

Reconnaissance des formes

TP9 : ARBRES DE DÉCISION

Par

Douaille Erwan & Francois Remy

Introduction

Le but de ce TP est de nous faire comprendre comment déterminer au mieux les sous ensembles à éliminer/conservé en fonction d'un indice nous permettant d'optimiser au minimum le nombre de proposition que nous devons faire pour trouver un élément parmi une liste. Concrètement avec le jeu du pendu, nous ferons une première analyse des mots dans laquelle nous choisirons quelle lettre proposer pour trouver le plus rapidement (donc avec un nombre de proposition minimal) le mot sélectionné par l'adversaire.

Question de bon sens

La valeur N vaut 16. Pour chaque essais on a 2 réponses possibles si ce n'est pas la bonne valeur (Supérieur ou Inférieur), de plus A indique qu'il faut exactement 4 propositions donc $2^4=16$.

Quand A dit qu'il faut exactement 4 propositions, cela signifie que (si l'on raisonne par dichotomie) la valeur n'est pas un nombre pair (hormis 16), si le nombre choisi est 8 par B, ce qui permet d'éliminer ces valeurs et de ne choisir que des valeurs impairs ce qui permettra de trouver plus rapidement la bonne valeur.

Variante du jeu du pendu

Question 2

La relation qui existe entre p et n est : $n = 2^p$

Question 3

La fonction :

```
1 alphabetIndex <- function(x) { strtoi(charToRaw(x),16L)-96 }
```

Permet d'indiquer la place dans l'alphabet qu'occupe chaque lettre de la chaine passée en paramètre.

```
1 ## Renvoies le nombre de fois ou apparaissent les lettres dans l'ensemble names
2 nbLettres <- function(names)
3 {
4     nbLettre = matrix(rep(0,26*length(names)),nrow=1, ncol=26);
5     for (i in 1:length(names))
6     {
7         c = alphabetIndex(names[i]);
8         nbLettre[1,c] <- nbLettre[1,c]+1;
9     }
10    nbLettre;
11 }
```

La matrice mat contient une matrice indiquant pour chaque mot si une lettre est présente ou pas. Chaque ligne de la matrice correspond à une lettre de l'alphabet, chaque colonne un mot. Par exemple, si la lettre 'a' est présente dans le mot numero 3, alors dans la case (1,3) de la matrice nous aurons la valeur 1, 0 sinon.

Question 4

Nous obtenons le tableau ci-dessous qui représente le nombre d'occurrences des lettres (1-26) présentes dans l'ensemble des noms.

1	2	3	4	5	6	7	8	9	10
150	36	86	33	184	15	41	49	124	2

11	12	13	14	15	16	17	18	19	20
5	82	53	118	140	61	12	130	53	82

21	22	23	24	25	26
107	19	1	4	13	5

Nous observons donc que la lettre *i* (9) apparait 124 fois. La lettre apparaissant le plus de fois est la lettre *e* (5), avec 184 occurrences.

Question 5

```

1 #Calcul de l'entropie
2 #log2(t^t)
3 calculEntropie <- function(indexLettre , names)
4 {
5   n = length(names);
6   nl = nbLettres(names)[1,indexLettre]
7
8   entropie <- log2((nl/n)^(-nl/n)) - log2(((n-nl)/n)^((n-nl)/n));
9   entropie;
10 }
```

L'entropie max est obtenue pour la lettre *o* (15) avec une entropie de 0.9999189.

Question 6

La première question doit être le max des entropies pour avoir la lettre qui revient le plus souvent dans l'ensemble actuel.

```
1 #premiere question
2 firstQuestion <- function(names)
3 {
4   mat = containsLetter(names);
5   entropies = matrix(rep(0,26),nrow=1,ncol=26);
6   props = matrix(rep(0,26),nrow=1,ncol=26);
7   for (i in 0:26)
8   {
9     entropies[1,i] <- calculEntropie(i,names);
10  }
11  entropies;
12 }
13 #On recuperera le max dans d'autres fonctions avec, questionTab <- firstQuestion(
    names);indiceLettre <- which.max(questionTab);
```

En utilisant cette fonction on observe que la lettre ayant la plus grande entropie sur notre liste de noms est la lettre *o*. Le nombre d'occurrences d'une lettre parmi un ensemble de mots n'est donc pas l'information à considérer dans ce genre d'exercice. Nous aurons donc intérêt à utiliser l'entropie et non le nombre d'occurrences comme facteur discriminant pour minimiser le nombre de proposition pour trouver le mot rapidement.

Question 7

```

1 #partage
2 partage <- function(names)
3 {
4   questionTab <- firstQuestion(names);
5   indiceLettre <- which.max(questionTab);
6   partition <- sum(containsLetter(names)[indiceLettre,]);
7   partition2 <- length(names) - sum(containsLetter(names)[indiceLettre,]);
8   cpt <- 1;
9   cpt2 <- 1;
10  for (i in 1 : length(names))
11  {
12    if (identical(TRUE, (containsLetter(names)[indiceLettre, i]) == 1)) {
13      partition[cpt] <- names[i];
14      cpt <- cpt + 1 ;
15    } else {
16      partition2[cpt2] <- names[i];
17      cpt2 <- cpt2 + 1 ;
18    }
19  }
20  return(list('lettre'=indiceLettre, 'partition1'=partition, 'partition2'
    '=partition2));
  }

```

Dans cette fonction nous récupérons l'indice de la lettre qui a la plus grande entropie sur l'ensemble names. Nous créons deux partitions, nous parcourons l'ensemble et si la lettre avec l'entropie max est contenue dans un mot, nous rangeons ce mot dans la partition1 sinon dans la partition2. Cela nous permet de diviser l'ensemble. Autrement dit, c'est dans cette fonction que nous pouvons binariser un noeud de notre arbre. Nous retournons trois éléments, l'indice de la lettre avec l'entropie max, la partition contenant les mots avec la lettre entropie max et la deuxième partition contenant les autres mots.

Question 8

```

1 #partageIteratif
2 partageIteratif <- function(names, cpt)
3 {
4   ensemble <- partage(names);
5   ens <- names;
6   if(length(names)>1 && (calculEntropie(ensemble$lettre, names)>0.01)) {
7     part1 <- partageIteratif(ensemble$partition1, cpt+1);
8     part2 <- partageIteratif(ensemble$partition2, cpt+1);
9     ens <- list(part1, part2);
10  }
11  return(ens);
12 }

```

Cette fonction va créer notre arbre de décision. On donne en entrée un paramètre, un ensemble de base, et la fonction va récursivement découpé cet ensemble de base en sous ensemble jusqu'à ce que toute les feuilles de l'arbres, autrement dit un sous ensemble de un mot, soit atteint. Nous nous servirons du retour de cette fonction pour dessiner l'arbre.

Question 9

```

1 #interactif
2 interactif <- function(names)
3 {
4   print("Saisir un mot a devinner dans l'ensemble donne:");
5   mot <- scan("", what="character", nlines=1);
6   print("Mode auto: (y/n)");
7   mode <- scan("", what="character", nlines=1);
8   if(mode == "y"){
9     jeuAuto(names, mot);
10  } else {
11    jeuManuel(names);
12  }
13 }
14
15 jeuManuel <- function(names)
16 {
17   if(length(names)==1){
18     print("trouve: ");
19     print(names)
20   } else {
21     ensemble <- partage(names);
22     print("Ce mot contient-il la lettre: (y/n)");
23     print(ensemble$lettre);
24     confirm <- scan("", what="character", nlines=1);
25     if(confirm == "y"){
26       jeuManuel(ensemble$partition1);
27     } else {
28       jeuManuel(ensemble$partition2);
29     }
30   }
31 }
32
33 jeuAuto <- function(names, mot)
34 {
35   if(length(names)==1){
36     print("trouve: ");
37     print(names)
38   } else {
39     ensemble <- partage(names);
40     containLettre <- containsLetter(mot);
41     if(containLettre[ensemble$lettre,]==1){
42       jeuAuto(ensemble$partition1, mot);
43     } else {
44       jeuAuto(ensemble$partition2, mot);
45     }
46   }
47 }

```

Dans cette partie de code R nous avons réaliser un déroulement de jeux, l'utilisateur saisie un mot qui devra être deviné par l'ordinateur. Le mode est une simple option, si on choisit le mode manuel l'ordinateur nous demandera de saisir si oui ou non un caractère est présent dans notre mot, autrement l'ordinateur obtiendra cette réponse en faisant appel à la méthode *containsLetter(ensemble)* par lui même.

Conclusion

Nous avons observé et mis en pratique le calcul de l'entropie. L'entropie est un indice qui nous indique l'importance, le poids de cette donnée. Exemple, nous avons constaté que la lettre *e* est la plus utilisée dans notre liste de mots mais que la lettre *o* nous permet d'obtenir une information plus grande (discriminer au mieux) par rapport aux autres lettres, son entropie se rapproche de 1 et est la valeur la plus grande parmi les entropies des autres lettres. Nous avons ensuite utilisé l'entropie pour obtenir un nombre de proposition minimum pour trouver le mot voulus de manière rapide.