

# PrivateTeacher

## *Cours Privés de Science*

### Chi2 et test d'indépendance Exercice d'entraînement

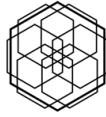
Julien RUPPEN

15 October, 2023

#### **Abstract**

Ce cours est une introduction au Chi2 et aux tests d'indépendance. Tout le matériel nécessaire est organisé de manière à ce qu'on puisse l'aborder avec comme seul prérequis des notions d'algèbre. Ce document s'adresse donc aux étudiants.es de première et deuxième année de bachelor de l'Université et des hautes écoles. Il contient le matériel nécessaire pour pratiquer des exercices type et sera donc utile aussi bien aux étudiants.es en sciences sociales et politiques, en psychologie, HEC ou encore aux étudiants.es de médecine.





Enoncé

## Enoncé

Une société de conseil est mandatée pour étudier une population. Il est possible de classer les individus en plusieurs groupes différents et on cherche à savoir si chacun de ces groupes sont semblables ou pas. La population est donc classée en différentes catégories que nous désignerons par Groupe-A Groupe-B, groupe-C, etc.

La société de conseil choisi un certain nombres de caractéristiques afin de comparer les groupes. Une fois ces critères choisi, elle se rend sur place parmi les individus, pour prendre des mesures et faire des observations. Nous désignerons chacune des valeurs observées par: Obs-1 Obs-2 Obs-3 etc.

Voici les données obtenue:

##		Obs-1	Obs-2
##	Groupe-A	50	36
##	Groupe-B	96	79
##	Groupe-C	103	137

L'objectif de l'étude est de déterminer si les valeurs observées permettent d'identifier le groupe dans lequel elles ont été prises ou si au contraire ces valeurs sont les même pour l'ensemble de la population.

## Réponse

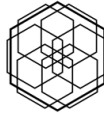
```
## [1] "La population est hétérogène, les deux groupes sont différents"  
## [1] "La variable groupe et les variables numériques sont dépendantes"  
## [1] "On peut expliquer l'appartenance à un groupe à l'aide des observations"
```

## Explication

Des données homogènes signifient que les observations ne sont pas significativement différentes d'un groupe à l'autre. Elles se ressemblent toutes, elles sont semblables et c'est pourquoi elles sont dites: "homogène" Lorsque les données sont homogènes, cela signifie que la variable "groupe" et les variables "observation" ne sont pas en relation. Ces deux variables sont donc indépendantes.

Au contraire si les données sont différentes d'un groupe à l'autre, la population est dite "hétérogène" Cela signifie que les groupe sont différents et qu'il est possible de les identifier à partir des observations. On peut alors expliquer l'appartenance à un groupe à l'aide des observation. La variable groupe et les variables observation sont dépendantes l'une de l'autre.





## Détail des calculs

Pour évaluer la force du lien entre le groupe et les observations, on calcul la valeur du Chi2:

### Étape 1: Calculer les effectifs théoriques.

```
##           Obs-1  Obs-2
## Groupe-A  42.743  43.257
## Groupe-B  86.976  88.024
## Groupe-C 119.281 120.719
```

### Étape 2: Calculer l'écart entre les valeurs attendues et les valeurs observées.

On calcule maintenant la valeur suivante pour chaque case du tableau:

$$\frac{(O_i - E_i)^2}{E_i}$$

```
##           Obs-1  Obs-2
## Groupe-A  1.232  1.218
## Groupe-B  0.936  0.925
## Groupe-C  2.222  2.196
```

### Étape 3: Additionner les valeurs.

La somme de chacun de ces nombre constitue la statistique Chi2

$$\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

```
## [1] "Chi2 calculé = 8.729"
```

### Situer cette valeur au sein de sa distribution

On commence par déterminer quel est le degrés de liberté de la distribution. Elle est donnée par le produit entre nombre de colonne -1 et nombre de ligne -1 de la table de contingence.

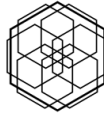
$$df = (n_{lgne} - 1) \times (n_{col} - 1)$$

```
## [1] "df = 2"
```

On choisit ensuite un seuil de tolérance alpha. Conventionnellement, on choisit alpha = 0.05. La table de distribution Chi2 enfin, nous donne la valeur suivante:

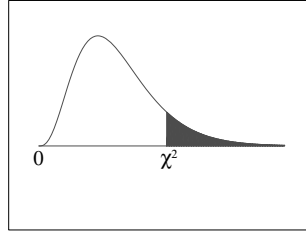
```
## [1] "Chi2 critique = 5.991"
```





# Table des quantiles

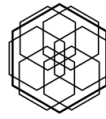
## Chi-Square Distribution Table



The shaded area is equal to  $\alpha$  for  $\chi^2 = \chi^2_{\alpha}$ .

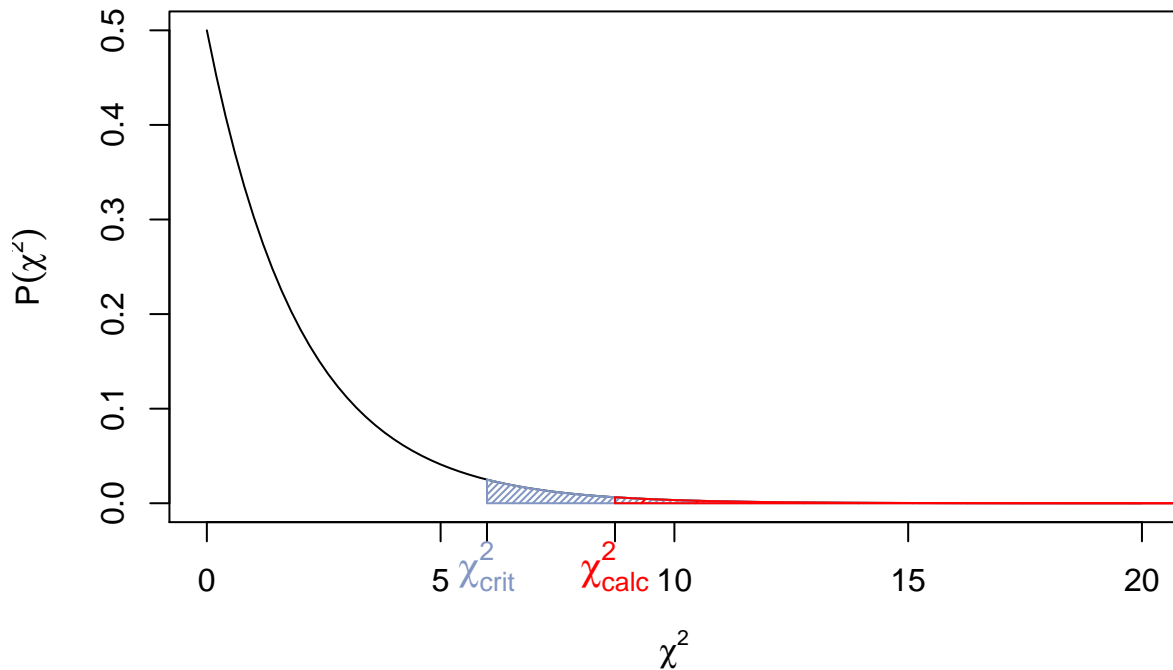
$df$	$\chi^2_{.995}$	$\chi^2_{.990}$	$\chi^2_{.975}$	$\chi^2_{.950}$	$\chi^2_{.900}$	$\chi^2_{.100}$	$\chi^2_{.050}$	$\chi^2_{.025}$	$\chi^2_{.010}$	$\chi^2_{.005}$
1	0.000	0.000	0.001	0.004	0.016	2.706	3.841	5.024	6.635	7.879
2	0.010	0.020	0.051	0.103	0.211	4.605	5.991	7.378	9.210	10.597
3	0.072	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.345	12.838
4	0.207	0.297	0.484	0.711	1.064	7.779	9.488	11.143	13.277	14.860
5	0.412	0.554	0.831	1.145	1.610	9.236	11.070	12.833	15.086	16.750
6	0.676	0.872	1.237	1.635	2.204	10.645	12.592	14.449	16.812	18.548
7	0.989	1.239	1.690	2.167	2.833	12.017	14.067	16.013	18.475	20.278
8	1.344	1.646	2.180	2.733	3.490	13.362	15.507	17.535	20.090	21.955
9	1.735	2.088	2.700	3.325	4.168	14.684	16.919	19.023	21.666	23.589
10	2.156	2.558	3.247	3.940	4.865	15.987	18.307	20.483	23.209	25.188
11	2.603	3.053	3.816	4.575	5.578	17.275	19.675	21.920	24.725	26.757
12	3.074	3.571	4.404	5.226	6.304	18.549	21.026	23.337	26.217	28.300
13	3.565	4.107	5.009	5.892	7.042	19.812	22.362	24.736	27.688	29.819
14	4.075	4.660	5.629	6.571	7.790	21.064	23.685	26.119	29.141	31.319
15	4.601	5.229	6.262	7.261	8.547	22.307	24.996	27.488	30.578	32.801
16	5.142	5.812	6.908	7.962	9.312	23.542	26.296	28.845	32.000	34.267
17	5.697	6.408	7.564	8.672	10.085	24.769	27.587	30.191	33.409	35.718
18	6.265	7.015	8.231	9.390	10.865	25.989	28.869	31.526	34.805	37.156
19	6.844	7.633	8.907	10.117	11.651	27.204	30.144	32.852	36.191	38.582
20	7.434	8.260	9.591	10.851	12.443	28.412	31.410	34.170	37.566	39.997
21	8.034	8.897	10.283	11.591	13.240	29.615	32.671	35.479	38.932	41.401
22	8.643	9.542	10.982	12.338	14.041	30.813	33.924	36.781	40.289	42.796
23	9.260	10.196	11.689	13.091	14.848	32.007	35.172	38.076	41.638	44.181
24	9.886	10.856	12.401	13.848	15.659	33.196	36.415	39.364	42.980	45.559
25	10.520	11.524	13.120	14.611	16.473	34.382	37.652	40.646	44.314	46.928
26	11.160	12.198	13.844	15.379	17.292	35.563	38.885	41.923	45.642	48.290
27	11.808	12.879	14.573	16.151	18.114	36.741	40.113	43.195	46.963	49.645
28	12.461	13.565	15.308	16.928	18.939	37.916	41.337	44.461	48.278	50.993
29	13.121	14.256	16.047	17.708	19.768	39.087	42.557	45.722	49.588	52.336
30	13.787	14.953	16.791	18.493	20.599	40.256	43.773	46.979	50.892	53.672
40	20.707	22.164	24.433	26.509	29.051	51.805	55.758	59.342	63.691	66.766
50	27.991	29.707	32.357	34.764	37.689	63.167	67.505	71.420	76.154	79.490
60	35.534	37.485	40.482	43.188	46.459	74.397	79.082	83.298	88.379	91.952
70	43.275	45.442	48.758	51.739	55.329	85.527	90.531	95.023	100.425	104.215
80	51.172	53.540	57.153	60.391	64.278	96.578	101.879	106.629	112.329	116.321
90	59.196	61.754	65.647	69.126	73.291	107.565	113.145	118.136	124.116	128.299
100	67.328	70.065	74.222	77.929	82.358	118.498	124.342	129.561	135.807	140.169





## Représentation graphique

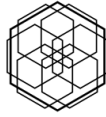
### Distribution de probabilité des valeur de Chi2 pour DF = 2



La statistique du Chi2 nous l'avons vu, est une mesure de l'écart entre les valeurs espérées et les valeurs observées. Plus cet écart est important, plus la différence entre les groupes est statistiquement significative.

```
## [1] "Ici la valeur de chi2 est au delà de la valeur critique."  
## [1] "Cela nous indique que les groupes sont différents."  
## [1] "La population est donc hétérogène."  
## [1] "La variable groupe et les variables observations sont dépendantes."
```





*Avec le logiciel R*

## **Avec le logiciel R**

On arrive au même résultat à l'aide des commandes suivantes:

### **1) entrer les données sous forme de matrix**

```
## data = matrix(c(50,96,103,36,79,137), ncol=2, byrow=FALSE)
```

### **2) Exécuter le test chi-squared**

```
## result = chisq.test(data)
```

### **3) Afficher les résultats**

```
## result
```

### **4) Sortie logicielle**

```
##  
## Pearson's Chi-squared test  
##  
## data: data  
## X-squared = 8.7295, df = 2, p-value = 0.01272
```

