

# Statistique avec SAS®: Tests du $\chi^2$ et alternatives

Nicolas Poulin



# Préparation des données.

- télécharger le jeu de données *mais.csv* sur <http://www-irma.u-strasbg.fr/~gardes/teaching.html>
- remplacer les NA par des espaces vides dans excel.
- sauvegarder le fichier au format *.xlsx* dans un répertoire.
- charger le fichier sur le serveur SAS.

# Import d'un fichier .x/sx en SAS®.

```
PROC IMPORT OUT=work.mais1  
            DATAFILE="folders/myfolders/UEdata/mais.xlsx"  
            DBMS= xlsx;  
            GETNAMES=YES;  
RUN;
```

# Exploration des données

1 Affichage des attributs des variables :

```
PROC CONTENTS DATA= work.mais1; RUN;
```

2 Impression des données :

```
PROC PRINT DATA= work.mais1; RUN;
```

# Statistiques descriptives pour des variables quantitatives

- 1 Sur l'ensemble du jeu de données :

```
PROC MEANS DATA= work.mais1;  
    VAR Hauteur Masse;  
RUN;
```

- 2 Statistiques descriptives dans les niveaux d'une variable qualitative :

```
PROC MEANS DATA= work.mais1;  
    VAR Hauteur Masse;  
    CLASS Parcelle;  
RUN;
```

# Statistiques descriptives pour des variables quantitatives

- 3 Statistiques descriptives dans les niveaux de plusieurs variables qualitatives nichées :

```
PROC MEANS DATA= work.mais1;  
  VAR Hauteur Masse;  
  CLASS Parcelle Couleur;  
RUN;
```

## Remarque

*L'ordre dans lequel les variables qualitatives sont écrites est important. Ici il s'agit de Couleur dans Parcelle.*

# PROC FREQ de SAS®.

```
PROC FREQ <options> ;  
  BY variables ;  
  EXACT statistic-options </ computation-options> ;  
  OUTPUT <OUT=SAS-data-set> options ;  
  TABLES requests </ options> ;  
  TEST options ;  
  WEIGHT variable </ option> ;  
RUN;
```

# PROC FREQ de SAS®.

- **BY** : pour des analyses par groupe
- **EXACT** : pour faire des tests exacts
- **OUTPUT** : création d'un dataset en sortie
- **TABLES** : spécification des tableaux et analyses
- **TEST** : tests de mesure d'association
- **WEIGHT** : définition d'une variable de pondération.

## Remarque

*Pour pouvoir utiliser le paramètre **BY** il faut généralement trier le tableau de données en fonction des niveaux de la variable qualitative.*



# Statistiques descriptives pour des variables qualitatives

1 variables qualitatives :

```
PROC FREQ DATA= work.mais1 nlevels;  
    TABLES Parcelle Couleur;  
RUN;
```

## Remarque

*L'option **nlevels** permet d'afficher le tableau des niveaux du facteur des variables spécifiées dans le paramètre **TABLES**.*

# Statistiques descriptives pour des variables qualitatives

## 2 variables qualitatives emboîtées :

- Il faut tout d'abord trier le jeu de données.
- Si on veut les fréquences des parcelles pour les différentes couleurs :

```
PROC SORT DATA= work.mais1;  
    BY Couleur Parcelle;  
RUN;
```

- On peut ensuite utiliser la PROC FREQ

```
PROC FREQ DATA= work.mais1 nlevels;  
    TABLES Parcelle;  
    BY Couleur;  
RUN;
```

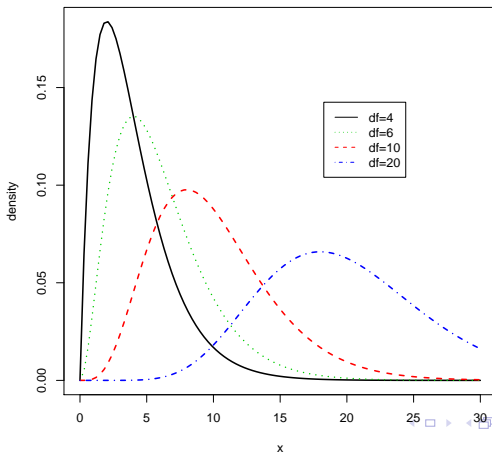
# La loi du $\chi^2$ .

En probabilité et statistique des variables aléatoires sont distribuées selon certaines lois de probabilité. Parmi elles, la loi du  $\chi^2$  est une des plus connues.

Le paramètre  $\nu$  d'une loi du  $\chi^2$  est appelé nombre de degrés de liberté.

La loi du  $\chi^2$  à  $\nu$  degrés de liberté correspond à la somme de  $\nu$  carrés de lois normales  $\mathcal{N}(0,1)$ .

# La loi du $\chi^2$ .



# Tests du $\chi^2$ .

La loi du  $\chi^2$  est souvent utilisée pour les tests d'hypothèses.

Le nom de tests du  $\chi^2$  correspond aux tests du  $\chi^2$  de Pearson.

Il existe 2 types de tests du  $\chi^2$  de Pearson :

- le test d'ajustement à une loi théorique : teste si un échantillon a pu être tiré selon une loi spécifique.
- le test d'indépendance : teste si 2 variables qualitatives observées sur 1 échantillon sont indépendantes.


# Test du $\chi^2$ d'ajustement à une loi.

**But :** tester si une variable qualitative observée sur 1 échantillon est issue d'une loi de probabilité  $\mathcal{L}$  spécifiée.

Les hypothèses testées sont :

- $H_0$  : la variable suit la loi  $\mathcal{L}$ .
- $H_1$  : la variable ne suit pas la loi  $\mathcal{L}$ .

Le test est basé sur le tableau des effectifs (resp. de fréquences) qui présentent les occurrences (resp. les fréquences) de la variable qualitative sur l'échantillon.

 **Anglais :** la frequency table est le tableau des effectifs.

# Test du $\chi^2$ pour l'ajustement à une loi.

- On dispose d'un dé à 6 faces dont on ne sait pas s'il est équilibré.
- On suppose que ce dé a été lancé  $N=100$  fois et que l'on a obtenu :

Face 1	Face 2	Face 3	Face 4	Face 5	Face 6
17	14	15	18	16	20

- Est-ce que ce dé est équilibré?

# Test du $\chi^2$ pour l'ajustement à une loi.

Les hypothèses testées sont :

- $H_0$  : le dé est équilibré.
- $H_1$  : le dé n'est pas équilibré.

C'est-à-dire :

- $H_0$  : le résultat du lancé de dé suit une loi uniforme.
- $H_1$  : le résultat du lancé de dé ne suit pas une loi uniforme.



# Conditions d'application du test du $\chi^2$ .

- les individus composant l'échantillon ont été choisis aléatoirement (i.e. indépendance des observations).
- les classes des variables sont exclusives.
- Règle de Cochran : au moins 80% des effectifs théoriques sont au moins égaux à 5.
- La taille de l'échantillon doit être assez grande.

Sous  $H_0$ , chaque face a une probabilité de réalisation de  $\frac{1}{6}$ .

## Remarque

*Le nombre suffisant d'observations est assez subjectif : pour certains 25 observations sont suffisantes, pour d'autres cela peut être 30 ou 50.*



# Tableau de contingence théorique.

Si  $H_0$  est vraie, sur 100 essais nous nous attendons à obtenir le tableau :

Face 1	Face 2	Face 3	Face 4	Face 5	Face 6
16.667	16.667	16.667	16.667	16.667	16.667

# Statistique de décision du test du $\chi^2$ .

La statistique de décision est une mesure de la distance entre le tableau observé et le tableau théorique :

$$\chi_{obs}^2 = \sum_{i=1}^k \frac{(n_i - t_i)^2}{t_i}$$

où

- $k$  est le nombre de classes.
- $n_i$  sont les effectifs observés de la classe  $i$ .
- $t_i$  sont les effectifs théoriques de la classe  $i$ .

Sous  $H_0$  la statistique de décision suit une loi du  $\chi^2$  à  $\nu = k - 1$  où  $k$  est le nombre de classes.

# Test du $\chi^2$ d'ajustement avec SAS®.

- 1 Grâce à l'étape DATA, créer le jeu de données.

```
DATA dice;  
    INPUT face $ effectif;  
    CARDS;  
    Face1 17  
    Face2 14  
    Face3 15  
    Face4 18  
    Face5 16  
    Face6 20  
    ;
```

# Test du $\chi^2$ d'ajustement avec SAS®.

2 Grâce à **PROC FREQ**, procéder au test du  $\chi^2$  d'ajustement.

Options de **TABLES** :

- mises après le symbole /
- **CHISQ** : indique que l'on veut un test du  $\chi^2$
- **TESTP** : proportions théoriques testées
- **TESTF** : effectifs théoriques testés
- **NOCUM** : effectifs cumulés non affichés dans le tableau

Paramètre de **WEIGHT** : donne le nom de la variable contenant les effectifs.

# Test du $\chi^2$ d'ajustement avec SAS®.

```
PROC FREQ DATA =work.dice ORDER=DATA;  
  TABLES face / CHISQ NOCUM  
  TESTF=(16.667 16.667 16.667 16.667 16.667 16.667);  
  WEIGHT effectif;  
RUN;
```

## Remarque

*L'option **ORDER=DATA** permet que les niveaux des variables qualitatives soient ordonnés comme dans le jeu de données.*

# Test du $\chi^2$ d'ajustement avec SAS®.

- **BUT** : tester si la répartition des couleurs dans le jeu de données *mais* est uniforme.
- Il faut d'abord trier le jeu de données :

```
PROC SORT DATA= work.mais1;  
    BY Couleur;  
RUN;
```

# Test du $\chi^2$ d'ajustement avec SAS®.

- On peut ensuite utiliser la PROC FREQ :

```
PROC FREQ DATA= work.mais1 ORDER=DATA;  
    TABLES Couleur / CHISQ TESTP=(0.334 0.334 0.334);  
RUN;
```



# Exercice.

- On suppose que 30% des parcelles sont orientées au nord, 20% au sud, 17% à l'est et 33% à l'ouest.
- Tester si la supposition ci-dessus peut être considérée comme vraie sur l'ensemble de la population.

# Test du $\chi^2$ d'indépendance.

**But :** tester si 2 variables qualitatives observées sur 1 échantillon sont indépendantes.

Les hypothèses testées sont :

- $H_0$  : les variables sont indépendantes.
- $H_1$  : les variables ne sont pas indépendantes.

Le test est basé sur les tableaux de contingence qui présentent les occurrences croisées de 2 variables qualitatives.

Exemple de tableau de contingence :

	<b>BLOND</b>	<b>CHATAIN</b>	<b>BRUN</b>	<b>ROUX</b>
<b>BLEU</b>	25	9	3	7
<b>GRIS ou VERT</b>	13	17	10	7
<b>MARRON</b>	7	13	8	5

# Conditions d'applications.

- les individus composant l'échantillon ont été choisis aléatoirement (i.e. indépendance des observations).
- les classes des variables sont exclusives.
- Règle de Cochran : au moins 80% des effectifs théoriques sont au moins égaux à 5.
- La taille de l'échantillon doit être assez grande.

## Remarque

*Le nombre suffisant d'observations est assez subjectif : pour certains 25 observations sont suffisantes, pour d'autres cela peut être 30 ou 50.*

# Tableau des effectifs théoriques.

- Le test consiste à mesurer si la différence entre ce qu'on observe et ce qui devrait idéalement se passer en cas d'indépendance est statistiquement significatif.
- Il faut construire un tableau d'effectifs théoriques basé sur les marges du tableau de contingence observé :

$$t_{ij} = \frac{n_{i.} \times n_{.j}}{n}$$

où

- ▶  $t_{ij}$  est l'effectif théorique pour la ligne  $i$  et la colonne  $j$ .
- ▶  $n_{i.}$  la somme des effectifs observés pour la ligne  $i$ .
- ▶  $n_{.j}$  la somme des effectifs observés pour la colonne  $j$ .

## Remarque

*Le tableau des effectifs théoriques peut comporter des effectifs qui ne sont pas des entiers naturels.*



# Tableau des effectifs théoriques.

- Tableau de contingence observé ( $n_{i,j}$ ) :

	<b>BLOND</b>	<b>CHATAIN</b>	<b>BRUN</b>	<b>ROUX</b>
<b>BLEU</b>	25	9	3	7
<b>GRIS ou VERT</b>	13	17	10	7
<b>MARRON</b>	7	13	8	5

- Tableau des effectifs théoriques ( $t_{i,j}$ ) :

	<b>BLOND</b>	<b>CHATAIN</b>	<b>BRUN</b>	<b>ROUX</b>
<b>BLEU</b>	15.96774	13.83871	7.451613	6.741935
<b>GRIS ou VERT</b>	17.05645	14.78226	7.959677	7.201613
<b>MARRON</b>	11.97581	10.37903	5.588710	5.056452

# Prise de décision.

- Statistique de décision :  $\chi_{obs}^2 = \sum_{i=1}^k \sum_{j=1}^c \frac{(n_{ij} - t_{ij})^2}{t_{ij}}$ .
- Loi sous  $H_0$  : loi du  $\chi^2$  à  $\nu = (k-1)(c-1)$  ddl où  $k$  et  $c$  sont les nombres de classes des deux variables.
- La décision est prise grâce à la p-value.

# Test du $\chi^2$ d'indépendance avec SAS®.

## 1 Grâce à l'étape DATA, créer le jeu de données.

```
DATA COULEUR;  
  INPUT YEUX $ CHEVEUX $ EFFECTIF;  
  CARDS;  
  BLEUS BLONDS 25  
  BLEUS BRUNS 9  
  BLEUS NOIRS 3  
  BLEUS ROUX 7  
  VERTS BLONDS 13  
  VERTS BRUNS 17  
  VERTS NOIRS 10  
  VERTS ROUX 7  
  MARRONS BLONDS 7  
  MARRONS BRUNS 13  
  MARRONS NOIRS 8  
  MARRONS ROUX 5
```

```
;
```

# Test du $\chi^2$ d'indépendance avec SAS®.

- 2 Grâce à **PROC FREQ** , créer le tableau de contingence observé.

```
PROC FREQ DATA=COULEUR ORDER=DATA;  
  TABLES YEUX*CHEVEUX ;  
  WEIGHT EFFECTIF ;  
  TITLE1 "Tableau de contingence observé";  
  TITLE2 "- - - - -";  
RUN;
```



# Test du $\chi^2$ d'indépendance avec SAS®.

Le quadrant supérieur gauche du tableau indique (en anglais) le contenu de chaque case  $(i,j)$ , à savoir :

- l'effectif  $n_{i,j}$  (Fréquence)
- le Pourcentage correspondant à  $f_{i,j} = n_{i,j}/N$
- le Pourcentage en ligne correspondant à  $n_{i,j}/n_{i,.}$
- le Pourcentage en colonne correspondant à  $n_{i,j}/n_{.,j}$

# Test du $\chi^2$ d'indépendance avec SAS®.

Sur la ligne Total on peut lire :

- les effectifs  $n_{.,j}$  des modalités de la variable colonne,
- les pourcentages ligne correspondant aux proportions  
 $f_{.,j} = n_{.,j}/N$

C'est la ligne marginale donnant la distribution (le tri-à-plat) de la variable CHEVEUX sans distinction de la couleur des yeux.

Sur la colonne Total, colonne marginale, on lit de même la distribution de la variable YEUX dans l'ensemble de la population (effectifs  $n_{i,.}$  et pourcentages colonne correspondant à  $f_{i,.} = n_{i,.}/N$ ).

# Test du $\chi^2$ d'indépendance avec SAS®.

- 3 Grâce à **PROC FREQ** , créer le tableau de contingence théorique.

```
PROC FREQ DATA=COULEUR ORDER=DATA;  
  TABLES YEUX*CHEVEUX / EXPECTED;  
  WEIGHT EFFECTIF ;  
RUN;
```

# Test du $\chi^2$ d'indépendance avec SAS®.

- 4 Grâce à **PROC FREQ** , tester si les variables YEUX et CHEVEUX sont indépendantes.

```
PROC FREQ DATA=COULEUR ORDER=DATA;  
  TABLES YEUX*CHEVEUX / CHISQ;  
  WEIGHT EFFECTIF ;  
RUN;
```

# Exercice.

Vérifier, sur l'exemple *mais* si la parcelle et la couleur sont des variables aléatoires indépendantes.

## Exercice.

- Il faut tout d'abord trier le jeu de données :

```
PROC SORT DATA= work.mais1;  
    BY Parcelle Couleur;  
RUN;
```

- On peut ensuite obtenir le tableau observé :

```
PROC FREQ DATA= work.mais1 ORDER= DATA;  
    TABLES Couleur*Parcelle;  
RUN;
```

### Remarque

*Pas besoin de paramètre **WEIGHT** : les poids sont calculés automatiquement car il s'agit du tableau complet.*



## Exercice.

- On peut ensuite obtenir le tableau théorique :

```
PROC FREQ DATA= work.mais1 ORDER= DATA;  
    TABLES Couleur*Parcelle / EXPECTED;  
RUN;
```

### Remarque

*Les conditions d'application du test du  $\chi^2$  d'indépendance ne sont pas vérifiées (trop d'effectifs théoriques inférieurs à 5). Il faut donc soit :*

- *fusionner les niveaux Jaune.rouge et Rouge pour la variable Couleur*
- *utiliser un test exact de Fisher*

# Exercice.

- Si on voulait tout de même faire le test du  $\chi^2$  :

```
PROC FREQ DATA= work.mais1 ORDER= DATA;  
    TABLES Couleur*Parcelle / CHISQ EXPECTED;  
RUN;
```



# Alternatives : test exact de Fisher

Il peut être utilisé pour tester l'indépendance entre deux variables qualitatives quand la condition de Cochran n'est pas vérifiée ou pour de petits échantillons.

Hypothèses :

- $H_0$ : les variables sont indépendantes.
- $H_1$ : les variables ne sont pas indépendantes.

## Remarque

*Le test est dit exact car il calcule en fait la probabilité qu'un tableau de contingence observé le soit sous une hypothèse d'indépendance (par des techniques de dénombrement).*

# Test exact de Fisher

Conditions d'application :

- les individus composant l'échantillon ont été choisis aléatoirement (i.e. indépendance des observations).
- les classes des variables sont exclusives.

Test exact de Fisher avec SAS®.

```
PROC FREQ DATA= work.mais1 ORDER= DATA;  
  TABLES Couleur*Parcelle / EXPECTED;  
  EXACT Fisher;  
RUN;
```

# Alternatives : G-test (Likelihood ratio Chi-Square)

- Le test du  $\chi^2$  de Pearson est basé sur une approximation d'un ratio de log-vraisemblance.
- L'utilisation du G-test est recommandé dans le livre de Robert R. Sokal et F. James Rohlf (1981), *Biometry : the principles and practice of statistics in biological research*, New-York: Freeman.
- Le G-test n'utilise pas l'approximation mais calcule le vrai rapport de log-vraisemblance ce qui permet d'obtenir des résultats plus fiables.
- Les hypothèses et conditions d'application du G-test sont les mêmes que celles du test du  $\chi^2$  mais sans la limite de la taille d'échantillon.

## Paramètres :

- Loi sous  $H_0$  : loi du  $\chi^2$  à  $\nu = (k-1)(c-1)$  ddl où  $k$  et  $c$  sont les nombres de classes des deux variables.
- Résultat donné dans le tableau des tests lorsqu'on utilise l'option **CHISQ** du paramètre **TABLES**

# Correction de continuité pour les tableaux $2 \times 2$

Lorsque l'on traite un tableau de contingence pour 2 variables qualitatives comportant chacune 2 niveaux, il est nécessaire d'appliquer la correction de continuité de Yates.

En effet, le test du  $\chi^2$  suppose que des probabilités discrètes (loi binomiale) peut être approchée par une distribution du  $\chi^2$  qui est continue.

# Correction de continuité pour les tableaux $2 \times 2$

Cette approximation n'est plus valable dans le cas de tableaux de contingences  $2 \times 2$ . Yates a proposé la correction :

- Statistique de décision :  $\chi^2_{\text{Yates}} = \sum_{i=1}^k \sum_{j=1}^c \frac{(|n_{ij} - t_{ij}| - 0.5)^2}{t_{ij}}$ .
- Loi sous  $H_0$  : loi du  $\chi^2$  à  $\nu = (k-1)(c-1)$  ddl où  $k$  et  $c$  sont les nombres de classes des deux variables.
- La décision est prise grâce à la p-value.
- Résultat donné dans le tableau des tests lorsqu'on utilise l'option **CHISQ** du paramètre **TABLES**

# Correction de continuité pour les tableaux $2 \times 2$

SAS<sup>®</sup> donne aussi automatiquement le résultat du test exact de Fisher pour les tableaux  $2 \times 2$  car certains statisticiens considèrent que l'approximation de Yates est obsolète.

Réaliser le test du  $\chi^2$  avec correction de Yates sur le tableau de contingence des variables *Attaque* et *Verse\_Traitement*.