

2 Rappel de Probabilité

2.1 Probabilité

Proposition 2 Propriétés des probabilités :

- $P(\emptyset) = 0$;
- $P(A) = 1 - P(A^c)$;
- $A \subset B \Rightarrow P(A) \leq P(B)$;
- $A \cap B \Rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$;
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$;

Proposition 3 A et B sont indépendants si $P(A \cap B) = P(A)P(B)$, n événements A_1, \dots, A_n sont indépendants si pour toute partie I de $\{1, \dots, n\}$: $P(\cap_{i \in I} A_i) = \prod_{i \in I} P(A_i)$.

2.2 Variables Aléatoires

Proposition 4 Propriétés de la densité

- $f_A(x) \geq 0$
- $\int_{\mathbb{R}} f_A(x) dx = 1$
- $F_A(x) = \int_{-\infty}^x f_A(t) dt$

2.3 Fonction de répartition

Définition 8 $F : \mathbb{R} \rightarrow [0, 1]$ définie par $\forall x \in \mathbb{R} : F(x) = P(X \leq x)$

Proposition 5 Propriété de la fonction de répartition

- F est croissante et continue à droite.
- $\lim_{x \rightarrow -\infty} F(x) = 0$
- $\lim_{x \rightarrow +\infty} F(x) = 1$
- $F(a) \leq F(b) = F(\theta) - F(a)$
- Si X est continue, alors F est dérivable et on a $F'(x) = f(x)$

2.4 Espérance mathématique

Proposition 6 Propriétés de l'espérance

- Si $X \geq 0$, alors $E(X) \geq 0$
- $E(a) = a$
- $E(aX + bY) = aE(X) + bE(Y)$

2.5 Variance

Proposition 7 Propriétés de la variance

- $E[(X - a)^2] = Var(X) + E[X(X - a)^2]$
- $Var(aX + b) = a^2 Var(X)$
- $E[(X - E(X))(Y - E(Y))] = Cov(X, Y)$

2.6 Autres moments

Définition 9

- Moment centré d'ordre k : $m_k = E[(X - E(X))^k]$
- Moment non centré d'ordre k : $m_k = E[X^k]$

2.7 Convergence Stochastique

Définition 10 La suite (X_n) converge en probabilité vers a si $\forall \epsilon > 0$ et n_0 tel que $\forall n > n_0$

$$P(|X_n - a| < \epsilon) > 1 - \eta$$

On note $(X_n) \xrightarrow{P} a$. La convergence en probabilité de X_n vers a revient à étudier la convergence de $X_n - X_n$ vers 0. Condition : Si $E(X_n) \xrightarrow{P} a$ et $Var(X_n) \rightarrow 0$, alors $X_n \xrightarrow{P} a$.

1 Statistiques Descriptives

1.1 Variable Discrète

Définition 1 Soit x une variable discrète, à valeurs dans $\mathcal{K} = \{\xi_1, \dots, \xi_K\}$ avec $\xi_1 < \dots < \xi_K$. Une distribution de n observations de x peut être représentée sous forme d'un tableau de fréquences où figure pour chaque modalité de ξ_k de x, le nombre n_k d'observations ayant la valeur ξ_k . La fréquence relative est $f_k = \frac{n_k}{n}$. La fréquence cumulée est $F_k = \sum_{j=1}^k f_j$

1.2 Variable Continue

Définition 2 Si x est continue, on partitionne le domaine de définition de x en K classes. On prendra la règle de Sturges pour calculer K. $K = 1 + \frac{3}{10} \log_{10} n$.

Pour une variable discrète, on va utiliser un diagramme en bâton, alors que pour une variable continue, on va utiliser un histogramme, où l'aire de chaque rectangle est proportionnelle à l'effectif.

1.3 Fonction de répartition empirique

Définition 3 La fonction de répartition empirique est :

$$F : \mathbb{R} \rightarrow [0, 1]$$

1.4 Indicateur de tendance centrale

Définition 4 Moyenne empirique : $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. Dans le cas des classes, on a : $\bar{x} = \frac{1}{n} \sum_{k=1}^K n_k \xi_k$

1.5 Indicateur de dispersion

Définition 6 La variance empirique : $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$

La variance empirique corrigée : $s^{*2} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n}{n-1} s^2$

1.6 Boîte à Moustaches

Définition 7 Il s'agit d'un graphique formé d'une boîte définissant la médiane par le premier et le troisième quartile, où figure aussi le médiane ($f_{0.5}$). Des segments de droites s'étendent de part et d'autre de la boîte jusqu'au point le plus extrême à une distance inférieure à 1.5H. Les autres points sont représentés individuellement. H est la hauteur de la boîte.

2 Intervalles de confiance

2.1 Définition

X_1, \dots, X_n étant un échantillon i.i.d dont la loi parente X a une distribution dépendante du paramètre réel θ , α un réel quelconque fixé appartenant à $]0, 1[$, on appelle intervalle de confiance pour θ au niveau de confiance $1 - \alpha$ tout intervalle $[T_1, T_2]$ où T_1 et T_2 sont deux statistiques vérifiant $P(T_1 \leq \theta \leq T_2) = 1 - \alpha$.

2.2 Construction

- Choisir un estimateur T de θ dont on connaît la loi de probabilité en fonction de θ ;
- Déterminer $f(T, \theta)$ dont la loi de probabilité ne dépend plus de θ : la fonction *pivotal*;
- Si c'est possible, en déduire $P(T_1 \leq \theta \leq T_2) = 1 - \alpha$ où T_1 et T_2 sont des statistiques fonctions de T.

2.3 Variable normale

Hypothèse : $X_1, \dots, X_n, X_i \sim \mathcal{N}(\mu, \sigma^2)$, niveau de confiance $1 - \alpha$, intervalle bilatéral.

2.3.1 Moyenne

1. Si σ est connue, on a $T = \bar{X}$ (estimateur) et $\theta = \mu$, on sait que $f(T, \theta) = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim \mathcal{N}(0, 1)$.

Comme $P(u_{\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq u_{1 - \frac{\alpha}{2}}) = 1 - \alpha$, on obtient comme I.C. $[\bar{X} - \frac{\sigma}{\sqrt{n}} u_{1 - \frac{\alpha}{2}}, \bar{X} + \frac{\sigma}{\sqrt{n}} u_{1 - \frac{\alpha}{2}}]$. En unilatéral, on écrit $P(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq u_{1 - \alpha}) = 1 - \alpha$ et on obtiendrait $[\bar{X} - \frac{\sigma}{\sqrt{n}} u_{1 - \alpha}, +\infty[$.

2. Sinon, on remplace σ par S^* pour remarquer que

$$\frac{\bar{X} - \mu}{\frac{S^*}{\sqrt{n}}} = \frac{\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}}{\frac{S^*}{\frac{\sigma}{\sqrt{n}}}} \sim \tau_{n-1} \text{ (pleure pas).}$$

On obtient alors $P(-t_{n-1; 1 - \frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{\frac{S^*}{\sqrt{n}}} \leq t_{n-1; 1 - \frac{\alpha}{2}}) = 1 - \alpha$ d'où I.C. $[\bar{X} - \frac{S^*}{\sqrt{n}} t_{n-1; 1 - \frac{\alpha}{2}}, \bar{X} + \frac{S^*}{\sqrt{n}} t_{n-1; 1 - \frac{\alpha}{2}}]$.

2.3.2 Variance

1. Si μ est connue, l'estimateur T est l'estimateur de maximum de vraisemblance $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$, d'où $f(T, \theta) = \frac{\frac{n\hat{\sigma}^2}{\sigma^2}}{\chi_n^2} \sim \chi_n^2$ et on utilise la relation $P(\frac{n\hat{\sigma}^2}{\sigma^2} < \sigma^2 < \frac{n\hat{\sigma}^2}{\chi_{n-1}^2})$

3 Tests d'hypothèse : comparaison de 2 échantillons

3.1 Cas particuliers

3.1.1 Hyp. simple vs Hyp. simple

On veut tester

$$\begin{cases} H_0 : \theta = \theta_0 \\ H_1 : \theta = \theta_1 (\neq \theta_0) \end{cases}$$

Théorème 6 (Neyman et Pearson) La région critique optimale vérifie une relation de la forme $\frac{L(\theta_1; x_1, \dots, x_n)}{L(\theta_0; x_1, \dots, x_n)} > \lambda$, où λ se calcule par $P_{H_0}((X_1, \dots, X_n) \in W) = \alpha^*$

3.1.2 Hyp. simple vs Hyp. composite

1. On veut résoudre

$$\begin{cases} H_0 : \theta = \theta_0 \\ H_1 : \theta \in E(\theta_0 \notin E) \end{cases} \Leftrightarrow \begin{cases} H_0 : \theta = \theta_0 \\ H_1 : \theta = \theta_1(\theta_1 \in E) \end{cases}$$

Si la région critique W_i est indépendante de i et toujours égale à W , elle est optimale pour le test $\theta = \theta_0$ contre $\theta \in E$. L'erreur de seconde espèce sera une fonction de θ , mais sera toujours minimale pour α^* donné. Le test est donc dit *UPP* (Uniformément Plus Puissant).

2. Cas non *UPP* : si W_1 n'est pas constant, la stratégie de Neyman-Pearson ne peut s'appliquer. On conserve la contrainte $\alpha = P_{H_0}((X_1, \dots, X_n) \in W) \leq \alpha^*$ fixé et on recherche une région critique en utilisant une fonction *pivotal*.

3.1.3 Hyp. composite vs Hyp. composite

1. Test *UPP*

Définition 29 Une famille de lois L_θ est une famille à rapport de vraisemblance monotone s'il existe une statistique T telle que $\frac{L(\theta_2; x_1, \dots, x_n)}{L(\theta_1; x_1, \dots, x_n)}$ soit une fonction croissante de t pour $\theta' > \theta$.

Théorème 7 (Lehman) Si la famille des lois L_θ est à rapport de vraisemblance monotone, alors le problème de test

$$\begin{cases} H_0 : \theta \leq \theta_0 \\ H_1 : \theta > \theta_0 \end{cases}$$

admet un test *UPP* défini par une région critique de la forme $T > A$ où A est défini par $P_{\theta_0}(T > A) = \alpha^*$.

2. Test non *UPP*

Forme générale : X suit une loi dépendant de plusieurs paramètres dont le paramètre réel θ sur lequel porte le test :

$$\begin{cases} H_0 : \theta = \theta_0 & \text{autres paramètres inconnus} \\ H_1 : \theta > \theta_0, \neq \theta_0 & \text{autres paramètres inconnus} \end{cases}$$

Stratégie :

- Trouver une fonction *pivotal* de θ
- Proposer une forme de région critique raisonnable
- Déterminer cette région critique en prenant $\alpha = P_{H_0}(X \in W) = \alpha^*$.

4 Statistiques Descriptives

4.1 Variable Discrète

Définition 1 Soit x une variable discrète, à valeurs dans $\mathcal{K} = \{\xi_1, \dots, \xi_K\}$ avec $\xi_1 < \dots < \xi_K$. Une distribution de n observations de x peut être représentée sous forme d'un tableau de fréquences où figure pour chaque modalité de ξ_k de x, le nombre n_k d'observations ayant la valeur ξ_k . La fréquence relative est $f_k = \frac{n_k}{n}$. La fréquence cumulée est $F_k = \sum_{j=1}^k f_j$

4.2 Variable Continue

Définition 2 Si x est continue, on partitionne le domaine de définition de x en K classes. On prendra la règle de Sturges pour calculer K. $K = 1 + \frac{3}{10} \log_{10} n$.

Pour une variable discrète, on va utiliser un diagramme en bâton, alors que pour une variable continue, on va utiliser un histogramme, où l'aire de chaque rectangle est proportionnelle à l'effectif.

4.3 Fonction de répartition empirique

Définition 3 La fonction de répartition empirique est :

$$F : \mathbb{R} \rightarrow [0, 1]$$

4.4 Indicateur de tendance centrale

Définition 4 Moyenne empirique : $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. Dans le cas des classes, on a : $\bar{x} = \frac{1}{n} \sum_{k=1}^K n_k \xi_k$

4.5 Indicateur de dispersion

Définition 6 La variance empirique : $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$

La variance empirique corrigée : $s^{*2} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n}{n-1} s^2$

4.6 Boîte à Moustaches

Définition 7 Il s'agit d'un graphique formé d'une boîte définissant la médiane par le premier et le troisième quartile, où figure aussi le médiane ($f_{0.5}$). Des segments de droites s'étendent de part et d'autre de la boîte jusqu'au point le plus extrême à une distance inférieure à 1.5H. Les autres points sont représentés individuellement. H est la hauteur de la boîte.

2 Tests d'hypothèse : comparaison de 2 échantillons

2.1 Cas particuliers

2.1.1 Hyp. simple vs Hyp. simple

On veut tester

$$\begin{cases} H_0 : \theta = \theta_0 \\ H_1 : \theta = \theta_1 (\neq \theta_0) \end{cases}$$

Théorème 6 (Neyman et Pearson) La région critique optimale vérifie une relation de la forme $\frac{L(\theta_1; x_1, \dots, x_n)}{L(\theta_0; x_1, \dots, x_n)} > \lambda$, où λ se calcule par $P_{H_0}((X_1, \dots, X_n) \in W) = \alpha^*$

2.1.2 Hyp. simple vs Hyp. composite

1. On veut résoudre

$$\begin{cases} H_0 : \theta = \theta_0 \\ H_1 : \theta \in E(\theta_0 \notin E) \end{cases} \Leftrightarrow \begin{cases} H_0 : \theta = \theta_0 \\ H_1 : \theta = \theta_1(\theta_1 \in E) \end{cases}$$

Si la région critique W_i est indépendante de i et toujours égale à W , elle est optimale pour le test $\theta = \theta_0$ contre $\theta \in E$. L'erreur de seconde espèce sera une fonction de θ , mais sera toujours minimale pour α^* donné. Le test est donc dit *UPP* (Uniformément Plus Puissant).

2. Cas non *UPP* : si W_1 n'est pas constant, la stratégie de Neyman-Pearson ne peut s'appliquer. On conserve la contrainte $\alpha = P_{H_0}((X_1, \dots, X_n) \in W) \leq \alpha^*$ fixé et on recherche une région critique en utilisant une fonction *pivotal*.

2.1.3 Hyp. composite vs Hyp. composite

1. Test *UPP*

Définition 29 Une famille de lois L_θ est une famille à rapport de vraisemblance monotone s'il existe une statistique T telle que $\frac{L(\theta_2; x_1, \dots, x_n)}{L(\theta_1; x_1, \dots, x_n)}$ soit une fonction croissante de t pour $\theta' > \theta$.

Théorème 7 (Lehman) Si la famille des lois L_θ est à rapport de vraisemblance monotone, alors le problème de test

$$\begin{cases} H_0 : \theta \leq \theta_0 \\ H_1 : \theta > \theta_0 \end{cases}$$

admet un test *UPP* défini par une région critique de la forme $T > A$ où A est défini par $P_{\theta_0}(T > A) = \alpha^*$.

2. Test non *UPP*

Forme générale : X suit une loi dépendant de plusieurs paramètres dont le paramètre réel θ sur lequel porte le test :

$$\begin{cases} H_0 : \theta = \theta_0 & \text{autres paramètres inconnus} \\ H_1 : \theta > \theta_0, \neq \theta_0 & \text{autres paramètres inconnus} \end{cases}$$

Stratégie :

- Trouver une fonction *pivotal* de θ
- Proposer une forme de région critique raisonnable
- Déterminer cette région critique en prenant $\alpha = P_{H_0}(X \in W) = \alpha^*$.

2.2 Moyennes indépendantes de loi de proba respectives $\chi_{\nu_1}^2$ et $\chi_{\nu_2}^2$, alors la variable $\frac{X_1/\nu_1}{Y_2/\nu_2}$ suit une loi de Fisher à ν_1 et ν_2 degrés de liberté.

2.3 Variances : test de Fisher

Définition 30 (loi de Fisher) Si X et Y sont 2 variables aléatoires indépendantes de loi de proba respectives $\chi_{\nu_1}^2$ et $\chi_{\nu_2}^2$, alors la variable $\frac{X/\nu_1}{Y/\nu_2}$ suit une loi de Fisher à ν_1 et ν_2 degrés de liberté.

$$F_{\nu_1, \nu_2, \alpha} = \frac{1}{F_{\nu_2, \nu_1, 1-\alpha}} \quad (\text{voir les tables})$$

2.4 Moyennes : test de Student

Estimateur de la variance commune à 2 populations gaussiennes ($N = n + m$) : $S^{*2} = \frac{1}{N-2} ((n-1)S_X^2 + (m-1)S_Y^2)$.

Remarque : $\frac{(N-2)S^{*2}}{\sigma^2} = \frac{(n-1)S_X^2}{\sigma^2} + \frac{(m-1)S_Y^2}{\sigma^2} \sim \chi_{N-2}^2$

$X \sim \mathcal{N}(\mu_X, \sigma^2)$ et $Y \sim \mathcal{N}(\mu_Y, \sigma^2)$:

$$\begin{cases} H_0 : \mu_X = \mu_Y & (HC) \\ H_1 : \mu_X \neq \mu_Y & (HC) \end{cases}$$

Région critique :

- $\frac{\bar{X} - \bar{Y}}{S^* \sqrt{\frac{1}{n} + \frac{1}{m}}} < t_{N-2, \alpha/2}$
- $\frac{\bar{X} - \bar{Y}}{S^* \sqrt{\frac{1}{n} + \frac{1}{m}}} > t_{N-2, 1-\alpha/2}$

2.5 Tests non paramétriques

2.5.1 Test de Wilcoxon

Suite W_X : fusion des réalisations des deux échantillons et classement par ordre croissant.

Si les deux échantillons on la même distribution ($F_X = F_Y$) alors $E(W_X) = \frac{n(n+m+1)}{2}$ et $Var(W_X) = \frac{nm(n+m+1)}{12}$

Région critique : $W_X < \frac{n(n+m+1)}{2} - u_{1-\alpha} \sqrt{\frac{nm(n+m+1)}{12}}$

2.5.2 Tests non paramétriques

2.5.2.1 Test de Wilcoxon

Suite W_X : fusion des réalisations des deux échantillons et classement par ordre croissant.

Si les deux échantillons on la même distribution ($F_X = F_Y$) alors $E(W_X) = \frac{n(n+m+1)}{2}$ et $Var(W_X) = \frac{nm(n+m+1)}{12}$

Région critique : $W_X < \frac{n(n+m+1)}{2} - u_{1-\alpha} \sqrt{\frac{nm(n+m+1)}{12}}$

2.5.3 Échantillons appariés

Pour l'étude de couples, en règle général, on étudie la différence : $D_i = Y_i - X_i$.

$$\begin{cases} H_0 : m_e = 0 \\ H_1 : m_e > 0 \text{ ou } < 0 \text{ ou } \neq 0 \end{cases}$$