

TD10 : Risques et monuments (DS de l'an dernier)

1 Tableaux de contingence « dilués »

On se donne un tableau d'effectifs n_{ij} ($1 \leq i \leq r, 1 \leq j \leq s$) de deux variables, ainsi que leurs marges en ligne $n_{i.}$ et en colonne $n_{.j}$. Pour un nombre réel $0 \leq \alpha \leq 1$ donné, on définit une version « diluée » des données par

$$\hat{n}_{ij} = \alpha n_{ij} + (1 - \alpha) \frac{n_{i.} n_{.j}}{n},$$

où n est l'effectif total. On voit facilement que $\hat{n}_{ij} = n_{ij}$ quand $\alpha = 1$, et qu'il est égal aux effectifs sous hypothèse d'indépendance $n_{i.} n_{.j} / n$ quand $\alpha = 0$.

Question 1 Montrer que la marge en ligne $\hat{n}_{i.}$ de \hat{n}_{ij} est égale à la marge en ligne $n_{i.}$ de n_{ij} pour tout α . Même question pour la marge en colonne $\hat{n}_{.j}$.

Question 2 Calculer les profils lignes centrés associés à \hat{n}_{ij} en fonction des profils ligne centrés associés à n_{ij} . On rappelle que le tableau des profils ligne est formé des grandeurs $n_{ij}/n_{i.}$ et que son point moyen est le profil marginal des colonnes (de coordonnées $\hat{n}_{.j}/n$).

Note ce dernier résultat permet (avec un peu de travail) de montrer que le tableau dilué conduit aux mêmes axes propres que le tableau d'origine, mais liés à des valeurs propres multipliées par α^2 .

2 ACM : risques médicaux et âge

Une compagnie d'assurance a compilé à propos de ses assurés des données sur leur taux de risque (0=normal, 1=fort) pour le système cardio-vasculaire (CVas, cœur), le système locomoteur (Loco, risque de paralysie), le système neurologique (Neuro, cerveau) et le diabète (Diab).

2.1 Les données

On obtient le tableau de Burt suivant, dans lequel une paire de données a été cachée (NA) :

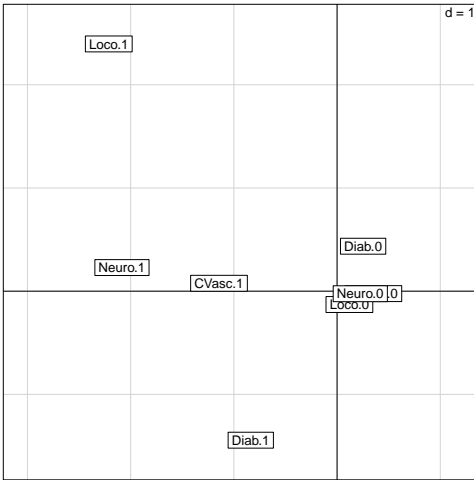
	CVasc.0	CVasc.1	Loco.0	Loco.1	Neuro.0	Neuro.1	Diab.0	Diab.1
CVasc.0	28464	0	27344	1120	26571	1893	22458	6006
CVasc.1	0	8742	7957	NA	7013	1729	6125	2617
Loco.0	27344	7957	35301	0	32186	3115	27312	7989
Loco.1	1120	NA	0	1905	1398	507	1271	634
Neuro.0	26571	7013	32186	1398	33584	0	26303	7281
Neuro.1	1893	1729	3115	507	0	3622	2280	1342
Diab.0	22458	6125	27312	1271	26303	2280	28583	0
Diab.1	6006	2617	7989	634	7281	1342	0	8623

Question 3 Expliquer comment on peut calculer les valeur manquantes du tableau de Burt, indiquées par NA, de six manières différentes.

Question 4 Les personnes ayant un risque locomoteur élevé ont-elles un risque de diabète plus grand ou plus petit que la moyenne ?

2.2 Analyse des correspondances multiples

On réalise une ACM sur les données ci-dessus. On fournit ci-dessous les variances des coordonnées des catégories, leurs poids et leurs contributions (en % pour ces deux derniers) pour tous les axes.



Variances		poids		Axis1	Axis2	Axis3	Axis4	
1	0.34	CVasc.0	19.1	CVasc.0	7.0	0.0	7.9	8.5
2	0.24	CVasc.1	5.9	CVasc.1	22.9	0.1	25.8	27.7
3	0.23	Loco.0	23.7	Loco.0	1.0	1.7	2.2	0.2
4	0.20	Loco.1	1.3	Loco.1	18.6	31.1	41.6	3.6
		Neuro.0	22.6	Neuro.0	3.4	0.1	0.5	5.8
		Neuro.1	2.4	Neuro.1	31.4	0.5	5.0	53.3
		Diab.0	19.2	Diab.0	3.6	15.4	3.9	0.2
		Diab.1	5.8	Diab.1	12.0	51.1	13.0	0.7

Question 5 Expliquez pourquoi il n'y a que 4 axes. Combien faut-il en conserver pour l'analyse ? Que peut-on dire alors de la qualité globale de la représentation ?

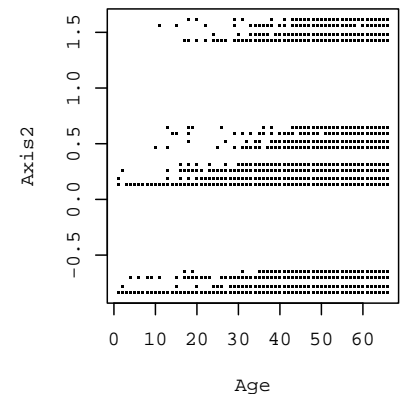
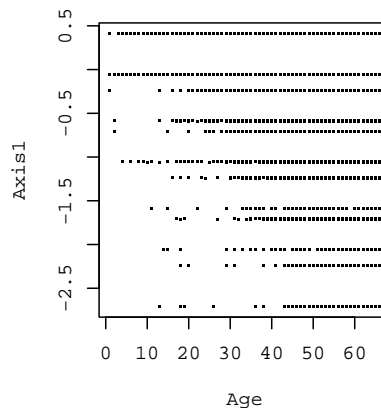
Question 6 Quelles sont les catégories qui déterminent les deux premiers axes principaux ? (on détaillera les critères et on cherchera à être précis dans la réponse).

Question 7 Calculez à partir des poids la contribution à l'inertie totale de chacune des catégories de risque fort (Xxx.1). Calculez la part d'inertie totale qu'elles représentent ensemble et montrez en quoi cela explique en partie la forme du nuage sur le premier plan principal.

2.3 Une variable supplémentaire : l'âge des assurés

On ajoute à l'analyse une nouvelle variable quantitative : l'âge des assurés. On calcule donc la corrélation de cette variable avec les deux premiers axes, que l'on donne ci-dessous accompagnée d'une représentation des couples (âge, coordonnées factorielles) pour les deux premiers axes.

	Axis1	Axis2	Axis3	Axis4
Age	-0.14	0.12	0.15	-0.09



Question 8 Expliquez pourquoi les points dans les graphiques ci-dessus sont regroupés par lignes.

Question 9 Que peut-on dire du lien entre la variable Age et les deux premiers axes ? La forme des nuages de points semble-t-elle donner des informations ?

En faisant l'hypothèse que l'âge n'a en fait pas une relation linéaire avec les axes, on regroupe les individus par tranches d'âge de façon à traiter l'âge comme une variable qualitative. On regroupe les individus en 4 groupes : moins de 19 ans (Age.0_19), de 20 à 39 ans (Age.20_39), de 40 à 59 ans (Age.40_59) et plus de 60 ans (Age.60plus). On donne ci-dessous les coordonnées des nouvelles catégories sur les axes, leur effectif et enfin la valeur test correspondante.

	Axis1	Axis2	Axis3	Axis4	suppl1.eff	Axis1	Axis2	Axis3	Axis4		
Age.0_19	0.06	-0.88	-0.63	0.14	Age.0_19	933	Age.0_19	1.93	-27.23	-19.49	4.24
Age.20_39	0.29	-0.11	-0.15	0.15	Age.20_39	5116	Age.20_39	22.27	-8.22	-11.62	11.44
Age.40_59	0.00	0.04	0.00	0.00	Age.40_59	23538	Age.40_59	0.81	10.50	0.32	0.32
Age.60plus	-0.21	0.05	0.17	-0.12	Age.60plus	7619	Age.60plus	-20.72	5.03	17.08	-11.79

Question 10 Quelles sont les catégories d'âge qui sont liées aux deux premiers axes ? On expliquera ce que sont les valeurs test et pourquoi on peut les utiliser. Quelles interprétation des axes peut-on en déduire ?

3 AFC : monuments historiques

3.1 Les données

On considère la répartition de 12387 monuments historiques en fonction de deux variables :

- leur propriétaire : COMM (municipalité), PRIV (privé), ETAT (état), DEPA (département), ETPU (établissement public) ;
- leur type : anti (antiquités), chat (châteaux), mili (architecture militaire), reli (monuments religieux), civi (architecture civile), dive (divers).

On donne ci-dessous le tableau de contingence, les profils marginaux des lignes et des colonnes (en %). Le χ^2 d'écart à l'indépendance vaut 4550.06.

	anti	chat	mili	reli	civi	dive		COMM	PRIV	ETAT	DEPA	ETPU
COMM	490	289	351	5022	563	967	62.0	29.2	5.7	1.7	1.3	
PRIV	956	964	76	426	956	242						
ETAT	161	82	59	160	138	109						
DEPA	32	58	7	49	26	40	13.4	11.6	4.0	45.9	14.1	11.0
ETPU	23	40	2	30	59	10						

Question 11 Avec une erreur inférieure à 1%, montrez que les variables type et propriétaire sont liées. On pourra s'aider de la table de χ^2 ci-dessous.

3.2 Analyse factorielle des correspondances

On réalise une analyse factorielle des correspondances sur ces données où on se limite aux deux premiers axes factoriels. On fournit ci-dessous, pour les lignes puis pour les colonnes : les coordonnées des modalités, leurs contribution aux axes (en %), la qualité de leur représentation par les deux premiers axes principaux (en % encore).

	Axis1	Axis2	Axis1	Axis2	Axis1	Axis2	Comp1	Comp2	Axis1	Axis2	Axis1	Axis2		
COMM	0.46	0.01	COMM	36.5	1.2	COMM	99.9	0.1	anti	17.4	10.5	anti	96.8	1.6
PRIV	-0.84	0.04	PRIV	58.8	4.6	PRIV	99.7	0.2	chat	26.3	16.6	chat	97.0	1.7
ETAT	-0.34	-0.38	ETAT	1.9	87.9	ETAT	44.0	55.2	mili	0.7	38.5	mili	38.5	60.8
DEPA	-0.42	-0.03	DEPA	0.9	0.2	DEPA	45.5	0.3	reli	38.4	11.7	reli	99.1	0.8
ETPU	-0.71	0.21	ETPU	1.9	6.1	ETPU	73.7	6.4	civi	15.7	0.6	civi	96.5	0.1
									dive	1.5	22.1	dive	59.9	24.8

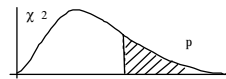
Question 12 La valeur propre associée au premier axe principal est 0.352. Montrez que la deuxième vaut à peu près 0.0096 à partir des données fournies.

Question 13 Calculez l'inertie totale ; quelle est la part d'inertie expliquée par le premier plan principal ?

Question 14 Quelles sont les modalités qui déterminent les deux premiers axes principaux ?

Question 15 Quels sont les types de monuments et de propriétaires qui sont mal représentés par le premier plan principal ?

TABLE DU CHI-DEUX : $\chi^2(n)$



n P	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.02	0.01
1	0,0158	0,0642	0,148	0,455	1,074	1,642	2,706	3,841	5,412	6,635
2	0,211	0,446	0,713	1,386	2,408	3,219	4,605	5,991	7,824	9,210
3	0,584	1,005	1,424	2,366	3,665	4,642	6,251	7,815	9,837	11,341
4	1,064	1,649	2,195	3,357	4,878	5,989	7,779	9,488	11,668	13,277
5	1,610	2,343	3,000	4,351	6,064	7,289	9,236	11,070	13,388	15,086
6	2,204	3,070	3,828	5,348	7,231	8,558	10,645	12,592	15,033	16,812
7	2,833	3,822	4,671	6,346	8,383	9,803	12,017	14,067	16,622	18,475
8	3,490	4,594	5,527	7,344	9,524	11,030	13,362	15,507	18,168	20,090
9	4,168	5,380	6,393	8,343	10,656	12,242	14,684	16,919	19,679	21,666
10	4,865	6,179	7,267	9,342	11,781	13,442	15,987	18,307	21,161	23,209
11	5,578	6,989	8,148	10,341	12,899	14,631	17,275	19,675	22,618	24,725
12	6,304	7,807	9,034	11,340	14,011	15,812	18,549	21,026	24,054	26,217
13	7,042	8,634	9,926	12,340	15,119	16,985	19,812	22,362	25,472	27,688
14	7,790	9,467	10,821	13,339	16,222	18,151	21,064	23,685	26,873	29,141
15	8,547	10,307	11,721	14,339	17,322	19,311	22,307	24,996	28,259	30,578
16	9,312	11,152	12,624	15,338	18,418	20,465	23,542	26,296	29,633	32,000
17	10,085	12,002	13,531	16,338	19,511	21,615	24,769	27,587	30,995	33,409
18	10,865	12,857	14,440	17,338	20,601	22,760	25,989	28,869	32,346	34,805
19	11,651	13,716	15,352	18,338	21,689	23,900	27,204	30,144	33,687	36,191
20	12,443	14,578	16,266	19,337	22,775	25,038	28,412	31,410	35,020	37,566
21	13,240	15,445	17,182	20,337	23,858	26,171	29,615	32,671	36,343	38,932
22	14,041	16,314	18,101	21,337	24,939	27,301	30,813	33,924	37,659	40,289
23	14,848	17,187	19,021	22,337	26,018	28,429	32,007	35,172	38,968	41,638
24	15,659	18,062	19,943	23,337	27,096	29,553	33,196	36,415	40,270	42,980
25	16,473	18,940	20,867	24,337	28,172	30,675	34,382	37,652	41,566	44,314
26	17,292	19,820	21,792	25,336	29,246	31,795	35,563	38,885	42,856	45,642
27	18,114	20,703	22,719	26,336	30,319	32,912	36,741	40,113	44,140	46,963
28	18,939	21,588	23,647	27,336	31,391	34,027	37,916	41,337	45,419	48,278
29	19,768	22,475	24,577	28,336	32,461	35,139	39,087	42,557	46,693	49,588
30	20,599	23,364	25,508	29,336	33,530	36,250	40,256	43,773	47,962	50,892

Pour $n > 30$, on peut admettre que $\sqrt{2\chi^2} - \sqrt{2n-1} = N(0,1)$