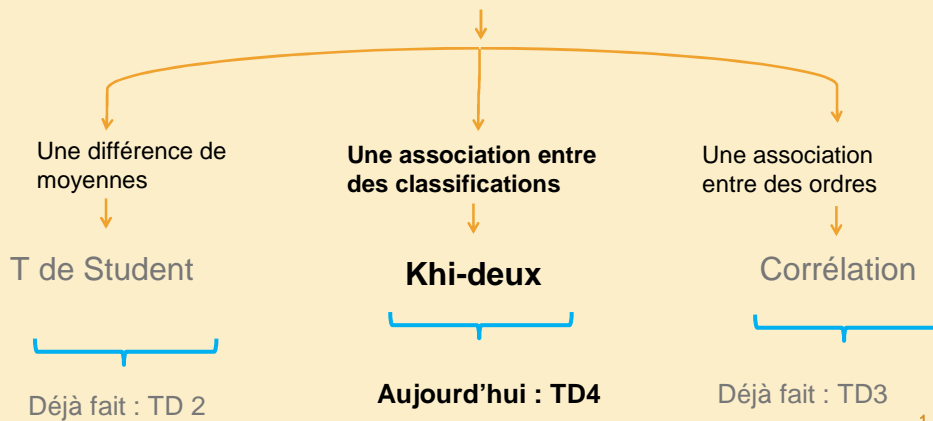


6. TD4 – Khi-deux

Les tests d'inférence statistiques permettent d'estimer le risque d'inférer un résultat d'un échantillon à une population et de décider si on « prend le risque » (si ≤ 0.05 ou 5 %)



6. TD4 – Khi-deux

6.1 Définition

- ✗ Ici, on travaille à partir de fréquences ou d'effectifs issus de variables nominales ou catégorielles.

Deux types de X^2 (= khi/chi deux/(au) carré).

✗ X^2 de conformité

- + Une seule variable
- + Permet de savoir si la distribution est conforme à une distribution qu'on connaît (par exemple, celle qu'on aurait pu obtenir par hasard).

✗ X^2 d'indépendance

- + Mesure le lien entre 2 variables nominales (ex: pays de naissance et couleur des yeux)
- + Permet, à partir de l'observation de l'échantillon, de décider si 2 variables nominales sont indépendantes ou non et de déterminer le risque qu'on prend de généraliser à tort de l'échantillon à la population.

6. TD4 – Khi-deux

6.1 Définition

- × Principe : le khi-deux permet de décider si une série d'**effectifs observés** diffère significativement d'une série d'**effectifs théoriques** (ou **attendus**)
- × La distribution (ensemble des effectifs) théorique est celle qui correspond à l'**hypothèse nulle**.

Hypothèse nulle (H_0) : notion centrale des tests d'inférence statistique.

C'est l'hypothèse selon laquelle il n'y a **pas de lien entre les variables**.

C'est à partir de cette hypothèse qu'est calculée la probabilité p .

Hypothèse alternative (H_1): Hypothèse selon laquelle il y a un **lien entre les variables**

3

6. TD4 – Khi-deux

6.1 Définition

6.2 Khi-deux de conformité

- × Enquête menée sur 110 sujets
Parmi ces 4 viennoiseries, laquelle préférez-vous?

	Croissant	Pain au chocolat	Pain aux raisins	Brioche
Effectifs observés	20	60	10	20
Effectifs théoriques/attendus	27,5	27,5	27,5	27,5

Posons H_0 : Toutes les viennoiseries sont aussi populaires les unes que les autres.

Si l'hypothèse nulle est vérifiée, on devrait avoir à peu près le même nbre de personnes / catégorie, soit: $110/4 = 27,5$ (**Effectifs théoriques**)

⇒ Le test de khi-deux va comparer les effectifs observés et les effectifs théoriques.

⇒ Si les 2 types d'effectifs sont proches, l'hypothèse nulle est validée, s'ils sont éloignés, alors l'hypothèse nulle est rejetée.

4

6. TD4 – Khi-deux
6.1 Définition
6.2 Khi-deux de conformité

- × 1/ On soustrait les effectifs théoriques aux effectifs observés

Effectifs observés	Effectifs théoriques	Effectifs observés - Effectifs théoriques
20	27.5	-7.5
60	27.5	32.5
10	27.5	-17.5
20	27.5	-7.5

- × 2/ On élève au carré les différences

$$(-7.5)^2 = 56.25$$

$$32.5^2 = 1056.25$$

$$(-17.5)^2 = 306.25$$

$$(-7.5)^2 = 56.25$$

5

6. TD4 – Khi-deux
6.1 Définition
6.2 Khi-deux de conformité

- × 3/ On divise ces différences au carré par la fréquence attendue (effectif théorique)

$$56.25/27.5 = 2.05$$

$$1056.25/27.5 = 38.41$$

$$306.25/27.5 = 11.14$$

$$56.25/27.5 = 2.05$$

- × 4/ On en fait la somme

$$2.05+38.41+11.14+2.05 = 53.65$$

$$X^2 = 53.65$$

- × 5/ Calcul du degré de liberté (ddl) [nombre de catégories -1]

$$\text{Ici, } 4 \text{ (viennoiseries) } - 1 = 3$$

$$\text{ddl} = 3$$

6

6. TD4 – Khi-deux
 6.1 Définition
 6.2 Khi-deux de conformité

✘ 6/ On se reporte à une table de Khi-Deux pour avoir la valeur de p

p	0.999	0.995	0.99	0.98	0.95	0.9	0.8	0.2	0.1	0.05	0.02	0.01	0.005	0.001
ddl														
1	0.0000	0.0000	0.0002	0.0005	0.0009	0.0158	0.0642	1.6424	2.7055	3.8415	5.4119	6.6349	7.8794	10.8276
2	0.0020	0.0100	0.0201	0.0404	0.1026	0.2107	0.4453	3.2189	4.6052	5.9915	7.3780	9.2103	10.5966	13.8155
3	0.0243	0.0717	0.1148	0.1848	0.3518	0.5844	1.0052	4.6416	6.2514	7.8147	9.3478	11.3449	12.8382	16.2662
4	0.0908	0.2070	0.2971	0.4294	0.7107	1.0636	1.6488	5.9886	7.7794	9.4877	11.6678	13.2767	14.8603	18.4668
5	0.2102	0.4117	0.5543	0.7513	1.1455	1.6103	2.3425	7.2893	9.2364	11.0705	13.3882	15.0863	16.7496	20.5150
6	0.3811	0.6757	0.8721	1.1344	1.6354	2.2041	3.0701	8.5581	10.6446	12.5916	15.0332	16.8119	18.5476	22.4577
7	0.5985	0.9893	1.2390	1.5643	2.1673	2.8331	3.8223	9.8032	12.0170	14.0671	16.6224	18.4753	20.2777	24.3219
8	0.8571	1.3444	1.6465	2.0325	2.7326	3.4695	4.5936	11.0301	13.3616	15.5073	18.1682	20.0902	21.9550	26.1245

$\chi^2 = 53.65$; ddl = 3

Dans la ligne de ddl = 3, on cherche dans la table une valeur de chi-deux immédiatement inférieure à celle qu'on a calculée (53.65) et on en déduit que le p associé au khi-deux calculé est inférieur à celui qui figure dans la table.

Valeur de p ? $p < 0.001$

Que peut-on conclure ? Globalement, les personnes interrogées n'aiment pas toutes les viennoiseries de la même façon : préférence pour les pains au chocolat

6. TD4 – Khi-deux
 6.1 Définition
 6.2 Khi-deux de conformité

Remarque: parallèle entre...

- ✘ Le khi-deux de conformité compare la distribution d'une **variable nominale** à une valeur fixe (ici la valeur de hasard, déterminée à partir du nombre de catégories de la variable)
- ✘ Test t univarié (cf. TD Comparaison de moyennes) compare la distribution d'une **variable quantitative** à une valeur fixe (qui peut aussi être la valeur de hasard)

6. TD4 – Khi-deux

6.1 Définition

6.2 Khi-deux de conformité

Mini-TD

On veut caractériser des écoles par leur degré de mixité sociale

école mixte (milieux sociaux équilibrés); **école non mixte** (les familles d'un certain milieu – Fav ou Défav.- sont plus nombreuses).

Pour chaque élève de classes de CM de chaque école, on prend en considération la profession des parents et chaque élève est ensuite affecté à un milieu social

4 milieux sociaux: Défavorisé; Intermédiaire -; Intermédiaire +; favorisé

Nbre d'élèves (effectifs)

École	Défav	Inter-	Inter+	Fav	Tot.
du Bois	6 13,9%	6 13,9%	8 18,6%	23 53,40%	43
Mille Chemins	7 17,5%	12 30%	12 30%	9 22,5%	40
Village	7 14%	19 38%	17 34%	7 14%	50
Portail Rouge	17 45,9%	6 16,2%	6 16,2%	8 21,6%	37
Renaudie+ Vercors	21 80,8%	2 7,7%	2 7,7%	1 3,8%	26
Total	58	45	45	48	196

Tableau 6 : Profil social des écoles

Buson, L. (2009). *Variation stylistique entre 5 et 11 ans et réseaux de socialisation scolaire: usages, représentations, acquisition et prise en compte éducative*. Thèse de doctorat, Université Stendhal, Grenoble.

6. TD4 – Khi-deux

6.1 Définition

6.2 Khi-deux de conformité

École	Défav	Inter-	Inter+	Fav	Tot.
du Bois	6 13,9%	6 13,9%	8 18,6%	23 53,40%	43
Mille Chemins	7 17,5%	12 30%	12 30%	9 22,5%	40
Village	7 14%	19 38%	17 34%	7 14%	50
Portail Rouge	17 45,9%	6 16,2%	6 16,2%	8 21,6%	37
Renaudie+ Vercors	21 80,8%	2 7,7%	2 7,7%	1 3,8%	26
Total	58	45	45	48	196

Tableau 6 : Profil social des écoles

Mini-TD

- L'auteur a pris une répartition équilibrée comme effectif théorique pour calculer le khi-deux (43/4 pour l'école du Bois).
- Recalculer les khi-deux et trouver les p associés dans la table distribuée (5 groupes d'étudiants, une école par groupe)
- Au vu des khi-deux, quelles écoles peuvent être qualifiées de socialement mixtes ou non mixtes ?

6. TD4 – Khi-deux
 6.1 Définition
 6.2 Khi-deux de conformité

43/4

Calcul du χ^2 pour l'école du Bois

Effectifs observés	Effectifs théoriques	Effectifs observés - Effectifs théoriques
6 (defav)	10.75	-4.75
6 (inter-)	10.75	-4.75
8 (inter+)	10.75	-2.75
23 (fav)	10.75	12.25

(Effectifs observés - Effectifs théoriques) ²	(Effectifs observés - Effectifs théoriques) ² / fréq. attendue	χ^2 (somme de la colonne précédente)
22.56	2.10	18.86
22.56	2.10	
7.56	0.70	
150.06	13.96	

p	0.999	0.995	0.99	0.98	0.95	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1	0.05	0.02	0.01	0.005	0.001
1	0.0000	0.0000	0.0002	0.0006	0.0039	0.0158	0.0642	0.1624	0.2705	0.3845	0.5119	0.6349	0.7894	1.0276	1.3501	1.7537	2.3029	2.8781	3.8415
2	0.0020	0.0100	0.0201	0.0404	0.1026	0.2107	0.4463	0.7189	1.0724	1.4861	1.9672	2.4454	3.0008	3.5795	4.1917	4.6052	5.0240	5.4119	5.9915
3	0.0243	0.0717	0.1148	0.1848	0.3128	0.5044	0.7682	1.0092	1.3288	1.7139	2.1637	2.6768	3.2528	3.8915	4.5937	5.2691	6.0258	6.7794	7.8794
4	0.0908	0.2070	0.2971	0.4294	0.7107	1.0636	1.5788	2.1425	2.7425	3.3882	4.0777	4.8082	5.5777	6.3882	7.2446	8.1446	9.0882	10.0682	11.1446
5	0.2102	0.4117	0.5543	0.7519	1.1455	1.6103	2.1425	2.7425	3.3882	4.0777	4.8082	5.5777	6.3882	7.2446	8.1446	9.0882	10.0682	11.1446	12.2662
6	0.3811	0.6757	0.8721	1.1344	1.6354	2.2041	2.8701	3.5381	4.2082	4.8777	5.5477	6.2177	6.8877	7.5577	8.2277	8.8977	9.5677	10.2377	10.9077
7	0.5985	0.9893	1.2390	1.5643	2.1673	2.8331	3.5023	4.1723	4.8423	5.5123	6.1823	6.8523	7.5223	8.1923	8.8623	9.5323	10.2023	10.8723	11.5423
8	0.8571	1.3444	1.6463	2.0323	2.7326	3.4893	4.2936	5.0479	5.8022	6.5565	7.3108	8.0651	8.8194	9.5737	10.3280	11.0823	11.8366	12.5909	13.3452

	École du Bois	École Mille Chemins	École Village	École Portail Rouge	Écoles Renaudie + Vercors
Khi-deux	18,86	1,80	9,84	8,95	43,23
p < 0.001		Non significatif	p < 0.02	p < 0.05	p < 0.001
degré de liberté	3	3	3	3	3

On cherche dans la table une valeur de khi-deux immédiatement inférieure à celle qu'on a calculée et on en déduit que le p associé au khi-deux calculé est inférieur à celui qui figure dans la table.

6. TD4 – Khi-deux

6.1 Définition

6.2 Khi-deux de conformité

École	Défav		Inter-		Inter+		Fav		Chi 2	p	conclusion
du Bois	6	13,9%	6	13,9%	8	18,6%	23	53,40%	18,86	<0.001	Non mixte favorisée
Mille Chemins	7	17,5%	12	30%	12	30%	9	22,5%	1,80	ns	mixte
Village	7	14%	19	38%	17	34%	7	14%	9,84	< 0.02	Non mixte intermédiaire
Portail Rouge	17	45,9%	6	16,2%	6	16,2%	8	21,6%	8,95	< 0.05	Non mixte défavorisé
Renaudie+Vercors	21	80,8%	2	7,7%	2	7,7%	1	3,8%	43,23	<0.001	Non mixte défavorisé
Total	58		45		45		48				

Tableau 6 : Profil social des écoles

Au sens strict, une seule école – Mille Chemins - est vraiment mixte. Le khi-deux est NS et tous les milieux sociaux y sont équitablement représentés.

Dans les 4 autres écoles, un milieu social est plus représenté : favorisé (Bois), intermédiaire (Village) ou défavorisé (Portail Rouge et Renaudie+Vercors).

NB – Au total, l'auteure trouve que les enfants de familles défavorisées scolarisées au contact d'enfants d'autres milieux ont plus de flexibilité stylistique : ils adaptent mieux leur façon de parler à l'interlocuteur.

Buson, L. (2009). *Variation stylistique entre 5 et 11 ans et réseaux de socialisation scolaire: usages, représentations, acquisition et prise en compte éducative*. Thèse de doctorat, Université Stendhal, Grenoble. 13

EXAMEN

× **Lundi 9 janvier à 10H30 à 12H30 en Amphi 4**

× 1/3 temps supplémentaire: RDV à 09H50 en C213

× Rappel:

+ **Documents papiers autorisés**

× **Mais PAS d'ordinateur ou téléphone portable**

+ Supports de cours mis en ligne sur *Alfresco Share*

<https://espaces-collaboratifs.grenet.fr/share/page/site-index>

!!! La semaine prochaine !!!

× Lundi 12 décembre

+ Cours avec T. Mout en F111

(8h30-10h30: groupe1; 10h30-12h30: groupe 2)

+ Séance sur l'examen de l'an passé (envoyé par mail pendant les vacances de la Toussaint)

6. TD4 – Khi-deux

6.1 Définition

6.2 Khi-deux de conformité

6.2 Khi-deux d'indépendance

- × Mesure le lien entre 2 variables nominales (ex: pays de naissance et couleur des yeux)
- × Permet, à partir de l'observation de l'échantillon, de décider si 2 variables nominales sont indépendantes ou non.
- × On veut savoir s'il y a un lien entre le fait de boire et de fumer chez les étudiants
 - + Enquête : On demande à 110 étudiants :
 - > fumez-vous régulièrement (plus d'une cigarette / semaine) ?
 - > buvez-vous de l'alcool régulièrement (au moins deux verres/ semaine) ?

Quatre catégories de réponses

Combien de personnes boivent et fument ?

Combien de personnes boivent mais ne fument pas ?

Combien de personnes ne boivent pas mais fument ?

Combien de personnes ne boivent pas et ne fument pas ?

- 6. TD4 – Khi-deux
- 6.1 Définition
- 6.2 Khi-deux de conformité
- 6.2 Khi-deux d'indépendance

1/ On résume les données observées dans une table 2X2 (2 lignes et 2 colonnes)

2/ On calcule les totaux des lignes et des colonnes

	Fume	Ne fume pas	Total
Boit	50	15	65
Ne boit pas	20	25	45
Total	70	40	110

La 1ère étape du test consiste à déterminer combien on peut s'attendre à trouver d'étudiants dans chaque catégorie en supposant qu'il n'y a aucune relation entre les 2 variables (= hypothèse nulle) → Calcul des effectifs théoriques (méthode assez similaire à celle utilisée par le khi-deux de conformité SAUF qu'on ne s'attend pas à avoir 4 cellules identiques. Les effectifs théoriques ne sont donc pas tous égaux).

Logique

Si pas de lien entre FUME et BOIT, la proportion de buveurs et de non buveurs devrait être à peu près la même chez les fumeurs et les non fumeurs : 65/110 et 45/110

- 6. TD4 – Khi-deux
- 6.1 Définition
- 6.2 Khi-deux de conformité
- 6.2 Khi-deux d'indépendance

3/ Calcul des effectifs théoriques:

On se demande ce qu'on aurait dans chaque case si la proportion restait la même que dans les totaux marginaux : combien y aurait-il de buveurs parmi les 70 fumeurs si il y en avait la même proportion que 65/110 ?

On multiplie le total de chaque ligne par le total de chaque colonne en divisant par le total général (110).

	Fume	Ne fume pas	Total	
Effectifs observés	Boit	50	15	65
	Ne boit pas	20	25	45
	Total	70	40	110

Pour FUME et BOIT : $(70 \cdot 65) / 110 = 41.4$

Pour NE FUME PAS et BOIT: $(40 \cdot 65) / 110 = 23.6$

	Fume	Ne fume pas	Total	
Effectifs théoriques	Boit	41.4	23.6	65
	Ne boit pas	28.6	16.4	45
	Total	70	40	110

- 6. TD4 – Khi-deux
- 6.1 Définition
- 6.2 Khi-deux de conformité
- 6.2 Khi-deux d'indépendance

4/ On soustrait les effectifs théoriques aux effectifs observés

Effectifs observés - Effectifs théoriques

	Fume	Ne fume pas
Boit	50 - 41.4 = 8.6	15 - 23.6 = -8.6
Ne boit pas	20 - 28.6 = -8.6	25 - 16.4 = 8.6

- * 5/ On élève au carré la valeur absolue des différences
- * 6/ On divise ces nombres par l'effectif théorique de chaque cellule

	Fume	Ne fume pas
Boit	8.6²/41.4 1.8	8.6²/23.6 3.1
Ne boit pas	8.6²/28.6 2.6	8.6²/16.4 4.5

* 7/ On additionne tous ces nombres

* $1.8 + 3.1 + 2.6 + 4.5 = 12$
KHI-DEUX = 12

* On associe une valeur de p à ce khi-deux de 12: on regarde dans la table Khi-deux au bon ddl

!!! Dans le cas de khi-deux à plus d'une variable :
ddl = (nbre de lignes -1) x (nbre de colonnes -1)
Ici 2 lignes et 2 colonnes, donc ddl = (2-1) X (2-1) = 1

p	0.999	0.995	0.99	0.98	0.95	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1	0.05	0.02	0.01	0.005	0.001
1	0.0000	0.0000	0.0002	0.0008	0.0039	0.0158	0.0642	0.16424	0.27055	0.38415	0.54119	0.6349	0.7794	0.8276	0.9052	0.95996	0.9849	0.99501	0.99999
2	0.0020	0.0100	0.0201	0.0404	0.1026	0.2107	0.4463	0.7189	1.0652	1.3815	1.6759	1.9431	2.1788	2.3381	2.4453	2.5191	2.5658	2.5964	2.6149
3	0.0243	0.0717	0.1433	0.2843	0.5518	1.0644	1.9052	2.9189	4.1652	5.3815	6.5411	7.6349	8.6794	9.6759	10.6253	11.5351	12.4019	13.2299	14.0149
4	0.0908	0.2070	0.3971	0.6743	1.2126	2.1788	3.5815	5.4119	7.7794	10.6253	13.9351	17.7249	21.9431	26.5964	31.5658	36.7815	42.2599	47.9964	53.9849
5	0.2102	0.4117	0.7543	1.3453	2.3425	3.8223	5.7881	8.2831	11.3381	14.9652	19.1759	23.9431	29.2599	35.1253	41.5351	48.4815	55.9599	63.9664	72.4949
6	0.3811	0.6757	1.1344	1.9354	3.2041	4.9701	7.2831	10.1446	13.5918	17.6332	22.2619	27.4815	33.2899	39.6753	46.6415	54.1819	62.2999	70.9864	80.2349
7	0.5985	0.9893	1.7390	2.9643	4.8331	7.3223	10.4831	14.2619	18.6759	23.6253	29.1159	35.1431	41.7099	48.8153	56.4515	64.6199	73.3164	82.5349	92.2649
8	0.8571	1.3444	2.3325	3.9326	6.1695	9.1431	13.0301	17.7419	23.2831	29.6515	36.7431	44.5599	53.0999	62.3615	72.3315	82.9999	94.3664	106.4349	119.1949

Avec un ddl de 1, le khi-deux calculé de 12 est supérieur à 10.8276. On en conclut donc que le p associé est inférieur à 0.001.
La probabilité d'obtenir un khi-deux aussi élevé que 12 dans l'échantillon en supposant que les 2 variables FUME et BOIT sont indépendantes dans la population (hypothèse nulle) est inférieure à 1 chance sur 1000.

Conclusion :
Chez les étudiants interrogés, il y a un lien entre le fait de boire et le fait de fumer

6. TD4 – Khi-deux
 6.1 Définition
 6.2 Khi-deux de conformité
 6.2 Khi-deux d'indépendance

Mini-TD

Version simplifiée du fichier sur les erreurs de liaison recueillies chez une fillette entre 2 ans 1 mois et 6 ans 4 mois

Mot2	Erreurs en N	Erreurs en Z	Orientation singulier / pluriel
Ami	1	11	pluriel
Arbre	10	7	pluriel
Enfant	11	14	pluriel
Escalier	6	0	pluriel
Habit	9	10	pluriel
Oiseau	23	32	pluriel
Ane	27	1	singulier
Anorak	6	0	singulier
Arc-en-ciel	11	0	singulier
Avion	12	1	singulier
Eléphant	15	0	singulier
Orage	11	2	singulier
Ours	20	0	singulier

Étape 1 : construisez le tableau à double entrée ci-dessous

	Nombre total d'erreurs /n/	Nombre total d'erreurs /z/
Noms orientés pluriel		
Noms orientés singulier		

Étape 2 : Calculez le khi-deux d'indépendance et concluez

21

Mot2	Erreurs en N	Erreurs en Z	Orientation singulier / pluriel
Ami	1	11	pluriel
Arbre	10	7	pluriel
Enfant	11	14	pluriel
Escalier	6	0	pluriel
Habit	9	10	pluriel
Oiseau	23	32	pluriel
Ane	27	1	singulier
Anorak	6	0	singulier
Arc-en-ciel	11	0	singulier
Avion	12	1	singulier
Eléphant	15	0	singulier
Orage	11	2	singulier
Ours	20	0	singulier

Étape 1 : construisez le tableau à double entrée ci-dessous

	Nombre total d'erreurs /n/	Nombre total d'erreurs /z/
Noms orientés pluriel	60	74
Noms orientés singulier	102	4

22

Etape 2 : Calculez le khi-deux d'indépendance et concluez

	Nombre total d'erreurs /n/	Nombre total d'erreurs /z/	
Noms orientés pluriel	60 $ET_1 = (162 \times 134) / 240 = 90,45$ $D_1 = (EO_1 - ET_1)^2 / ET_1$ $D_1 = 10,250995$	74 $ET_2 = (78 \times 134) / 240 = 43,55$ $D_2 = (EO_2 - ET_2)^2 / ET_2$ $D_2 = 21,2905281$	134
Noms orientés singulier	102 $ET_3 = (162 \times 106) / 240 = 71,55$ $D_3 = (EO_3 - ET_3)^2 / ET_3$ $D_3 = 12,958805$	4 $ET_4 = (78 \times 106) / 240 = 34,45$ $D_4 = (EO_4 - ET_4)^2 / ET_4$ $D_4 = 26,9144412$	106
	162	78	240

$X^2 = D_1 + D_2 + D_3 + D_4$

Donc $X^2 \approx 71,41$

DDL = 1

$p < 0.001$

Conclusion

Les erreurs en /n/ et /z/ sont plus ou moins fréquentes selon que les noms sont orientés au pluriel et au singulier.

- Que se passe-t-il exactement pour chaque type de noms ?