

## Chapitre 11 – Régression linéaire simple

### Analyse de corrélation linéaire (R)

- permet de déterminer l'intensité de la liaison linéaire entre 2 VAR
- utilise R linéaire pour mesure l'intensité de la liaison,

$$R = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} = \frac{n \sum X_i Y_i - (\sum X_i)(\sum Y_i)}{\sqrt{n \sum X_i^2 - (\sum X_i)^2} \sqrt{n \sum Y_i^2 - (\sum Y_i)^2}} = \frac{\sum_{i=1}^n X_i Y_i - \frac{\sum_{i=1}^n X_i \sum_{i=1}^n Y_i}{n}}{\sqrt{\sum_{i=1}^n X_i^2 - \frac{(\sum_{i=1}^n X_i)^2}{n}} \sqrt{\sum_{i=1}^n Y_i^2 - \frac{(\sum_{i=1}^n Y_i)^2}{n}}}$$

### Coefficient de corrélation (LINÉAIRE)

Nombre compris entre  $-1 \leq R \leq 1$

#### Analyse de régression

- On veut déterminer la relation statistique entre deux variables Y et X (c'est l'analyse de régression)
- En général, la variable X est contrôlable
- La 1<sup>ère</sup> étape d'une analyse de régression consiste à faire un graphique de Y en fonction de X

#### Variables dépendante et indépendante

Variable dépendante Y: variable expliquée

Variable indépendante X: variable explicative

#### Régression linéaire simple

Permet de déterminer une équation reliant la variable dépendante Y à la variable indépendante X, de déterminer la validité du modèle obtenu, de faire des prévisions

#### La droite de régression

$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ ,  $\beta_0$ ,  $\beta_1$  et  $\varepsilon_i$  sont les paramètres à estimer

**ROUGE** : Variable expliquée (variable aléatoire)

Terme constant pour une valeur donnée de  $X_i$  (variable explicative)

**ORANGE** : Facteur aléatoire

#### Composantes du modèle

$Y_i$  = composante non aléatoire + composante aléatoire

Composante non aléatoire :  $E(Y_i/X_i) = \beta_0 + \beta_1 X_i$

Composante aléatoire :  $\varepsilon_i = Y_i - (\beta_0 + \beta_1 X_i)$

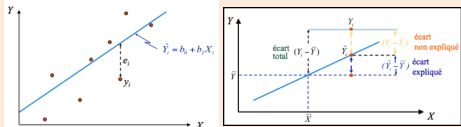
#### Hypothèses fondamentales du modèle

Les  $\varepsilon_i$  sont des VA ind. (normalement distribuées, d'esp math  $E(\varepsilon_i) = 0$  et VAR constante  $Var(\varepsilon_i) = \sigma_\varepsilon^2$ ) pour toutes valeurs de  $X_i$

#### Estimation de la droite de régression

Plusieurs droites peuvent être ajustées aux données obtenues

Par la méthode des moindres carrés, on peut déterminer celle qui minimise la somme des carrés des erreurs ---  $E(Y_i) = \beta_0 + \beta_1 X_i$  avec un estimateurs,  $\hat{Y}_i = b_0 + b_1 X_i$



#### Détermination de la droite de régression

- Les résidus :  $e_i = y_i - \hat{y}_i = y_i - (b_0 + b_1 X_i)$
- La somme des carrés des erreurs :  $SC_{rés} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$
- $b_0$  et  $b_1$  sont les coefficients de régression ( $b_0$  = l'ordonnée à l'origine &  $b_1$  = la pente de la droite de régression empirique)

#### Méthode des moindres carrés

- On veut minimiser la somme des carrés des erreurs,  $SC_{rés}$
- Min  $SC_{rés} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 X_i)^2$
- En éliminant les équations précédentes, on obtient :
- $b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$  &  $b_0 = \bar{y} - b_1 \bar{x}$

#### Droite de régression empirique

$\hat{Y} = b_0 + b_1 X_i$ , puisque  $b_0 = \bar{Y} - b_1 \bar{X}$  - alors en substituant  $b_0$ , on obtient  $\hat{Y}_i = \bar{Y} + b_1 (X_i - \bar{X})$

#### Analyse de la variance

Dans quelle mesure la droite de régression est-elle utile pour expliquer la variation des observations  $Y_i$

- Variation expliquée par la régression
- Variation non expliquée ou résiduelle

#### Décomposition de la variation totale (voir graphique)

$(Y_i - \bar{Y}) = (\hat{Y}_i - \bar{Y}) + (Y_i - \hat{Y}_i)$

Écart total = écart expliqué par la droite + écart non expliqué

#### Somme des carrés et calcul de somme des carrés

Variation total = variation ex. par droite + variation non expliquée

$$SCT = SCR + SC_{rés} = \sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum (Y_i - \hat{Y}_i)^2$$

- $SCR = \sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - \frac{(\sum Y_i)^2}{n}$
- $SC_{rés} = \sum (Y_i - \hat{Y}_i)^2 = b_1^2 \left[ \sum X_i^2 - \frac{(\sum X_i)^2}{n} \right]$
- $SC_{rés} = \sum (Y_i - \hat{Y}_i)^2 = \sum Y_i^2 - b_0 \sum Y_i - b_1 \sum X_i Y_i$

#### Coefficient de détermination simple $r^2$

La qualité de l'ajustement du modèle aux données peut être mesurée grâce au coefficient de détermination simple  $r^2$

Ce coefficient donne contribution de variable explicative ds l'équation de régression pour expliquer les variations de la variable dépendante

$$r^2 = \frac{\text{variation expliquée}}{\text{variation totale}} = \frac{SCR}{SCT} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2}, 0 \leq r^2 \leq 1$$

#### Coefficient de corrélation linéaire simple

Le signe du coefficient de corrélation est le même que celui du coefficient de régression  $b_1$ ,  $r = \pm \sqrt{r^2}$ ,  $-1 \leq r \leq 1$

#### Variation totale inexpliquée

$$(1 - r^2) = 1 - \frac{\text{variation expliquée}}{\text{variation totale}} = \frac{\text{variation résiduelle}}{\text{variation totale}} = 1 - \frac{SCR}{SCT} = \frac{SC_{rés}}{SCT}$$

#### Les sommes des carrés et leurs degrés de liberté

Variation	Somme des carrés	Degré de liberté
Régression	$SCR = \sum (\hat{Y}_i - \bar{Y})^2$	1
Résiduelle	$SC_{rés} = \sum (Y_i - \hat{Y}_i)^2$	(n-2)
Totale	$SCT = \sum (Y_i - \bar{Y})^2$	(n-1)

#### Carrés moyens ou variances

$$CMR = \frac{SCR}{1} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{1} \quad CM_{rés} = \frac{SC_{rés}}{n-2} = \frac{\sum (Y_i - \hat{Y}_i)^2}{n-2} = S_\varepsilon^2$$

#### Espérance des carrés moyens

$E(CM_{rés}) = \sigma_\varepsilon^2$   $E(CMR) = \sigma_\varepsilon^2 + \beta_1^2 \sum (X_i - \bar{X})^2$  si  $\beta_1 = 0$ , alors

$E(CMR) = E(CM_{rés})$  et si  $\beta_1 \neq 0$ , alors  $E(CMR) > E(CM_{rés})$

#### Tableau d'analyse de variance

Source de variation	Somme des carrés	Degrés de liberté	Carrés moyens	F
Due à la régression	SCR	1	CMR	$\frac{CMR}{S_\varepsilon^2}$
Résiduelle	$SC_{rés}$	n-2	$S_\varepsilon^2$	
Total	SCT	n-1		

**Test sur  $\beta_1$  : test de Fisher** \*\*les valeurs critique sont tjrs ds tables

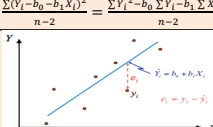
- Si  $\beta_1 = 0$ , alors  $CMR/CM_{rés} = 1$
- $F = \frac{CMR}{CM_{rés}} = \frac{SCR/\sigma_\varepsilon^2/1}{SC_{rés}/\sigma_\varepsilon^2/n-2} = \frac{SCR/1}{SC_{rés}/n-2} = F_\alpha(1, n-2)$
- $H_0 : \beta_1 = 0$ ,  $[E(CMR) = E(CM_{rés})]$
- $H_1 : \beta_1 \neq 0$ ,  $[E(CMR) > E(CM_{rés})]$
- Seuil de signification  $\alpha$
- Si  $H_0$  est vraie alors  $F = \frac{CMR}{CM_{rés}} \in F(1, n-2)$
- Règle de décision : rejeter  $H_0$  si  $F > F_{\alpha;1;n-2}$

#### Résidus et mesure de distribution $\sigma_\varepsilon^2$

- Veut obtenir estimation de dist. des  $Y_i$  autour de droite régression
- $Var(Y_i/X_i) = \sigma_\varepsilon^2 = Var(\varepsilon_i)$
- Les résidus :  $e_i = y_i - \hat{y}_i$  = permettent de vérifier l'hypothèse de normalité des erreurs et l'homogénéité de la variance des erreurs

#### Estimation de la variance $\sigma_\varepsilon^2$

$$S_\varepsilon^2 = \frac{\sum (Y_i - \hat{Y}_i)^2}{n-2} = \frac{SC_{rés}}{n-2} = \frac{\sum \varepsilon_i^2}{n-2} = \frac{\sum (Y_i - b_0 - b_1 X_i)^2}{n-2} = \frac{\sum Y_i^2 - b_0 \sum Y_i - b_1 \sum X_i Y_i}{n-2}$$



#### Résidus et droite de régression

### Analyse graphique des résidus

L'analyse des résidus permet de vérifier que :  $E(\varepsilon_i) = 0$ ,  $Var(\varepsilon_i) = \sigma_\varepsilon^2$

- les  $\varepsilon_i$  sont normalement distribués & ne sont pas autos corrélés
- L'analyse des résidus peut se faire graphiquement
- Les résidus devraient se distribuer aléatoirement autour de zéro(0)
- Si l'analyse du graphique des résidus en fonction de la variable dépendante révèle une forme non aléatoire, alors le modèle n'est pas linéaire, les erreurs ne sont pas normalement distribuées, n'ont pas une variance constante (hétéroélastique) et sont autos corrélées

#### Inférence sur $\beta_1$

$b_1$  est un estimateur efficace de  $\beta_1$ ,  $b_1$  étant un estimateur, il possède une distribution d'échantillonnage avec une moyenne et une variance

#### Distribution d'échantillonnage de $b_1$

si  $\varepsilon \in N(0, \sigma_\varepsilon^2)$  et  $Cov(\varepsilon_i, \varepsilon_j) = 0$  et  $Cov(\varepsilon_i, X_i) = 0$ , alors  $b_1 \in$

$$N\left(\beta_1, \sigma_\varepsilon^2 \left(\frac{1}{\sum (X_i - \bar{X})^2}\right)\right)$$

#### Estimation de la variance $\sigma^2(b_1)$

$b_1$  distribué normalement, taille d'échantillon petite,  $\sigma^2(b_1)$  inconnue

$$S^2(b_1) = \frac{S_\varepsilon^2}{\sum (X_i - \bar{X})^2}, S(b_1) = \frac{S_\varepsilon}{\sqrt{\sum (X_i - \bar{X})^2}} \text{ et } t = \frac{b_1 - \beta_1}{s(b_1)} \in t(n-2)$$

#### Intervalle de confiance pour $\beta_1$

$$b_1 - t_{\alpha/2;n-2} S(b_1) \leq \beta_1 \leq b_1 + t_{\alpha/2;n-2} S(b_1), S(b_1) = \frac{S_\varepsilon}{\sqrt{\sum (X_i - \bar{X})^2}}$$

#### Test d'hypothèse pour $\beta_1$

Le paramètre  $\beta_1$  mesure la dépendance linéaire entre Y et X

Si cette dépendance est significative, alors  $\beta_1 \neq 0$

Si  $\beta_1 = 0$ , le modèle devient :  $Y_i = \beta_0 + \varepsilon_i$  ou  $Y_i = E(Y_i) + \varepsilon_i$

**H<sub>0</sub> :  $\beta_1 = 0$     H<sub>1</sub> :  $\beta_1 \neq 0$**

**Conditions du test**: le modèle est valide,  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ ,  $b_1$  distribué selon la loi normale,  $\sigma^2(b_1)$  inconnue,  $\alpha$  fixé

- La statistique :  $t = \frac{b_1 - \beta_1}{s(b_1)} = \frac{b_1}{s(b_1)} \in t(n-2)$

**La règle de décision** : **H<sub>1</sub> :  $\beta_1 \neq 0$ , Rejeter  $H_0$  si  $t < -t_{\alpha/2;n-2}$  ou si  $t > t_{\alpha/2;n-2}$**

#### Inférence sur $\beta_0$

$b_0$  est un estimateur efficace de  $\beta_0$ ,  $b_0$  étant un estimateur, il possède une distribution d'échantillonnage avec une moyenne et une variance

#### Distribution d'échantillonnage de $b_0$

si  $\varepsilon \in N(0, \sigma_\varepsilon^2)$  et  $Cov(\varepsilon_i, \varepsilon_j) = 0$  et  $Cov(\varepsilon_i, X_i) = 0$ , alors  $b_0 \in$

$$N\left(\beta_0, \sigma_\varepsilon^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (X_i - \bar{X})^2}\right)\right)$$

#### Estimation de la variance $\sigma^2(b_0)$

$$S^2(b_0) = S_\varepsilon^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (X_i - \bar{X})^2}\right) S = S_\varepsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum (X_i - \bar{X})^2}} \text{ et } t = \frac{b_0 - \beta_0}{s(b_0)} \in t(n-2)$$

#### Intervalle de confiance pour $\beta_0$

- $b_0 - t_{\alpha/2;n-2} S(b_0) \leq \beta_0 \leq b_0 + t_{\alpha/2;n-2} S(b_0)$
- $S(b_0) = S_\varepsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum (X_i - \bar{X})^2}}$

#### Test d'hypothèse pour $\beta_0$

Si le test d'hypothèse sur  $\beta_0$  favorise l'hypothèse nulle  $H_0: \beta_0 = 0$ , le modèle devient :  $Y_i = \beta_1 X_i + \varepsilon_i$

Dans ce cas, la droite de régression passe par l'origine

**H<sub>0</sub> :  $\beta_0 = 0$     H<sub>1</sub> :  $\beta_0 \neq 0$**

**Conditions du test**: le modèle est valide,  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ ,  $b_0$  distribué selon la loi normale,  $\sigma^2(b_0)$  inconnue,  $\alpha$  fixé

- La statistique :  $t = \frac{b_0}{s(b_0)} \in t(n-2)$

**La règle de décision** : **H<sub>1</sub> :  $\beta_0 \neq 0$ , Rejeter  $H_0$  si  $t < -t_{\alpha/2;n-2}$  ou si  $t > t_{\alpha/2;n-2}$**

#### Intervalle de confiance pour $E(Y_n)$

Pour  $X = X_n$ , l'estimation ponctuelle de  $E(Y_n)$  est  $\hat{Y}_n = b_0 + b_1 X_n$

La distribution d'échantillonnage de  $E(Y_n)$  est  $\frac{Y_n - E(Y_n)}{s(Y_n)} \in t(n-2)$

L'intervalle de conf :  $\hat{Y}_n - t_{\alpha/2;n-2} S(\hat{Y}_n) \leq E(Y_n) \leq \hat{Y}_n + t_{\alpha/2;n-2} S(\hat{Y}_n)$

#### Calcul de $S(\hat{Y}_n)$

$$S^2(\hat{Y}_n) = S^2 \left[ \frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right], S(\hat{Y}_n) = S \sqrt{\frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum (X_i - \bar{X})^2}}$$

#### Intervalle de confiance pour $Y_n$

Pour  $X = X_n$ , l'estimation ponctuelle de  $Y_n$  est  $\hat{Y}_n = b_0 + b_1 X_n$

La distribution d'échantillonnage de  $Y_n$  est  $t = \frac{Y_n - \hat{Y}_n}{s \sqrt{1 + \frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum (X_i - \bar{X})^2}}} \in t(n-2)$

L'intervalle de conf :  $\hat{Y}_n - t_{\alpha/2;n-2} S(d_n) \leq Y_n \leq \hat{Y}_n + t_{\alpha/2;n-2} S(d_n)$

#### Calcul de $S(d_n)$

$$S^2(d_n) = S^2 + \frac{S^2}{n} + S^2 \frac{(X_n - \bar{X})^2}{\sum (X_i - \bar{X})^2}, S(d_n) = S \sqrt{1 + \frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum (X_i - \bar{X})^2}}$$

#### Calcul de $S(d_n)$

$$S \sqrt{1 + \frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum (X_i - \bar{X})^2}}$$

## Chapitre 8

Test d'égalité des espérances pour 2 pop de variances égales		
CAS 3 - test de 2 moyenne	Classe 1	Classe 2
Moyenne	$\bar{x}_1$	$\bar{x}_2$
Variance	$S_1^2$	$S_2^2$
Observations	$n_1$	$n_2$
Variance pondérée	$S_c^2 = \text{Var}(\bar{X}_1 - \bar{X}_2)$	
Différence hypothétique...		-
Degré de liberté	n-2	-
Statistique t	Résultat empi.	-
P(T<=t) unilatéral		-
Valeur critique de t (unilatéral)	t calculé	-
P(T<=t) bilatéral		-
Valeur critique de t (bilatéral)		-
Test d'égalité des espérances : observations paires		
Moyenne	$\bar{x}_1$	$\bar{x}_2$
Variance	$S_1^2$	$S_2^2$
Observations	$n_1$	$n_2$
Coefficient de cor. de Pearson		-
Différence hypothétique...		-
Degré de liberté	n-1 (1 n slmt)	-
Statistique t	Résultat empi.	-
P(T<=t) unilatéral		-
Valeur critique de t (unilatéral)	t calculé	-
P(T<=t) bilatéral		-
Valeur critique de t (bilatéral)		-

Moyenne	$\bar{x}$
Erreur-type	$S_\varepsilon$
Écart-type	$\sigma$
Variance de l'échantillon	$\sigma^2$
Plage	Max - min
Min	# min par observation
Max	# max par observation
Somme	Total d'échantillon
Nombre d'observations	n
Test d'égalité des variances (F-test)	
Moyenne	$\bar{x}_1$
Variance	$S_1^2$
Observation	$n_1$
Degré de liberté	$n_1 - 1$
F	Résultat empi. ( $S_1^2/S_2^2$ )
P(F<=f) unilatéral	-
Valeur critique F (unilatéral)	t calculé

## Chapitre 9

Résultat empirique : calcul de $\chi^2$						
	Fobs.	Fthéor.	Calcul de $\chi^2$			
Résultats	$f_{oi}$	$P_i$	$f_{ti}$	$f_{oi} \cdot f_{ti}$	$(f_{oi} \cdot f_{ti})^2$	$(f_{oi} \cdot f_{ti})^2 / f_{ti}$
1						

CHAPITRE 8 - Test d'hypothèse sur deux paramètres

Différence de deux moyennes

- Soit  $X_1$  et  $X_2$  deux variables aléatoires de moyenne  $\mu_1$  et  $\mu_2$  et de variance  $\sigma_1^2$  et  $\sigma_2^2$  respectivement
- La statistique  $(\bar{X}_1 - \bar{X}_2)$  est utilisée pour comparer l'espérance mathématique  $(\mu_1 - \mu_2)$  des deux populations

Cas 1 : Pop. normales, et  $\sigma_1^2$  et  $\sigma_2^2$  connues

Si on prélève indépendamment des échantillons aléatoires de taille  $n_1$  et  $n_2$  respectivement de 2 pop normales tel que  $X_1 \in N(\mu_1, \sigma_1^2)$  et  $X_2 \in N(\mu_2, \sigma_2^2)$  alors la distribution d'échantillonnage de l'estimateur  $(\bar{X}_1 - \bar{X}_2)$  est une loi nor. avec :

$$E(\bar{X}_1 - \bar{X}_2) = \mu_1 - \mu_2 \text{ et } \text{Var}(\bar{X}_1 - \bar{X}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

Soit  $X_1$  et  $X_2$  2 V.A. indépendantes, tel que  $X_1 \in N(\mu_1, \sigma_1^2)$  et  $X_2 \in N(\mu_2, \sigma_2^2)$  alors :  $(\bar{X}_1 - \bar{X}_2) \in N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$  et

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \in N(0,1)$$

Cas 2 : Grands échantillons,  $n_1 \geq 30$  et  $n_2 \geq 30$

Si on prélève indépendamment des échantillons aléatoires de taille  $n_1$  et  $n_2$  suffisamment grandes, alors grâce au théorème central limite, la différence  $(\bar{X}_1 - \bar{X}_2)$  suit approximativement une loi normale avec :

$$E(\bar{X}_1 - \bar{X}_2) = \mu_1 - \mu_2 \text{ et } \text{Var}(\bar{X}_1 - \bar{X}_2) = \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}$$

Si  $X_1$  et  $X_2$  sont 2 V.A ind. et que  $n_1$  et  $n_2$  sont grands, alors on a approx.

$$(\bar{X}_1 - \bar{X}_2) \in N(\mu_1 - \mu_2, \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}) \text{ et } Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \in N(0,1)$$

Cas 3 : Pop. normales,  $\sigma_1^2$  et  $\sigma_2^2$  inconnues mais supposées =, et  $n_1$  et/ou  $n_2$  petits (si les var sont  $\neq$ , faut utiliser le test de Satterthwaite, cas 4)

Si on prélève indépendamment des échantillons aléatoires de petites tailles ( $n_1$  et/ou  $n_2 < 30$ ) de populations normales de variances inconnues mais supposées égales à  $\sigma^2$ , on obtient :

$$E(\bar{X}_1 - \bar{X}_2) = \mu_1 - \mu_2 \text{ et } \text{Var}(\bar{X}_1 - \bar{X}_2) = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)$$

Si  $X_1 \in N(\mu_1, \sigma^2)$  et  $X_2 \in N(\mu_2, \sigma^2)$ ,  $\sigma_1^2$  et  $\sigma_2^2$  sont inconnues mais supposées égales à  $\sigma^2$ , et  $n_1$  et/ou  $n_2$  sont petits, alors on a :

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \in t(n_1 + n_2 - 2)$$

Estimation de  $\sigma^2$  par  $S_c^2$  & Calcul de  $\text{Var}(\bar{X}_1 - \bar{X}_2)$

$$S_c^2 = \frac{\sum(X_{1i} - \bar{X}_1)^2 + \sum(X_{2i} - \bar{X}_2)^2}{(n_1-1) + (n_2-1)} = \frac{\sum(X_{1i} - \bar{X}_1)^2 + \sum(X_{2i} - \bar{X}_2)^2}{(n_1-1)s_1^2 + (n_2-1)s_2^2}$$

$$\text{Var}(\bar{X}_1 - \bar{X}_2) = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right) = \frac{S_c^2}{n_1} + \frac{S_c^2}{n_2} = S_c^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)$$

Test sur une différence de moyennes  $(\mu_1 - \mu_2)$

- $H_0 : \mu_1 - \mu_2 = 0$ ,  $H_1 : \mu_1 - \mu_2 \neq 0$ ;  $\mu_1 - \mu_2 < 0$ ;  $\mu_1 - \mu_2 > 0$ ,  $\alpha$  fixé
- Conditions du test :
  - o CAS 1; 2 E.A. ind. pop. Norm. et  $\sigma_1^2$  et  $\sigma_2^2$  connues
  - o CAS 2; 2 E.A. ind. de grande taille (CHANGE  $\sigma$  POUR  $S$ )
- La statistique pour :  $Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \in N(0,1)$

La règle de décision : i)  $H_1 : \mu_1 - \mu_2 \neq 0$ , Rejeter  $H_0$  si

$$|\bar{X}_1 - \bar{X}_2| < (\mu_1 - \mu_2) - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \text{ ou si}$$

$$|\bar{X}_1 - \bar{X}_2| > (\mu_1 - \mu_2) + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

ii)  $H_1 : \mu_1 - \mu_2 < 0$ , Rejeter  $H_0$  si  $(\bar{X}_1 - \bar{X}_2) < (\mu_1 - \mu_2) - z_{\alpha} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

iii)  $H_1 : \mu_1 - \mu_2 > 0$ , Rejeter  $H_0$  si  $(\bar{X}_1 - \bar{X}_2) > (\mu_1 - \mu_2) + z_{\alpha} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

CAS 3; 2 E.A. ind.,  $\sigma_1^2$  et  $\sigma_2^2$  inconnues, et  $n_1$  et/ou  $n_2$  petits

La statistique pour :  $t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \in t(n_1 + n_2 - 2)$

La règle de décision : i)  $H_1 : \mu_1 - \mu_2 \neq 0$ , Rejeter  $H_0$  si

$$|\bar{X}_1 - \bar{X}_2| < (\mu_1 - \mu_2) - t_{\alpha/2} \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} \text{ ou si}$$

$$|\bar{X}_1 - \bar{X}_2| > (\mu_1 - \mu_2) + t_{\alpha/2} \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

La règle de décision : ii)  $H_1 : \mu_1 - \mu_2 < 0$

$$\text{Rejeter } H_0 \text{ si } (\bar{X}_1 - \bar{X}_2) < (\mu_1 - \mu_2) - t_{\alpha} \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

La règle de décision : iii)  $H_1 : \mu_1 - \mu_2 > 0$

$$\text{Rejeter } H_0 \text{ si } (\bar{X}_1 - \bar{X}_2) > (\mu_1 - \mu_2) + t_{\alpha} \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

Test d'hypothèse pour observations couplées

- On veut comparer les moyennes des 2 caractéristiques provenant de  $m$  unités stats d'une pop. --- Les variables  $X_1$  et  $X_2$  ne sont plus indépendantes --- Les observations sont couplées

Observations couplées

- Soit les observations pour  $X_1 : X_{1,1}, X_{1,2}, \dots, X_{1,n}$
- Soit les observations pour  $X_2 : X_{2,1}, X_{2,2}, \dots, X_{2,n}$
- Ces observations sont couplées :  $(X_{1,1}, X_{2,1}), (X_{1,2}, X_{2,2}), \dots, (X_{1,n}, X_{2,n})$
- On calc la dif pour chak couple  $i : D_i = X_{1,i} - X_{2,i}$  pour  $i = 1, 2, \dots, n$
- Si  $D_i$  provient d'une pop. d'espérance  $\mu_D$  et de variance  $\sigma_D^2$  alors :  $\bar{D} = \frac{\sum D_i}{n}$  est un estim de  $\mu_D$ ,  $S_D^2 = \frac{\sum (D_i - \bar{D})^2}{n-1}$  est un estim de  $\sigma_D^2$

$\bar{D}$  est la valeur empirique

Inférence pour observations couplées

- Cas 1 : pop. normale, et  $\sigma_D^2$  connue
- Cas 2 : grand échantillon,  $n \geq 30$
- Cas 3 : pop. normale,  $\sigma_D^2$  inconnue, et  $n$  est petit

Test pour  $\mu_D$

- $H_0 : \mu_D = 0$ ,  $H_1 : \mu_D \neq 0$ ;  $\mu_D < 0$ ;  $\mu_D > 0$ ,  $\alpha$  fixé
- Conditions du test : 2 E.A. avec observ. couplées [cas 1, 2 ou 3]
- La statistique [cas 1, cas 2 =  $\sigma \rightarrow S$ , cas 3 =  $\sigma \rightarrow S, Z \rightarrow t$ ] :  $Z = \frac{\bar{D} - \mu_D}{\sigma_D} \in N(0,1)$  [cas 3 :  $t \in t(n-1)$ ]

La règle de décision : i)  $H_1 : \mu_D \neq 0$

$$\text{Rejeter } H_0 \text{ si } \bar{D} < \mu_D - z_{\alpha/2} \sqrt{\frac{\sigma_D^2}{n}} \text{ ou si } \bar{D} > \mu_D + z_{\alpha/2} \sqrt{\frac{\sigma_D^2}{n}} \text{ ou } = \left( \frac{\sigma_D^2}{n} \right)$$

La règle de déc : ii)  $H_1 : \mu_D < 0$ ; Rejeter  $H_0$  si  $\bar{D} < \mu_D - z_{\alpha} \sqrt{\frac{\sigma_D^2}{n}}$

La règle de déc : iii)  $H_1 : \mu_D > 0$ ; Rejeter  $H_0$  si  $\bar{D} > \mu_D + z_{\alpha} \sqrt{\frac{\sigma_D^2}{n}}$

Différence entre deux proportions

- Soit  $X_1$  et  $X_2$  2 V.A. ind distribués selon une loi binomiale de paramètre  $p_1$  et  $p_2$  respectivement et provenant de 2 pop. dif.
- On veut comparer la proportion de succès  $(p_1 - p_2)$  entre les 2 pop
- Soit des EA ind taille  $n_1$  et  $n_2$  provenant 2 pop bin. de  $p_1$  et  $p_2$  alors :  $(\hat{p}_1 - \hat{p}_2) \in N(p_1 - p_2, \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2})$ , autant  $n_1$  et  $n_2$  sont grands
- Donc en autant que  $n_1$  et  $n_2$  sont  $t$ , on a :  $E(\hat{p}_1 - \hat{p}_2) = p_1 - p_2$ ,  $\text{Var}(\hat{p}_1 - \hat{p}_2) = \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}$

Test d'hypothèse pour  $(p_1 - p_2)$  -  $\alpha$  fixé (p est un %)

$H_0 : p_1 - p_2 = 0$  ( $p_1 = p_2 = p$ )  $H_1 : p_1 - p_2 \neq 0$ ;  $p_1 - p_2 < 0$ ;  $p_1 - p_2 > 0$

Conditions du test : deux échantillons aléatoires indépendants,  $n_1$  et  $n_2$  grands ( $n_1 p_1, n_2 p_2, n_1(1-p_1), n_2(1-p_2)$  soient tous  $\geq 5$ )

Puisque  $p_1 = p_2 = p$ , on doit estimer  $p : \hat{p} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$

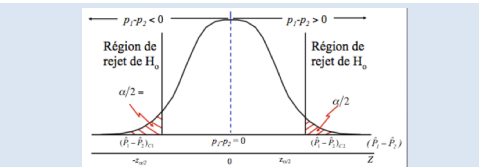
La statistique à utiliser est  $Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \in N(0,1)$

Puisque selon  $H_0$  on  $p_1 = p_2 = p$ , stat devin  $Z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \in N(0,1)$

La règle de décision : i)  $H_1 : p_1 - p_2 \neq 0$ , Rejeter  $H_0$  si  $|\hat{p}_1 - \hat{p}_2| > z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$

ii)  $H_1 : p_1 - p_2 < 0$ , Rejeter  $H_0$  si  $(\hat{p}_1 - \hat{p}_2) < -z_{\alpha} \sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$

iii)  $H_1 : p_1 - p_2 > 0$ , Rejeter  $H_0$  si  $(\hat{p}_1 - \hat{p}_2) > z_{\alpha} \sqrt{\hat{p}(1-\hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$



Distribution F de Fisher

Pour utiliser loi de Fisher, les obs doivent provenir de lois normales La loi de Fisher est formée par le rapport entre 2 V.A. indépendantes suivant une distribution de KHi-deux divisées par leurs dl  $0 \leq F \leq \infty$

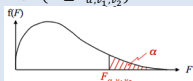
Propriétés de loi Fisher

- La courber de la loi de Fisher a une asymétrie positive
- La distribution de Fisher dépend uniquement de  $v_1$  et  $v_2$

Tableau du F de Fisher

La table de la loi de Fisher donne les probabilités,  $\alpha$ , en fonction du nombre de dl  $v_1$ , et  $v_2$ , et de la valeur du F

La table donne  $P(F \geq F_{\alpha; v_1, v_2})$



Relation complémentaire pour le F de Fisher - si on ne peut pas l'utiliser on se sert de  $F_{1-\alpha; v_1, v_2} = \frac{1}{F_{\alpha; v_2, v_1}}$

Distribution du quotient de deux variances  $S_1^2/S_2^2$

On veut comparer la variance de deux populations On tire un échantillon de taille  $n_1$  et  $n_2$  respectivement de 2 pop norm Pour comparer leur variance on se sert du quotient  $S_1^2/S_2^2$  Si la variable analysée pour les deux pop suit une distribution de probabilité normale alors  $\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \in F(n_1 - 1, n_2 - 1)$ ,  $0 \leq F \leq \infty$

Test d'hypothèse pour deux variances

$H_0 : \frac{\sigma_1^2}{\sigma_2^2} = 1 - H_1 : \sigma_1^2/\sigma_2^2 \neq 1$ ;  $\sigma_1^2/\sigma_2^2 < 1$ ;  $\sigma_1^2/\sigma_2^2 > 1 - \alpha$  fixé

Conditions du test : 2 échantillons aléatoires indépendantes, pop norm

La stat.  $\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \in F(n_1 - 1, n_2 - 1)$  sous  $H_0$ ,  $F = \frac{S_1^2}{S_2^2} \in F(n_1 - 1, n_2 - 1)$

La règle de décision : i)  $H_1 : \sigma_1^2/\sigma_2^2 \neq 1$ , Rejeter  $H_0$  si  $\frac{S_1^2}{S_2^2} <$

$$\frac{1}{F_{\alpha/2}(n_2-1, n_1-1)} \text{ ou si } \frac{S_1^2}{S_2^2} > \frac{1}{F_{\alpha/2}(n_1-1, n_2-1)}$$

ii)  $H_1 : \sigma_1^2/\sigma_2^2 < 1$ , Rejeter  $H_0$  si  $\frac{S_1^2}{S_2^2} < \frac{1}{F_{\alpha}(n_2-1, n_1-1)}$

iii)  $H_1 : \sigma_1^2/\sigma_2^2 > 1$ , Rejeter  $H_0$  si  $\frac{S_1^2}{S_2^2} > \frac{1}{1 - \alpha}$

Chapitre 9 - KHi-deux

Tableau croisé et test d'indépendance

- On veut tester l'hypothèse selon laquelle deux variables X et Y mesurées sur une échelle nominale ou ordinaire sont ind
- Les obs sont réparties selon les modalités croisées des 2 variables
- On obtient un tableau à 2 dimensions appelé tableau croisé ou tableau de contingence

Structure de tableau croisé

	Y				Total de la ligne
X	B <sub>1</sub>	B <sub>2</sub>	...	B <sub>k</sub>	
A <sub>1</sub>	f <sub>011</sub>	f <sub>012</sub>	...	f <sub>01k</sub>	L <sub>1</sub>
A <sub>2</sub>	f <sub>021</sub>	f <sub>022</sub>	...	f <sub>02k</sub>	L <sub>2</sub>
...	...	...	...	...	...
A <sub>r</sub>	f <sub>0r1</sub>	f <sub>0r2</sub>	...	f <sub>0rk</sub>	L <sub>r</sub>
Total de la colonne	C <sub>1</sub>	C <sub>2</sub>	...	C <sub>k</sub>	n

Les fréquences observées

- $f_{0ij}$  : fréquence observée des données ayant la modalité  $A_i$  de la variable X et la modalité  $B_j$  de la variable Y
- $C_j = \sum_{i=1}^r f_{0ij}$ ,  $j = 1, 2, \dots, k$
- $L_i = \sum_{j=1}^k f_{0ij}$ ,  $i = 1, 2, \dots, r$
- $\sum_{i=1}^r L_i = \sum_{j=1}^k C_j = \sum_{i=1}^r \sum_{j=1}^k f_{0ij} = n$
- $H_0$  : les deux caractères X et Y sont indépendants
- $H_1$  : les deux caractères X et Y ne sont pas indépendants
- Conditions du test : Seuil de signification du test,  $\alpha$ , est fixé, Échantillon aléatoire de taille  $n$ , Les fréquences théoriques sont supérieures ou égales à 5 (les fréquences observées sont tjrs données)

Calcul du  $\chi^2$  pour le test d'indépendance (calcul chak donnée)

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^k \frac{(f_{0ij} - f_{t_{ij}})^2}{f_{t_{ij}}} \in \chi^2((r-1)(k-1)) \text{ où } f_{t_{ij}} = \frac{L_i C_j}{n}$$

1. On rejette  $H_0$  si  $\chi^2 > \chi^2_{\alpha; (r-1)(k-1)}$  ( $\chi^2$  théorique  $>$   $\chi^2$  calculée)

Mesure d'association entre deux caractères

- Coefficient de contingence  $v$  de Cramer,  $v = \sqrt{\frac{\chi^2}{m(k-1)}}$ , où  $k$  est le min entre le # de lignes et le nombre de colonnes et  $0 \leq v \leq 1$

Rappel :  $P(A \cap B) = P(A) \times P(B)$  seulement si indépendant

Test d'ajustement analytique (test du  $\chi^2$  de Pearson)

- On veut vérifier si une distribution de fréquence empirique provient d'une distribution théorique connue
- Le test du  $\chi^2$  permet de vérifier s'il y a une différence significative entre les  $f_{0i}$  et les  $f_{t_{i}}$  selon la distribution théorique

Distribution de fréquences observées

- Les observations de l'échantillon sont réparties en classe (distribution de fréquences)
  - $L_1 \leq X \leq L_{S_1}, F_{01}$
  - $L_2 \leq X \leq L_{S_2}, F_{02}$  ...
  - $L_k \leq X \leq L_{S_k}, F_{0k}$   $\rightarrow$  fréquence observée dans la classe  $k$
  - $\sum_{i=1}^k F_{0i} = n$ , où  $n$  : taille de l'échantillon
- Distribution de fréquences théoriques
  - Calculer les fréquences théoriques selon  $H_0$ 
    - $p_1 = P(L_1 \leq X \leq L_{S_1}), F_{t1} = p_1 n$
    - $p_2 = P(L_2 \leq X \leq L_{S_2}), F_{t2} = p_2 n$  ...
    - $p_k = P(L_k \leq X \leq L_{S_k}), F_{tk} = p_k n$
    - $\sum_{i=1}^k p_i = 1, \sum_{i=1}^k F_{ti} = n$
- Test d'ajustement analytique
  - $H_0$  : les obsv proviennent de la distribution théorique  $f(x; \theta)$
  - lois discrètes (binomiale, Poisson, ...)
  - lois continues (exponentielle, uniforme, normale, ...)
- $H_1$  : les obsv ne proviennent pas de la distribution théorique  $f(x; \theta)$
- $\alpha$  : le seuil de signification est fixé
- Conditions du test (même que le test d'indépendance + Échantillon aléatoire,  $n$  est grand ( $n \geq 50$ ), tel que les  $F_{t_{ij}} \geq 5$ , (on peut regrouper les classes aux extrémités de la distribution pour respecter cette condition)

- La statistique :  $\chi^2 = \sum_{i=1}^k \frac{(f_{0i} - f_{t_i})^2}{f_{t_i}} \in \chi^2(k-1-r)$ , où  $r$  est le # de paramètres à estimer de la distribution théorique et  $k$  le # de classes
- Région critique : rejeter  $H_0$  si  $\chi^2 > \chi^2_{\alpha; k-1-r}$
- Résultat empirique : calculer  $\chi^2 = \sum_{i=1}^k \frac{(f_{0i} - f_{t_i})^2}{f_{t_i}}$

Le rejet de  $H_0$  peut provenir de trois causes :

- la durée de vie ne suit pas une loi normale
- la moyenne de la durée de vie n'est pas  $m$
- l'écart-type de la durée de vie n'est pas  $s$

Les causes 2) et 3) peuvent