

# Statistique

Vincent Charvillat et Jean-Yves Tourneret<sup>(1)</sup>

(1) Université de Toulouse, ENSEEIHT-IRIT

[Vincent.Charvillat@enseeiht.fr](mailto:Vincent.Charvillat@enseeiht.fr) et [jyt@n7.fr](mailto:jyt@n7.fr)

# Plan du cours

## ● Chapitre 1 : Estimation

- Modèle statistique, qualités d'un estimateur, exemples
- Inégalité de Cramér Rao
- Maximum de vraisemblance
- Méthode des moments
- Estimation Bayésienne
- Intervalles de confiance

## ● Chapitre 2 : Tests Statistiques

# Bibliographie

- B. Lacaze, M. Maubourguet, C. Mailhes et J.-Y. Tourneret, Probabilités et Statistique appliquées, Cépadues, 1997.
- Athanasios Papoulis and S. Unnikrishna Pillai, Probability, Random Variable and Stochastic Processes, McGraw Hill Higher Education, 4th edition, 2002.

# Modèle Statistique

- Observations

$$x_1, \dots, x_n$$

- Échantillon

$$X_1, \dots, X_n$$

$n$  va iid associées aux observations

- Estimateur

$$\hat{\theta}(X_1, \dots, X_n) \text{ ou } \hat{\theta}_n \text{ ou } \hat{\theta}$$

# Qualités d'un estimateur

•  $\theta \in \mathbb{R}$

• **Biais** (erreur systématique) :  $b_n(\theta) = E(\hat{\theta}_n) - \theta$

• **Variance**

$$v_n(\theta) = E\left[\left(\hat{\theta}_n - E(\hat{\theta}_n)\right)^2\right] = E\left[\hat{\theta}_n^2\right] - E(\hat{\theta}_n)^2$$

• **Erreur quadratique moyenne** (précision)

$$e_n(\theta) = E\left[\left(\hat{\theta}_n - \theta\right)^2\right] = v_n(\theta) + b_n^2(\theta)$$

CS de **convergence** :  $\hat{\theta}_n$  est un estimateur  
convergent si  $\lim_{n \rightarrow +\infty} b_n(\theta) = \lim_{n \rightarrow +\infty} v_n(\theta) = 0$

# Qualités d'un estimateur

- $\boldsymbol{\theta} \in \mathbb{R}^p$

- Biais

$$b_n(\boldsymbol{\theta}) = E \left( \widehat{\boldsymbol{\theta}}_n \right) - \boldsymbol{\theta} \in \mathbb{R}^p$$

- Matrice de covariance

$$E \left[ \left( \widehat{\boldsymbol{\theta}}_n - E \left( \widehat{\boldsymbol{\theta}}_n \right) \right) \left( \widehat{\boldsymbol{\theta}}_n - E \left( \widehat{\boldsymbol{\theta}}_n \right) \right)^T \right]$$

# Exemples

- **Exemple 1** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\theta = m$  et  $\sigma^2$  connue

- Moyenne empirique

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

- Autre estimateur

$$\tilde{\theta}_n = \frac{2}{n(n+1)} \sum_{i=1}^n i X_i$$

- **Exemple 2** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\theta = \sigma^2$ ,  $m$  connue ou inconnue

- **Exemple 3** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\boldsymbol{\theta} = (m, \sigma^2)^T$

# Plan du cours

- **Chapitre 1 : Estimation**

- Modèle statistique, qualités d'un estimateur, exemples

- Inégalité de Cramér Rao

- Maximum de vraisemblance

- Méthode des moments

- Estimation Bayésienne

- Intervalles de confiance

- **Chapitre 2 : Tests Statistiques**



# Inégalité de Cramér Rao

- **Vraisemblance**

$$L(x_1, \dots, x_n; \theta) = \begin{cases} X_i \text{ va discrète : } P[X_1 = x_1, \dots, X_n = x_n; \theta] \\ X_i \text{ va continue : } p(x_1, \dots, x_n; \theta) \end{cases}$$

- **Inégalité pour  $\theta \in \mathbb{R}$**

- **Définition**

$$\text{Var}(\hat{\theta}_n) \geq \frac{[1 + b'_n(\theta)]^2}{-E\left[\frac{\partial^2 \ln L(X_1, \dots, X_n; \theta)}{\partial \theta^2}\right]} = \text{BCR}(\theta)$$

BCR( $\theta$ ) est appelée **Borne de Cramér Rao** de  $\theta$

- **Hypothèses**

Log-vraisemblance deux fois dérivable et support de la loi indépendant de  $\theta$  (contre-exemple : loi  $\mathcal{U}[0, \theta]$ )

# Inégalité de Cramér Rao

- **Estimateur Efficace** : estimateur sans biais tel que

$$\text{Var} \left( \hat{\theta}_n \right) = \text{BCR}(\theta) \text{ (Il est unique !)}$$

**Exemple** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\theta = m$  et  $\sigma^2$  connue

- **Cas où  $(X_1, \dots, X_n)$  est un échantillon**

$$\text{Var} \left( \hat{\theta}_n \right) \geq \frac{[1 + b'_n(\theta)]^2}{-nE \left[ \frac{\partial^2 \ln L(X_1; \theta)}{\partial \theta^2} \right]} = \text{BCR}(\theta)$$

# Cas multivarié

- Inégalité pour un estimateur non biaisé de  $\boldsymbol{\theta} \in \mathbb{R}^p$

$$\text{Cov}(\widehat{\boldsymbol{\theta}}) \geq I_n^{-1}(\boldsymbol{\theta})$$

avec  $I_{ij} = E \left[ -\frac{\partial^2 \ln L(X_1, \dots, X_n; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right]$  pour  $i, j = 1, \dots, p$  et  
 $A \geq B$  signifie  $A - B$  matrice semi définie positive

$$x^T (A - B)x \geq 0, \quad \forall x \in \mathbb{R}^p$$

On en déduit

$$\text{Var}(\widehat{\theta}_i) \geq [I_n^{-1}(\boldsymbol{\theta})]_{ii}$$

**Exemple** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\boldsymbol{\theta} = (m, \sigma^2)^T$ .

# Plan du cours

## ● Chapitre 1 : Estimation

- Modèle statistique, qualités d'un estimateur, exemples
- Inégalité de Cramér Rao
- Maximum de vraisemblance
- Méthode des moments
- Estimation Bayésienne
- Intervalles de confiance

## ● Chapitre 2 : Tests Statistiques

# Méthode du Maximum de Vraisemblance

## • Définition

$$\hat{\theta}_{MV} = \arg \max_{\theta} L(X_1, \dots, X_n; \theta)$$

## • Recherche du maximum pour $\theta \in \mathbb{R}$

Si  $L(X_1, \dots, X_n; \theta)$  est deux fois dérivable et si les bornes de la loi de  $X_i$  sont indépendantes de  $\theta$

$$\frac{\partial L(X_1, \dots, X_n; \theta)}{\partial \theta} = 0 \text{ ou } \frac{\partial \ln L(X_1, \dots, X_n; \theta)}{\partial \theta} = 0$$

On vérifie qu'on a bien un maximum en étudiant

$$\frac{\partial \ln L(X_1, \dots, X_n; \theta)}{\partial \theta} \geq 0 \text{ ou } \frac{\partial^2 \ln L(X_1, \dots, X_n; \hat{\theta}_{MV})}{\partial \theta^2} < 0$$

# Méthode du Maximum de Vraisemblance

- Recherche du maximum pour  $\theta \in \mathbb{R}^p$

$$\frac{\partial L(X_1, \dots, X_n; \theta)}{\partial \theta_i} = 0 \text{ ou } \frac{\partial \ln L(X_1, \dots, X_n; \theta)}{\partial \theta_i} = 0$$

pour  $i = 1, \dots, p$

- Exemples

- Exemple 1 :  $X_i \sim \mathcal{P}(\lambda)$ ,  $\theta = \lambda$

- Exemple 2 :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\boldsymbol{\theta} = (m, \sigma^2)^T$

# Propriétés

- Estimateur **asymptotiquement non biaisé**

$$\lim_{n \rightarrow +\infty} E \left[ \hat{\boldsymbol{\theta}}_{MV} \right] - \boldsymbol{\theta} = 0$$

- Estimateur **convergent**
- Estimateur **asymptotiquement efficace**

$$\lim_{n \rightarrow +\infty} \frac{\text{Var} \left( \hat{\theta}_i \right)}{\left[ I_n^{-1}(\boldsymbol{\theta}) \right]_{ii}} = 1$$

- Normalité Asymptotique**

# Propriétés

- **Invariance Fonctionnelle**

Si  $\mu = h(\theta)$ , où  $h$  est une fonction bijective d'un ouvert  $O \subset \mathbb{R}^p$  dans un ouvert  $V \subset \mathbb{R}^p$ , alors

$$\hat{\mu}_{MV} = h\left(\hat{\theta}_{MV}\right)$$

- **Conclusions**

L'estimateur du maximum de vraisemblance  $\hat{\theta}_{MV}$  possède **beaucoup de bonnes propriétés asymptotiques** mais peut être **difficile à étudier** car il est la solution d'un problème d'optimisation.



# Plan du cours

## ● Chapitre 1 : Estimation

- Modèle statistique, qualités d'un estimateur, exemples

- Inégalité de Cramér Rao

- Maximum de vraisemblance

- Méthode des moments

- Estimation Bayésienne

- Intervalles de confiance

## ● Chapitre 2 : Tests Statistiques

# Méthode des moments

- **Définition** : supposons que  $X_1, \dots, X_n$  ont la même loi de paramètre inconnu  $\theta \in \mathbb{R}^p$ . En général, le vecteur paramètre à estimer  $\theta$  est lié aux premiers moments de la loi commune des va  $X_i$  par une relation notée

$$\theta = h(m_1, \dots, m_q)$$

avec  $m_k = E[X_i^k]$  et  $q \geq p$ . Un estimateur des moments de  $\theta$  est défini par

$$\hat{\theta}_{\text{Mo}} = h(\hat{m}_1, \dots, \hat{m}_q) \text{ avec } \hat{m}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

- **Exemple** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\theta = (m, \sigma^2)^T$

# Propriétés

- Estimateur convergent
- Normalité Asymptotique
- Conclusion : peu de propriétés mais cet estimateur est généralement facile à étudier.

# Plan du cours

## ● Chapitre 1 : Estimation

- Modèle statistique, qualités d'un estimateur, exemples
- Inégalité de Cramér Rao
- Maximum de vraisemblance
- Méthode des moments
- Estimation Bayésienne
- Intervalles de confiance

## ● Chapitre 2 : Tests Statistiques

# Estimation Bayésienne

- **Principe** : l'estimation Bayésienne consiste à estimer un vecteur paramètre inconnu  $\theta \in \mathbb{R}^p$  à l'aide de la **vraisemblance** de  $X_1, \dots, X_n$  (paramétrée par  $\theta$ ) et d'une **loi a priori**  $p(\theta)$ . Pour cela, on minimise une fonction de coût  $c(\theta, \hat{\theta})$  qui représente l'erreur entre  $\theta$  et  $\hat{\theta}$ .

- **Estimateur MMSE** : l'estimateur qui minimise l'erreur quadratique moyenne  $c(\theta, \hat{\theta}) = E \left[ (\theta - \hat{\theta})^2 \right]$  est

$$\hat{\theta}_{\text{MMSE}} = E(\theta | X_1, \dots, X_n)$$

- **Remarque** :  $p(\theta | x_1, \dots, x_n)$  est la **loi a posteriori** de  $\theta$

- **Preuve** : voir cours

# Estimation Bayésienne

- **Estimateur MAP** : l'estimateur du maximum a posteriori (MAP) est défini par

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta | X_1, \dots, X_n)$$

Cet estimateur minimise la fonction de coût  $E [c(\theta, \hat{\theta})]$  avec

$$c(\theta, \hat{\theta}) = \begin{cases} 1 & \text{si } \left\| \theta - \hat{\theta} \right\| > \Delta \\ 0 & \text{si } \left\| \theta - \hat{\theta} \right\| < \Delta \end{cases}$$

avec  $\Delta$  arbitrairement petit.

- **Preuve** : voir livre de H. Van Trees

# Exemple

- **Vraisemblance**

$$X_i \sim \mathcal{N}(\theta, \sigma^2)$$

- **Loi a priori**

$$\theta \sim \mathcal{N}(\mu, \nu^2)$$

- **Loi a posteriori**

$$\theta | X_1, \dots, X_n \sim \mathcal{N}(m_p, \sigma_p^2)$$

- **Estimateurs Bayésiens**

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{MMSE}} = m_p = \bar{X} \left( \frac{n\nu^2}{n\nu^2 + \sigma^2} \right) + \mu \left( \frac{\sigma^2}{\sigma^2 + n\nu^2} \right)$$

# Plan du cours

## ● Chapitre 1 : Estimation

- Modèle statistique, qualités d'un estimateur, exemples

- Inégalité de Cramér Rao

- Maximum de vraisemblance

- Méthode des moments

- Estimation Bayésienne

- Intervalles de confiance

## ● Chapitre 2 : Tests Statistiques



# Intervalles de confiance

- **Principe** : un intervalle de confiance  $[a, b]$  pour le paramètre  $\theta \in \mathbb{R}$  est un intervalle tel que  $P[a < \theta < b] = \alpha$ , où  $\alpha$  est le **paramètre de confiance** (en général  $\alpha = 0.99$  ou  $\alpha = 0.95$ ).
- **Détermination pratique de l'intervalle** : on cherche un **estimateur** de  $\theta$  noté  $\hat{\theta}$  (par la méthode des moments, du maximum de vraisemblance, ...), on en déduit une **statistique**  $T(X_1, \dots, X_n)$  qui dépend de  $\theta$  de loi connue, on cherche  $c(\theta)$  et  $d(\theta)$  tels que

$$P[c(\theta) < T(X_1, \dots, X_n) < d(\theta)] = \alpha$$

On en déduit l'intervalle  $[a, b]$ .

# Exemples

- **Exemple 1** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $m$  inconnue,  $\sigma^2$  connue.

$$T = \frac{\frac{1}{n} \sum_{i=1}^n X_i - m}{\sigma / \sqrt{n}} \sim \mathcal{N}(0, 1)$$

- **Exemple 3** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ , IC pour  $m$ ,  $\sigma^2$  inconnue.

$$T \sim \mathcal{N}(0, 1) \quad \text{et} \quad U = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi_{n-1}^2$$

donc

$$\frac{T}{\sqrt{\frac{U}{n-1}}} \sim t_{n-1}$$

suit une loi de **Student** à  $n - 1$  degrés de liberté.

# Que faut-il savoir ?

- Notions de **biais**, **variance** et **convergence** d'un estimateur
- Calcul d'une **borne de Cramér-Rao** et notion d'**efficacité**
- Détermination de l'estimateur du **maximum de vraisemblance** (MV)
- Propriétés de l'estimateur MV
- Principe et application de la **méthode des moments**
- Principe et application de l'**estimation Bayésienne**
- Détermination des **intervalles de confiance**

# Plan du cours

- Chapitre 1 : Estimation
- Chapitre 2 : Tests Statistiques
  - Généralités, exemple
  - Courbes COR
  - Théorème de Neyman Pearson
  - Test du rapport de vraisemblance généralisé
  - Test du  $\chi^2$
  - Test de Kolmogorov

# Généralités

- **Principe** : un test statistique est un mécanisme qui permet de décider entre plusieurs **hypothèses**  $H_0, H_1, \dots$  à partir de  $n$  observations  $x_1, \dots, x_n$ . On se limitera dans ce cours à deux hypothèses  $H_0$  et  $H_1$ . Effectuer un test, c'est déterminer une **statistique de test**  $T(X_1, \dots, X_n)$  et un **ensemble**  $\Delta$  tel que

$$\begin{aligned} \mathcal{H}_0 \text{ rejetée si } T(X_1, \dots, X_n) \in \Delta \\ \mathcal{H}_0 \text{ acceptée si } T(X_1, \dots, X_n) \notin \Delta. \end{aligned} \tag{1}$$

- **Vocabulaire**
  - $H_0$  est l'hypothèse **nulle**
  - $H_1$  est l'hypothèse **alternative**
  - $\{(x_1, \dots, x_n) | T(x_1, \dots, x_n) \in \Delta\}$  : **région critique**

# Définitions

- Tests **paramétriques** et **non paramétriques**
- Hypothèses **simples** et hypothèses **composites**
- **Risque de première espèce** = probabilité de fausse alarme

$$\alpha = \text{PFA} = P[\text{Rejeter } H_0 | H_0 \text{ vraie}]$$

- **Risque de seconde espèce** = probabilité de non-détection

$$\beta = \text{PND} = P[\text{Rejeter } H_1 | H_1 \text{ vraie}]$$

- **Puissance du test** = probabilité de détection :  $\pi = 1 - \beta$

# Exemple

$$X_i \sim \mathcal{N}(m, \sigma^2), \sigma^2 \text{ connue}$$

## • Hypothèses

$$H_0 : m = m_0, H_1 : m = m_1 > m_0$$

## • Stratégie du test

$$\text{Rejet de } H_0 \text{ si } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i > S_\alpha$$

## • Problèmes

Déterminer le seuil  $S_\alpha$ , le risque  $\beta$  et la puissance du test  $\pi$ .

# Plan du cours

- Chapitre 1 : Estimation
- Chapitre 2 : Tests Statistiques
  - Généralités, exemple
  - Courbes COR
  - Théorème de Neyman Pearson
  - Test du rapport de vraisemblance généralisé
  - Test du  $\chi^2$
  - Test de Kolmogorov



# Courbes COR

- Caractéristiques opérationnelles du récepteur

$$PD = h(\text{PFA})$$

- Exemple :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\sigma^2$  connue

$$H_0 : m = m_0, H_1 : m = m_1 > m_0$$

- Probabilité de fausse alarme

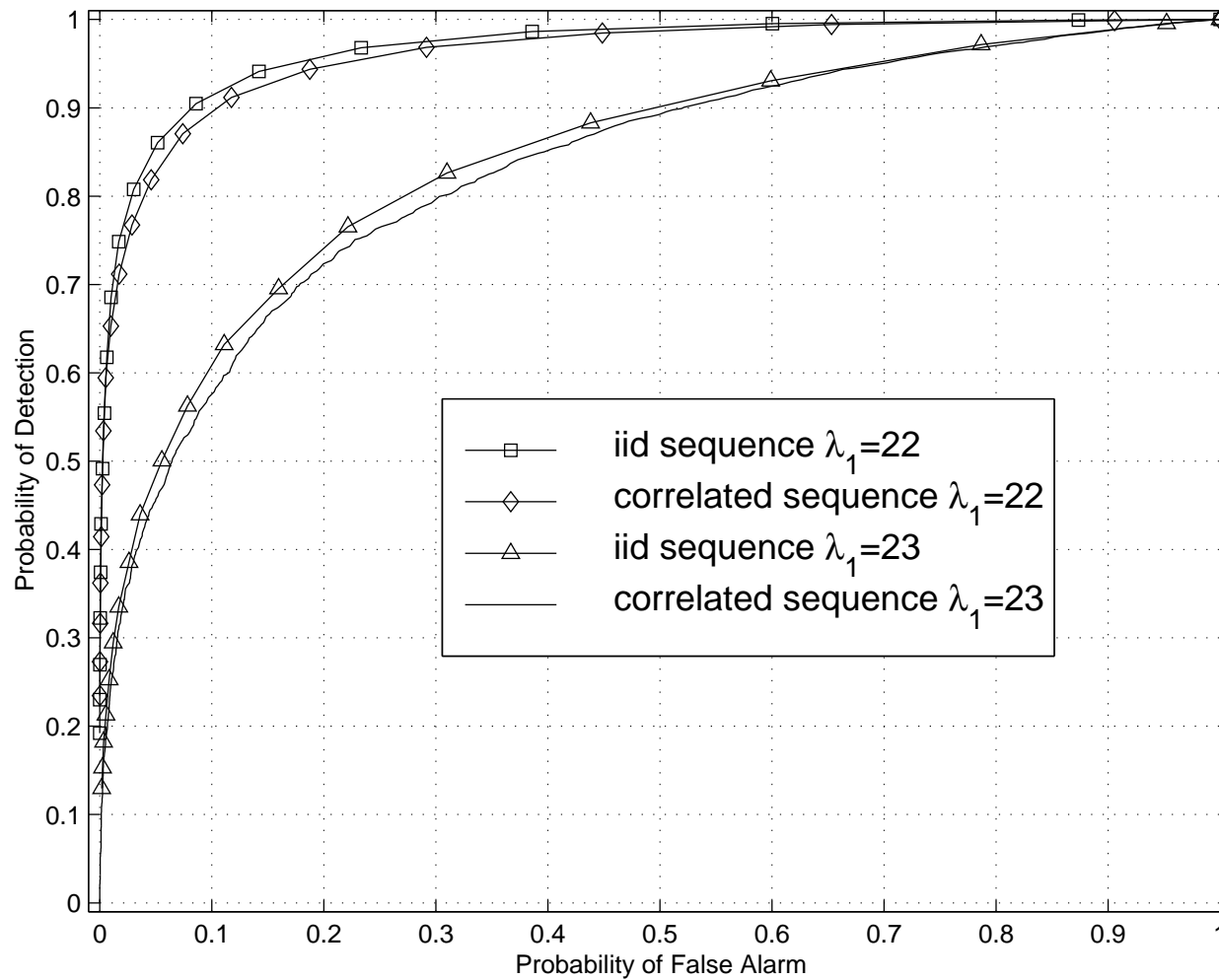
$$S_\alpha = m_0 + \frac{\sigma}{\sqrt{n}} F^{-1}(1 - \alpha)$$

- Probabilité de détection

$$PD = \pi = 1 - F\left(\frac{S_\alpha - m_1}{\frac{\sigma}{\sqrt{n}}}\right)$$

# Représentation graphique

*ROC's for iid and correlated sequences*



# Plan du cours

- Chapitre 1 : Estimation
- Chapitre 2 : Tests Statistiques
  - Généralités, exemple
  - Courbes COR
  - Théorème de Neyman Pearson
  - Test du rapport de vraisemblance généralisé
  - Test du  $\chi^2$
  - Test de Kolmogorov

# Théorème de Neyman-Pearson

## Test paramétrique à hypothèses simples

$$H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0 \text{ et } H_1 : \boldsymbol{\theta} = \boldsymbol{\theta}_1 \quad (2)$$

- Variables aléatoires continues

- **Théorème** : à  $\alpha$  fixé, le test qui minimise  $\beta$  (ou maximise  $\pi$ ) est défini par

$$\text{Rejet de } H_0 \text{ si } \frac{L(x_1, \dots, x_n | H_1)}{L(x_1, \dots, x_n | H_0)} > S_\alpha$$

- **Remarque** :  $L(x_1, \dots, x_n | H_i) = f(x_1, \dots, x_n | \boldsymbol{\theta}_i)$

- **Exemple** :  $X_i \sim \mathcal{N}(m, \sigma^2)$ ,  $\sigma^2$  connue

$$H_0 : m = m_0, \quad H_1 : m = m_1 > m_0$$

# Résumé

Effectuer un test de Neyman-Pearson, c'est

- 1) Déterminer la **statistique** et la **région critique** du test
- 2) Déterminer la relation entre le **seuil**  $S_\alpha$  et le risque  $\alpha$
- 3) Calculer le risque  $\beta$  et la **puissance**  $\pi$  du test (ou la courbe COR)
- 4) **Application numérique** : on accepte ou rejette l'hypothèse  $H_0$  en précisant le risque  $\alpha$  donné

# Théorème de Neyman-Pearson

- Variables aléatoires discrètes

- **Théorème** : parmi tous les tests de risque de première espèce  $\leq \alpha$  fixé, le test de puissance maximale est défini par

$$\text{Rejet de } H_0 \text{ si } \frac{L(x_1, \dots, x_n | H_1)}{L(x_1, \dots, x_n | H_0)} > S_\alpha$$

- **Remarque** :

$$L(x_1, \dots, x_n | H_i) = P[X_1 = x_1, \dots, X_n = x_n | \theta_i]$$

- **Exemple** :  $X_i \sim \mathcal{P}(\lambda)$ ,  $H_0 : \lambda = \lambda_0$ ,  $H_1 : \lambda = \lambda_1 > \lambda_0$
- **Loi asymptotique** : quand  $n$  est suffisamment grand, utilisation du théorème de la limite centrale

# Plan du cours

- Chapitre 1 : Estimation
- Chapitre 2 : Tests Statistiques
  - Généralités, exemple
  - Courbes COR
  - Théorème de Neyman Pearson
  - Test du rapport de vraisemblance généralisé
  - Test du  $\chi^2$
  - Test de Kolmogorov

# Test du rapport de vraisemblance généralisé

Test paramétrique à hypothèses composites

$$H_0 : \boldsymbol{\theta} \in \Theta_0 \text{ et } H_1 : \boldsymbol{\theta} \in \Theta_1 \quad (3)$$

## • Définition (Test GLR)

$$\text{Rejet de } H_0 \text{ si } \frac{L\left(x_1, \dots, x_n \mid \hat{\boldsymbol{\theta}}_1^{\text{MV}}\right)}{L\left(x_1, \dots, x_n \mid \hat{\boldsymbol{\theta}}_0^{\text{MV}}\right)} > S_\alpha$$

où  $\hat{\boldsymbol{\theta}}_0^{\text{MV}}$  et  $\hat{\boldsymbol{\theta}}_1^{\text{MV}}$  sont les estimateurs du maximum de vraisemblance de  $\boldsymbol{\theta}$  sous les hypothèses  $H_0$  et  $H_1$ .

## • Remarque

$$L\left(x_1, \dots, x_n \mid \hat{\boldsymbol{\theta}}_i^{\text{MV}}\right) = \sup_{\boldsymbol{\theta} \in \Theta_i} L\left(x_1, \dots, x_n \mid \boldsymbol{\theta}\right)$$



# Plan du cours

- Chapitre 1 : Estimation
- Chapitre 2 : Tests Statistiques
  - Généralités, exemple
  - Courbes COR
  - Théorème de Neyman Pearson
  - Test du rapport de vraisemblance généralisé
  - Test du  $\chi^2$
  - Test de Kolmogorov

# Test du $\chi^2$

Le test du  $\chi^2$  est un test **non paramétrique d'ajustement** (ou d'adéquation) qui permet de tester les deux hypothèses suivantes

$$H_0 : L = L_0, \quad H_1 : L \neq L_0$$

où  $L_0$  est une loi donnée. Le test consiste à déterminer si  $(x_1, \dots, x_n)$  est de loi  $L_0$  ou non. On se limitera dans ce cours au cas simple où  $x_i \in \mathbb{R}$ .

- **Définition**

$$\text{Rejet de } H_0 \text{ si } \phi_n = \sum_{k=1}^K \frac{(Z_k - np_k)^2}{np_k} > S_\alpha$$

- **Remarque** :  $L_0$  peut être une loi discrète ou continue

# Test du $\chi^2$

## • Statistique de test

- $Z_k$  : nombre d'observations  $x_i$  appartenant à la classe  $C_k$ ,  $k = 1, \dots, K$
- $p_k$  : probabilité qu'une observation  $x_i$  appartienne à la classe  $C_k$  sachant  $X_i \sim L_0$

$$P[X_i \in C_k | X_i \sim L_0]$$

- $n$  : nombre total d'observations

## • Loi (asymptotique) de la statistique de test sous $H_0$

$$\phi_n \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \chi_{K-1}^2$$

# Remarques

- **Interprétation de  $\phi_n$**

$$\phi_n = \sum_{k=1}^K \frac{n}{p_k} \left( \frac{Z_k}{n} - p_k \right)^2$$

Distance entre probabilités théoriques et empiriques

- **Loi asymptotique de  $\phi_n$**  : voir notes de cours ou livres

- **Nombre d'observations fini**

Une heuristique dit que la loi asymptotique de  $\phi_n$  est une bonne approximation pour  $n$  fini si 80% des classes vérifient  $np_k \geq 5$  et si  $p_k > 0, \forall k = 1, \dots, K$

☞ **Classes équiprobables**

# Remarques

- **Correction**

Lorsque les paramètres de la loi  $L_0$  sont **inconnus**

$$\phi_n \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \chi_{K-1-n_p}^2$$

où  $n_p$  est le nombre de paramètres inconnus estimés par la méthode du maximum de vraisemblance

- **Constitution des classes dans le cas d'une loi discrète**

- **Puissance du test**

Non calculable

# Exemple

4.13	1.41	-1.16	-0.75	1.96	2.46	0.197	0.24	0.42	2.00
2.08	1.48	1.73	0.82	0.33	-0.76	0.42	4.60	-2.83	0.197
2.59	0.54	4.06	-0.69	4.99	0.67	2.45	5.61	2.13	1.76
5.03	0.85	1.29	0.17	-0.38	2.76	-1.03	1.87	4.48	0.73

Est-il raisonnable de penser que ces observations sont issues d'une population de loi  $\mathcal{N}(1, 4)$  ?

## Solution

### • Classes

$$C_1 : ]-\infty, -0.34], C_2 : ]-0.34, 1], C_3 : ]1, 2.34], C_4 : ]2.34, \infty[$$

### • Nombres d'observations

$$Z_1 = 7, Z_2 = 12, Z_3 = 10, Z_4 = 11$$

# Exemple

- Statistique de test

$$\phi_n = 1.4$$

- Seuils

	$\chi_2^2$	$\chi_3^2$
$S_{0.05}$	5.991	7.815
$S_{0.01}$	9.210	11.345

donc on accepte l'hypothèse  $H_0$  avec les risques  $\alpha = 0.01$  et  $\alpha = 0.05$ .

# Plan du cours

- Chapitre 1 : Estimation
- Chapitre 2 : Tests Statistiques
  - Généralités, exemple
  - Courbes COR
  - Théorème de Neyman Pearson
  - Test du rapport de vraisemblance généralisé
  - Test du  $\chi^2$
  - Test de Kolmogorov



# Test de Kolmogorov

Le test de Kolmogorov est un test **non paramétrique d'ajustement** (ou d'adéquation) qui permet de tester les deux hypothèses suivantes

$$H_0 : L = L_0, \quad H_1 : L \neq L_0$$

où  $L_0$  est une loi donnée. Le test consiste à déterminer si  $(x_1, \dots, x_n)$  est de loi  $L_0$  ou non. On se limitera dans ce cours au cas simple où  $x_i \in \mathbb{R}$ .

## • Définition

$$\text{Rejet de } H_0 \text{ si } D_n = \sup_{x \in \mathbb{R}} |\hat{F}(x) - F_0(x)| > S_\alpha$$

• **Remarque** :  $L_0$  doit être une loi continue

# Remarques

- **Statistique de test**

$F_0(x)$  est la fonction de répartition théorique associée à  $L_0$  et  $\hat{F}(x)$  est la fonction de répartition empirique de  $(x_1, \dots, x_n)$

- **Loi asymptotique de  $D_n$  sous  $H_0$  : voir livres**

$$P[\sqrt{n}D_n < y] \xrightarrow[n \rightarrow \infty]{} \sum_{k=-\infty}^{+\infty} (-1)^k \exp(-2k^2 y) = K(y)$$

- **Détermination du seuil  $S_\alpha$  :  $S_\alpha = \frac{1}{\sqrt{n}} K^{-1}(1 - \alpha)$**

Le seuil dépend de  $\alpha$  et de  $n$ .

# Remarques

- Calcul de  $D_n$

$$D_n = \max_{i \in \{1, \dots, n\}} \max\{E_i^+, E_i^-\}$$

$$E_i^+ = \left| \widehat{F}(x_i^{*+}) - F_0(x_i^*) \right|, \quad E_i^- = \left| \widehat{F}(x_i^{*-}) - F_0(x_i^*) \right|$$

- Rq :  $x_1^*, \dots, x_n^*$  est la statistique d'ordre de  $x_1, \dots, x_n$ .

- Rq :  $\widehat{F}(x_i^{*+}) = i/n$  et  $\widehat{F}(x_i^{*-}) = (i-1)/n$ .

- Puissance du test

Non calculable

# Exemple

Est-il raisonnable de penser que ces observations sont issues d'une population de loi uniforme sur  $[0, 1]$  ?

$x_i$	0.0078	0.063	0.10	0.25	0.32	0.39	0.40	0.48	0.49	0.53
$E_i^-$	0.0078	0.013	0.00	0.10	0.07	0.14	0.05	0.008	0.04	0.03
$E_i^+$	0.0422	0.037	0.05	0.05	0.12	0.09	0.10	0.13	0.09	0.08
$\text{Max}(E_i^+, E_i^-)$	0.0422	0.037	0.05	0.1	0.12	0.14	0.10	0.13	0.09	0.08

$x_i$	0.67	0.68	0.69	0.73	0.79	0.80	0.87	0.88	0.90	0.996
$E_i^-$	0.17	0.13	0.04	0.03	0.04	0.05	0.07	0.03	0.05	0.046
$E_i^+$	0.12	0.08	0.09	0.08	0.09	0.00	0.02	0.02	0.00	$4e - 3$
$\text{Max}(E_i^+, E_i^-)$	<b>0.17</b>	0.13	0.09	0.08	0.09	0.05	0.07	0.03	0.05	0.046

# Exemple

- Statistique de test

$$D_n = 0.17$$

- Seuils pour  $n = 20$

$S_{0.05}$	0.294
$S_{0.01}$	0.352

donc on accepte l'hypothèse  $H_0$  avec les risques  $\alpha = 0.01$  et  $\alpha = 0.05$ .

# Que faut-il savoir ?

- Définition et calcul des **risques de première et seconde espèce** et de la **puissance** d'un test binaire
- Définition et détermination des courbes **COR**
- Appliquer le théorème de **Neyman-Pearson** dans le cas de variables aléatoires discrètes et continues
- Test du **rapport de vraisemblance généralisé**
- Principe et mise en oeuvre d'un **test du  $\chi^2$**
- Principe et mise en oeuvre d'un **test de Kolmogorov**