

# Aléatoire

## MAP 311 - X2013

### Leçon 8

*Mots clés :*

1. Statistique gaussienne : estimation de la moyenne
  - Variance inconnue
  - Estimateur du maximum de vraisemblance
2. Estimation non-paramétrique
3. Tests d'hypothèses dans le cas gaussien
4. Test du  $\chi^2$
5. Pour en savoir plus : estimation d'un jeu avec observation partielle

# STATISTIQUE SUR ÉCHANTILLON GAUSSIEN, ESTIMATION DE MOYENNE

Modèle :  $(Y_1, \dots, Y_n)$  v.a. i.i.d. de loi  $\mathcal{N}(\mu, \sigma^2)$  (moyenne et variance inconnues)

Moyenne empirique :  $\bar{Y}_n := \frac{1}{n} \sum_{i=1}^n Y_i \stackrel{\text{loi}}{=} \mathcal{N}(\mu, \frac{\sigma^2}{n})$

Variance empirique :  $\sigma_n^2 := \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2$

**Théorème.**  $\bar{Y}_n$  et  $\sigma_n^2$  sont des estimateurs convergents et sans biais de  $(\mu, \sigma^2)$ .

PREUVE. Pour  $\bar{Y}_n$ , voir leçon 7. Il reste à démontrer

- $\sigma_n^2 \xrightarrow{p.s.} \sigma^2$ ,
- $\mathbb{E}(\sigma_n^2) = \sigma^2$ .

On a  $\sigma_n^2 = \frac{n}{n-1} \left( \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}_n^2 \right)$ .

- Par la LGN,  $\sigma_n^2 \xrightarrow{p.s.} \mathbb{E}(Y^2) - (\mathbb{E}(Y))^2 = \sigma^2$ .
- $$\begin{aligned} \mathbb{E}(\sigma_n^2) &= \frac{n}{n-1} \mathbb{E}(Y^2) - \frac{n}{n-1} \left[ \frac{\text{Var}(Y)}{n} + (\mathbb{E}(Y))^2 \right] \\ &= \frac{n}{n-1} \text{Var}(Y) - \frac{n}{n-1} \frac{\text{Var}(Y)}{n} = \text{Var}(Y). \end{aligned}$$

□

## Vitesse de convergence

**Théorème.**  $\sqrt{n} \frac{(\bar{Y}_n - \mu)}{\sigma}$  et  $(n - 1) \frac{\sigma_n^2}{\sigma^2}$  sont indépendants, respectivement de loi  $\mathcal{N}(0, 1)$  et de loi  $\chi^2(n - 1)$ .

**Corollaire.**  $\sqrt{n} \frac{(\bar{Y}_n - \mu)}{\sigma_n} \stackrel{\text{loi}}{=} \frac{\mathcal{N}(0, 1)}{\sqrt{\frac{\chi^2(n-1)}{n-1}}}$  avec numérateur et dénominateur indépendants, c.à-d. une v.a. de loi de Student avec paramètre  $n - 1$ .

**Intervalles de confiance.** Avec probabilité 95%, on a (pour chaque  $n$ )

$$\bar{Y}_n - t_{n-1}(97.5\%) \frac{\sigma_n}{\sqrt{n}} \leq \mu \leq \bar{Y}_n + t_{n-1}(97.5\%) \frac{\sigma_n}{\sqrt{n}}.$$

$n$	10	20	100	$+\infty$
$t_{n-1}(97.5\%)$	2.262	2.093	1.984	1.96

# ESTIMATEUR DU MAXIMUM DE VRAISEMBLANCE (E.M.V.)

Etant donnés des paramètres  $(\mu, \sigma^2)$ , le vecteur aléatoire  $(Y_1, \dots, Y_n)$  i.i.d. de loi  $\mathcal{N}(\mu, \sigma^2)$  a une loi de densité

$$p_{\mu, \sigma^2}(y_1, \dots, y_n) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_j - \mu)^2}{2\sigma^2}}.$$

**Définition.** L'e.m.v. est défini par

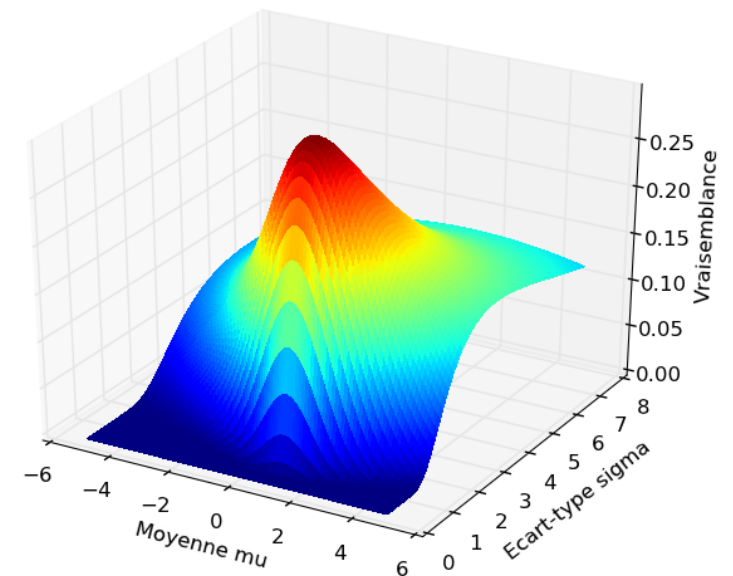
$$(\hat{\mu}_n, \hat{\sigma}_n^2) = \arg \max_{\mu, \sigma^2} p_{\mu, \sigma^2}(\mathbf{Y}_1, \dots, \mathbf{Y}_n)$$

(paramètres qui rendent l'observation la plus "probable").

**Proposition.** Dans ce cas gaussien, on trouve  $\hat{\mu}_n = \bar{X}_n$  et  $\hat{\sigma}_n^2 = \frac{n-1}{n} \sigma_n^2$ .

**Remarque.** En général, l'e.m.v. a de très bonnes propriétés de convergence et d'optimalité (voir cours de Statistique de 2A).

Vraisemblance (renormalisée): 100 données i.i.d.  $\mathcal{N}(1, 2^2)$



## ESTIMATION NON PARAMÉTRIQUE

$(X_n)_n$  : variables aléatoires indépendantes, de même loi.

**Proposition (LGN).** Convergence (ponctuelle) des fonctions de répartition :

$$\mathbf{F}_n(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\mathbf{X}_i(\omega) \leq \mathbf{x}} \rightarrow \mathbb{P}(\mathbf{X} \leq \mathbf{x}) := \mathbf{F}(\mathbf{x}), \quad \forall \mathbf{x} \quad \text{p.s.}$$

**Proposition (TCL).**  $\sqrt{n} \left( F_n(x) - F(x) \right) \xrightarrow{\text{loi}} \mathcal{N} \left( 0, F(x)(1 - F(x)) \right)$ .

**Théorème (Glivenko-Cantelli, admis).** Convergence (uniforme) :

$$\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0 \quad \text{p.s.}$$

▣ Validation de la représentation en histogramme :

$$\frac{1}{n} \#\{\mathbf{i} : \mathbf{1} \leq \mathbf{i} \leq n, \mathbf{a} < \mathbf{X}_i \leq \mathbf{b}\} = \mathbf{F}_n(\mathbf{b}) - \mathbf{F}_n(\mathbf{a}) \xrightarrow{\text{p.s.}} \mathbf{F}(\mathbf{b}) - \mathbf{F}(\mathbf{a}).$$

On suppose maintenant que la loi de  $X$  a une densité  $f > 0$ .

**Théorème.**  $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$  a une loi qui ne dépend pas de la loi  $F$  :

$$\sup_{\mathbf{x} \in \mathbb{R}} |\mathbf{F}_n(\mathbf{x}) - \mathbf{F}(\mathbf{x})| \stackrel{\text{loi}}{=} \sup_{\mathbf{u} \in [0,1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq \mathbf{u}} - \mathbf{u} \right|$$

où  $(U_i)_i$  sont des v.a. i.i.d. de loi uniforme sur  $[0, 1]$ .

PREUVE.  $U_i = F(X_i)$  définit une suite de v.a. i.i.d. de loi uniforme sur  $[0, 1]$ , et  $F$  est bijective de  $\mathbb{R}$  dans  $[0, 1]$  :

$$\sup_{x \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq x} - F(x) \right| = \sup_{u \in [0,1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq F^{-1}(u)} - u \right| = \sup_{u \in [0,1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq u} - u \right|. \quad \square$$

**Théorème (Kolmogorov-Smirnov, admis).**

La v.a.  $\sqrt{n} \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$  converge en loi vers une variable aléatoire  $B$  de fonction de répartition

$$\mathbb{P}(B \leq x) = 1 + 2 \sum_{k \geq 1} (-1)^k e^{-2k^2 x}, \quad x \geq 0.$$

## Comment tester si $F$ (inconnue) est égale à $F_0$ donnée ?

Posons  $Q_0 = F_0^{-1}$ . Si  $F = F_0$ ,  $u \mapsto F_n(Q_0(u)) \xrightarrow{p.s.} F(Q_0(u)) = u$  : alignement asymptotique sur la diagonale.

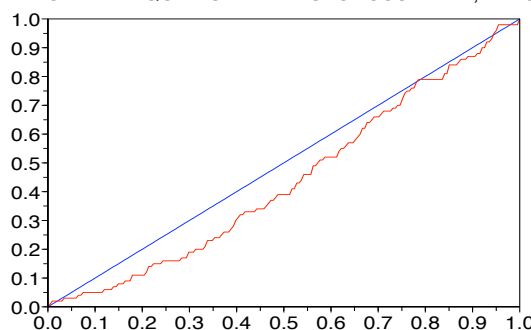
**Exemple.** Données de loi normale  $\mathcal{N}(0, 1)$ .

Test d'adéquation à

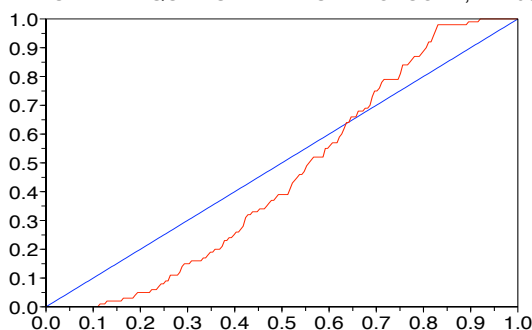
- la loi normale,
- la loi de Cauchy.

Représentation dite en Probabilité/Probabilité.

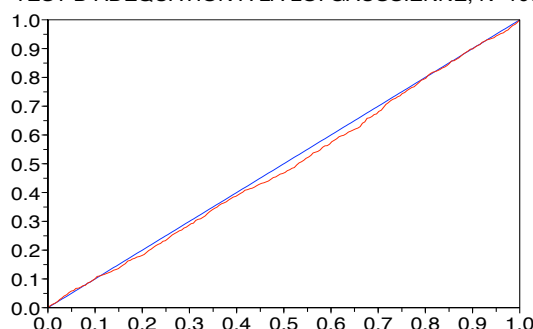
TEST D ADEQUATION A LA LOI GAUSSIENNE, N=100



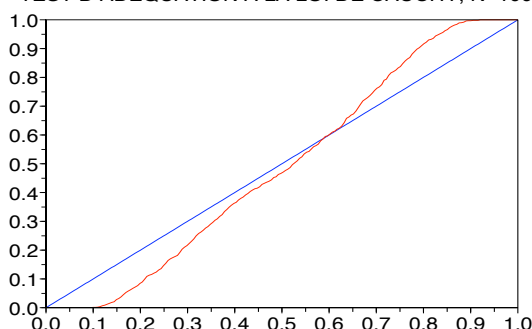
TEST D ADEQUATION A LA LOI DE CAUCHY, N=100



TEST D ADEQUATION A LA LOI GAUSSIENNE, N=1000



TEST D ADEQUATION A LA LOI DE CAUCHY, N=1000



Utiliser le résultat de Kolmogorov-Smirnov permet d'obtenir des intervalles de confiance sur les fluctuations uniformes.

## POUR EN SAVOIR PLUS : CONVERGENCE DES QUANTILES EMPIRIQUES

A partir des observations  $(X_i)_{1 \leq i \leq n}$  i.i.d., on forme la statistique d'ordre :

$$X_{(1)} \leq X_{(2)} \leq \cdots \leq X_{(i)} \leq \cdots \leq X_{(n)}.$$

Pour simplifier la discussion, on suppose que *la loi de  $X$  est à densité* : alors les  $\leq$  sont des  $<$  avec probabilité 1.

La fonction de répartition empirique  $F_n$  associée est alors telle que :

$$\mathbf{F}_n(\mathbf{x}) = \frac{\mathbf{i}}{\mathbf{n}} \quad \text{si } \mathbf{x} \in [\mathbf{X}_{(\mathbf{i})}, \mathbf{X}_{(\mathbf{i}+1)}[ \quad (\text{avec } X_{(0)} := -\infty, X_{(n+1)} = +\infty).$$

Le quantile empirique est alors  $\mathbf{Q}_n(\alpha) = \mathbf{X}_{(\lceil n\alpha \rceil)}$ .

**Exemple.** La 11ème donnée la plus grande sur 1000 données donne  $Q(99\%)$ .



## Convergence - Loi des grands nombres et Théorème de la Limite Centrale

**Théorème.**  $Q_n(\alpha) \rightarrow Q(\alpha)$ .

PREUVE. D'une part  $F_n(Q_n(\alpha)) = \frac{[n\alpha]}{n} \rightarrow \alpha$ .

D'autre part, par Glivenko-Cantelli,  $F_n(Q_n(\alpha)) - F(Q_n(\alpha)) \rightarrow 0$ . □

**Théorème (fluctuations, admis).** Supposons que  $X$  a une densité  $f > 0$ .

$$\sqrt{n}(Q_n(\alpha) - Q(\alpha)) \xrightarrow{\text{loi}} \mathcal{N}\left(0, \frac{\alpha(1-\alpha)}{f^2(Q(\alpha))}\right).$$

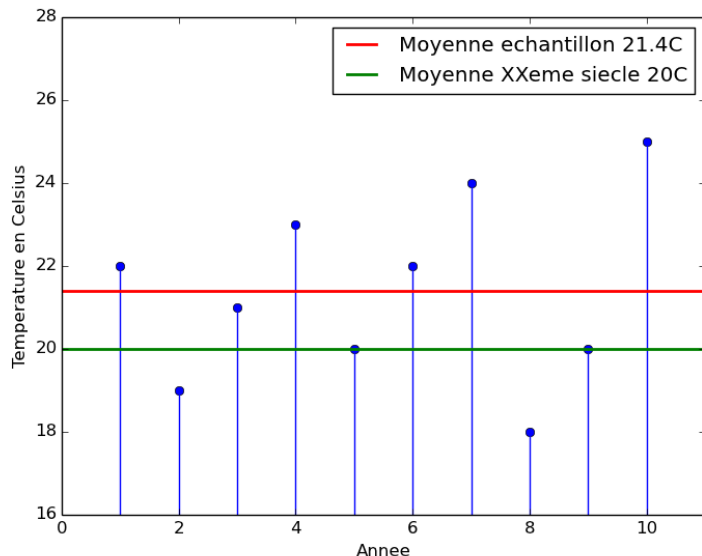
## TESTS D'HYPOTHÈSES

On considère une modélisation paramétrique décrit par une famille de probabilités  $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$  pour  $\Theta \subset \mathbb{R}^d$ .

**Exemple.** Modèle gaussien  $\mathcal{N}(\mu, \sigma^2)$  avec  $\Theta = \{(\mu, \sigma) : \mu \in \{\mu_0, \mu_1\}, \sigma \in [\sigma_0, \sigma_1]\}$

QUESTION : au vu de l'observation  $(X_1, \dots, X_n)$ , le paramètre  $\theta$  du modèle est-il ou non dans un sous-ensemble  $\mathbf{H}_0$  (appelé *hypothèse nulle*) ?

**Exemple.** Température moyenne à Paris au XXème siècle : modélisé par  $\mathcal{N}(20, 1)$ .



Températures observées sur les 10 dernières années :

- ▷ peut-on décider que le modèle sur ces dernières années (supposé gaussien) a une moyenne inférieure ou supérieure à  $20^\circ\text{C}$  ?
- ▷ réchauffement ?


## DÉFINITIONS USUELLES

▷ **Test d'hypothèse** : règle de décision qui, au vu de l'observation de  $X$ , permet de décider si  $\theta \in \mathbf{H}_0$  (*hypothèse nulle*) ou si  $\theta \in \mathbf{H}_1$  (*hypothèse alternative*) ?

Ici  $\mathbf{H}_0$  et  $\mathbf{H}_1$  partitionnent l'espace  $\Theta$  des paramètres.

▷ **Région critique  $W$**  :

- si  $\mathbf{X} := (X_1, \dots, X_n) \in W$ , on rejette  $\mathbf{H}_0$  et on accepte  $\mathbf{H}_1$
- si  $\mathbf{X} \notin W$ , on accepte  $\mathbf{H}_0$  et on rejette  $\mathbf{H}_1$

 **Notation** : pour simplifier, on note de la même manière la région critique  $W$  et l'événement critique  $\{\mathbf{X} \in W\}$  associé à la règle de décision.

▷ **Erreur/risque de 1ère espèce** : on rejette  $\mathbf{H}_0$  à tort, avec probabilité

$$\mathbb{P}_\theta(W) \text{ pour } \theta \in \mathbf{H}_0 .$$

**Niveau du test** :  $\sup_{\theta \in \mathbf{H}_0} \mathbb{P}_\theta(W)$  (typiquement 1%, 5%...)

▷ **Erreur/risque de 2nde espèce** : on rejette  $\mathbf{H}_1$  à tort, avec probabilité :

$$\mathbb{P}_\theta(W^c) \text{ pour } \theta \in \mathbf{H}_1 .$$

On accepte à raison  $\mathbf{H}_1$  avec  $\mathbb{P}_\theta(W)$

**Puissance du test** :  $\mathbb{P}_\theta(W)$  ( $\theta \in \mathbf{H}_1$  )

▷ **Objectif** : à un niveau de test fixé, trouver le test de puissance maximale



Dissymétrie entre  $\mathbf{H}_0$  et  $\mathbf{H}_1$

## EXEMPLE MODÈLE GAUSSIEN $\mathcal{N}(\mu, \sigma^2)$ : TEST DE 2 MOYENNES AVEC $\sigma^2$ CONNU

$\mathbf{H}_0 = \{\mu = \mu_0\}$ ,  $\mathbf{H}_1 = \{\mu = \mu_1\}$ ,  $\mu_0 > \mu_1$ ,  $(X_i)_i$  de loi  $\mathcal{N}(\mu, \sigma^2)$

On propose un test qui accepte  $\mathbf{H}_0$  si la moyenne empirique  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  est suffisamment grande.

▸ région critique  $W = \{\bar{X}_n \leq a\}$

▷ **Niveau du test** : pour  $\theta \in \mathbf{H}_0$

$$\mathbb{P}_\theta(W) = \mathbb{P}\left(\mathcal{N}\left(\mu_0, \frac{\sigma^2}{n}\right) \leq a\right) = \mathbb{P}\left(\mathcal{N}(0, 1) \leq \frac{a - \mu_0}{\sigma/\sqrt{n}}\right)$$

Pour un niveau de test  $\alpha$ , prendre  $a = \mu_0 + \mathcal{N}^{-1}(\alpha) \frac{\sigma}{\sqrt{n}} \stackrel{\alpha=5\%}{=} \mu_0 - 1.65 \frac{\sigma}{\sqrt{n}}$

▷ **Puissance** : pour  $\theta \in \mathbf{H}_1$

$$\mathbb{P}_\theta(W) = \mathbb{P}\left(\mathcal{N}(0, 1) \leq \frac{a - \mu_1}{\sigma/\sqrt{n}}\right) = \mathbb{P}\left(\mathcal{N}(0, 1) \leq \mathcal{N}^{-1}(\alpha) + \frac{\mu_0 - \mu_1}{\sigma/\sqrt{n}}\right) \xrightarrow{n \rightarrow +\infty} 1$$

La puissance du test dépend de la distance entre  $\mu_1$  et  $\mu_0$ , et la taille de l'échantillon.

**ENCORE UN PEU DE TERMINOLOGIE**

**Définition.** On se donne  $(W_n)_n$  une suite de régions critiques (une suite de tests) où  $n$  est la taille de l'échantillon.

La suite de tests  $(W_n)_n$  est dite

- convergente si  $\mathbb{P}_\theta(W_n) \xrightarrow{n \rightarrow +\infty} 1$ , pour  $\theta \in \mathbf{H}_1$
- de niveau asymptotique  $\alpha$  si  $\sup_{\theta \in \mathbf{H}_0} \mathbb{P}_\theta(W_n) \xrightarrow{n \rightarrow +\infty} \alpha$

**Exemple (modèle gaussien).** Test de 2 moyennes avec

$$W_n = \left\{ \bar{X}_n \leq \mu_0 + \mathcal{N}^{-1}(\alpha) \frac{\sigma}{\sqrt{n}} \right\}.$$

Ce test est *exactement* de niveau  $\alpha$ , il est convergent.

## MODÈLE GAUSSIEN : TEST DE MOYENNE

$$\mathbf{H}_0 = \{\mu = \mu_0\} \text{ VERSUS } \mathbf{H}_1 = \{\mu \neq \mu_0\}$$

Rappels :

- sous  $\mathbf{H}_0$  ,  $\zeta_n = \sqrt{n} \frac{(\bar{X}_n - \mu_0)}{\sigma_n} \stackrel{\text{loi}}{=} \text{Student } t(n-1)$
- sous  $\mathbf{H}_0$  ,  $\bar{X}_n \xrightarrow{p.s.} \mu_0$
- ▣ Région critique :  $W_n = \{|\zeta_n| > a\}$

**Théorème.** En prenant  $a_n = t_{n-1}^{-1}(1 - \alpha/2)$ , la suite de tests associés à  $(W_n = \{|\zeta_n| > a_n\})_n$  est convergente et de niveau  $\alpha$ .

PREUVE.

- Niveau du test :  $\mathbb{P}_{\mu_0}(|\zeta_n| > a_n) = 2t_{n-1}(-a_n) = 2(1 - t_{n-1}(a_n)) = \alpha$
- Convergence : pour tout  $\mu \neq \mu_0$ ,  $\mathbb{P}_{\mu}(|\zeta_n| > a_n) \rightarrow 1$  car
  - ▶  $|\zeta_n| \underset{n \rightarrow +\infty}{\sim} \left| \sqrt{n} \frac{(\mu - \mu_0)}{\sigma} \right| \rightarrow +\infty$  p.s. et  $a_n \rightarrow \mathcal{N}^{-1}(1 - \alpha/2)$
  - ▶  $\mathbf{1}_{|\zeta_n| > a_n} \xrightarrow{p.s.} 1$  et par convergence dominée  $\mathbb{E}(\mathbf{1}_{|\zeta_n| > a_n}) \rightarrow 1$ .

□

## RÉCHAUFFEMENT À PARIS ?

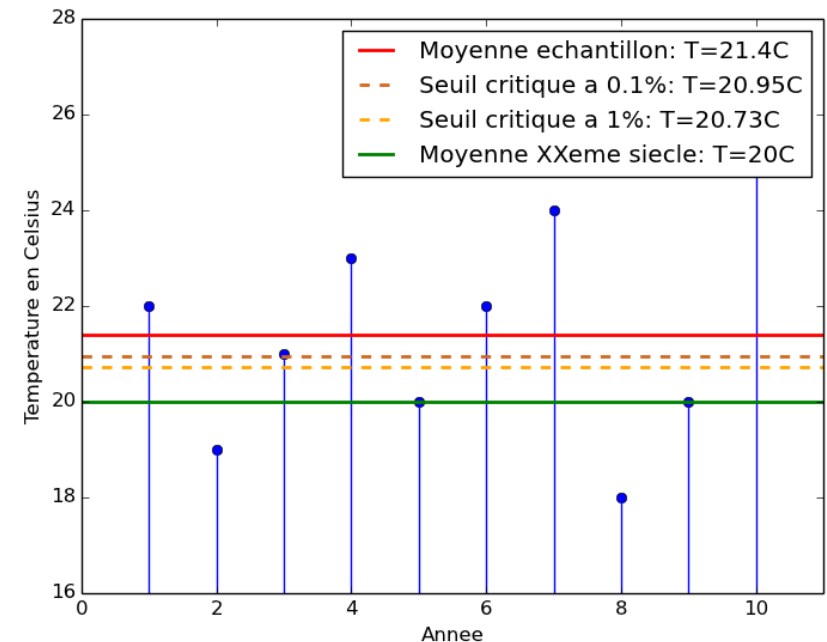
On teste  $\mathbf{H}_0 = \{\mu \leq 20\}$  versus  $\mathbf{H}_1 = \{\mu > 20\}$ ,  $\sigma^2 = 1$  connue.

On pose le test  $W_n = \{\sqrt{n}(\bar{X}_n - 20) \geq \mathcal{N}^{-1}(1 - \alpha)\}$ .

**Proposition (à faire en exercice).** C'est un test de niveau  $\alpha$  et convergent.

**Exemple.** Sur 10 données et au niveau  $\alpha = 1\%$ ,  $W_n = \{\bar{X}_n \geq 20 + \frac{2.33}{\sqrt{10}} = 20.73\}$ .

$x_n^{\text{obs}} = 21.4 > 20.73 \implies$  on rejette  $\mathbf{H}_0$  et on accepte  $\mathbf{H}_1$



**Définition.** La  $p$ -valeur du test est la probabilité que la statistique de test prenne des valeurs plus défavorables que la valeur observée, c'ad  $\mathbb{P}_\mu(\bar{X}_n \geq x_n^{\text{obs}})$ ,  $\mu \in \mathbf{H}_0$ . Ici  $p = 4.7 \times 10^{-6}$  : on rejette  $\mathbf{H}_0$  à tous les niveaux  $\alpha \geq p$ .



## TEST DU $\chi^2$

### CONTEXTE :

- on pose une question avec 3 réponses possibles : 1, 2, 3
- on interroge  $n$  personnes indépendamment : les réponses sont  $X_i$ .

$$\mathbb{P}(X_i = j) = p_j, \quad j = 1, 2, 3.$$

- $\Theta = \{p = (p_1, p_2, p_3) \in ]0, 1[^3 : p_1 + p_2 + p_3 = 1\}$
- On veut tester l'hypothèse nulle  $\mathbf{H}_0 = \{p = p^0\}$  contre l'hypothèse alternative  $\mathbf{H}_1 = \{p \neq p^0\}$

Plus généralement,  $k$  réponses possibles. Espace fini.

### FRÉQUENCE EMPIRIQUE :

$$\hat{p}_j = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i=j}, \quad 1 \leq j \leq k.$$

**Définition.** La distance du  $\chi^2$  est définie par  $\zeta_n := n \sum_{j=1}^k \frac{(\hat{p}_j - p_j^0)^2}{p_j^0}$ .

### Intuition :

- Quantifie la proximité entre  $\hat{p}$  et  $p^0$
- Normalisation par  $n$  pour tenir compte des fluctuations (TCL) en  $1/\sqrt{n}$  pour les grands échantillons

**Théorème.** 1) La loi de  $\zeta_n$  converge sous  $\mathbf{H}_0$  vers la loi du  $\chi^2(k-1)$ .  
2) Sous  $\mathbf{H}_1$ ,  $\zeta_n \xrightarrow{p.s.} +\infty$ .

PREUVE. 1) (dans le cas  $k=2$ ) Admis dans le cas général (basé sur le TCL multidimensionnel).

- $\sqrt{n}(\hat{p}_1 - p_1^0) \xrightarrow{\mathcal{L}} \mathcal{N}(0, p_1^0(1-p_1^0))$
- $\zeta_n = n \left[ \frac{(\hat{p}_1 - p_1^0)^2}{p_1^0} + \frac{(1-\hat{p}_1 - (1-p_1^0))^2}{1-p_1^0} \right] = n \frac{(\hat{p}_1 - p_1^0)^2}{p_1^0(1-p_1^0)} \xrightarrow{\mathcal{L}} \text{carré d'une gaussienne standard} \stackrel{\text{loi}}{=} \chi^2(1)$

$$2) \zeta_n \underset{n \rightarrow +\infty}{\sim} n \sum_{j=1}^k \frac{(p_j - p_j^0)^2}{p_j^0} \xrightarrow{p.s.} +\infty. \quad \square$$

**Corollaire.**  $W_n = \{\zeta_n > \chi_{k-1}^2(1 - \alpha)\}$  définit une suite de tests convergents et de niveau asymptotique  $\alpha$ .

**Exemple.** 100 lancers d'un dé à 6 faces.

$j$	1	2	3	4	5	6
$\hat{p}_j$	0.2	0.13	0.17	0.12	0.23	0.15

**Le dé est-il pipé ?**

- $p^0 = (\frac{1}{6}, \dots, \frac{1}{6})$ ,
  - $k = 6, n = 100$
  - $\zeta_n^{\text{obs}} = 5.36$
  - $\chi_{k-1}^2(1 - 5\%) = 11.07$
- ▣ au niveau 5%, on accepte l'hypothèse que le dé est équilibré.

## POUR EN SAVOIR PLUS

### ESTIMATION D'UN JEU AVEC OBSERVATION PARTIELLE

#### Deux pièces $A$ et $B$ .

- La pièce  $A$  donne P avec probabilité  $\alpha_A$
- La pièce  $B$  donne P avec probabilité  $\alpha_B$

#### Deux individus : joueur, observateur

- $n$  tours de jeu répétés :
  1. Le joueur choisit une pièce :  $Z_i \in \{A, B\}$ . Il choisit  $A$  avec proba  $\beta$  et  $B$  avec proba  $1 - \beta$ .
  2. Il la lance 3 fois.
  3. Il transmet le nombre de P obtenu à un observateur :  $X_i \in \{0, 1, 2, 3\}$ .

**Objectif :** l'observateur cherche les paramètres  $\theta = (\alpha_A, \alpha_B, \beta)$  à partir des seules observations  $X = (X_1, \dots, X_n)$  (et pas  $Z = (Z_1, \dots, Z_n)$ )

**Maximum de vraisemblance :**  $\theta^* = \arg \max_{\theta} p_{\theta}(X)$ .

Ici, optimisation directe délicate :  $p_{\theta}(X) = \sum_{(z_1, \dots, z_n) \in \{A, B\}^n} \prod_{i=1}^n p_{\theta}(X_i, z_i)$

## Algorithme EM (*Expectation-Maximization*)

**PRINCIPE DE L'ALGORITHME** : remplacer  $p_\theta(X)$  par une fonction plus facile à maximiser et qui a le même arg max.

### Propriétés.

1.  $\log p_\theta(X) - \log p_{\theta'}(X) \geq C(\theta, \theta') := \mathbb{E}_{\theta'} \left( \log \left( \frac{p_\theta(X, Z)}{p_{\theta'}(X, Z)} \right) \mid X \right)$ .
2.  $C(\theta, \theta) = 0$ .

PREUVE. 2) est clair. Pour 1), par Bayes,

$$\begin{aligned}
 \log p_\theta(X) - \log p_{\theta'}(X) &= \log \left( \sum_z p_\theta(X, z) \right) - \log p_{\theta'}(X) \\
 &= \log \left( \sum_z p_{\theta'}(z \mid X) \frac{p_\theta(X \mid z) p_\theta(z)}{p_{\theta'}(z \mid X)} \right) - \sum_z p_{\theta'}(z \mid X) \log p_{\theta'}(X) \\
 &\geq \sum_z p_{\theta'}(z \mid X) \log \left( \frac{p_\theta(X \mid z) p_\theta(z)}{p_{\theta'}(z \mid X) p_{\theta'}(X)} \right) \quad (\text{inégalité de Jensen}) \\
 &= \sum_z p_{\theta'}(z \mid X) \log \left( \frac{p_\theta(X, z)}{p_{\theta'}(X, z)} \right) = \mathbb{E}_{\theta'} \left[ \log \left( \frac{p_\theta(X, Z)}{p_{\theta'}(X, Z)} \right) \mid X \right].
 \end{aligned}$$

□

## Algorithme EM en deux étapes

- **Expectation (E)** : calculer

$$\theta \mapsto \hat{C}(\theta, \theta_k) = \mathbb{E}_{\theta_k} [\log p_\theta(X, Z) \mid X]$$

car  $C(\theta, \theta_k) - \hat{C}(\theta, \theta_k) = \mathbb{E}_{\theta_k} [\log p_{\theta_k}(X, Z) \mid X]$  est indépendant de  $\theta$ .

- **Maximisation (M)** :

$$\theta_{k+1} := \arg \max_{\theta} \hat{C}(\theta, \theta_k) = \arg \max_{\theta} C(\theta, \theta_k)$$

**Théorème.** Etant donné  $\theta_0$ , le choix  $\theta_{k+1} := \arg \max_{\theta} C(\theta, \theta_k)$  assure que  $(p_{\theta_k}(X))_k$  est une suite croissante en  $k$ .

PREUVE.  $\log p_{\theta_{k+1}}(X) - \log p_{\theta_k}(X) \geq C(\theta_{k+1}, \theta_k) \geq C(\theta_k, \theta_k) \geq 0$ . □

▮▮▮  $(\theta_k)_k$  converge vers un maximum local ? global ?

## Mise en œuvre sur l'exemple du jeu

### Etape E :

$$\begin{aligned}
 \hat{C}(\theta, \theta_k) &:= \mathbb{E}_{\theta_k} [\log p_{\theta}((X_1, \dots, X_n), (Z_1, \dots, Z_n)) \mid X_1, \dots, X_n] \\
 &= \sum_{i=1}^n \mathbb{E}_{\theta_k} [\log p_{\theta}(X_i, Z_i) \mid X_i] \quad (\text{indépendance des } (X_i, Z_i)_i) \\
 &= \sum_{i=1}^n p_{\theta_k}(A \mid X_i) \log p_{\theta}(X_i, A) + p_{\theta_k}(B \mid X_i) \log p_{\theta}(X_i, B), \\
 &= \sum_{i=1}^n p_{\theta_k}(A \mid X_i) \log \left( \beta \binom{3}{X_i} \alpha_A^{X_i} (1 - \alpha_A)^{3-X_i} \right) + p_{\theta_k}(B \mid X_i) \log \left( (1 - \beta) \binom{3}{X_i} \alpha_B^{X_i} (1 - \alpha_B)^{3-X_i} \right).
 \end{aligned}$$

avec

$$p_{\theta_k}(A \mid X_i) = \frac{\beta_k p_{\theta_k}(X_i \mid A)}{\beta_k p_{\theta_k}(X_i \mid A) + (1 - \beta_k) p_{\theta_k}(X_i \mid B)} \quad (\text{Bayes}),$$

$$p_{\theta_k}(B \mid X_i) = \frac{(1 - \beta_k) p_{\theta_k}(X_i \mid B)}{\beta_k p_{\theta_k}(X_i \mid A) + (1 - \beta_k) p_{\theta_k}(X_i \mid B)} \quad (\text{Bayes}).$$

**Etape M :** le maximum de  $\hat{C}(\theta, \theta_k)$  est  $\theta_{k+1} = (\alpha_{A,k+1}, \alpha_{B,k+1}, \beta_{k+1})$  donné par

- $\alpha_{\mathbf{A},k+1} = \frac{\sum_{i=1}^n \frac{X_i}{3} \mathbf{p}_{\theta_k}(\mathbf{A} | \mathbf{X}_i)}{\sum_{i=1}^n \mathbf{p}_{\theta_k}(\mathbf{A} | \mathbf{X}_i)}$  : moyenne des fréquences de P pondérée par la probabilité a posteriori d'avoir choisi la pièce A,
- $\alpha_{\mathbf{B},k+1} = \frac{\sum_{i=1}^n \frac{X_i}{3} \mathbf{p}_{\theta_k}(\mathbf{B} | \mathbf{X}_i)}{\sum_{i=1}^n \mathbf{p}_{\theta_k}(\mathbf{B} | \mathbf{X}_i)}$ ,
- $\beta_{k+1} = \frac{1}{n} \sum_{i=1}^n \mathbf{p}_{\theta_k}(\mathbf{A} | \mathbf{X}_i)$ .

### A la limite $k \rightarrow +\infty$ (si convergence vers le max global)

▷ On retrouve  $p_{\theta^*}(A | X_i)$  la probabilité a posteriori du choix de la pièce A au tirage  $i$  sachant l'observation  $X_i$ , pour le paramètre  $\theta^*$  le plus vraisemblable.

▷ Classification des pièces choisies  $Z_i = A$  ou  $B$  avec certaines probabilités.



## Exemple numérique. Valeurs exactes : $\alpha_A = 0.3$ , $\alpha_B = 0.5$ , $\beta = 0.7$ , $n = 503$

Initialisation : alphaA=0.450, alphaB=0.550, beta=0.500

Iteration 10: alphaA=0.285, alphaB=0.404, beta=0.548

.....

Iteration 990: alphaA=0.289, alphaB=0.467, beta=0.720

Iteration 1000: alphaA=0.289, alphaB=0.467, beta=0.722

=====  
 Probabilite a posteriori (exacte) de piece A sachant Xi:

Xi=0 -> 86.5% | Xi=1 -> 73.3% | Xi=2 -> 54.0% | Xi=3 -> 33.5% |

Probabilite a posteriori (estimee) de piece A sachant Xi:

Xi=0 -> 86.1% | Xi=1 -> 74.1% | Xi=2 -> 57.1% | Xi=3 -> 38.2% |

Probabilite a posteriori (aveugle) de piece A sachant Xi:

Xi=0 -> 50.0% | Xi=1 -> 50.0% | Xi=2 -> 50.0% | Xi=3 -> 50.0% |

=====

