

# TP1 TID

SCAIA Matteo, MARIAC Damien

October 27, 2024



# Contents

<b>0</b>	<b>Préambule</b>	<b>3</b>
<b>1</b>	<b>Exercice 1</b>	<b>3</b>
1.1	Modélisation et variable . . . . .	3
1.2	Arbre de segmentation binaire . . . . .	5
<b>2</b>	<b>Exercice 2</b>	<b>6</b>
2.1	Recodages . . . . .	6
2.1.1	Recodage en 2 variables . . . . .	6
2.1.2	Recodage en 3 variables . . . . .	7
2.2	Le meilleur recodage de X pour prédire Y . . . . .	7
<b>3</b>	<b>Exercice 3</b>	<b>8</b>
3.1	Description . . . . .	8
3.2	Calculs . . . . .	8
3.2.1	Première étape . . . . .	8
3.2.2	Deuxième étape : champignon jaune . . . . .	9
3.3	Arbre de discrimination . . . . .	10
<b>4</b>	<b>ANNEXE</b>	<b>11</b>

## 0 Préambule

Pour les calculs, nous avons choisi de faire l'exercice 1 (1) en utilisant le logarithme en base 2, et les exercices 2 et 3 (2 et 3) en logarithme népérien.

En effet, calculer avec le logarithme népérien ou avec le logarithme en base 2 revient au même. Les deux logarithmes sont les mêmes à un facteur près.

$$\log_2(x) = \frac{\ln(x)}{\ln(2)}$$

## 1 Exercice 1

On considère le tableau ci-dessous, répartissant la population active occupée selon l'âge ( $A$ ), le sexe ( $S$ ) et la catégorie socioprofessionnelle ( $C$ ) (source: INSEE, enquête emploi 2016).

Catégorie socioprofessionnelle des actifs occupés selon le sexe et l'âge						
Âge	De 15 à 29 ans	De 30 à 49 ans	De 30 à 39 ans	De 40 à 49 ans	De 50 à 59 ans	60 ans ou plus
	Effectifs (en milliers)	Effectifs (en milliers)	Effectifs (en milliers)	Effectifs (en milliers)	Effectifs (en milliers)	Effectifs (en milliers)
SEXE : Femmes						
<b>Agriculteurs</b>	27,8	189,7	70,0	119,6	187,1	76,9
<b>Artisans, con</b>	117,4	914,0	357,9	556,1	525,8	184,8
<b>Cadres et pro</b>	564,9	2 638,5	1 209,0	1 429,5	1 161,6	360,0
<b>Professions i</b>	1 353,7	3 735,7	1 840,7	1 895,0	1 507,8	256,2
<b>Employés</b>	1 570,9	3 486,4	1 605,6	1 880,9	1 819,6	397,0
<b>Ouvriers</b>	1 271,6	2 648,6	1 285,9	1 362,7	1 300,4	180,5
SEXE : Hommes						
<b>Agriculteurs</b>	24,2	146,0	56,2	89,8	138,0	43,4
<b>Artisans, con</b>	79,2	645,6	258,4	387,2	378,7	128,7
<b>Cadres et pro</b>	315,7	1 538,5	685,3	853,2	719,0	240,1
<b>Professions i</b>	613,3	1 750,4	834,7	915,7	755,9	123,7
<b>Employés</b>	476,0	865,5	449,5	416,0	329,8	58,3
<b>Ouvriers</b>	1 085,8	2 133,9	1 068,4	1 065,5	987,9	130,1

Figure 1: Tableau répartissant la population active occupée selon des catégories

### 1.1 Modélisation et variable

Tout d'abord de manière intuitive, nous avons envie de modéliser la variable socioprofessionnelle avec les deux autres. Cependant, nous devons le montrer de manière formelle. Grâce au code fourni dans la partie 4, nous calculons l'information mutuelle de chacune des variables.

Premièrement, calculons l'entropie de chacune de ces variables. Pour la variable  $A$ , nous avons le tableau suivant (en fréquence).

Table 1: Distribution par âge ( $A$ )

	15-29 ans	30-39 ans	40-49 ans	50-59 ans	60+ ans
Age	0.1866	0.2419	0.2730	0.2441	0.0542

Nous pouvons calculer l'entropie de  $A$ .

$$H(A) = - \sum_{n=1}^5 p_i \log_2(p_i) = 2,1833$$

De la même manière, nous calculons l'entropie de  $C$  et  $S$ . Nous utilisons les tableaux ci-dessous.

**Table 2:** Distribution par sexe ( $S$ )

	Femme	Homme
Proportion	0,6589	0,3411

**Table 3:** Distribution par catégorie socioprofessionnelle ( $C$ )

	Agriculteur	Artisans	Cadres	Profession In	Employes	Ouvrier
Proportion	0,0207	0,0740	0,1875	0,2512	0,2240	0,2423

Nous obtenons.

$$H(S) = 0,9258 \quad H(C) = 2,3266$$

A cette étape, nous pouvons interpréter les résultats. L'entropie d'une variable mesure l'incertitude liée à celle-ci. Donc la variable  $C$  est la variable avec l'incertitude la plus élevée.

Nous allons maintenant calculer les valeurs suivantes :  $H(A, S)$ ,  $H(A, C)$  et  $H(S, C)$ . Pour ce faire, nous déterminerons les tableaux joints correspondants pour chacune de ces paires de variables. Puis nous calculerons l'entropie de ces paires.

$$H(X) = - \sum_{x \in \text{Tableau}} P(x) \log_2(P(x))$$

Nous trouvons.

**Table 4:** Distribution jointe sexe ( $S$ ) et âge ( $A$ )

	15-29 ans	30-39 ans	40-49 ans	50-59 ans	60+
Femme	0,1220	0,1584	0,1803	0,1618	0,0362
Homme	0,0645	0,0834	0,0927	0,0823	0,1802

**Table 5:** Distribution jointe ( $C$ ) et âge ( $A$ )

	15-29 ans	30-39 ans	40-49 ans	50-59 ans	60+
Agriculteur	0,0013	0,0031	0,0052	0,0081	0,0030
Artisans	0,0049	0,0153	0,0234	0,0225	0,0080
Cadres	0,0219	0,0471	0,0568	0,0468	0,0149
Professions In	0,0489	0,0665	0,0699	0,0563	0,0094
Employes	0,0509	0,0511	0,0571	0,0535	0,0113
Ouvrier	0,0586	0,0586	0,0604	0,0569	0,0077

**Table 6:** Distribution jointe ( $C$ ) et ( $S$ )

	Agriculteur	Artisans	Cadres	Profession IN	Employes	Ouvrier
Femmes	0.0119	0.0433	0.1175	0.1705	0.1810	0.1344
Homme	0.0087	0.0307	0.0700	0.0807	0.0430	0.1079

Nous obtenons les valeurs suivantes.

$$H(A, S) = 3,1092 \quad H(A, C) = 4,4817 \quad H(C, S) = 3,2242$$

De plus, nous obtenons pour  $H(A, S, C)$  la valeur suivante.

$$H(A, S, C) = - \sum_{n=1}^{72} p_i \log_2(p_i) = 5,3778$$

Nous pouvons calculer les informations mutuelles.

$$I(C, (AS)) = H(C) + H(A, S) - H(A, S, C) = 0,0580$$

$$I(A, (SC)) = 0,029803$$

$$I(S, (CA)) = 0,029813$$

L'information mutuelle mesure la quantité d'information partagée entre deux variables aléatoires, indiquant dans quelle mesure la connaissance de l'une des variables réduit l'incertitude concernant l'autre. Avec les calculs précédents, nous observons que la connaissance de la catégorie socioprofessionnelle permet de mieux réduire l'incertitude sur l'âge et le sexe.

De plus, une information mutuelle proche de 0 suggère une indépendance.

Calculons le rapport le plus élevé entre l'information mutuelle et la variable conditionnée.

$$R_1 = \frac{I(C, (AS))}{H(A, C)} = 0,0187$$

$$R_2 = \frac{I(A, (SC))}{H(S, C)} = 0,0092$$

$$R_3 = \frac{I(S, (CA))}{H(A, C)} = 0,0066$$

Le rapport  $R_1$  est le plus élevé. Donc la variable  $C$  modélisée par les deux autres ( $A$  et  $S$ ) nous donne le plus d'information. En d'autres termes, la connaissance de la catégorie socioprofessionnelle et du couple âge, sexe apporte un gain d'information sur le couple âge, catégorie socioprofessionnelle. Et ce gain est de 1,8%.

## 1.2 Arbre de segmentation binaire

Grâce à la partie 1.1, nous pouvons calculer les informations mutuelles suivantes.

$$I(A, S) = H(A) + H(S) - H(A, S) = 5,1104 * 10^{-5}$$

$$I(A, C) = 0,02830$$

$$I(S, C) = 0,02831$$

Les résultats ci-dessus, permettent de faire les interprétations suivantes. Les variables âge et sexe sont proches de l'indépendance. Donc l'une n'influe presque pas sur l'autre. Et l'information mutuelle entre  $A$  et  $C$  est très proche de celle entre  $S$  et  $C$ .

Calculons maintenant la somme des informations mutuelles de chaque variable avec les deux autres. En pratique, nous souhaitons identifier la variable qui fournit le plus d'informations sur les autres. Cela revient à calculer :

$$\bar{I}_k = \sum_{n \neq k} I(X_n, X_k) \quad (1)$$

Ensuite, nous choisirons la variable pour laquelle  $\bar{I}_k$  est maximale.

$$\bar{I}_A = I(A, S) + I(A, C) = 0,02835$$

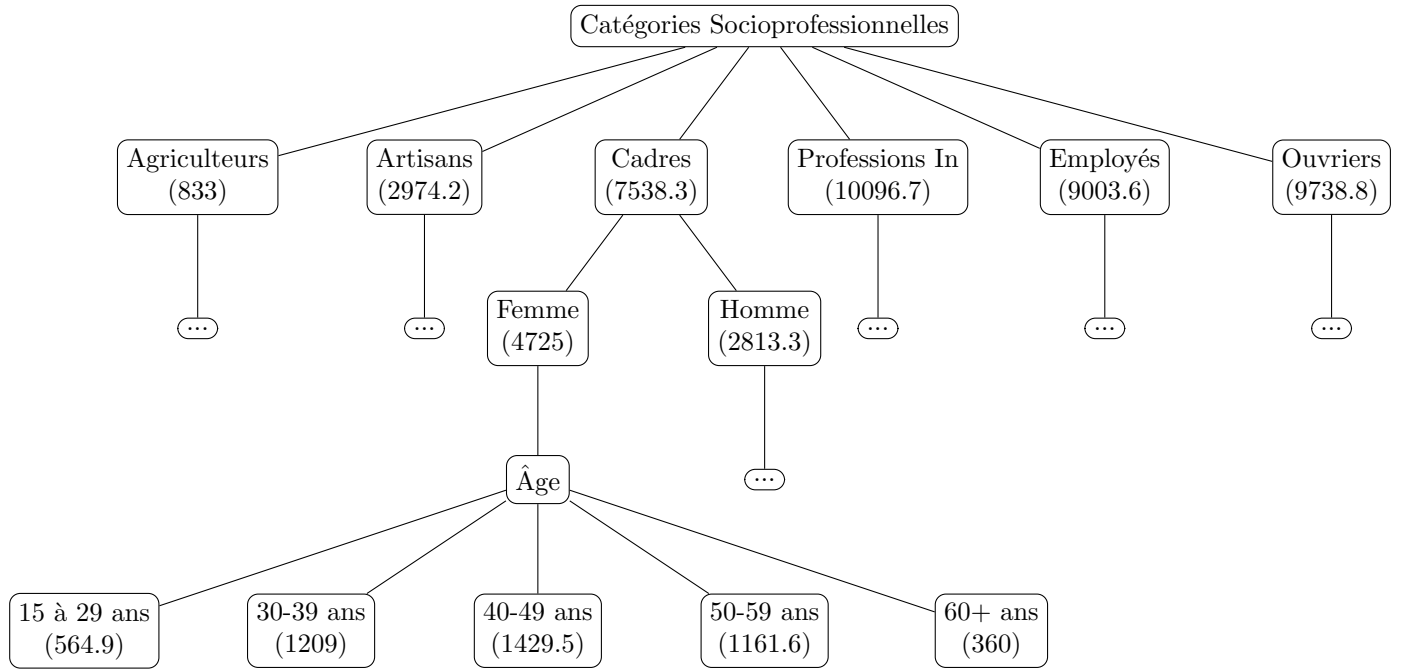
$$\bar{I}_S = 0,02836$$

$$\bar{I}_C = 0,0566$$

Nous avons donc  $\bar{I}_C$  qui est le plus élevé. Ainsi, l'arbre de segmentation commencera avec la variable  $C$ , car elle fournit le plus d'informations sur  $A$  et  $S$ .

Ce résultat était attendu, il correspond avec ce qu'on a trouvé dans la partie 1.1. Cela nous conforte dans l'idée que la variables  $C$  nous apporte plus d'informations sur les deux autres.

De plus, lors de la construction de l'arbre, le choix de la deuxième variable de segmentation n'a pas d'importance, étant donné qu'il nous reste deux variables de segmentation.



**Figure 2:** Arbre de segmentation binaire

## 2 Exercice 2

### 2.1 Recodages

#### 2.1.1 Recodage en 2 variables

En agrégeant seulement les classes contiguës, nous avons 4 possibilités de regroupement binaire de  $X$ .

$$Z_1 = \{\{0\%\}, \{0 - 0.5\%, 0.5 - 1\%, 1 - 3\%, > 3\%\}\}$$

$$Z_2 = \{\{0\%, 0 - 0.5\%\}, \{0.5 - 1\%, 1 - 3\%, > 3\%\}\}$$

$$Z_3 = \{\{0\%, 0 - 0.5\%, 0.5 - 1\%\}, \{1 - 3\%, > 3\%\}\}$$

$$Z_4 = \{\{0\%, 0 - 0.5\%, 0.5 - 1\%, 1 - 3\%\}, \{> 3\%\}\}$$

L'entropie de chacun de ses recodages se calcule numériquement.

Détaillons le premier calcul. Nous cherchons le tableau associé à  $Z_1$ , puis nous calculons  $H(Z_1)$ .

**Table 7:** Tableau entre  $Y$  et  $Z_1$

$Y \setminus Z_1$	0%	$\{0 - 0.5\%, 0.5 - 1\%, 1 - 3\%, \geq 3\%\}$
O	2	20
N	27	23

$$H(Z_1) = - \left( \frac{29}{72} \ln\left(\frac{29}{72}\right) + \frac{43}{72} \ln\left(\frac{43}{72}\right) \right) = 0.674$$

De la même manière, nous trouvons.

$$H(Z_2) = 0.661 \quad H(Z_3) = 0.427 \quad H(Z_4) = 0.073$$

Le meilleur recodage est donc le premiers c'est à dire :  $Z_1 = \{\{0\%\}, \{0 - 0.5\%, 0.5 - 1\%, 1 - 3\%, > 3\%\}\}$ . En effet  $H(Z_1)$  a l'entropie la plus élevée.

### 2.1.2 Recodage en 3 variables

En procédant de la même façon, on considère alors 6 cas :

$$\begin{aligned} Z_1 &= \{\{0\%\}, \{0 - 0.5\%\}, \{0.5 - 1\%, 1 - 3\%, > 3\%\}\} \\ Z_2 &= \{\{0\%\}, \{0 - 0.5\%, 0.5 - 1\%\}, \{1 - 3\%, > 3\%\}\} \\ Z_3 &= \{\{0\%\}, \{0 - 0.5\%, 0.5 - 1\%, 1 - 3\%\}, \{> 3\%\}\} \\ Z_4 &= \{\{0\%, 0 - 0.5\%\}, \{0.5 - 1\%, 1 - 3\%\}, \{> 3\%\}\} \\ Z_5 &= \{\{0\%, 0 - 0.5\%\}, \{0.5 - 1\%\}, \{1 - 3\%, > 3\%\}\} \\ Z_6 &= \{\{0\%, 0 - 0.5\%, 0.5 - 1\%\}, \{1 - 3\%\}, \{> 3\%\}\} \end{aligned}$$

L'entropie est calculée numériquement :

$$H(Z_1) = 1,068 \quad H(Z_2) = 1,014 \quad H(Z_3) = 0,740 \quad H(Z_4) = 0,721 \quad H(Z_5) = 0,915 \quad H(Z_6) = 0,474$$

Nous remarquons que le meilleur recodage en trois variables est  $Z_1$ .

## 2.2 Le meilleur recodage de X pour prédire Y

Il s'agit ici de recoder  $X$  en réduisant l'incertitude sur  $Y$ . Nous devons donc trouver le  $Z_k$  qui maximise l'information mutuelle entre  $Z_k$  et  $Y$ .

C'est-à-dire, trouvons le recodage  $Z_k$  qui maximise :

$$I(Z_k, Y) = H(Z_k) + H(Y) - H(Z_k, Y).$$

Détaillons le calcul pour le premier recodage :

	0%	$\neq 0\%$
OUI	2/72	20/72
NON	27/72	23/72

On a :

$$\begin{aligned} H(Z_1) &= - \left( \frac{29}{72} \ln\left(\frac{29}{72}\right) + \frac{43}{72} \ln\left(\frac{43}{72}\right) \right) \\ H(Y) &= - \left( \frac{22}{72} \ln\left(\frac{22}{72}\right) + \frac{50}{72} \ln\left(\frac{50}{72}\right) \right) \\ H(Z_1, Y) &= - \left( \frac{2}{72} \ln\left(\frac{2}{72}\right) + \frac{27}{72} \ln\left(\frac{27}{72}\right) + \frac{20}{72} \ln\left(\frac{20}{72}\right) + \frac{23}{72} \ln\left(\frac{23}{72}\right) \right) \end{aligned}$$

Nous trouvons alors:

$$I(Z_1, Y) = 0.102$$

Faisons le même cheminement, mais avec les autres variables. Nous obtenons alors les informations mutuelles suivantes.

$$I(Z_2, Y) = 0.063$$

$$I(Z_3, Y) = 0.023$$

$$I(Z_4, Y) = 0.016$$

On observe que l'information mutuelle la plus élevée est  $I(Z_1, Y)$ . Donc le recodage avec la variable  $Z_1$  est le meilleur permettant de prédire  $Y$ . En effet, une information mutuelle élevée indique que les deux variables partagent de l'information en commun. De plus lors d'un recodage de variable, choisir l'information mutuelle la plus élevée nous donne plus d'information sur  $Y$ .

### 3 Exercice 3

#### 3.1 Description

On considère le tableau recensant des espèces de champignons suivant 4 variables qualitatives.

**Table 8:** Caractéristiques des espèces de champignons

Espèce	Comestible	Chapeau	Tige	Couleur
a	o	a	e	b
b	o	a	e	j
c	o	a	e	b
d	o	pl	f	j
e	o	pl	f	b
f	n	po	f	r
g	n	po	f	j
h	n	po	e	r
i	n	a	f	j
j	n	pl	f	j

On cherche à faire un arbre de discrimination qui permet de prédire la comestibilité à partir des autres caractéristiques, tout en étant le plus court possible.

Pour ce faire, on calcule les informations mutuelles de  $X_1$  avec chaque autre variable. En effet, nous allons chercher la variable qui informe le plus sur la comestibilité. Nous effectuerons nos calculs avec le logarithme népérien.

Nous posons pour la suite.

$$X_1 = \{Comestible\} \quad X_2 = \{Chapeau\} \quad X_3 = \{Tige\} \quad X_4 = \{Couleur\}$$

#### 3.2 Calculs

##### 3.2.1 Première étape

On calcule l'information mutuelle avec la formule suivante,  $I(X_1, X_2) = H(X_1) + H(X_2) - H(X_1, X_2)$ .

**Table 9:** Croisement entre comestibilité ( $X_1$ ) et forme du chapeau ( $X_2$ )

$X_1 \backslash X_2$	Po	a	Pl
O	0	3	2
N	3	1	1

Nous trouvons grâce à ce tableau,  $I(X_1, X_2) = H(X_1) + H(X_2) - H(X_1, X_2) = 0.277$ .

Pour les variables suivantes, nous avons ces tableaux.

**Table 10:** Croisement entre comestibilité ( $X_1$ ) et la tige ( $X_3$ )

$X_1 \backslash X_3$	e	f
O	3	2
N	1	4

**Table 11:** Croisement entre comestibilité ( $X_1$ ) et la couleur ( $X_4$ )

$X_1 \backslash X_4$	b	j	r
O	3	2	0
N	0	3	2



De la même manière, nous pouvons calculer l'information mutuelle et nous trouvons :

$$I(X_1, X_3) = 0.106 \quad I(X_1, X_4) = 0.357$$

Nous pouvons conclure à cette étape que la variable qui apporte le plus d'informations sur la comestibilité du champignon est sa couleur. En effet, l'information mutuelle  $I(X_1, X_4)$  est la plus élevée.

De plus, on se rend compte dans le tableau 8 que les champignons bruns sont comestibles alors que les rouges ne le sont pas.

Traisons le cas des champignons jaunes.

### 3.2.2 Deuxieme étape : champignon jaune

En considérant seulement les champignons jaunes, on calcule les informations mutuelles suivant les autres variables ( $X_2$  et  $X_3$ ). Nous obtenons les tableaux suivants.

**Table 12:** Tableau de  $X_1, X_2$  avec champignon jaune

$X_1 \backslash X_2$	a	pl	po
O	1	1	0
N	1	1	1

**Table 13:** Tableau de  $X_1, X_3$  avec champignon jaune

$X_1 \backslash X_3$	e	f
O	1	1
N	0	2

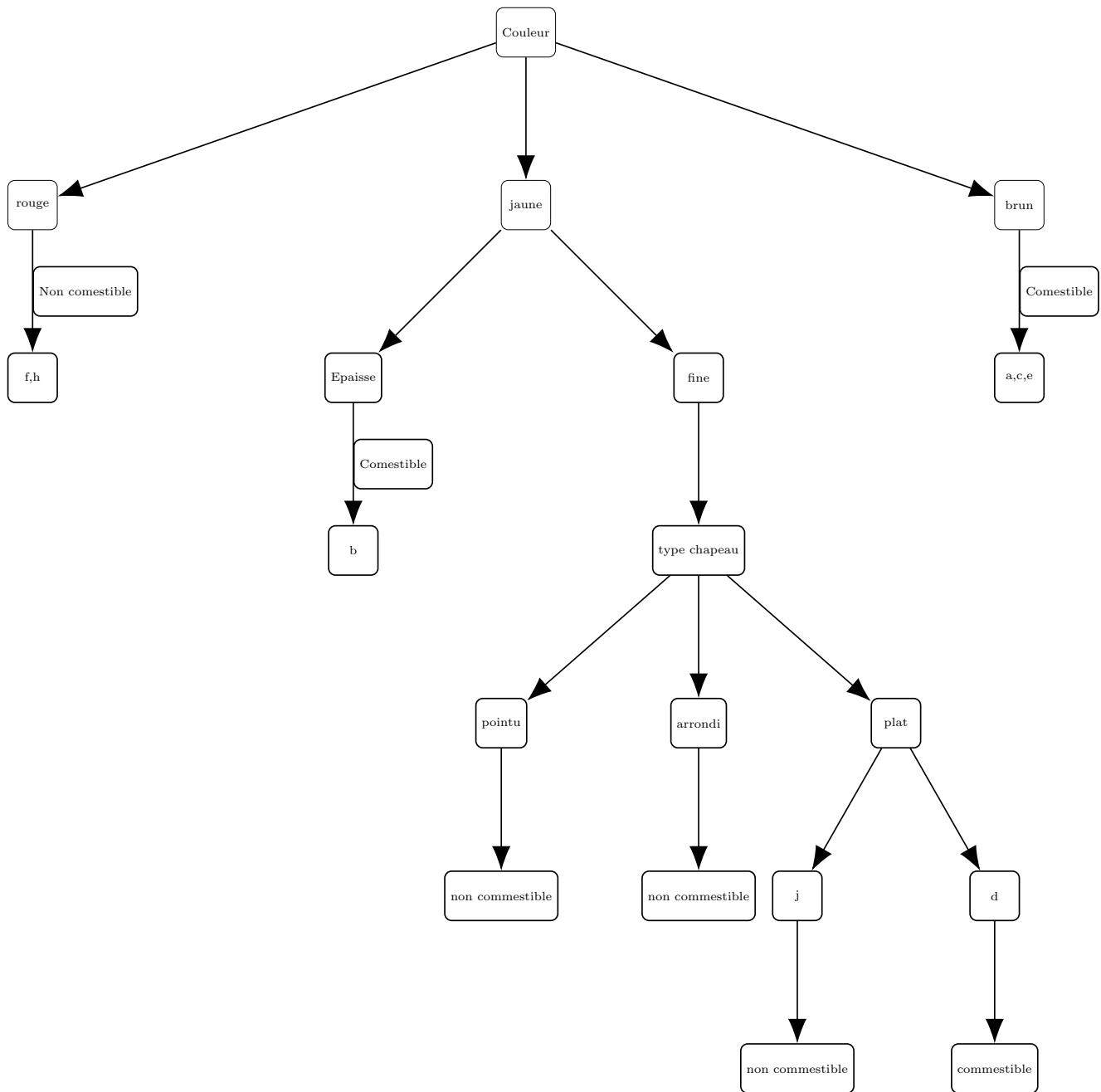
Nous calculons les informations mutuelles et nous obtenons les résultats suivants.

$$I(X_1, X_2) = 0.118 \quad I(X_1, X_3) = 0.223.$$

On considère alors la variable  $X_3$ , celle qui correspond à la tige du champignon. Si le champignon est jaune et que la tige est épaisse, alors il est comestible. Sinon il faut distinguer entre le champignon d (comestible) et les champignons g,i,j. Mais seuls les champignons d et j partagent la même forme du chapeau (plat), alors que les champignons g et i sont respectivement pointu, arrondi et tout deux non comestibles.

### 3.3 Arbre de discrimination

Avec toutes les informations mutuelles calculées, on obtient alors l'arbre suivant :



**Figure 3:** Arbre de discrimination pour la comestibilité d'un champignon

## 4 ANNEXE

```

1 # Question 1
2 # Création d'une matrice des effectifs
3 effectifs <- matrix(c(
4   27.8, 70.0, 119.6, 187.1, 76.9, # Femmes - Agriculteurs
5   117.4, 357.9, 556.1, 525.8, 184.8, # Femmes - Artisans
6   564.9, 1209.0, 1429.5, 1161.6, 360.0, # Femmes - Cadres
7   1353.7, 1840.7, 1895.0, 1507.8, 256.2, # Femmes - Professions intermédiaires
8   1570.9, 1605.6, 1880.9, 1819.6, 397.0, # Femmes - Employés
9   1271.6, 1285.9, 1362.7, 1300.4, 180.5, # Femmes - Ouvriers
10  24.2, 56.2, 89.8, 138.0, 43.4, # Hommes - Agriculteurs
11  79.2, 258.4, 387.2, 378.7, 128.7, # Hommes - Artisans
12  315.7, 685.3, 853.2, 719.0, 240.1, # Hommes - Cadres
13  613.3, 834.7, 915.7, 755.9, 123.7, # Hommes - Professions intermédiaires
14  476.0, 449.5, 416.0, 329.8, 58.3, # Hommes - Employés
15  1085.8, 1068.4, 1065.5, 987.9, 130.1 # Hommes - Ouvriers
16 ), nrow = 12, byrow = TRUE)
17
18 # Noms des colonnes (groupes d'âge)
19 colnames(effectifs) <- c("15-29", "30-39", "40-49", "50-59", "60+")
20
21 # Noms des lignes (combinaison sexe et catégorie socioprofessionnelle)
22 rownames(effectifs) <- c("Femmes_Agriculteurs", "Femmes_Artisans", "Femmes_Cadres",
23   "Femmes_Professions_inter", "Femmes_Employes", "Femmes_Ouvriers",
24   "Hommes_Agriculteurs", "Hommes_Artisans", "Hommes_Cadres",
25   "Hommes_Professions_inter", "Hommes_Employes", "Hommes_Ouvriers")
26
27
28 tableau <- effectifs/sum(effectifs) # On met notre tableau sous forme de pourcentages
29
30 #On calcule H(A)
31 age_effectif <- colSums(tableau)
32 H.A = 0
33 for(i in 1:length(age_effectif)){
34   H.A <- H.A + (-age_effectif[i]*log2(age_effectif[i]))
35 }
36 H.A <- as.numeric(H.A)
37 #On calcule H(S)
38 femmes_effectif <- sum(rowSums(tableau[1:6,]))
39 hommes_effectif <- sum(rowSums(tableau[7:12,]))
40 sexe_effectif <- c(femmes_effectif, hommes_effectif)
41 H.S = 0
42 for(i in 1:length(sexe_effectif)){
43   H.S <- H.S + (-sexe_effectif[i]*log2(sexe_effectif[i]))
44 }
45 #On calcule H(C)
46 # Calcul de la somme des effectifs pour chaque catégorie socioprofessionnelle
47 categorie_effectif <- list(
48   Agriculteurs = sum(rowSums(tableau[c(1, 7), ])), # Agriculteurs (Femmes + Hommes)
49   Artisans = sum(rowSums(tableau[c(2, 8), ])), # Artisans (Femmes + Hommes)
50   Cadres = sum(rowSums(tableau[c(3, 9), ])), # Cadres (Femmes + Hommes)
51   Professions_inter = sum(rowSums(tableau[c(4, 10), ])), # Professions intermédiaires (
52     Femmes + Hommes)
53   Employes = sum(rowSums(tableau[c(5, 11), ])), # Employés (Femmes + Hommes)
54   Ouvriers = sum(rowSums(tableau[c(6, 12), ])) # Ouvriers (Femmes + Hommes)
55 )
56 H.C = 0

```

```

56 for(i in 1:length(categorie_effectif)){
57   H.C <- H.C + (-categorie_effectif[[i]]*log2(categorie_effectif[[i]]))
58 }
59 #Nous venons de calculer les entropies de chaque variable.
60
61 # Calculons les entropies H(A,S), H(A,C) et H(S,C)
62 # Pour H(A,S)
63 age_femme <- colSums(tableau[1:6,])
64 age_homme <- colSums(tableau[7:12,])
65 age_sexe <- c(age_femme, age_homme)
66 H.A.S = 0
67 for(i in 1:length(age_sexe)){
68   H.A.S <- H.A.S + (-age_sexe[i]*log2(age_sexe[i]))
69 }
70 H.A.S <- as.numeric(H.A.S)
71 # Pour H(A,C)
72 Agriculteurs_age = colSums(tableau[c(1, 7), ])
73 Artisans_age = colSums(tableau[c(2, 8), ])
74 Cadres_age = colSums(tableau[c(3, 9), ])
75 Professions_inter_age = colSums(tableau[c(4, 10), ])
76 Employes_age = colSums(tableau[c(5, 11), ])
77 Ouvriers_age = colSums(tableau[c(6, 12), ])
78 categorie_age <- list(Agriculteurs_age, Artisans_age, Cadres_age, Professions_inter_age,
79   Employes_age, Ouvriers_age)
80 H.A.C = 0
81 for(i in 1:length(categorie_age)){
82   for(j in 1:length(categorie_age[[i]])){
83     H.A.C <- H.A.C + (-categorie_age[[i]][j]*log2(categorie_age[[i]][j]))
84   }
85 }
86 H.A.C <- as.numeric(H.A.C)
87 #Pour H(S,C)
88 femmes_catesociop <- rowSums(tableau[1:6,])
89 hommes_catesociop <- rowSums(tableau[7:12,])
90 sexe_catesociop <-list(femmes_catesociop, hommes_catesociop)
91 H.S.C = 0
92 for(i in 1:length(sexe_catesociop)){
93   for(j in 1:length(sexe_catesociop[[i]])){
94     H.S.C <- H.S.C + (-sexe_catesociop[[i]][j]*log2(sexe_catesociop[[i]][j]))
95   }
96 }
97 H.S.C <- as.numeric(H.S.C)
98 # Calculons maintenant l'entropie suivante H(A,S,C)
99 H.A.S.C = 0
100 for (i in 1:nrow(tableau)) {
101   for (j in 1:ncol(tableau)) {
102     H.A.S.C <- H.A.S.C + (-tableau[i, j] * log2(tableau[i, j]))
103   }
104 }
105 # Calculons les informations mutuelles I(A,SC) I(S,AC) I(C,AS)
106 I.A.SC = H.A + H.S.C - H.A.S.C
107 I.S.AC = H.S + H.A.C - H.A.S.C
108 I.C.AS = H.C + H.A.S - H.A.S.C
109 # On calcule les rapports
110 R1 <- I.A.SC / H.S.C
111 R2 <- I.S.AC / H.A.C
112 R3 <- I.C.AS / H.A.S
113 #Le rapport le plus eleve est R3, cela veut dire que l'environnement professionnel est
114   mieux expliqué par
115 # le sexe et l age.

```

```
116 # Question 2
117 # En reprenant les calculs plus haut.
118 I.A.S = H.A + H.S - H.A.S
119 I.A.C = H.A + H.C - H.A.C
120 I.S.C = H.S + H.C - H.S.C
121 Ib_A = I.A.S + I.A.C
122 Ib_S = I.A.S + I.S.C
123 Ib_C = I.A.C + I.S.C
```