leftarrow \{\{\text{Puits}\}, Q \setminus \{\{\text{Puits}\}\}\}$ ;  $Pt(k) \leftarrow \{\{\{\text{Puits}\}\}\}$ ;
  tant que  $0 \equiv \bar{0}[R(k)]$  et  $Pt(k) \not\emptyset$  faire
    début  $R(k+1) \leftarrow R(k) \cap \left( \bigcap_{\substack{B \in Pt(k) \\ a \in A}} [Ba^{-1}, Q \setminus Ba^{-1}] \right)$ ;
     $k \leftarrow k+1$ 
fin;
si  $0 \not\equiv \bar{0}[R(k)]$  alors retour  $(k-1)$ 
  sinon retour  $(\infty)$ 
fin
```

PROPOSITION 8 : L'algorithme 7 effectue le calcul de $d(u, v)$ en temps $O(|A|(|u|+|v|)\log(|u|+|v|))$.

Démonstration : On peut mettre en œuvre la construction de la suite $R(k)$ ($k \in \mathbb{N}$) de façon que le temps d'exécution soit proportionnel au nombre d'opérations qa^{-1} effectuées ($q \in Q$, $a \in A$) [CC 82]. Soit α ce nombre pour q et a fixés. L'opération qa^{-1} est effectuée à chaque fois que l'état q se trouve dans une petite classe. Ceci se produit au plus $\log_2 |Q|$ fois — en effet, pour toute petite classe B de $R(k)$ contenue dans une classe B' de $R(k-1)$ on a : $|B| \leq |B'|/2$. Par conséquent $\alpha \leq \log_2 |Q|$, d'où le résultat. ■

Exemple :

$$A = \{a, b, c\}, \quad u = cabacb, \quad v = bacabc.$$

$$\mathcal{S}(u) : \textcircled{0} \quad c \quad \textcircled{1} \quad a \quad \textcircled{2} \quad b \quad \textcircled{3} \quad a \quad \textcircled{4} \quad c \quad \textcircled{5} \quad b \quad \textcircled{6}$$

Puits

$$\mathcal{S}(v) : \textcircled{0} \quad b \quad \textcircled{1} \quad a \quad \textcircled{2} \quad c \quad \textcircled{3} \quad a \quad \textcircled{4} \quad b \quad \textcircled{5} \quad c \quad \textcircled{6}$$

$$P(0) = \{ \{0, 1, 2, 3, 4, 5, 6, \bar{0}, \bar{1}, \bar{2}, \bar{3}, \bar{4}, \bar{5}, \bar{6}\}, \{\text{Puits}\} \},$$

$$Pt(0) = \{ \{\text{Puits}\} \},$$

$$\{\text{Puits}\} a^{-1} = \{4, 5, 6, \bar{4}, \bar{5}, \bar{6}, \text{Puits}\}, \quad \{\text{Puits}\} b^{-1} = \{6, \bar{5}, \bar{6}, \text{Puits}\},$$

$$\{\text{Puits}\} c^{-1} = \{5, 6, \bar{6}, \text{Puits}\}.$$

$$P(1) = \{ \{0, 1, 2, 3, \bar{0}, \bar{1}, \bar{2}, \bar{3}\}, \{4, \bar{4}\}, \{5\}, \{\bar{5}\}, \{6, \bar{6}\}, \{\text{Puits}\} \}$$

$$Pt(1) = \{ \{4, \bar{4}\}, \{5\}, \{\bar{5}\}, \{6, \bar{6}\} \}$$

$$\{4, \bar{4}\} a^{-1} = \{2, 3, \bar{2}, \bar{3}\}, \quad \{4, \bar{4}\} b^{-1} = \emptyset, \quad \{4, \bar{4}\} c^{-1} = \emptyset,$$

$$\{5\} a^{-1} = \emptyset, \quad \{5\} b^{-1} = \emptyset, \quad \{5\} c^{-1} = \{1, 2, 3, 4\},$$

$$\{\bar{5}\} a^{-1} = \emptyset, \quad \{\bar{5}\} b^{-1} = \{\bar{1}, \bar{2}, \bar{3}, \bar{4}\}, \quad \{\bar{5}\} c^{-1} = \emptyset,$$

$$\{6, \bar{6}\} a^{-1} = \emptyset, \quad \{6, \bar{6}\} b^{-1} = \{3, 4, 5\}, \quad \{6, \bar{6}\} c^{-1} = \{\bar{3}, \bar{4}, \bar{5}\}.$$

$$P(2) = \{ \{0, \bar{0}\}, \{1\}, \{\bar{1}\}, \{2\}, \{\bar{2}\}, \{3\}, \{\bar{3}\},$$

$$\{4\}, \{\bar{4}\}, \{5\}, \{\bar{5}\}, \{6, \bar{6}\}, \{\text{Puits}\} \}$$

$$Pt(2) = \{ \{1\}, \{\bar{1}\}, \{2\}, \{\bar{2}\}, \{3\}, \{\bar{3}\}, \{4\} \}$$

$$\{1\} a^{-1} = \emptyset, \quad \{1\} b^{-1} = \emptyset, \quad \{1\} c^{-1} = \{0\}$$

$$\{\bar{1}\} a^{-1} = \emptyset, \quad \{\bar{1}\} b^{-1} = \{\bar{0}\}, \quad \{\bar{1}\} c^{-1} = \emptyset$$

$$\{2\} a^{-1} = \{0, 1\}, \quad \{2\} b^{-1} = \emptyset, \quad \{2\} c^{-1} = \emptyset$$

$$\{\bar{2}\} a^{-1} = \{\bar{0}, \bar{1}\}, \quad \{\bar{2}\} b^{-1} = \emptyset, \quad \{\bar{2}\} c^{-1} = \emptyset,$$

$$\{3\} a^{-1} = \emptyset, \quad \{3\} b^{-1} = \{0, 1, 2\}, \quad \{3\} c^{-1} = \emptyset$$

$$\{\bar{3}\} a^{-1} = \emptyset, \quad \{\bar{3}\} b^{-1} = \emptyset, \quad \{\bar{3}\} c^{-1} = \{\bar{0}, \bar{1}, \bar{2}\}$$

$$\{4\} a^{-1} = \{2, 3\}, \quad \{4\} b^{-1} = \emptyset, \quad \{4\} c^{-1} = \emptyset.$$

$$P(3) = \{\{0\}, \{\bar{0}\}, \{1\}, \{\bar{1}\}, \{2\}, \{\bar{2}\}, \{3\}, \{\bar{3}\},$$

$$\{4\}, \{\bar{4}\}, \{5\}, \{\bar{5}\}, \{\bar{6}\}, \{\text{Puits}\}$$

$$0 \neq \bar{0} [R(3)] \text{ et } d(u, v) = 2.$$

3. UN ALGORITHME QUASI LINEAIRE

Le calcul de la distance par les sous-mots est un problème voisin de celui de l'équivalence de deux automates [AHU 74]. En effet, le test d'équivalence de $\mathcal{S}(u)$ et $\mathcal{S}(v)$, et le calcul de $d(u, v)$, conduisent tous les deux à l'étude de la différence symétrique Δ de $|\mathcal{S}(u)|$ et $|\mathcal{S}(v)|$. Dans un cas on teste si Δ est vide, dans l'autre on calcule la longueur minimale des mots de Δ . En fait, la méthode que nous présentons est plus précise encore : il s'agit de déterminer le mot h de Δ le plus petit pour l'ordre généalogique; on a alors $d(u, v) = |h| - 1$. Ceci est obtenu en générant une suite de mots finie, croissante et contenant nécessairement le mot h s'il existe. Naturellement, le calcul est d'autant plus rapide que la suite comporte moins d'éléments. Nous considérons successivement deux suites. La première, qui donne un algorithme quadratique, correspond au parcours en largeur du graphe de l'automate produit $\mathcal{S}(u) \times (v)$. La construction de la seconde, sous-suite de la précédente, se ramène à un problème de type « UNION-FIND » pour lequel il existe une solution simple et de complexité optimale [AHU 74]; elle permet d'effectuer le calcul de $d(u, v)$ en temps quasi linéaire.

Considérons les deux mots $u = u_1 \dots u_m$ ($u_i \in A$), $v = v_1 \dots v_n$ ($v_i \in A$) et les automates

$$\mathcal{S}(u) = (\{0, \dots, m+1\}, 0, \{0, \dots, m\}, t_u),$$

$$\mathcal{S}(v) = (\{\bar{0}, \dots, \bar{n}+1\}, \bar{0}, \{\bar{0}, \dots, \bar{n}\}, t_v).$$

Soient

$$Z = \{w \in A^* / \forall x \in A^*, (0.w, \bar{0}.w) = (0.x, \bar{0}.x) \Rightarrow w \leq x\},$$

et

$$l = \text{card}(Z) (\leq (m+2)(n+2)).$$

Soit $z(i)$ ($1 \leq i \leq l$) la suite obtenue en ordonnant, de façon croissante pour l'ordre généalogique les éléments de Z .

On remarque que Z est clos par préfixe :

$$\forall x, y \in A^*, \quad xy \in Z \Rightarrow x \in Z.$$

Il est alors facile de construire z à l'aide d'une file f . La construction fait appel aux opérations suivantes :

INITIALISER (f) : procédure qui vide f .

AJOUTER (w, f) : procédure qui ajoute l'élément w à la fin de f .

TÊTE (f) : fonction dont le résultat est le premier élément de f .

SUPPRIMER (f) : procédure qui supprime le premier élément de f .

MARQUER (p, q) : procédure qui permet de marquer un élément (p, q) de $\{0, \dots, m+1\} \times \{\bar{0}, \dots, \bar{n}+1\}$.

Construction de z

```

début
   $i \leftarrow 1$ ;  $z(i) \leftarrow \epsilon$ ; INITIALISER ( $f$ ); AJOUTER ( $\epsilon, f$ );
  tant que  $f \neq$  file vide faire
    début
       $w \leftarrow$  TÊTE ( $f$ ); SUPPRIMER ( $f$ );
      pour  $a$  variant de la première à la dernière lettre de  $A$  faire
        si  $(0, wa, \bar{0}, wa)$  n'est pas marqué alors
          début  $i \leftarrow i + 1$ ;  $z(i) \leftarrow wa$ ; AJOUTER ( $wa, f$ );
            MARQUER ( $0, wa, \bar{0}, wa$ )
          fin
        fin {tant que}
      fin {construction}
    fin
  fin

```

La suite z permet de calculer $d(u, v)$. En effet, si $u \neq v$, il est clair que Z contient le mot le plus petit pour l'ordre généalogique qui distingue u et v .

Pratiquement on utilise la suite t d'éléments de $\{0, \dots, m+1\} \times \{\bar{0}, \dots, \bar{n}+1\}$ définie par $t(i) = (0, z(i), \bar{0}, z(i))$ ($1 \leq i \leq l$). Cette suite est celle que l'on obtiendrait en effectuant un parcours en largeur du graphe de l'automate produit $\mathcal{S}(u) \times \mathcal{S}(v)$, les sommets adjacents d'un sommet s donné étant examinés dans l'ordre des étiquettes des flèches issues de s .

Soit $D = \{0, \dots, m\} \times \{\bar{n}+1\} \cup \{m+1\} \times \{\bar{0}, \dots, \bar{n}\}$. $z(i)$ distingue u et v si et seulement si $t(i) \in D$. L'algorithme 9 calcule $d(u, v)$, il utilise une file F d'éléments de la forme (p, \bar{q}, k) où $(p, \bar{q}) = t(i)$, $k = |z(i)|$ ($1 \leq i \leq l$).

ALGORITHME 9 : Fonction $d(u, v)$.

```

Début
  INITIALISER ( $F$ ); AJOUTER ( $((0, \bar{0}, 0), F)$ );
  tant que  $F \neq$  file vide faire
    début
       $(p, \bar{q}, k) \leftarrow$  TÊTE ( $F$ ); SUPPRIMER ( $F$ );
      pour  $a$  variant de la première à la dernière lettre de  $A$  faire

```

$\underline{\text{si}} (p.a, \bar{q}.a) \text{ n'est pas marqué alors}$
 $\quad \underline{\text{début}}$
 $\quad \quad \underline{\text{si}} (p.a, \bar{q}.a) \in D \text{ alors retour } (k);$
 $\quad \quad \text{AJOUTER } ((p.a, \bar{q}.a, k+1), F);$
 $\quad \quad \text{MARQUER } (p.a, \bar{q}.a)$
 $\quad \underline{\text{fin}} \{\text{si}\}$
 $\quad \underline{\text{fin}}; \{\text{tant que}\}$
 $\quad \underline{\text{retour}} (\infty)$
 $\underline{\text{fin}} \{\text{fonction}\}$

PROPOSITION 10 : La complexité maximale de l'algorithme 9 est $O(|A||u||v|)$.

Démonstration : Le test $(p.a, \bar{q}.a) \in D$ et les opérations AJOUTER et MARQUER se font en temps constant. Le temps d'exécution est donc proportionnel au nombre d'éléments placés dans file; celui-ci est inférieur à $(|u|+2)(|v|+2)$. ■

Améliorer l'algorithme 9 : $d(u, v)$ est plus rapidement calculée en utilisant une sous-suite de z .

Soient $z'(i)$ ($1 \leq i \leq l'$) une suite de mots et $G(i) = (S, A(i))$ ($1 \leq i \leq l'$) une suite de graphes non orientés, définies de la façon suivante : $S = \{0, \dots, m+1\} \cup \{\bar{0}, \dots, \bar{n}+1\}$ est l'ensemble des sommets de $G(i)$, $A(i)$ l'ensemble des arêtes; $z'(1) = \varepsilon$ et $A(1) = \emptyset$; supposons connus les éléments

$$z'(1), \dots, z'(i), \quad A(1), \dots, A(i),$$

et soit

$$E(i) = \{x \in A^*/z'(i) \leq x \quad \text{et} \quad C_i(0.x) \neq C_i(\bar{0}.x)\}$$

où $C_i(p)$ ($p \in S$) est la composante connexe de p dans $G(i)$; si $E(i) = \emptyset$ alors $l' = i$, sinon

$$z'(i+1) = \min E(i)$$

et

$$A(i+1) = A(i) \cup \{\{0, z'(i+1), \bar{0}, z'(i+1)\}\}.$$

L'ensemble $Z' = \{z'(i) / 1 \leq i \leq l'\}$ est clos par préfixe :

$$\forall x, y \in A^*, \quad xy \in Z' \Rightarrow x \in Z'.$$

Le schéma suivant décrit la construction de la suite z' . COMPOSANTE (p) ($p \in S$) désigne la composante connexe de p dans le graphe courant et UNION (p, q) ($p, q \in S$) est une procédure qui effectue l'ajout de l'arête $\{p, q\}$; f est une file de mots.

Construction de z'

```

débüt
   $i \leftarrow 1$ ;  $z'(i) \leftarrow \varepsilon$ ; INITIALISER ( $f$ ); AJOUTER ( $\varepsilon, f$ );
  tant que  $f \neq$  file vide faire
    débüt
       $w \leftarrow$  TÊTE ( $f$ ); SUPPRIMER ( $f$ );
      pour  $a$  variant de la première à la dernière lettre de  $A$  faire
        si COMPOSANTE ( $0.w$ )  $\neq$  COMPOSANTE ( $\bar{0}.w$ ) alors
          débüt  $i \leftarrow i + 1$ ;  $z'(i) \leftarrow wa$ ; AJOUTER ( $wa, f$ );
          UNION ( $0.w, \bar{0}.w$ )
          fin
        fin {tant que}
    fin {tant que}
fin (construction)

```

La proposition suivante montre que la suite z' permet de calculer $d(u, v)$.

PROPOSITION 11 : Si $u \neq v$ et si h est le plus petit mot (pour l'ordre généalogique) qui distingue u et v , alors il existe $i \in \{1, \dots, l'\}$ tel que $z'(i) = h$.

Démonstration : Si $|h| = 1$, c'est évident. Supposons $|h| \geq 2$. Soit $k = \max \{i | z'(i) \leq h \text{ et } C_i(0.h) \neq C_i(\bar{0}.h)\}$; k existe car $|h| \geq 2$. $E(k) \neq \emptyset$ car $h \in E(k)$; donc $k < l'$. Si $h \neq \min E(k)$ alors $z'(k+1) < h$, et $C_{k+1}(0.h) \neq C_{k+1}(\bar{0}.h)$ [en effet, l'un des deux états $0.h$ et $\bar{0}.h$ est final et l'autre pas, par conséquent, l'égalité $C_{k+1}(0.h) = C_{k+1}(\bar{0}.h)$ impliquerait l'existence d'un entier $j \leq k+1$ tel que l'un des deux états $0.z'(j)$ et $\bar{0}.z'(j)$ soit final et l'autre pas, $z'(j)$ distinguerait alors u et v , ce qui contredirait la minimalité de h]. Ceci est impossible, k étant maximal, donc $h = z'(k+1)$ et $k+1 \leq l'$. ■

On en déduit l'algorithme 12. Il effectue le calcul de $d(u, v)$ en utilisant la suite $t'(i) = (0.z'(i), \bar{0}.z'(i))$ ($1 \leq i \leq l'$). F est une file d'éléments de la forme (p, \bar{q}, k) , où $(p, \bar{q}) = t'(i)$, $k = |z'(i)|$ ($1 \leq i \leq l'$).

ALGORITHME 12 : Fonction $d(u, v)$

```

Débüt
  INITIALISER ( $F$ ); AJOUTER ( $(0, \bar{0}, 0), F$ );
  tant que  $F \neq$  file vide faire
    débüt
       $(p, \bar{q}, k) \leftarrow$  TÊTE ( $F$ ); SUPPRIMER ( $F$ );
      pour  $a$  variant de la première à la dernière lettre de  $A$  faire
        si COMPOSANTE ( $p.a$ )  $\neq$  COMPOSANTE ( $\bar{q}.a$ ) alors
          débüt
            si  $(p.a, \bar{q}.a) \in D$  alors retour ( $k$ );
            AJOUTER ( $(p.a, \bar{q}.a, k+1), F$ );
            UNION ( $p.a, \bar{q}.a$ )
          fin {si}
        fin {tant que}
    fin {tant que}

```

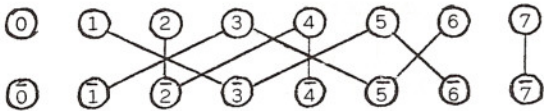
retour (∞)
fin {fonction}

REMARQUE : Si l'on souhaite obtenir effectivement h , le mot le plus petit qui distingue u et v , il suffit de définir un tableau PERE et d'effectuer PERE $[p.a, \bar{q}.a] \leftarrow (p, \bar{q}, a)$ à chaque fois que l'élément $(p.a, \bar{q}.a, k+1)$ est placé dans la file. On peut alors facilement construire h à partir de PERE.

Exemple :

$A = \{a, b, c\}, \quad u = cabacb, \quad v = bacabc$
 $\mathcal{S}(u) : \textcircled{0} \ c \ \textcircled{1} \ a \ \textcircled{2} \ b \ \textcircled{3} \ a \ \textcircled{4} \ c \ \textcircled{5} \ b \ \textcircled{6} \ \textcircled{7}$
 $\textcircled{0} \ b \ \textcircled{1} \ a \ \textcircled{2} \ c \ \textcircled{3} \ a \ \textcircled{4} \ b \ \textcircled{5} \ c \ \textcircled{6} \ \textcircled{7}$

Graphe



File

$(0, \bar{0}, 0)$
 $(2, \bar{2}, 1) (3, \bar{1}, 1) (1, \bar{3}, 1)$
 $(3, \bar{1}, 1) (1, \bar{3}, 1) (4, \bar{4}, 2) (3, \bar{5}, 2) (5, \bar{3}, 2)$
 $(1, \bar{3}, 1) (4, \bar{4}, 2) (3, \bar{5}, 2) (5, \bar{3}, 2) (4, \bar{2}, 2) (6, \bar{5}, 2)$
 $(4, \bar{4}, 2) (3, \bar{5}, 2) (5, \bar{3}, 2) (4, \bar{2}, 2) (6, \bar{5}, 2) (5, \bar{6}, 2)$
 $(3, \bar{5}, 2) (5, \bar{3}, 2) (4, \bar{2}, 2) (6, \bar{5}, 2) (5, \bar{6}, 2) (7, \bar{7}, 3)$
 $(5, \bar{3}, 2) (4, \bar{2}, 2) (6, \bar{5}, 2) (5, \bar{6}, 2) (7, \bar{7}, 3)$
 $(p, \bar{q}, k) = (3, \bar{5}, 2), \quad 3.a = 4, \quad \bar{5}.a = \bar{7}$
 $(4, \bar{7}) \in D \Rightarrow \underline{d(u, v) = 2}$

PÈRE

$\begin{matrix} + \\ 2, \bar{2} & 3, \bar{1} & 1, \bar{3} & 4, \bar{4} \\ 0, \bar{0}, a & 0, \bar{0}, b & 0, \bar{0}, c & 2, \bar{2}, a \end{matrix}$
 $\begin{matrix} + \\ 3, \bar{5} & 5, \bar{3} & 4, \bar{2} & 6, \bar{5} & 5, \bar{6} & 7, \bar{7} \\ 2, \bar{2}, b & 2, \bar{2}, c & 3, \bar{1}, a & 3, \bar{1}, b & 1, \bar{3}, c & 4, \bar{4}, a \end{matrix}$

$h = aba$

On remarque que, par exemple, $(2, \bar{4}, 2)$ n'a pas été placé dans la file car 2 et $\bar{4}$ étaient déjà dans la même composante connexe.

Mise en œuvre : On connaît une méthode optimale pour mettre en œuvre les opérations UNION et COMPOSANTE (cf. [AHU 74]). Elle permet d'effectuer une suite de $O(|Z|)$ opérations UNION et COMPOSANTE en temps $O(|S|\alpha(|S|))$, où α est une fonction dont la croissance est extrêmement faible : pour tout entier n , $\alpha(n)$ est le plus petit entier k tel que $F(k) \geq n$, avec

$$F(0)=1 \quad \text{et} \quad F(i)=2^{F(i-1)}$$

pour $i > 0$ ($\alpha(n) \leq 5$ pour tout $n \leq 2^{65536}$). En outre, la complexité spatiale de la méthode est linéaire.

PROPOSITION 13. — Soient u et v deux mots de A^* . Le calcul de $d(u, v)$ par l'algorithme 12 est effectué en temps

$$O(|A|(|u|+|v|)\alpha(|u|+|v|)),$$

sa complexité en espace est $O(|A|(|u|+|v|))$.

Démonstration : Soit N_1 et N_2 respectivement le nombre d'opérations UNION et COMPOSANTE effectuées par l'algorithme. L'opération UNION provoque l'ajout d'une arête entre deux sommets dont les composantes connexes sont distinctes, elle est donc réalisée au plus $|S|-1$ fois. Par conséquent, le nombre d'éléments placés dans la file est inférieur ou égal à $|S|-1$. On a

$$N_1 \leq |S|-1 \quad \text{et} \quad N_2 \leq 2|A|(|S|-1). \quad |S| = |u| + |v| + 4,$$

d'où le résultat. ■

BIBLIOGRAPHIE

- AHU 74. A. V. AHO, J. E. HOPCROFT et J. D. ULLMAN, *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, Mass., 1974.
- CC 82. A. CARDON et M. CROCHEMORE, *Partitioning a graph in $O(|A|\log_2|V|)$* , Theor. Comput. Sci. vol. 19, 1982, p. 85-98.
- He 84. J. J. HEBRARD, *Distances sur les mots. Application à la recherche de motifs*, Thèse de 3^e cycle, Université de Haute-Normandie, 1984.
- Hi 77. D. S. HIRSCHBERG, *Algorithms for the Longest Common Subsequence Problem*, J. Assoc. Comput. Mach., vol. 24, 1977, p. 664-675.
- HS 77. J. W. HUNT et T. G. SZYMANSKI, *A Fast Algorithm for Computing Longest Common Subsequences*, Comm. ACM., vol. 20, 1977, p. 350-353.

- Lo 82. LOTHAIRE, *Combinatorics on Words*, Addition-Wesley, Reading, Mass., 1982.
- Mo 70. H. L. MORGAN, *Spelling Correction in Systems Programs*, Comm. ACM., vol. 13, 1970, p. 90-94.
- MP 80. W. J. MASEK et M. S. PATERSON, *A Faster Algorithm Computing String Edit Distances*, J. Comput. and Sys. Sci., vol. 20, 1980, p. 18-31.
- NKY 82. N. NAKATSU, Y. KAMBAYASHI et S. YAJIMA, *A Longest Common Subsequence Algorithm Suitable for Similar Test Strings*, Acta Informatica, vol. 18, 1982, p. 171-179.
- Se 74. P. H. SELLERS, *An Algorithm for the Distance between two Finite Sequences*, J. Combinatorial Theory, Series A, vol. 16, 1974, p. 253-258.
- Si 84. I. SIMON, *An Algorithm to Distingsh Words Efficiently by their Subwords* Communication at "Combinatorial Algorithms on words" conference Maratea (1984).
- SK 83. D. SANKOFF et J. B. KRUSKAL, *Time Warps, String Edits, and Macromolecules: the Theory and Practice of Sequence Comparison*, Addison-Wesley, Reading, Mass., 1983.
- WF 74. R. A. WAGNER et M. J. FISCHER, *The String to String Correction Problem*, J. Assoc. Comput. Mach., vol. 21, 1974, p. 168-173.