

Épreuve d'informatique de l'X - 2010 - PT/PSI

Rechercher un mot dans un texte et compter ses occurrences

Question 1. On récupère les tailles du mot et du texte (ligne 3 et 4). Si le suffixe du texte est plus court que le mot, on retourne **faux** (ligne 5). Dans la boucle des lignes 6 à 8, on parcourt simultanément le mot et le suffixe. Si on rencontre deux lettres distinctes, on arrête le programme en retournant la valeur **faux**. Si la boucle s'exécute complètement, c'est que le mot est en tête du suffixe étudié, on retourne alors la valeur **vrai** (ligne 9).

```

1  enTeteDeSuffixe := proc(mot::array, tab::array, k::posint)
2  local m, n, i;
3  m := taille(mot);
4  n := taille(tab);
5  if k+m > n+1 then return(false) fi;
6  for i from 1 to m do
7    if tab[k+i-1]<>mot[i] then return(false) fi;
8  od;
9  return(true);
10 end;
```

Question 2. : On regarde si le mot est en tête de l'un des suffixes du texte en étudiant tous les suffixes du texte (boucle des lignes 4 à 6). Si c'est le cas, on arrête le programme en retournant la valeur **vrai** (ligne 5). Si la boucle s'exécute complètement, c'est que le mot n'est pas dans le texte, on retourne alors la valeur **faux** (ligne 7).

```

1  rechercherMot := proc(mot::array, tab::array)
2  local k, m, n;
3  m := taille(mot); n := taille(tab);
4  for k from 1 to n do
5    if enTeteDeSuffixe(mot, tab, k) then return(true) fi;
6  od;
7  return(false);
8  end;
```

Question 3 : On commence par récupérer les tailles du mot et du texte (ligne 3 et 4). Puis, pour compter les occurrences d'un mot dans un texte, on compte le nombre de fois que le mot est en tête d'un suffixe du texte en étudiant tous les suffixes du texte (boucle des lignes 6 à 8). On retourne le résultat à la ligne 9.

```

1  compterOccurrences := proc(mot::array, tab::array)
2  local k, m, n, compt;
3  m := taille(mot);
4  n := taille(tab);
5  compt := 0;
6  for k from 1 to n do
7    if enTeteDeSuffixe(mot, tab, k) then compt := compt+1 fi;
8  od;
9  return(compt)
10 end;
```

Question 4 : Pour compter les nombres d'apparitions de chacune des lettres de l'alphabet, on commence par créer et initialiser un tableau **freq** de 26 cases (lignes 3 et 4). Ensuite, on parcourt le texte et, pour chaque lettre rencontrée, on incrémente le compteur associé dans le tableau **freq** (boucles des lignes 6 à 9). À la sortie de la boucle, on retourne le tableau des fréquences (ligne 10).

```

1  frequenceLettre := proc(tab::array)
2  local freq, i, n, lettre;
3  freq := allouer(26);
4  for i from 1 to 26 do freq[i] := 0 od;
5  n:= taille(tab);
6  for i from 1 to n do
7    lettre := tab[i];
8    freq[lettre] := freq[lettre]+1;
9  od;
10 return(freq)
11 end;
```

Question 5 : Pour afficher les fréquences des bigrammes d'un texte, on étudie chacun des bigrammes du texte (ce sont les séquences de deux lettres consécutives du texte). Il y a au plus $N = 26 \times 26 = 676$ bigrammes possibles. On stocke les fréquences d'apparition de chaque bigramme dans le tableau «**freq**» dans le seul but d'éviter d'afficher plusieurs fois la fréquence d'apparition d'un même bigramme et d'éviter de recalculer la fréquence d'apparition d'un bigramme déjà étudié. La boucle des lignes 8 à 20 permet d'étudier tous les bigrammes possibles du texte. Chaque bigramme correspond à une et une seule des cases du tableau **freq**. L'indice de la case est calculé ligne 11. Si la fréquence d'apparition n'est pas nulle, c'est que le bigramme a déjà été étudié donc on ne fait rien, sinon on calcule et on affiche la fréquence d'apparition du bigramme (ligne 16 à 18).

```

1 afficherFrequenceBigramme := proc(tab::array)
2 local i, freq, N, n, bigramme, a, b, c;
3 bigramme := allouer(2);
4 N := 26*26;
5 freq := array(0..N);
6 for i from 1 to N do freq[i] := 0 od;
7 n := taille(tab);
8 for i from 1 to n-1 do
9   a := tab[i];
10  b := tab[i+1];
11  c := a+26*b;
12  bigramme[1] := a;
13  bigramme[2] := b;
14  if freq[c] = 0
15    then
16      freq[c] := compterOccurences(bigramme, tab);
17      afficherMot(tab, i, 2);
18      printf(" est un bigramme de fréquence %d \n", freq[c]);
19    fi;
20 od;
21 end;

```

II. Tableau des suffixes

Question 6 : Pour comparer deux suffixes, on parcourt simultanément les deux suffixes. Dès qu'ils ont des lettres différentes, on retourne le résultat de la différence entre ces deux lettres (ligne 10). Si on a fini de parcourir le suffixe le plus court, sans quitter la boucle, c'est que l'un des suffixes est en tête de l'autre. On retourne alors la différence des longueurs des deux suffixes (ligne 13).

```

1 maxi := proc(a, b)
2   if a > b then return(a) else return(b) fi;
3 end;
4
5 comparerSuffixes:=proc(tab::array, k1::posint, k2::posint)
6 local i, n;
7 n := taille(tab);
8 for i from 0 to n-maxi(k1, k2) do
9   if tab[k2+i]<> tab[k1+i]
10     then return(tab[k2+i] - tab[k1+i])
11   fi;
12 od;
13 return(k1-k2)
14 end;

```

Question 7 : Pour calculer le tableau des suffixes, on commence par créer le tableau **tabS** = [1, 2, 3, ..., n]. (où n est la taille du texte) (lignes 9 et 10). On effectue ensuite un tri par bulles sur le tableau **tabS** en utilisant comme relation d'ordre sur l'ensemble {1, ..., n}, la relation d'ordre induite par la fonction «**comparerSuffixes**» (double boucle des lignes 11 à 17). On retourne enfin le tableau des suffixes à la ligne 18.

```

1 permute := proc(tab :: array, i :: posint, j :: posint)
2 local tmp;
3   tmp := tab[i]; tab[i] := tab[j]; tab[j] := tmp
4 end;
5
6 calculerSuffixes := proc(tab :: array)
7 local tabS, i, j, n;
8 n := taille(tab);
9 tabS := allouer(n);
10 for i from 1 to n do tabS[i] := i od;
11 for i from 1 to n-1 do
12   for j from 1 to n-i do
13     if comparerSuffixes(tab, tabS[j], tabS[j+1]) < 0
14     then permute(tabS, j, j+1)
15     fi;
16   od;
17 od;
18 return(tabS);
19 end;

```

III. Exploitation du tableau des suffixes

Question 8 : Pour comparer un mot et un suffixe, on adapte la fonction «comparerSuffixes». On parcourt simultanément le mot et le suffixe. Dès qu'ils ont des lettres différentes, on retourne la différence entre ces deux lettres (ligne 10). Si on sort de la boucle, on arrive à la ligne 12 et alors

- soit le mot est en tête du suffixe et alors on retourne 0,
- soit le suffixe est strictement en tête du mot, et alors on retourne -1.

```

1 mini := proc(a,b)
2   if a < b then return(a) else return(b) fi;
3 end;
4
5 comparerMotSuffixe := proc(mot :: array, tab :: array, k :: posint)
6 local i, n, m;
7 n := taille(tab);
8 m := taille(mot);
9 for i from 1 to mini(m, n+1-k) do
10   if mot[i] <> tab[k+i-1] then return(mot[i] - tab[k+i-1]) fi;
11 od;
12 if m <= n+1-k then return(0) else return(-1) fi
13 end;

```

Question 9 : Pour rechercher un mot dans un texte, on procède par recherche dichotomique du mot dans le tableau des suffixes du texte en appliquant la méthode proposée par l'énoncé. Les variables a et b représentent respectivement le plus petit et le plus grand des indices du tas dans lequel on effectue la recherche. Initialement, on effectue la recherche sur tout le tableau, ce qui explique les initialisations de a et b aux lignes 3 et 4. On effectue la boucle des lignes 5 à 13 tant que $b-a \geq 0$. Dans cette boucle, on calcule l'indice m de l'élément médian du sous-tas dans lequel on effectue notre recherche. On compare le mot avec cet élément médian. S'il y a égalité, on arrête le programme et on retourne la valeur **vrai** (ligne 8). Sinon on restreint le domaine de recherche suivant la position de l'élément médian par rapport au mot cherché. Si on sort de la boucle sans avoir quitté le programme, c'est que le mot n'est pas présent dans le texte, on retourne alors la valeur **faux** (ligne 14).

```

1 rechercherMot2 := proc(mot :: array, tab :: array, tabS :: array)
2 local m, n, a, b, comp;
3 a := 1;
4 b := taille(tabS);
5 while b-a >= 0 do
6   m := floor((a+b)/2);
7   comp := comparerMotSuffixe(mot, tab, tabS[m]);
8   if comp = 0 then return(true) fi;
9   if comp > 0
10     then a := m+1
11     else b := m-1
12   fi;
13 od;
14 return(false)
15 end;

```

Question 10 : Pour la fonction «rechercherMot», l'ordre de grandeur des comparaisons est en $O(n.m)$. Pour la fonction «rechercheMot2», l'ordre de grandeur du nombre de comparaisons est en $O(\log(n).m)$.

L'intérêt du tableau des suffixes est de réduire drastiquement le temps de recherche, surtout dans le cas d'un moteur de recherche internet où la taille du «texte» est très grande en comparaison de la taille du ou des mots recherchés.

Question 11 : Pour rechercher le premier suffixe contenant un mot donné, on adapte le principe de la recherche dichotomique mise en œuvre dans la fonction «rechercherMot2». Les variables a et b représentent respectivement le plus petit et le plus grand des indices du tas dans lequel on effectue la recherche. À chaque étape, on divise par 2 la taille du tas dans lequel on effectue la recherche. Contrairement à la fonction «rechercherMot2», on ne s'arrête pas dès que le mot est un préfixe du suffixe désigné par l'élément médian m . Les instructions conditionnelles des lignes 9 à 15 gèrent les différents cas de figure. Si mot est en tête du suffixe d'indice $\text{tabS}[m]$, alors c'est que l'indice cherché est dans l'intervalle $\llbracket a, m \rrbracket$ et b reçoit m (ligne 10). Si mot est après le suffixe d'indice $\text{tabS}[m]$, alors c'est que l'indice cherché est dans l'intervalle $\llbracket m+1, b \rrbracket$ et a reçoit $m+1$ (ligne 12). Sinon, c'est que l'indice cherché est dans l'intervalle $\llbracket a, m-1 \rrbracket$ et b reçoit $m-1$ (ligne 13). Si mot est après le suffixe d'indice $\text{tabS}[m]$, alors c'est que l'indice cherché est dans l'intervalle $\llbracket m+1, b \rrbracket$ et a reçoit $m+1$ (ligne 12). On quitte la boucle des lignes 6 à 16 quand le tas dans lequel on effectue la recherche est réduit à un élément. À la sortie de la boucle, on teste pour savoir si le mot est un préfixe du suffixe correspondant. Si ce n'est pas le cas on retourne 0 (ligne 19), sinon on retourne b qui est, en sortie de boucle le plus petit indice i de tabS tel que mot apparaît en tête du suffixe numéro $\text{tabS}[i]$ du texte tab .

```

1  rechercherPremierSuffixe:=proc(mot::array, tab::array, tabS::array)
2  local i;
3  local m,n,a,b,comp;
4  a := 1;
5  b := taille(tabS);
6  while b-a > 0 do
7    m := floor((a+b)/2);
8    comp := comparerMotSuffixe(mot, tab, tabS[m]);
9    if comp = 0
10     then b:= m;
11     else if comp > 0
12         then a := m+1
13         else b := m-1;
14     fi;
15 fi;
16 od;
17 if comparerMotSuffixe(mot, tab, tabS[b]) = 0
18     then return(b)
19     else return(0)
20 fi;
21 end;
```

Question 12 : Pour compter le nombre des occurrences d'un mot dans un texte, on détermine les indices m et M des premier et dernier suffixes du tableau des suffixes du texte tabS contenant le mot (ligne 3 et 4). Si le mot est présent dans le texte, on retourne $(M - m + 1)$ qui correspond au nombre d'apparition du mot dans le texte (ligne 5). Si ce n'est pas le cas, on retourne 0.

```

1  compterOccurrences2:=proc(mot::array, tab::array, tabS::array)
2  local m, M;
3  m := rechercherPremierSuffixe(mot, tab, tabS);
4  M := rechercherDernierSuffixe(mot, tab, tabS);
5  if m <> 0 then return(M-m+1)
6     else return(0)
7  fi;
8  end;
```

Question 13 : Pour afficher les K-grammes et leurs fréquences d'apparition, on étudie tous les K-grammes possibles (boucles des lignes 5 à 16). Dans la double boucle des lignes 5 à 16, on commence par créer le $i^{\text{ème}}$ K-gramme (boucle des lignes 6 à 8). Pour éviter d'étudier plusieurs fois le même K-gramme, on vérifie que l'indice du premier suffixe contenant le K-gramme dans le tableau des suffixes tabS coïncide avec i (ligne 9). Si c'est le cas, on affiche le K-gramme et le nombre d'occurrence du K-gramme dans le texte.

```

1  afficherFrequenceKgramme:=proc(tab::array, tabS::array, k::posint)
2  local i, j, m, n, mot;
3  mot := allouer(k);
4  n := taille(tab);
5  for i from 1 to n-k+1 do
6    for j from 1 to k do
7      mot[j] := tab[i+j-1];
8    od;
9    if tabS[rechercherPremierSuffixe(mot, tab, tabS)] = i
10     then
11       afficherMot(tab, i, k);
12       printf(" est de fréquence %d\n",
13             compterOccurrences2(mot, tab, tabS))
14     fi;
15 fi;
16 od;
17 end;
```