



UNIVERSITÉ DE LIÈGE

Recherche de motifs dans des séquences d'ADN

INFO0902-1 – Structure de données et algorithmes

François ROZET (s161024)

Bachelier Ingénieur civil
Année académique 2018-2019

Table des matières

1	Analyse théorique	1
1.1	Recherche dichotomique par force brute ou table de hachage	1
1.2	Fonction de coût	3
1.3	Programmation dynamique	3
1.4	Complexité	3

1 Analyse théorique

1.1 Recherche dichotomique par force brute ou table de hachage

Pseudo-code 1 Recherche dichotomique par force brute

```

1 function DICHOTOMIC-SEARCH( $X, Y$ )
2    $t = NIL$ 
3    $p = 1; r = \min(X.length, Y.length)$ 
4   while  $p \leq r$  do
5      $q = \lfloor \frac{p+r}{2} \rfloor$ 
6      $s = \text{BRUTE-FORCE-AUX}(X, Y, q)$ 
7     if  $s \neq NIL$  then
8        $p = q + 1$ 
9     else
10       $r = q - 1$ 
11       $t = s$ 
12  return  $t$ 

```

Pseudo-code 2 Recherche dichotomique par table de hachage

```

1 function HASH-TABLE-SEARCH( $X, Y$ )
2    $t = NIL$ 
3    $p = 1; r = \min(X.length, Y.length)$ 
4   while  $p \leq r$  do
5      $q = \lfloor \frac{p+r}{2} \rfloor$ 
6      $s = \text{HASH-TABLE-AUX}(X, Y, q)$ 
7     if  $s \neq NIL$  then
8        $p = q + 1$ 
9     else
10       $r = q - 1$ 
11       $t = s$ 
12  return  $t$ 

```

Pseudo-code 3 Recherche auxiliaire par force brute

```

1 function BRUTE-FORCE-AUX( $X, Y, l$ )
2   for  $i = 1$  to  $X.length - l + 1$  do
3     for  $j = 1$  to  $Y.length - l + 1$  do
4       for  $k = 1$  to  $l$  do
5         if  $X[i + k - 1] \neq Y[j + k - 1]$  then
6            $k = 0$ 
7           break
8       if  $k \neq 0$  then
9         return  $(l, i, j)$ 
10  return  $NIL$ 

```

Pseudo-code 4 Recherche auxiliaire par table de hachage sans calcul incrémental

```

1 function HASH-TABLE-AUX1( $X, Y, l$ )
2   Let  $H$  be a new hash table
3   for  $i = 1$  to  $X.length - l + 1$  do
4      $f = \text{ENCODE}(X, i, l)$ 
5     HASH-INSERT( $H, f, i$ )
6   for  $j = 1$  to  $Y.length - l + 1$  do
7      $f = \text{ENCODE}(Y, j, l)$ 
8      $i = \text{HASH-SEARCH}(H, f)$ 
9     if  $i \neq \text{NIL}$  then
10      return  $(l, i, j)$ 
11  return  $\text{NIL}$ 

```

Pseudo-code 5 Recherche auxiliaire par table de hachage avec calcul incrémental

```

1 function HASH-TABLE-AUX2( $X, Y, l$ )
2   Let  $H$  be a new hash table
3    $b = 4; p = b^{l-1}$ 
4    $f = \text{ENCODE}(X, 1, l)$ 
5   HASH-INSERT( $H, f, 1$ )
6   for  $i = 2$  to  $X.length - l + 1$  do
7      $f = (f - X[i - 1] \cdot p) \cdot b + X[i + l - 1]$ 
8     HASH-INSERT( $H, f, i$ )
9    $f = \text{ENCODE}(Y, 1, l)$ 
10   $i = \text{HASH-SEARCH}(H, f)$ 
11  if  $i \neq \text{NIL}$  then
12    return  $(l, i, 1)$ 
13  for  $j = 2$  to  $Y.length - l + 1$  do
14     $f = (f - Y[j - 1] \cdot p) \cdot b + Y[j + l - 1]$ 
15     $i = \text{HASH-SEARCH}(H, f)$ 
16    if  $i \neq \text{NIL}$  then
17      return  $(l, i, j)$ 
18  return  $\text{NIL}$ 

```

Pseudo-code 6 Fonction d'encodage

```

1 function ENCODE( $X, i, l$ )
2    $b = 4$ 
3    $f = 0$ 
4   for  $j = 1$  to  $l$  do
5      $f = f \cdot b + X[i + j - 1]$ 
6   return  $f$ 

```

1.2 Fonction de coût

$$C[i, j] = \begin{cases} C[i-1, j-1] + 1 & \text{si } i, j > 0 \text{ et } X[i] = Y[j] \\ 0 & \text{sinon} \end{cases} \quad (1)$$

1.3 Programmation dynamique

Pseudo-code 7 Recherche par programmation dynamique

```

1 function DYNAMIC-SEARCH( $X, Y$ )
2   Let  $C[0 \dots X.length, 0 \dots Y.length]$  be a new table
3    $l = 0$ 
4   for  $i = 0$  to  $X.length$  do
5     for  $j = 0$  to  $Y.length$  do
6       if  $i > 0$  and  $j > 0$  and  $X[i] == Y[j]$  then
7          $C[i, j] = C[i-1, j-1] + 1$ 
8         if  $C[i, j] > l$  then
9            $l = C[i, j]$ 
10           $i_{end} = i$ 
11           $j_{end} = j$ 
12       else  $C[i, j] = 0$ 
13   if  $l \neq 0$  then
14     return  $(l, i_{end} - l + 1, j_{end} - l + 1)$ 
15   return  $NIL$ 

```

1.4 Complexité

Soit n la longueur des deux séquences et αn la longueur de leur sous-séquence contiguë commune la plus longue.

- (a) La solution par force brute est composée de quatre boucles imbriquées. Les boucles initiale et intermédiaires exécutent toujours de l'ordre de $(1 - \alpha)n$ itérations. La boucle finale, quant à elle, en exécute de l'ordre de αn . Les opérations à l'intérieur de ces boucles étant toutes $\mathcal{O}(1)$, la complexité en temps est

$$\mathcal{O}(\alpha(1 - \alpha)^3 n^4) = \mathcal{O}(n^4) \quad (2)$$

Néanmoins, lorsque αn tend vers n ($\alpha \rightarrow 1$), les trois premières boucles voient leur nombre d'itérations tendre vers 1. Ainsi, la complexité en temps est $\Omega(n)$. On trouve aussi que $\alpha = \frac{1}{4}$ est le pire cas, bien qu'il soit équivalent au cas moyen en terme de complexité.

Concernant la complexité en espace, selon ce qu'il a été décidé de renvoyer, elle peut être $\Theta(1)$ ou $\Theta(\alpha n)$. Si la recherche renvoie un triplet (l, i, j) (l la longueur, i l'indice de départ dans $G1$ et j l'indice de départ dans $G2$), la complexité sera

$\Theta(1)$. À l'inverse, si la recherche renvoie la sous-séquence sous forme d'une table¹, la complexité sera $\Theta(\alpha n)$.

- (b) Comme la solution par force brute, la recherche dichotomique est composée de quatre boucles imbriquées. Cependant, ici, la boucle initiale s'arrête lorsque un interval de taille n a été divisé *suffisamment* de fois, c.-à-d. $\log_2(n)$ fois. De plus, les boucles intermédiaires sont maintenant en $\mathcal{O}(n)$. Ainsi, la complexité en temps est

$$\mathcal{O}(\alpha n^3 \log(n)) = \mathcal{O}(n^3 \log(n)) \quad (3)$$

Pour cette solution, lorsque αn tend vers n , les boucles intermédiaires ne voient pas leur nombre d'itérations tendre vers 1. Cependant, lorsque $\alpha n = 0$, la boucle la plus profonde n'effectue plus qu'une seule itération. Ainsi, la complexité en temps est $\Omega(n^2 \log(n))$. Le pire cas est équivalent au cas moyen.

Concernant la complexité en espace, cette solution est équivalent à la force brute.

- (c) La solution par table de hachage ne possède que deux boucles imbriquées qui correspondent aux deux premières de la recherche dichotomique. Néanmoins, sans calcul incrémental, l'opération d'encodage exécute $\mathcal{O}(n)$ opérations à chaque appel. On suppose aussi que l'insertion et la lecture dans la table de hachage sont faites en $\mathcal{O}(1)$. Dès lors, la complexité en temps est

$$\Theta(n^2 \log(n)) \quad (4)$$

Concernant la complexité en espace, à chaque appel de la fonction auxiliaire pour une taille l de sous-séquence, il est nécessaire d'initialiser une table de hachage sur un univers $U = \{0, 1, \dots, 4^l - 1\}$ afin d'y stocker $n - l + 1$ sous-séquences. Pour conserver un accès (insertion, recherche et suppression) aux éléments de la table en $\mathcal{O}(1)$, il est nécessaire d'avoir $m = \Omega(n - l)$ où m est sa taille². Or, la valeur de départ de l est $\frac{n}{2}$, dès lors, la complexité en espace de la solution par table de hachage est $\Theta(n)$.

- (d) Avec le calcul incrémental, l'encodage ne demande plus que $\mathcal{O}(1)$ opérations par appel. Ainsi, la complexité en temps devient

$$\Theta(n \log(n)) \quad (5)$$

La complexité en espace n'a, quant à elle, pas changé.

- (e) La solution par programmation dynamique n'utilise que deux boucles imbriquées itérant sur l'entière des séquences. Dès lors, la complexité en temps est

$$\Theta(n^2) \quad (6)$$

Concernant la complexité en espace, la matrice de coût possède $(n + 1)^2$ éléments. La complexité spatiale est donc aussi $\Theta(n^2)$.

1. Les pseudo-codes ont été rédigés selon cette politique.

2. Dans le cas d'une table de hachage par adressage direct, la condition $m \geq n - l + 1$ doit aussi être respectée.