Université de Picardie Jules Verne **IAE Amiens**

2017-2018

Licence mention Gestion parcours Management et Marketing Vente - Semestre 5 Statistiques appliquées

Les tests de khi-deux

1. Conformité à un modèle théorique

Dans une population donnée, on étudie un caractère X pouvant prendre r modalités et on cherche à savoir si on peut considérer que ce caractère est d'un type donné. Plus précisément, désignant par p_i la probabilité d'apparition dans la population de la $i^{ème}$ modalité du caractère, on se demande si les p_i correspondent à une certaine loi de probabilité.

On choisit alors une loi théorique : par exemple une distribution particulière (valeurs de p_i choisies arbitrairement avec $\sum p_i = 1$) ou une loi usuelle (loi de Poisson, loi Normale, ...). Dans ce dernier cas, il faut

choisir le(s) paramètre(s) de la loi : on procède alors par estimation ponctuelle (moyenne ou variance pour le paramètre de la loi de Poisson, moyenne et écart-type pour les paramètres de la loi Normale, ...).

Effectuant plusieurs échantillonnages de même taille n, on désigne par N_i la variable aléatoire égale à l'effectif observé de la $i^{ème}$ modalité du caractère ; l'effectif théorique étant égal à np_i .

Test de H_0 : X suit la loi théorique contre H_1 : X ne suit pas la loi théorique

Ce test s'appuie sur la distance
$$D$$
 entre les effectifs observés et théoriques :
$$D = \sum_{i=1}^{r} \frac{(N_i - np_i)^2}{np_i} = \sum_{i=1}^{r} \frac{N_i^2}{np_i} - n ,$$

En pratique, pour un échantillon, on observe un effectif n_i pour la $i^{\grave{e}me}$ modalité du caractère et on calcule $d = \sum_{i=1}^r \frac{(n_i - np_i)^2}{np_i} = \sum_{i=1}^r \frac{n_i^2}{np_i} - n.$

$$d = \sum_{i=1}^{r} \frac{(n_i - np_i)^2}{np_i} = \sum_{i=1}^{r} \frac{n_i^2}{np_i} - n$$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi-deux à r-1-k degrés de liberté, où k est le nombre de paramètres à estimer de la loi théorique choisie. On détermine b_{α} tel que $P(D \ge b_{\alpha}) = \alpha$ (table 4), et on décide que :

- si $d < b_{\alpha}$, alors on ne peut rejeter H_0 ;
- si $d \ge b_{\alpha}$, alors on rejette H_0 avec une probabilité α de se tromper.

La qualité de l'approximation de la loi de D est satisfaisante lorsque les effectifs théoriques vérifient tous la condition $np_i \ge 5$. Si ce n'est pas le cas, on peut regrouper certains effectifs de modalités voisines, r désignant alors le nombre de modalités après le(s) regroupement(s). Cependant, on peut ne pas faire de regroupement si les effectifs théoriques vérifient tous la condition $np_i \ge \frac{5s}{r}$, où s est égal au nombre de modalités ayant un effectif théorique $np_i < 5$.

Exemple 1 : test de conformité à une distribution théorique

Dans une population de consommateurs, on enregistre la présence de 5 catégories, notés A_1 à A_5 , et auxquels une théorie attribue les probabilités p_1 à p_5 données dans le tableau ci-dessous.

Sur un échantillon de n = 400 individus choisis au hasard dans la population, on désigne par n_i le nombre d'individus de catégorie A_i . Les n_i sont données dans le tableau ci-dessous.

Peut-on dire, au risque $\alpha = 0.05$, que la répartition des catégories dans l'échantillon est conforme à celle de la population ?

Population : celle qui est étudiée.

Caractère : la catégorie X, à r = 5 modalités de probabilité théorique p_i .

Echantillon $(X_1, ..., X_n)$ de taille n = 400.

Les p_i étant donnés, il n'y a pas de paramètre à estimer : k = 0.

Test de H_0 : X suit la loi théorique contre H_1 : X ne suit pas la loi théorique

				•	
x_i	n_i	p_i	np_i	$n_i - np_i$	$\frac{(n_i - np_i)^2}{np_i}$
A_1	200	0,40	160	+40	10
A_2	40	0,20	80	-40	20
A_3	96	0,20	80	+16	3,2
A_4	36	0,10	40	-4	0,4
A_5	28	0,10	40	-12	3,6
	400	1	400		37,2

On calcule
$$d = \sum_{i=1}^{r} \frac{(n_i - np_i)^2}{np_i} = 37, 2.$$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi-deux à r-1-k=4 degrés de liberté.

On détermine b_{α} tel que $P(D \ge b_{\alpha}) = \alpha$: pour $\alpha = 0,05$, on trouve $b_{\alpha} = 9,49$.

Comme $d \ge b_{\alpha}$, on rejette l'hypothèse H_0 , i.e. la conformité à la loi théorique : la répartition des catégories dans l'échantillon n'est pas conforme à celle de la population. En prenant cette décision de rejet de H_0 , on a une probabilité $\alpha = 0,05$ de se tromper.

Exemple 2 : test de conformité à une loi de Poisson $\mathcal{P}(\lambda)$

Une enquête effectuée auprès du comptoir de 150 coopératives agricoles a permis d'étudier l'arrivée dans le temps des usagers de ces coopératives. Pendant l'unité de temps, soit une heure, on a obtenu les résultats suivants :

nombre d'usagers arrivés	0	1	2	3	4	5	6
nombre de coopératives	37	46	39	19	5	3	1

Peut-on admettre que le nombre d'usagers arrivés dans cette population suit une loi de Poisson ?

Population : les coopératives. Caractère : nombre d'usagers arrivés X, à r=7 modalités.

Echantillon $(X_1, ..., X_n)$ de taille n = 150 de X. On cherche à ajuster à la distribution observée une loi théorique suivie par X (i.e. les probabilités p_i des modalités de X).

Test de H_0 : X suit une loi de Poisson contre H_1 : X ne suit pas une loi de Poisson

Rappel : X suit la loi de Poisson $\mathcal{P}(\lambda)$ si X est à valeurs dans \mathbb{N} et si, pour tout $k \in \mathbb{N}$, $P(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}$. On a $E(X) = Var(X) = \lambda$.

Les p_i devant être calculés à l'aide la loi de Poisson, il y a un paramètre à estimer : k = 1.

On calcule
$$\bar{x} = \frac{1}{n} \sum_{i=1}^{7} n_i x_i \simeq 1,48, s^2 = \frac{1}{n} \left(\sum_{i=1}^{7} n_i x_i^2 \right) - (\bar{x})^2 = 1,57 \text{ et } s_c^2 = \frac{n}{n-1} s^2 \simeq 1,58.$$

Comme \bar{x} et s_c^2 sont très proches, on pouvait effectivement penser à une loi de Poisson.

On peut alors estimer le paramètre λ à 1,48.

Sous l'hypothèse
$$H_0$$
, on a alors : $p_i = P(X = i) = e^{-1.48} \frac{(1.48)^i}{i!}$.

Voir le tableau en page suivante. On a, après regroupements, r = 5.

On calcule
$$d = \sum_{i=1}^{r} \frac{(n_i - np_i)^2}{np_i} \approx 0,75.$$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi-deux à r-1-k=3 degrés de liberté (après regroupements).

On détermine b_{α} tel que $P(D \ge b_{\alpha}) = \alpha$ (table 4) : pour $\alpha = 0,05$, on trouve $b_{\alpha} = 7,81$.

Comme $d < b_{\alpha}$, on ne peut rejeter l'hypothèse H_0 , i.e. la conformité à la loi théorique de Poisson : la répartition du nombre d'usagers arrivés est conforme à une loi de Poisson. En prenant cette décision de non-rejet de H_0 , on ne connait pas la probabilité de se tromper (erreur de deuxième espèce).

x_i	n_i	p_i	np_i	$n_i - np_i$	$\frac{(n_i - np_i)^2}{np_i}$
0	37	0,2276	34,15	+2,85	0,24
1	46	0,3369	50,54	-4,54	0,41
2	39	0,2493	37,40	+1,60	0,07
3	19	0,1230	18,45	+0,55	0,02
4	5	0,0455	6,82	-1,82	
5	$ 3\rangle 9$	0,0135 > 0,0632	2,02 > 9,46	+0,98 > -0,46	0,02
6 et +	1	0,0042	0,62	+0,38	
	150	1	150	0	0,76

On a regroupé les effectifs théoriques inférieurs à 5. On a maintenant r = 5.

Exemple 3 : test de conformité à une loi Normale $\mathcal{N}(\mu, \sigma)$

Lors d'une étude biologique portant sur une certaine espèce de mollusques, on a mesuré le taux de protéines de 36 individus appartenant à cette espèce. On a obtenu les résultats suivants.

ta	aux de protéine (en mg)]0;1,5]]1,5;3]]3;4,5]]4,5;6]]6;7,5]]7,5;9]]9;10,5]
	nombre d'individus	8	7	4	9	2	3	3

Peut-on admettre que le taux de protéines dans cette population suit une loi Normale?

Population : les mollusques. Caractère : taux de protéines X, à r=7 modalités.

Echantillon $(X_1, ..., X_n)$ de taille n = 36 de X.

On cherche à ajuster à la distribution observée une loi théorique suivie par X (i.e. les probabilités p_i des modalités de X). Lorsque la représentation graphique (histogramme) est plutôt symétrique et "en cloche", on peut penser à une loi de Normale. A noter que ce n'est pas tout à fait le cas ici!

Test de $H_0: X$ suit une loi Normale contre $H_1: X$ ne suit pas une loi Normale

Les p_i devant être calculés à l'aide la loi Normale, il y a deux paramètres à estimer : k=2.

Comme dans l'exemple 1, on calcule $\bar{x} \simeq 4,21$ et $s_c \simeq 2,86$.

On peut alors estimer les paramètres μ et σ par 4,21 et 2,86. Il s'agira donc de tester si X suit la loi Normale $\mathcal{N}(4,21;2,86)$, i.e. si $U=\frac{X-4,21}{2,86}$ suit la loi Normale $\mathcal{N}(0;1)$.

Voir le tableau en page suivante. On a, après regroupements,
$$r = 5$$
. On calcule $d = \sum_{i=1}^{r} \frac{(n_i - np_i)^2}{np_i} \approx 2,99$.

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi-deux à r-1-k=2 degrés de liberté.

On détermine b_{α} tel que $P(D \ge b_{\alpha}) = \alpha$ (table 4) : pour $\alpha = 0.05$, on trouve $b_{\alpha} = 5.99$.

Comme $d < b_{\alpha}$, on ne peut rejeter l'hypothèse H_0 , i.e. la conformité à la loi théorique Normale : la répartition du taux de protéines est conforme à une loi Normale. En prenant cette décision de non-rejet de H_0 , on ne connait pas la probabilité de se tromper (erreur de deuxième espèce).

Classes de X	n_i	Classes de $U:]u_i; u_{i+1}[$	$\phi(u_i)$	$p_i = \phi(u_{i+1}) - \phi(u_i)$	np_i	$n_i - np_i$	$\frac{(n_i - np_i)^2}{np_i}$
		$-\infty$	0				
]-∞;1,5]	8]-∞;-0,95]		0,1711	6,16	+1,84	0,55
		-0,95	0,1711				
]1,5;3]	7]-0,95;-0,42]		0,1661	5,98	+1,02	0,17
		-0,42	0,3372				
]3;4,5]	4]-0,42;0,10]		0,2026	7,29	-3,29	1,49
		0,10	0,5398				
]4,5;6]	9]0,10;0,63]		0,1959	7,05	+1,95	0,54
		0,63	0,7357				
]6;7,5]	2]0,63;1,15]		0,1392	5,01	-3,01	
		1,15	0,8749				
]7,5;9]	3]1,15;1,67]		0,0776	2,79	+0,21	0,24
		1,67	0,9525				
]9;+∞[3]1,67;+∞[0,0475	1,71	+1,29	
		+∞	1				
	36			1	36		2,99

On a regroupé les effectifs inférieurs à 5, c'est-à-dire les trois dernières classes, ce qui donne :

]6;+∞[8		9,51	-1,51	0,24	
----------	--	------	-------	------	--

2. Indépendance de 2 caractères

Dans une population donnée, on étudie deux caractères X et Y pouvant prendre respectivement r et s modalités. Effectuant plusieurs échantillonnages de même taille n, on désigne par $N_{i,j}$ la variable aléatoire égale à l'effectif observé du couple formé de la $i^{\grave{e}me}$ modalité du caractère X et de la $j^{\grave{e}me}$ modalité du caractère Y. En pratique, pour un échantillon, on observe des effectifs $n_{i,j}$.

Sous l'hypothèse d'indépendance de X et Y, l'effectif théorique est égal à $np_{i,j} = \frac{n_{i,\bullet}n_{\bullet,j}}{n}$, avec $n_{i,\bullet} = \sum_{i=1}^{s} n_{i,j}$ et $n_{\bullet,j} = \sum_{i=1}^{r} n_{i,j}$.

Test de $H_0: X$ et Y sont indépendantes contre $H_1: X$ et Y ne sont pas indépendantes

Ce test s'appuie sur la distance D entre les effectifs observés et théoriques :

$$D = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{(N_{i,j} - np_{i,j})^2}{np_{i,j}} = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{N_{i,j}^2}{np_{i,j}} - n.$$
On calcule $d = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{(n_{i,j} - np_{i,j})^2}{np_{i,j}} = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{n_{i,j}^2}{np_{i,j}} - n.$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi deux à (r-1)(s-1) degrés de liberté. On détermine le réel b_α tel que $P(D \ge b_\alpha) = \alpha$ (table 4), et on décide que :

- si $d < b_{\alpha}$, alors on ne peut rejeter H_0 ;
- si $d \ge b_{\alpha}$, alors on rejette H_0 avec une probabilité α de se tromper.

La qualité de l'approximation de la loi de D est satisfaisante lorsque les effectifs théoriques vérifient tous la condition $np_{i,j} \ge 5$. Si ce n'est pas le cas, on peut effectuer des regroupements de lignes ou de colonnes : r et s désignent alors le nombre de modalités après le(s) regroupement(s). Cependant, on peut ne pas faire de regroupement si les effectifs théoriques vérifient tous la condition $np_{i,j} \ge \frac{5t}{rs}$, où t est égal au nombre de couples de modalités ayant un effectif théorique $np_{i,j} < 5$.

Cas particulier : r = s = 2.

Dans ce cas, le test d'indépendance se confond strictement avec le test (bilatéral) d'égalité de deux proportions présenté dans le chapitre précédent. En effet, d est alors le carré de u et b_{α} le carré de u_{α} .

Exemple 4 : test d'indépendance

Une statistique effectuée sur 800 personnes donne la répartition suivante :

$n_{i,j}$	gros fumeurs	moyen fum.	petits fum.	non fum.	$n_{i\bullet}$
hypertension	74	116	68	82	340
pas d'hypert.	126	174	82	78	460
$n_{\bullet j}$	200	290	150	160	800

Tester au risque 10% l'indépendance entre l'hypertension et la consommation de tabac.

Les deux caractères sont *X* : hypertension et *Y* : consommation de tabac. On a r = 2 et s = 4.

Sous l'hypothèse d'indépendance de X et Y, les effectifs théoriques sont $np_{i,j} = \frac{n_{i\bullet} n_{\bullet j}}{n}$.

$np_{i,j}$	gros fumeurs	moyen fum.	petits fum.	non fum.	
hypertension	85	123,25	63,75	68	340
pas d'hypert.	115	166,75	86,25	92	460
	200	290	150	160	800

Par exemple,
$$np_{1,2} = \frac{n_{1\bullet} n_{\bullet 2}}{n} = \frac{340 \times 290}{800} = 123,25.$$

$\frac{(n_{ij}-np_{i,j})^2}{np_{i,j}}$	gros fumeurs	moyen fum.	petits fum.	non fum.
hypertension	1,424	0,426	0,283	2,882
pas d'hypert.	1,052	0,315	0,209	2,130

Par exemple,
$$\frac{(n_{12} - np_{1,2})^2}{np_{1,2}} = \frac{(116 - 123, 25)^2}{123, 25} = 0,426.$$
On obtient : $d = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{(n_{i,j} - np_{i,j})^2}{np_{i,j}} = 8,721.$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi deux à (r-1)(s-1)=3 degrés de liberté.

On détermine le réel b_{α} tel que $P(D \ge b_{\alpha}) = \alpha$ (table 4) : pour $\alpha = 0, 10$, on trouve $b_{\alpha} = 6, 25$.

Comme $d \ge b_{\alpha}$, on rejette l'hypothèse H_0 avec une probabilité α de se tromper : on rejette donc l'indépendance des deux caractères.

Remarque.

Si on teste au risque $\alpha=0,05$, on a $b_{\alpha}=7,81$, et donc $d \geq b_{\alpha}$: même décision qu'avec $\alpha=0,10$, et on a diminué la probabilité de se tromper.

Si on teste au risque $\alpha = 0.025$, on a $b_{\alpha} = 9.35$, et donc $d < b_{\alpha}$: on ne rejette pas H_0 mais on ne connait pas la probabilité de se tromper (erreur de deuxième espèce).

Exemple 5 : test d'indépendance de deux caractères à r=2 et s=2 modalités

Dans une même catégorie sociale, un échantillon de 40 hommes a fourni 8 fumeurs et un échantillon de 60 femmes a fourni 18 fumeuses.

On se demande si la proportion de fumeurs est la même pour les deux sexes.

On a déjà traité cette question dans le précédent chapitre par un test d'homogénéité (comparaison de deux proportions).

Population 1 : hommes. Variable X_1 de loi de Bernoulli $\mathcal{B}(p_1)$, où p_1 est la proportion d'hommes fumeurs. Echantillon de taille $n_1 = 40$ de X_1 . Estimation de $p_1 : f_1 = \frac{8}{40} = 0, 2$.

Population 2 : femmes. Variable X_2 de loi de Bernoulli $\mathcal{B}(p_2)$, où p_2 est la proportion de femmes fumeuses. Echantillon de taille $n_2 = 60$ de X_2 . Estimation de $p_2 : f_2 = \frac{18}{60} = 0, 3$.

Les échantillons sont indépendants.

Test (bilatéral) de H_0 : $p_1 = p_2 = p$ contre H_1 : $p_1 \neq p_2$.

On a
$$n_1f_1 = 8 \ge 5$$
, $n_1(1-f_1) = 32 \ge 5$, $n_2f_2 = 18 \ge 5$, $n_2(1-f_2) = 42 \ge 5$.

On a $n_1 f_1 = 8 \ge 5$, $n_1 (1 - f_1) = 32 \ge 5$, $n_2 f_2 = 18 \ge 5$, $n_2 (1 - f_2) = 42 \ge 5$. Sous l'hypothèse H_0 , $U = \frac{F_1 - F_2}{\sqrt{(\frac{1}{n_1} + \frac{1}{n_2})p(1 - p)}}$ suit approximativement la loi normale $\mathcal{N}(0; 1)$, et en

regroupant les deux échantillons, on peut estimer p par $f_{1,2}=\frac{n_1f_1+n_2f_2}{n_1+n_2}=\frac{8+18}{40+60}=0,26$. En

remplaçant
$$p$$
 par $f_{1,2}$, on ne modifie pas la loi approchée de U .

On calcule $u = \frac{f_1 - f_2}{\sqrt{(\frac{1}{n_1} + \frac{1}{n_2})f_{1,2}(1 - f_{1,2})}} = \frac{0,2 - 0,3}{\sqrt{(\frac{1}{40} + \frac{1}{60})0,26(1 - 0,26)}} \simeq -1,12.$

On détermine u_{α} tel que $P(-u_{\alpha} < U < u_{\alpha}) = 1 - \alpha$, i.e. $u_{\alpha} = \phi^{-1} \left(1 - \frac{\alpha}{2}\right)$ (table 2) : pour $\alpha = 0,05$, on

Comme $u \in]-u_{\alpha}, u_{\alpha}[$, on ne peut rejeter H_0 : la proportion de fumeurs ne diffère pas significativement entre les deux sexes. Pour cette décision de non-rejet, on ne connait pas la probabilité de se tromper (erreur de deuxième espèce).

On peut également traiter cette question par un **test d'indépendance** des deux caractères X: sexe, à r=2modalités (hommes, femmes), et Y: être fumeur, à s=2 modalités (fumeur, non fumeur).

Test de $H_0: X$ et Y sont indépendantes contre $H_1: X$ et Y ne sont pas indépendantes

Ce test s'appuie sur la distance D entre les effectifs observés et théoriques : $D = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{(N_{i,j} - np_{i,j})^2}{np_{i,j}}$

Sous l'hypothèse H_0 d'indépendance de X et Y, les effectifs théoriques sont $np_{i,j} = \frac{n_{i\bullet} n_{\bullet j}}{n}$.

$n_{i,j}$	fumeurs	non fum	$n_{i\bullet}$
hommes	8	32	40
femmes	18	42	60
$n_{\bullet j}$	26	74	100

$np_{i,j}$	fumeurs	non fum	
hommes	10,4	29,6	40
femmes	15,6	44,4	60
	26	74	100

$\frac{(n_{ij}-np_{i,j})^2}{np_{i,j}}$	fumeurs	non fum	
hommes	0,55	0,19	
femmes	0,37	0,13	
			1,24

On obtient :
$$d = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{(n_{ij} - np_{i,j})^2}{np_{i,j}} = 1,24.$$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi deux à (r-1)(s-1)=1 degré de liberté.

On détermine le réel b_{α} tel que $P(D \ge b_{\alpha}) = \alpha$ (table 4) : pour $\alpha = 0,05$, on trouve $b_{\alpha} = 3,84$.

Comme $d < b_{\alpha}$, on ne rejette pas l'hypothèse H_0 d'indépendance de X et Y: on peut donc considérer que les caractères "sexe" et "être fumeur" sont indépendants, ce qui signifie que les proportions de fumeurs chez les hommes et chez les femmes ne diffèrent pas significativement. Cela correspond aux résultats du test d'homogénéité précédent. En prenant cette décision de non-rejet de H_0 , on ne connait pas la probabilité de se tromper (erreur de deuxième espèce).

3. Homogénéité : comparaison de plusieurs échantillons

Dans une population donnée, on étudie un caractère X pouvant prendre s modalités.

On dispose de *r* échantillons pouvant provenir de cette population.

Effectuant plusieurs échantillonnages, on désigne par $N_{i,j}$ la variable aléatoire égale à l'effectif observé de la $j^{\grave{e}me}$ modalité du caractère X dans le $i^{\grave{e}me}$ échantillon. En pratique, pour un échantillonnage, on observe des effectifs $n_{i,j}$.

Sous l'hypothèse d'homogénéité des échantillons, l'effectif théorique est égal à $np_{i,j} = \frac{n_{i,\bullet}n_{\bullet,j}}{n}$, avec

$$n_{i,\bullet} = \sum_{j=1}^{s} n_{i,j} \text{ et } n_{\bullet,j} = \sum_{i=1}^{r} n_{i,j}.$$

Test de H_0 : les échantillons sont issus de la même population contre $H_1 = \overline{H_0}$

Ce test se déroule comme le test d'indépendance décrit au paragraphe 2, même si le problème posé est de nature différente.

Exemple 6

Dans deux échantillons de populations de deux villes, d'effectifs respectifs 100 et 400, on demande la radio préférée. Les résultats sont les suivants :

n_{ij}	R_1	R_2	R_3	R_4
e_1	10	30	50	10
e_2	60	120	180	40

Les deux populations présentent-elles les mêmes proportions de radio préférée ?

Autrement dit, les deux populations sont-elle identiques en termes de répartition de radio préférée ? Ce qui nous amène à tester si les deux échantillons proviennent de la même population.

On a r = 2 échantillons et s = 4 modalités pour le caractère radio préférée.

$n_{i,j}$	R_1	R_2	R_3	R_4	$n_{i\bullet}$
e_1	10	30	50	10	100
e_2	60	120	180	40	400
$n_{\bullet j}$	70	150	230	50	500

$np_{i,j}$	R_1	R_2	R_3	R_4	$n_{i\bullet}$
e_1	14	30	46	10	100
e_2	56	120	184	40	400
$n_{\bullet j}$	70	150	230	50	500

$\frac{(n_{i,j}-np_{i,j})^2}{np_{i,j}}$	R_1	R_2	R_3	R_4
e_1	1,14	0	0,35	0
e_2	0,29	0	0,09	0

On obtient :
$$d = \sum_{i=1}^{r} \sum_{j=1}^{s} \frac{(n_{i,j} - np_{i,j})^2}{np_{i,j}} = 1,87.$$

On sait que sous l'hypothèse H_0 , D suit approximativement la loi de khi deux à (r-1)(s-1)=3 degrés de liberté. On détermine le réel b_α tel que $P(D \ge b_\alpha)=\alpha$ (table 4) : pour $\alpha=0,05$, on trouve $b_\alpha=7,81$.

Comme $d < b_{\alpha}$, on ne peut rejeter l'hypothèse H_0 : les deux échantillons proviennent de la même population.

4. Exercices

Sauf mention explicite, les tests seront réalisés au risque 5%.

Exercice 1. (d'après examen de janvier 2016)

1) On a demandé à 160 étudiant(e)s de l'UPJV d'estimer le temps mensuel en heures qu'ils passent à cuisiner. On a obtenu les résultats suivants :

Heures	[0;5]]5,10]]10,15]]15,30]
Nombre d'étudiants	62	49	19	30

- a) Représenter graphiquement cette série statistique.
- b) Déterminer la médiane de cette série statistique et interpréter le résultat obtenu.

2) Des études antérieures sur l'ensemble de la population française ont permis d'établir la répartition suivante :

Heures	[0;5]]5,10]]10,15]]15,30]
Pourcentage	40%	35%	15%	10%

Effectuer un test statistique au risque 5% pour répondre à la question suivante : s'agissant du temps passé à cuisiner, l'échantillon d'étudiant(e)s de l'UPJV est-il représentatif de la population française ?

Exercice 2.

Dans une usine de production d'un laboratoire pharmaceutique, on a dénombré pendant deux mois, soit 50 jours d'activité, le nombre de pannes quotidiennes. On a consigné les résultats dans le tableau suivant :

x_i	0	1	2	3	4 et +
n_i	21	18	7	3	1

où n_i est le nombre de jours où l'on a observé x_i pannes.

- 1) Calculer la moyenne et la variance de cette distribution.
- 2) Tester l'ajustement à cette distribution d'une loi de Poisson.

Exercice 3.

A l'aide d'un programme informatique, on simule 100 lancers d'un dé à 6 faces numérotées de 1 à 6 On obtient les résultats suivants :

Face	1	2	3	4	5	6
Nombre de lancers	17	22	18	14	13	16

- 1) Préciser la(les) population(s) et le(s) caractère(s) étudié(s), ainsi que la(les) taille(s) d'échantillon.
- 2) Peut-on considérer, au risque 5%, que la simulation est bien celle d'un dé équilibré ?

Exercice 4. (d'après examen de janvier 2016)

Une chaîne d'agences immobilières a fixé un objectif de vente à ses agents de 1,5 biens en moyenne par mois. Pour savoir si cet objectif est atteint, elle a étudié le nombre de biens vendus par agent et par mois.

Elle a observé 50 agents pendant un mois dans la moitié nord de la France et obtenu la répartition suivante

Nombre de biens vendus	0	1	2	3	4	5
Nombre d'agents	14	18	10	5	2	1

- 1) Préciser la population et le caractère étudiés, la taille d'échantillon, le(s) estimateur(s) mis en jeu et leur loi.
 - 2) Représenter graphiquement cette série statistique.
- 3) Donner une estimation ponctuelle de la moyenne et de la variance du nombre de bien vendus par agent et par mois.
- 4) Effectuer un test statistique au risque 5% pour répondre à la question suivante : peut-on considérer que le nombre de biens vendus par agent par mois suit une loi de Poisson ? Présenter le détail des calculs permettant d'effectuer ce test.

Exercice 5.

A la suite du même traitement, on a observé 40 bons résultats chez 70 malades jeunes et 50 bons résultats chez 100 malades agés.

Peut-on dire qu'il y a indépendance entre l'âge du malade et l'effet du traitement ?

Exercice 6.

En novembre 2004, beaucoup d'étudiants ont déclaré être stressé par les changements consécutifs à la mise en place du LMD. C'est pourquoi l'Université leur a proposé de suivre un stage de relaxation proposant plusieurs méthodes différentes : méthode "be cool", méthode "be aware", méthode "be zen". A l'issue du stage, on leur a demandé comment ils se sentaient : moins, autant ou plus stressé qu'avant le stage. On a obtenu la répartition suivante des étudiants :

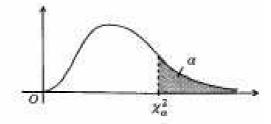
	moins	autant	plus
be cool	30	15	15
be aware 10		5	15
be zen	15	10	35

Effectuer un test statistique adéquat pour répondre à la question suivante : peut-on considérer que la méthode de relaxation choisie a une influence sur le niveau de stress après le stage ?

TABLE 4

Lois de Pearson ou lois du χ^2

Si Y^2 est une variable aléatoire qui suit la loi du χ^2 à v degrés de liberté, la table donne, pour α choisi, le nombre χ^2_s tel que $P(Y^2 \geqslant \chi^2_s) = \alpha$.



/x	0.99	0,975	0,95	0,90	0,10	0,05	0,025	10,0	0.001
1	0,000 2	0,001	0.004	0,016	2,71	3,84	5,02	6,63	10,83
2	0,02	0.05	0,10	0,21	4,61	5,99	7,38	9,21	13.82
3	0,12	0,22	0,35	0,58	6,25	7,81	9,35	11,34	16,27
4	0,30	0,48	0.71	1,06	7,78	9,49	11,14	13,28	18,47
5	0,55	0,83	1,15	1,61	9.24	11,07	12,83	15,09	20,52
6	0.87	1,24	1,64	2,20	10,64	12,59	14,45	16,81	22,46
7	1.24	1,69	2,17	2,83	12,02	14,07	16,01	18,47	24,32
8	1,65	2.18	2.73	3,49	13,36	15,51	17,53	20,09	26,13
9	2,09	2,70	3.33	4.17	14,68	16,92	19,02	21,67	27,88
10	2,56	3,25	3,94	4,87	15,99	18,31	20,48	23,21	29,59
11	3.05	3.82	4,57	5,58	17,27	19,67	21,92	24,72	31,26
12	3,57	4.40	5.23	6,30	18,55	21.03	23,34	26,22	32,91
1.3	4.11	5.01	5,89	7,04	19,81	22.36	24,74	27,69	34,53
14	4.66	5,63	6,57	7,79	21,06	23,68	26,12	29,14	36,12
15	5.23	6.26	7.26	8,55	22.31	25,00	27,49	30,58	37,70
16	5.81	6.91	7.96	9,31	23,54	26,30	28,84	32,00	39,25
17	6.41	7,56	8,67	10,08	24,77	27,59	30,19	33,41	40,79
18	7.01	8,23	9,39	10,86	25.99	28,87	31,53	34,80	42,31
19	7.63	8,91	10.12	11,65	27,20	30,14	32,85	36,19	43,82
20	8,26	9,59	10,85	12,44	28.41	31,41	34,17	37,57	45,32
21	8.90	10,28	11.59	13.24	29,61	32,67	35,48	38,93	46,80
22	9.54	10.98	12.34	14.04	30,81	33,92	36,78	40.29	48,27
22 23	10.20	11.69	13.09	14.85	32.01	35,17	38,08	41,64	49,73
24	10.86	12.40	13.85	15,66	33.20	36,41	39,37	42,98	51,18
25	11,52	13.12	14.61	16,47	34,38	37,65	40.65	44,31	52,62
26	12.20	13.84	15.38	17,29	35,56	38,88	41,92	45.64	54,05
27	12,88	14,57	16,15	18,11	36,74	40,11	43,19	46,96	55,48
28	13,57	15.31	16,93	18,94	37,92	41,34	44,46	48,28	56,89
29	14.26	16.05	17,71	19,77	39,09	42,56	45,72	49,59	58,30
30	14.95	16,79	18,49	20,60	40.26	43,77	46,98	50,89	59,70

Lorsque le degré de liberté v est tel que v > 30, la variable aléatoire :

$$U = \sqrt{2Y^2} - \sqrt{2v - 1}$$

suit à peu près la loi normale réduite.