

Recherche textuelle

DIU Enseignement de l'informatique au lycée
Bloc 5 : Algorithmique avancée

Bruno Grenet – Université de Montpellier

Juin 2021

Problème

Recherche de motif

Entrées : Un texte t de longueur n , un motif x de longueur $m < n$

Sortie : Toutes les positions de x dans t

Problème

Recherche de motif

Entrées : Un texte t de longueur n , un motif x de longueur $m < n$

Sortie : Toutes les positions de x dans t

abaa dans *acaabbabaaa* ?

Un premier algorithme

RECHERCHE NAÏVE

Entrées : t et x

1. $P \leftarrow$ liste vide
2. Pour $i = 0$ à $n - m$:
3. $j \leftarrow m - 1$
4. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
5. $j \leftarrow j - 1$
6. Si $j = -1$: Ajouter i à P
7. Renvoyer P

$x = abaa$

$a c a a b b a b a a a$

Un premier algorithme

RECHERCHE NAÏVE

Entrées : t et x

1. $P \leftarrow$ liste vide
2. Pour $i = 0$ à $n - m$:
3. $j \leftarrow m - 1$
4. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
5. $j \leftarrow j - 1$
6. Si $j = -1$: Ajouter i à P
7. Renvoyer P

$x = abaa$

$a c a a b b a b a a a$

Théorème

L'algorithme RechercheNaïve renvoie la liste des positions de x dans t en temps $O(mn)$.

Notations

- ▶ Longueur : $|w|$
- ▶ Lettres : $w_0, \dots, w_{|w|-1}$
- ▶ Fenêtre de taille m en position i : $w_{[i, i+m[} = w_i w_{i+1} \cdots w_{i+m-1}$
- ▶ x apparaît en position i dans t : $x = t_{[i, i+|x|[}$

La règle du mauvais caractère : exemple

$x = abaa$

$a c a a b b a b a a a$

Dernière occurrence non finale

$d_x(\ell)$: position de la dernière occurrence de ℓ dans x , non finale

Dernière occurrence non finale

$d_x(\ell)$: position de la dernière occurrence de ℓ dans x , non finale

Exemple : $x = abaabcab$

ℓ		a	b	c	d	\dots
$d_x(\ell)$		6	4	5	-1	-1

Règle du mauvais caractère

Théorème

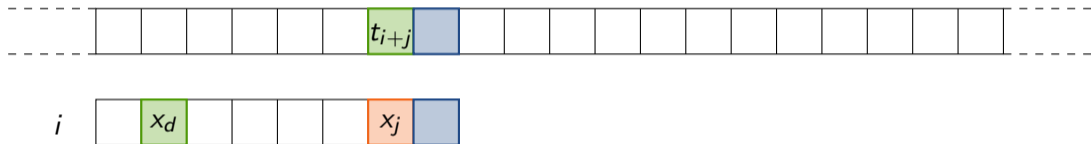
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.

Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.

Règle du mauvais caractère

Théorème

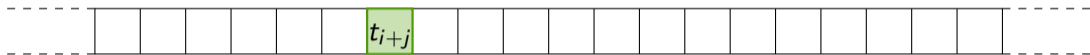
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.
Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.



Règle du mauvais caractère

Théorème

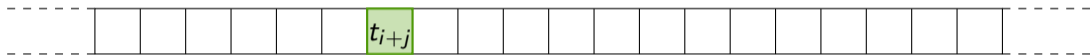
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.
Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.



Règle du mauvais caractère

Théorème

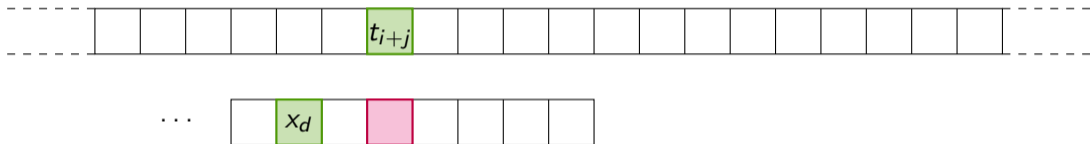
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.
Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.



Règle du mauvais caractère

Théorème

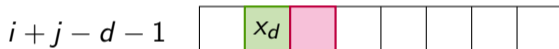
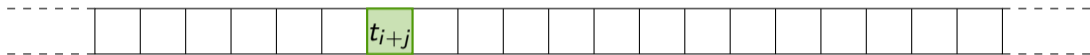
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.
Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.



Règle du mauvais caractère

Théorème

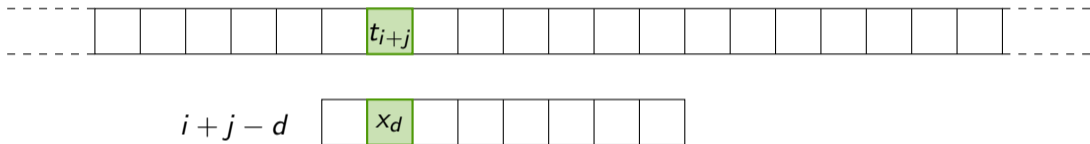
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.
Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.



Règle du mauvais caractère

Théorème

Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_k = t_{i+k}$ pour $k > j$ et $d = d_x(t_{i+j})$.
Alors si $j > d$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j - d$.



Algorithme de Boyer et Moore simplifié – version 1

BOYERMOORE-MC

Entrées : t , x et d_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + \max(1, j - d_x(t_{i+j}))$
10. Renvoyer P

Algorithme de Boyer et Moore simplifié – version 1

BOYERMOORE-MC

Entrées : t , x et d_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + \max(1, j - d_x(t_{i+j}))$
10. Renvoyer P

Théorème

L'algorithme BoyerMoore-MC renvoie la liste des positions de x dans t en temps $O(mn)$.

Algorithme de Boyer et Moore simplifié – version 1

BOYERMOORE-MC

Entrées : t , x et d_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + \max(1, j - d_x(t_{i+j}))$
10. Renvoyer P

Théorème

L'algorithme BoyerMoore-MC renvoie la liste des positions de x dans t en temps $O(mn)$.

- ▶ même pire cas que la recherche naïve :
 - ▶ motif $ba \cdots a$
 - ▶ texte $a \cdots \cdots a$

Algorithme de Boyer et Moore simplifié – version 1

BOYERMOORE-MC

Entrées : t , x et d_x

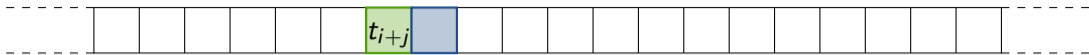
1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + \max(1, j - d_x(t_{i+j}))$
10. Renvoyer P

Théorème

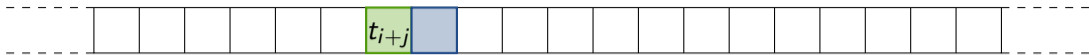
L'algorithme BoyerMoore-MC renvoie la liste des positions de x dans t en temps $O(mn)$.

- ▶ même pire cas que la recherche naïve :
 - ▶ motif $ba \cdots a$
 - ▶ texte $a \cdots \cdots a$
- ▶ *entrée supplémentaire* : tableau d des dernières occurrences non finales
 - ▶ exercice !

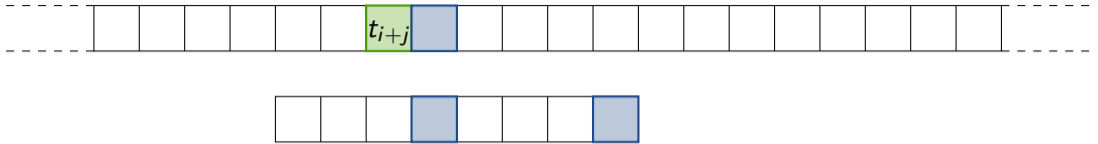
Variantes



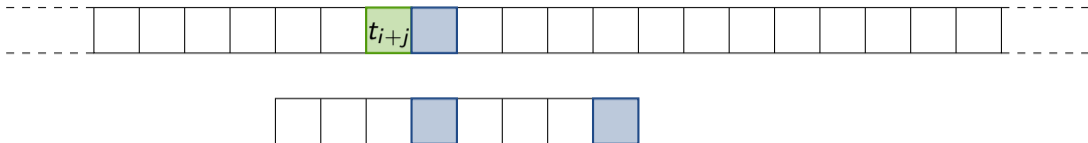
Variantes



Variantes

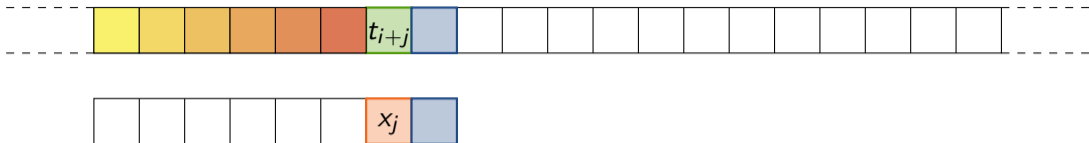


Variantes



- ▶ Algorithme de Horspool : alignement de la dernière lettre de la fenêtre $O(mn)$

Variantes



- ▶ Algorithme de Horspool : alignement de la dernière lettre de la fenêtre $O(mn)$
- ▶ Alignement le plus à droite possible, en tenant compte de toutes les lettres $O(mn)$

Règle du bon suffixe : exemple

$x = abaaaa$

a b b c a a c a a a a b a a a a

Règle du bon suffixe : p_x et s_x

- ▶ $p_x(j)$: longueur du plus long préfixe de x ($\neq x$) qui soit aussi suffixe de $x_{[j,m[}$
- ▶ $s_x(j)$: position la plus à droite d'une copie de $x_{[j,m[}$, non précédée de x_{j-1}

Règle du bon suffixe : p_x et s_x

- ▶ $p_x(j)$: longueur du plus long préfixe de x ($\neq x$) qui soit aussi suffixe de $x_{[j,m[}$
- ▶ $s_x(j)$: position la plus à droite d'une copie de $x_{[j,m[}$, non précédée de x_{j-1}

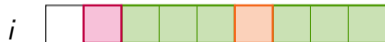
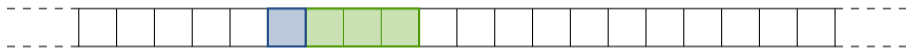
Exemple

x		b	a	b	a	b	a
j		0	1	2	3	4	5
$p_x(j)$		4	4	4	2	2	0
$s_x(j)$		-1	-1	0	-1	0	-1

Règle(s) du bon suffixe (1/2)

Théorème

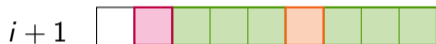
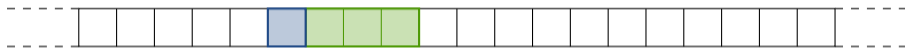
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) \geq 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j + 1 - s_x(j)$.



Règle(s) du bon suffixe (1/2)

Théorème

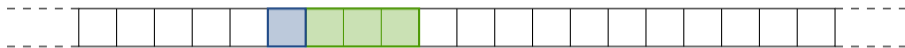
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) \geq 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j + 1 - s_x(j)$.



Règle(s) du bon suffixe (1/2)

Théorème

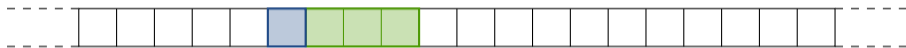
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) \geq 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j + 1 - s_x(j)$.



Règle(s) du bon suffixe (1/2)

Théorème

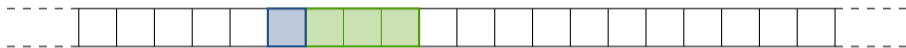
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) \geq 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j + 1 - s_x(j)$.



Règle(s) du bon suffixe (1/2)

Théorème

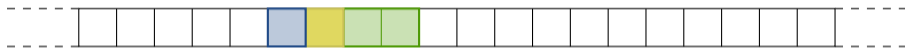
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) \geq 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + j + 1 - s_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

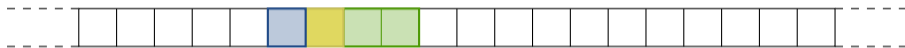
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

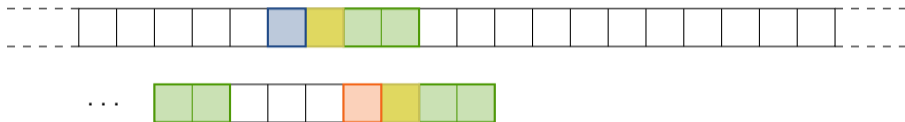
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

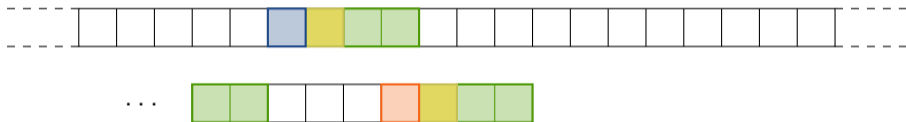
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

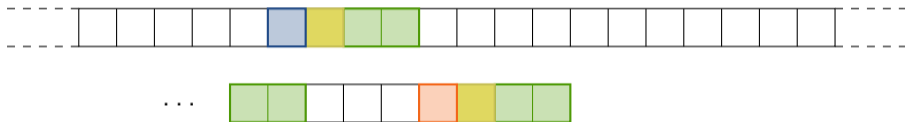
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

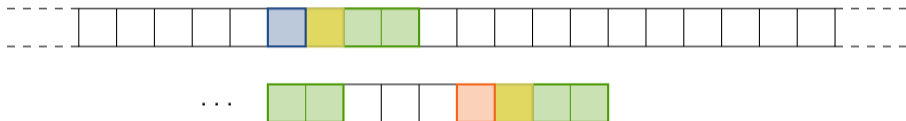
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

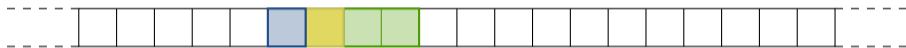
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

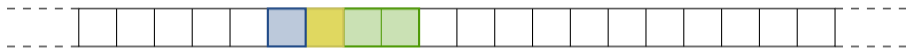
Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Règle(s) du bon suffixe (2/2)

Théorème

Soit $t_{[i, i+m[}$ la fenêtre courante, avec $x_j \neq t_{i+j}$ et $x_{[j+1, m[} = t_{[i+j+1, i+m[}$.
Si $s_x(j) < 0$, $x \neq t_{[i_0, i_0+m[}$ pour $i < i_0 < i + m - p_x(j)$.



Algorithme de Boyer et Moore simplifié – version 2

BOYERMOORE-BS

Entrées : t , x , p_x et s_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + m - p_x(1)$
10. Sinon si $s_x(j + 1) \geq 0$:
11. $i \leftarrow i + j + 1 - s_x(j + 1)$
12. Sinon : $i \leftarrow i + m - p_x(j)$
13. Renvoyer P

Algorithme de Boyer et Moore simplifié – version 2

BOYERMOORE-BS

Entrées : t , x , p_x et s_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + m - p_x(1)$
10. Sinon si $s_x(j + 1) \geq 0$:
11. $i \leftarrow i + j + 1 - s_x(j + 1)$
12. Sinon : $i \leftarrow i + m - p_x(j)$
13. Renvoyer P

Théorème

L'algorithme BoyerMoore-BS renvoie la liste des positions de x dans t en temps $O(mn)$.

Algorithme de Boyer et Moore simplifié – version 2

BOYERMOORE-BS

Entrées : t , x , p_x et s_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + m - p_x(1)$
10. Sinon si $s_x(j + 1) \geq 0$:
11. $i \leftarrow i + j + 1 - s_x(j + 1)$
12. Sinon : $i \leftarrow i + m - p_x(j)$
13. Renvoyer P

Théorème

L'algorithme BoyerMoore-BS renvoie la liste des positions de x dans t en temps $O(mn)$.

- même pire cas que la recherche naïve

Algorithme de Boyer et Moore simplifié – version 2

BOYERMOORE-BS

Entrées : t , x , p_x et s_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + m - p_x(1)$
10. Sinon si $s_x(j + 1) \geq 0$:
11. $i \leftarrow i + j + 1 - s_x(j + 1)$
12. Sinon : $i \leftarrow i + m - p_x(j)$
13. Renvoyer P

Théorème

L'algorithme BoyerMoore-BS renvoie la liste des positions de x dans t en temps $O(mn)$.

- ▶ même pire cas que la recherche naïve
- ▶ tableaux s et p
 - ▶ calcul en temps $O(m^2)$ (exercice !)
 - ▶ difficile : temps linéaire $O(m)$

Algorithme de Boyer et Moore – version complète !

BOYERMOORE

Entrées : t , x , d_x , p_x et s_x

1. $P \leftarrow$ liste vide
2. $i \leftarrow 0$
3. Tant que $i \leq n - m$:
4. $j \leftarrow m - 1$
5. Tant que $j \geq 0$ et $x_j = t_{i+j}$:
6. $j \leftarrow j - 1$
7. Si $j = -1$:
8. Ajouter i à P
9. $i \leftarrow i + m - p_x(1)$
10. Sinon :
11. Si $s_x(j + 1) \geq 0$: $i \leftarrow i + \max(j + 1 - s_x(j + 1), j - d_x(t_{i+j}))$
12. Sinon : $i \leftarrow i + \max(m - p_x(j), j - d_x(t_{i+j}))$
13. Renvoyer P

Propriétés de l'algorithme de Boyer et Moore

Théorème

L'algorithme BoyerMoore renvoie la liste de toutes les positions de x dans t en temps $O(mn)$.

Propriétés de l'algorithme de Boyer et Moore

Théorème

L'algorithme BoyerMoore renvoie la liste de toutes les positions de x dans t en temps $O(mn)$.

▶ **Tout ça pour ça !**

Propriétés de l'algorithme de Boyer et Moore

Théorème

L'algorithme BoyerMoore renvoie la liste de toutes les positions de x dans t en temps $O(mn)$.

- ▶ **Tout ça pour ça !**
- ▶ Très rapide en pratique
- ▶ Complexité est $O(m + n)$ dans bcp de situations
- ▶ Existe une variante en complexité $O(m + n)$ dans tous les cas

Conclusion

- ▶ Autres algorithmes de complexité $O(m + n)$:
 - ▶ Knuth-Morris-Pratt (automates finis)
 - ▶ Karp-Rabin (technique probabiliste)
- ▶ Extensions :
 - ▶ Expressions régulières (grep, etc.)
 - ▶ Mots bidimensionnels (plus connus sous le nom d'*images*)
 - ▶ Motifs non contigus
- ▶ Utilisation :
 - ▶ Ctrl + F dans bcp de logiciels
 - ▶ Détection de gènes dans une séquence ADN
 - ▶ ...

$a^*bc^*.^*a$

a _____ b _____ a