

Recherche d'occurrences

Pierre Duret, Stéphane Rodriguez, Jérôme Signouret

27/06/2019

```
jerome@Nextcloud: ~  
Fichier Édition Affichage Rechercher Terminal Aide  
root@Nextcloud: /var/log/apache2# ls -l  
total 3224  
-rw-r----- 1 root adm 1228825 juin 25 22:49 access.log  
-rw-r----- 1 root adm 1168977 juin 24 23:59 access.log.1  
-rw-r----- 1 root adm 39320 juin 16 00:08 access.log.10.gz  
-rw-r----- 1 root adm 23238 juin 15 00:07 access.log.11.gz  
-rw-r----- 1 root adm 27672 juin 13 23:50 access.log.12.gz  
-rw-r----- 1 root adm 19198 juin 13 00:09 access.log.13.gz  
-rw-r----- 1 root adm 28602 juin 11 23:17 access.log.14.gz  
-rw-r----- 1 root adm 94628 juin 24 00:06 access.log.2.gz  
-rw-r----- 1 root adm 74364 juin 22 23:39 access.log.3.gz  
-rw-r----- 1 root adm 83763 juin 22 00:00 access.log.4.gz  
-rw-r----- 1 root adm 61930 juin 21 00:08 access.log.5.gz  
-rw-r----- 1 root adm 72862 juin 19 23:46 access.log.6.gz  
-rw-r----- 1 root adm 46464 juin 19 00:04 access.log.7.gz  
-rw-r----- 1 root adm 66438 juin 17 23:49 access.log.8.gz  
-rw-r----- 1 root adm 46359 juin 16 23:05 access.log.9.gz  
-rw-r----- 1 root adm 88717 juin 25 22:31 error.log  
-rw-r----- 1 root adm 823 juin 25 00:08 error.log.1  
-rw-r----- 1 root adm 6942 juin 16 00:08 error.log.10.gz  
-rw-r----- 1 root adm 535 juin 15 00:07 error.log.11.gz  
-rw-r----- 1 root adm 7791 juin 14 00:07 error.log.12.gz  
-rw-r----- 1 root adm 501 juin 13 00:09 error.log.13.gz  
-rw-r----- 1 root adm 699 juin 12 00:07 error.log.14.gz  
-rw-r----- 1 root adm 7088 juin 24 00:07 error.log.2.gz  
-rw-r----- 1 root adm 333 juin 23 00:07 error.log.3.gz  
-rw-r----- 1 root adm 8179 juin 22 00:05 error.log.4.gz  
-rw-r----- 1 root adm 612 juin 21 00:08 error.log.5.gz  
-rw-r----- 1 root adm 422 juin 20 00:09 error.log.6.gz  
-rw-r----- 1 root adm 408 juin 19 00:05 error.log.7.gz  
-rw-r----- 1 root adm 336 juin 18 00:08 error.log.8.gz  
-rw-r----- 1 root adm 335 juin 17 00:06 error.log.9.gz  
-rw-r----- 1 root adm 0 juin 15 2018 other_vhosts_access.log  
-rw-r----- 1 root adm 7839 juin 4 2018 other_vhosts_access.log.1  
root@Nextcloud: /var/log/apache2#
```



L^AT_EX

Table des matières

I	Le travail à faire	3
1	Principe	3
2	Accroche possible	3
2.1	Ce que l'on sait faire en seconde	3
3	La partie du programme NSI première concernée	4
4	Activité : Rechercher un caractère	4
4.1	Trouvez le premier	4
4.2	Trouvez les tous	4
4.3	Comptez le nombre de lignes	4
4.4	Combien de caractères lus ?	4
4.5	Lien entre le nombre de caractères lus et la taille de la chaîne	4
5	Le cours	4
5.1	Objectifs de la partie cours	4
5.2	Présentation de l'objectif de l'activité globale	5
5.2.1	Courbes	5
5.2.2	Partie théorique	5
5.2.3	Partie pratique	7
6	Activité : Rechercher un motif	7
6.1	Trouvez le premier	7
6.2	Trouvez les tous	8
6.3	Combien de caractères lus ?	8
6.4	Lien entre le nombre de caractères lus et la taille de la chaîne	8
6.5	Lien entre le nombre de caractères lus et la taille du motif	8
7	Exercices	8
7.1	Compter le nombre de caractère "e"	8
7.2	Compteur de performance	8
7.3	Évolution de l'algorithme précédent	9
II	Le travail fait	9
1	Activité : Rechercher un caractère	9
1.1	Trouvez le premier	9
1.2	Trouvez les tous	9
1.3	Comptez le nombre de lignes	9
1.4	Combien de caractères lus ?	10
1.5	Lien entre le nombre de caractères lus et la taille de la chaîne	10
2	Le cours	10

3	Activité : Rechercher un motif	12
3.1	Trouvez le premier	12
3.2	Trouvez les tous	12
4	Exercices	14
4.1	Compter le nombre de caractère "e"	14
4.2	Compteur de performance	14
4.3	Évolution de l'algorithme précédent	15
III	Remarques, allons plus loin	15

Première partie

Le travail à faire

1 Principe

On cherche à trouver un motif dans un fichier comportant un grand nombre d'entrées.

On commencera par rechercher un motif constitué d'un seul élément puis on généralisera à un motif fait de plusieurs éléments.

On calculera puis on expérimentera le "nombre d'opérations" que l'on réalise lors de la recherche et sa variation en fonction de la taille des données.

On pourra essayer d'introduire la notion de performance d'un algorithme.

2 Accroche possible

On souhaite analyser un fichier de log d'un serveur pour savoir s'il y a eu des tentatives de connections illicites et si elles ont échoué ou réussi.

2.1 Ce que l'on sait faire en seconde

Un fichier *log* contient des messages enregistrés par le système ou les applications. Cela permet d'améliorer les configs, déboguer, comprendre le fonctionnement...

Un fichier *log* est un fichier texte dont les caractères sont encodés par défaut en UTF-8 sur nos systèmes. Il est comme des fichiers *csv*, une ligne par enregistrement (jusqu'au caractère <fin de ligne>) les données étant séparées par des champs spéciaux (virgule, tabulation, espace, point-virgule).

- Ouvrir le fichier *access.txt* avec le tableur, choisir le caractère de séparation "espace" à l'import.
- Utiliser les fonctionnalités du tableur pour dénombrer le nombre de fois où apparaît le motif "GET"

3 La partie du programme NSI première concernée

Contenus	Capacités attendues	Commentaires
Parcours séquentiel d'un tableau	Écrire un algorithme de recherche d'une occurrence sur des valeurs de type quelconque.	On montre que le coût est linéaire.

4 Activité : Rechercher un caractère

4.1 Trouvez le premier

Vous devez écrire un algorithme de fonction qui renvoie la position du premier caractère recherché.

On considère qu'une chaîne de caractères est un tableau dont les éléments sont chacun un des caractères constituant la chaîne.

La chaîne de caractère est contenue dans la variable *chaîne*.

4.2 Trouvez les tous

Vous devez écrire un algorithme de fonction qui renvoie la position de chaque caractère recherché.

4.3 Comptez le nombre de lignes

Vous devez écrire un algorithme de fonction qui renvoie le nombre de caractère recherché et trouvé.

En choisissant correctement le caractère recherché, on obtient le nombre de lignes dans le fichier.

4.4 Combien de caractères lus ?

En reprenant votre algorithme précédent, vous devez utiliser une variable qui compte le nombre de fois que vous avez effectué une lecture pour trouver le caractère recherché.

Cette fonction doit renvoyer le nombre de caractères lus.

4.5 Lien entre le nombre de caractères lus et la taille de la chaîne

On se pose la question de savoir si l'algorithme est plus ou moins efficace. C'est à dire : le nombre de lectures va-t-il doubler quand la taille de la chaîne va doubler ?

Proposez une méthode qui permet de répondre à cette question en TP informatique.

5 Le cours

5.1 Objectifs de la partie cours

- Présentation de l'objectif de l'activité globale
- maîtriser les 3 structures de bases (si, pour, tant que)

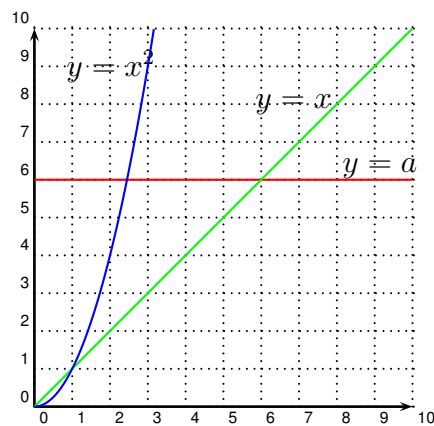
- Introduction de critères de comparaison d’algorithme (taille mémoire et/ou temps de calcul)
- Être capable reconnaître des évolutions en fonction des données à traiter.
- Être capable de choisir la (ou les) donnée(s) à analyser.
- Être capable d’analyser la complexité (théorique).
- Être capable d’analyser la complexité (pratique).
- Introduire la notion de cas (min, max, moyen).

5.2 Présentation de l’objectif de l’activité globale

Critères de comparaison d’algorithme (taille mémoire et/ou temps de calcul)

5.2.1 Courbes

- $y = a$: coût constant
- $y = x$: coût linéaire ou proportionnel
- $y = x^2$: coût quadratique



5.2.2 Partie théorique

P0

1. Combien de fois « Au revoir » va-t-il s’afficher ?
2. Combien de fois « On ne peut pas prendre la racine d’un réel négatif ! » va-t-il s’afficher ?
3. Combien de fois « La racine de ... est ... » va-t-il s’afficher ?

```

1 from math import sqrt
2 x = 2
3 if x >= 0:
4     y = sqrt(x)
5     print("La racine de {:.2f} est {:.3f}".format(x, y))
6 else:
7     print("On ne peut pas prendre la racine d'un reel negatf !")
8 print("\nAu revoir")

```

P1

1. Combien de fois « Au revoir » va-t-il s’afficher ?
2. Combien de fois « On ne peut pas prendre la racine d’un réel négatif ! » va-t-il s’afficher ?

3. Combien de fois « La racine de ... est va-t-il s'afficher ?

```
1 from math import sqrt
2 from random import *
3
4 x = float(input())
5 if x >= 0:
6     y = sqrt(x)
7     print("La racine de {:.2f} est {:.3f}".format(x, y))
8 else:
9     print("On ne peut pas prendre la racine d'un reel negatif !")
10 print("\nAu revoir")
```

P2

1. Combien de fois « Au revoir » va-t-il s'afficher ?
2. Combien de fois « On ne peut pas prendre la racine d'un réel négatif ! » va-t-il s'afficher ?
3. Combien de fois « La racine de ... est va-t-il s'afficher ?

```
1 from math import sqrt
2 from random import *
3
4 x = float(randint(-10,10))
5 if x >= 0:
6     y = sqrt(x)
7     print("La racine de {:.2f} est {:.3f}".format(x, y))
8 else:
9     print("On ne peut pas prendre la racine d'un reel negatif !")
10 print("\nAu revoir")
```

P3

1. Quelles sont les données ?
2. Combien de fois « L'élément cherche n'est pas dans la sequence. » va-t-il s'afficher ?
3. Combien de fois « Fin » va-t-il s'afficher ?

P4

1. Combien de fois « L'élément cherche n'est pas dans la sequence. » va-t-il s'afficher ?
2. Combien de fois « Fin » va-t-il s'afficher ?

```
1 from easygui import integerbox, msgbox
2
3 uneSequence = list(input())
4 print(uneSequence)
5
6 cible = integerbox("Entrez un entier :", "")
7 for i in uneSequence:
8     print(int(i))
9     if int(i) == cible:
10         break
11 else:
```

```
12     messagebox("L'element cherche n'est pas dans la sequence.")
13 messagebox("Fin")
```

P5

1. Combien de fois « facteur trouvé : » va-t-il s'afficher ?

```
1 # -*- coding: utf-8 -*-
2 from easygui import integerbox, messagebox
3 d=2
4 n=14
5 while n>1:
6     while n%d==0:
7         n = n//d
8         print("facteur trouvé:", d)
9     d += 1
```

P6

1. Combien de fois « facteur trouvé : » va-t-il s'afficher ?

```
1 from easygui import integerbox, messagebox
2 d=2
3 n=integerbox("Entrez un entier :", "")
4 while n>1:
5     while n%d==0:
6         n = n//d
7         print("facteur trouvé:", d)
8     d += 1
```

P0 bis

1. Sur le programme P0, modifier le programme afin de donner la réponse à la question :
Combien de fois « La racine de ... est va-t-il s'afficher ?

5.2.3 Partie pratique

Même question mais les élèves devront répondre aux questions en faisant des tests et devront insérer des compteurs pour avoir directement les réponses.

6 Activité : Rechercher un motif

6.1 Trouvez le premier

Un motif est une chaîne de caractères que l'on souhaite trouver dans un fichier. Vous devez écrire un algorithme de fonction qui renvoie la position du premier motif recherché. On considère qu'une chaîne de caractères est un tableau dont les éléments sont chacun un des caractères constituant la chaîne.

La chaîne de caractère est contenue dans la variable *chaine*.

6.2 Trouvez les tous

1. Vous devez écrire / compléter / analyser un algorithme de fonction qui renvoie le nombre d'occurrences du motif recherché dans un fichier.
2. Vous devez chercher à évaluer le coût en nombre d'opérations de l'exécution de la recherche.
3. Vous devez observer les variations éventuelles du coût :
 - en fonction de la taille du fichier
 - en fonction de la longueur du motif

6.3 Combien de caractères lus ?

En reprenant votre algorithme précédent, vous devez utiliser une variable qui compte le nombre de fois que vous avez effectué une lecture pour trouver le motif recherché. Cette fonction doit renvoyer le nombre de caractères lus.

6.4 Lien entre le nombre de caractères lus et la taille de la chaîne

On se pose la question de savoir si l'algorithme est plus ou moins efficace. C'est à dire : le nombre de lectures va-t-il doubler quand la taille de la chaîne va doubler ? Proposez une méthode qui permet de répondre à cette question en TP informatique.

6.5 Lien entre le nombre de caractères lus et la taille du motif

On se pose la question de savoir si l'algorithme est plus ou moins efficace. C'est à dire : le nombre de lectures va-t-il doubler quand la taille de la chaîne va doubler ? Proposez une méthode qui permet de répondre à cette question en TP informatique.

7 Exercices

7.1 Compter le nombre de caractère "e"

À vous d'écrire l'algorithme permettant de trouver le nombre de caractères "e" présent dans le fichier.

Vous pourrez essayer de coder en Python votre algorithme.

7.2 Compteur de performance

À vous d'ajouter un compteur qui s'incrémente à chaque affectation dans l'algorithme suivant.

```
1 fonction maFonction(laliste):
2     Pour i de 0 à len(laliste) - 1
3         x = laliste[i]
4         j = i
5         Tant que j > 0 et laliste[j - 1] > x
6             laliste[j] = laliste[j - 1]
7             j = j - 1
```

```
8     laliste[j] = x
9 liste = [2,6,8,1,0,5]
10 maFonction(liste)
```

7.3 Évolution de l'algorithme précédent

Comment se comporte l'algorithme précédent quand la longueur de liste s'allonge? Reprendre la méthode vue dans l'exercice 4.5.

Deuxième partie

Le travail fait

1 Activité : Rechercher un caractère

1.1 Trouvez le premier

Vous devez écrire un algorithme de fonction qui renvoie la position du premier caractère recherché.

On considère qu'une chaîne de caractères est un tableau dont les éléments sont chacun un des caractères constituant la chaîne.

La chaîne de caractère est contenue dans la variable *chaine*.

```
1 fonction rechercher(recherche, chaine):
2     Pour i variant de 0 à len(chaine):
3         si chaine[i]==recherche:
4             retourner i
5
6 rechercher('a', chaine)
```

1.2 Trouvez les tous

Vous devez écrire un algorithme de fonction qui renvoie la position de chaque caractère recherché.

```
1 fonction rechercher(recherche, chaine):
2     liste = []
3     Pour i variant de 0 à len(chaine):
4         si chaine[i]==recherche:
5             liste.ajouter(i)
6     retourner liste
7
8 rechercher('a', chaine)
```

1.3 Comptez le nombre de lignes

Vous devez écrire un algorithme de fonction qui renvoie le nombre de caractère recherché et trouvé.

En choisissant correctement le caractère recherché, on obtient le nombre de lignes dans le fichier.

```
1 fonction rechercher(recherche, chaine):
2     nombre = 0
3     Pour i variant de 0 à len(chaine):
4         si chaine[i]==recherche:
5             nombre = nombre + 1
6         retourner nombre
7
8 rechercher('\n', chaine)
```

1.4 Combien de caractères lus ?

En reprenant votre algorithme précédent, vous devez utiliser une variable qui compte le nombre de fois que vous avez effectué une lecture pour trouver le caractère recherché. Cette fonction doit renvoyer le nombre de caractères lus.

```
1 fonction rechercher(recherche, chaine):
2     nombre = 0
3     compteur = 0
4     Pour i variant de 0 à len(chaine):
5         compteur = compteur + 1
6         si chaine[i]==recherche:
7             nombre = nombre + 1
8     retourner compteur
9
10 rechercher('\n', chaine)
```

1.5 Lien entre le nombre de caractères lus et la taille de la chaîne

On se pose la question de savoir si l'algorithme est plus ou moins efficace. C'est à dire : le nombre de lectures va-t-il doubler quand la taille de la chaîne va doubler ? Proposez une méthode qui permet de répondre à cette question en TP informatique.

On espère que les élèves proposeront de compter le nombre de lecture pour différentes taille de chaîne.

Pour la mise en oeuvre, on peut penser à utiliser un fichier texte relativement long dont on ne copiera qu'une partie des données dans une chaîne pour créer des couples (taille,compteur) que l'on affiche avec *matplotlib*.

2 Le cours

P0

1. Combien de fois « Au revoir » va-t-il s'afficher ? *Une fois*
2. Combien de fois « On ne peut pas prendre la racine d'un réel négatif ! » va-t-il s'afficher ?
1
3. Combien de fois « La racine de ... est ... va-t-il s'afficher ?

0

P1

1. Combien de fois « Au revoir » va-t-il s'afficher ? *Une fois*
2. Combien de fois « On ne peut pas prendre la racine d'un réel négatif ! » va-t-il s'afficher ?
1
3. Combien de fois « La racine de ... est va-t-il s'afficher ? *0*

P2

1. Combien de fois « Au revoir » va-t-il s'afficher ? *Une fois*
2. Combien de fois « On ne peut pas prendre la racine d'un réel négatif ! » va-t-il s'afficher ?
moy=? ? min=0 max=1
3. Combien de fois « La racine de ... est va-t-il s'afficher ? *moy=? ? min=0 max=1*

P3

1. Quelles sont les données ? *Une fois*
2. Combien de fois « L'élément cherche n'est pas dans la sequence. » va-t-il s'afficher ?
moy= 0,5 min=0 max=1
3. Combien de fois « Fin » va-t-il s'afficher ? *moy= 0,5 min=0 max=1*

P4

1. Combien de fois « L'élément cherche n'est pas dans la sequence. » va-t-il s'afficher ?
Moy=? ? ? min=0 max=n avec n =nombre d'éléments dans la séquence
2. Combien de fois « Fin » va-t-il s'afficher ? *1*

P5

1. Combien de fois « facteur trouvé : » va-t-il s'afficher ? *Deux fois*

P6

1. Combien de fois « facteur trouvé : » va-t-il s'afficher ? *??? Tout dépend de n mini=0
moyen=? ? ? maxi=dépend de n*

P0 bis

1. Sur le programme P0, modifier le programme afin de donner la réponse à la question :
Combien de fois « La racine de ... est va-t-il s'afficher ? *Voir P0-1.py*

```
1 from math import sqrt
2 x = 2
3 compteur=0
4 if x >= 0:
5     y = sqrt(x)
6     print("La racine de {:.2f} est {:.3f}".format(x, y))
7 else:
8     print("On ne peut pas prendre la racine d'un reel negatf !")
9     compteur+=1
10 print("\nAu revoir")
11 print("réponse à la question ",compteur)
```

3 Activité : Rechercher un motif

3.1 Trouvez le premier

Un motif est une chaîne de caractères que l'on souhaite trouver dans un fichier. Vous devez écrire un algorithme de fonction qui renvoie la position du premier motif recherché. On considère qu'une chaîne de caractères est un tableau dont les éléments sont chacun un des caractères constituant la chaîne. La chaîne de caractère est contenue dans la variable *chaine*.

3.2 Trouvez les tous

1. Vous devez écrire / compléter / analyser un algorithme de fonction qui renvoie le nombre d'occurrences du motif recherché dans un fichier.
2. Vous devez chercher à évaluer le coût en nombre d'opérations de l'exécution de la recherche.
3. Vous devez observer les variations éventuelles du coût :
 - en fonction de la taille du fichier
 - en fonction de la longueur du motif

```
1  -- algorithme naif
2  fonction recherche(motif, file)
3  ouvrir file
4  -- (precondition) verifier file non vide
5  parcourir toutes les lettres du fichier:
6      si lettreCourante == premiere lettre du motif:
7          chercher si les lettres suivantes "valident" le motif
8          si oui :
9              occurrence +=1
10 -- (postcondition)
11 -- tous les lettres du fichier ont été parcourues
12 -- toutes les occurrences ont été trouvées...
```

```
1  fonction rechercher(recherche, chaine):
2      nombre = 0
3      Pour i variant de 0 à len(chaine):
4          si chaine[i]==recherche:
5              nombre = nombre + 1
6              retourner nombre
7
8  rechercher('e', chaine)
```

```
1  #!/usr/bin/python3
2
3  # fonction de recherche motif dans file
4  def chercheMotif(motif, file):
5
6  # lecture du fichier
7      try:
8          with open(file, "r") as f:
9              contenu=f.read()
10             assert{contenu != ''}
```

```

11 # si assertion fausse ecrire "fichier vide"
12     except AssertionError :
13         print("Le fichier est vide")
14
15 # initialisations
16     i = 0
17     j = 0
18     compteur = 0
19     cout = 0
20
21 # parcourir toutes les lettres du fichier
22     while i < len(contenu):
23         cout += 1
24         # si lettreCourante == premiere lettre du motif:
25         if contenu[i] == motif[j]:
26             # chercher si les lettres suivantes "valident" le motif
27             while j < len(motif):
28                 cout += 1
29                 if contenu[i+j] != motif[j]:
30                     j = 0
31                     break
32                 j += 1
33             if j == len(motif):
34                 compteur += 1
35                 j = 0
36         i += 1
37     return (compteur, cout)
38
39
40
41 # appel
42 chaine = input('Saisis la chaine que tu cherches : ')
43 fichier = 'access.txt'
44 compteur, cout = chercheMotif(chaine,fichier)
45 print ('nombre d\'occurences : ', compteur)
46 print ('cout : ', cout)

```

Deuxième version

```

1 #!/usr/bin/python3
2
3 # fonction de recherche de motif dans un fichier texte
4 def chercheMotif(motif, file):
5
6     # lecture du fichier
7     try:
8         with open(file,"r") as f:
9             contenu=f.read()
10            assert{contenu != ''}
11 # si assertion fausse ecrire "fichier vide"
12     except AssertionError :
13         print("Le fichier est vide")
14
15 # initialisations

```

```

16     i = 0
17     compteur = 0
18 # parcours du contenu
19     for i in range(len(contenu)-len(motif)):
20         j = 0
21         trouve = True
22         while trouve and j < len(motif):
23             if contenu[i+j] == motif[j]:
24                 j +=1
25             else:
26                 trouve = False
27         if j == len(motif):
28             compteur +=1
29     return (compteur)
30
31 # appel
32 chaine = input('Saisis la chaine que tu cherches : ')
33 fichier = 'access.txt'
34 print ('nombre d\'occurences : ', chercheMotif(chaine,fichier))

```

4 Exercices

4.1 Compter le nombre de caractère "e"

À vous d'écrire l'algorithme permettant de trouver le nombre de caractères 'à' présent dans le fichier.

Vous pourrez essayer de coder en Python votre algorithme.

4.2 Compteur de performance

À vous d'ajouter un compteur qui s'incrémente à chaque affectation dans l'algorithme suivant.

```

1 fonction maFonction(laliste):
2     Pour i de 0 à len(laliste) - 1
3         x = laliste[i]
4         j = i
5         Tant que j > 0 et laliste[j - 1] > x
6             laliste[j] = laliste[j - 1]
7             j = j - 1
8         laliste[j] = x
9 liste = [2,6,8,1,0,5]
10 maFonction(liste)

```

4.3 Évolution de l'algorithme précédent

Comment se comporte l'algorithme précédent quand la longueur de liste s'allonge? Reprendre la méthode vue dans l'exercice 4.5.

Lors du TP, on pourra fournir une fonction qui crée une liste de n entiers mélangés.

```
1 def listeLongueurVariable("nomdufichier", taille):
2     f = open(nomdufichier, 'rb')
3     chaine = f.read(taille)
4     f.close()
5     return chaine
6 machaine = listeLongueurVariable("access.log", 1000)
```

Troisième partie

Remarques, allons plus loin

Pour aller plus loin, on peut faire appel à d'autres algorithmes :

- Algorithme à l'aide d'automates finis
- Algorithme de Morris-Pratt

Un lien : <http://www.lsv.fr/haddad/coursalgorithmique.pdf>