

STATISTIQUES (STA01) – EXAMEN

ENSTA, Université Paris-Saclay
Année 2025 – 2026

Vendredi 14 novembre 2025
10h – 12h (tiers-temps 10h – 12h30)

Une feuille A4 recto manuscrite est autorisée en tant que support. La calculatrice est autorisée.

Avant de commencer :

- *Le sujet comporte trois exercices. Il est demandé de les traiter dans l'ordre sur la copie.*
- *Des résultats pourront être admis d'une question sur l'autre, à la condition de l'écrire clairement.*
- *Un barème indicatif est donné :*
 - Exercice 1 (questions de cours) : 6 points,*
 - Exercice 2 : 9 points,*
 - Exercice 3 (QCM) : 5 points.*

Exercice 1 – Questions de cours (6 points)

Ces questions sont indépendantes. Répondre à chaque question en faisant les démonstrations demandées. On prendra soin de bien rédiger.

1. Énoncer et redémontrer la décomposition biais-variance du risque quadratique, pour l'estimation d'un paramètre θ en dimension $k \geq 1$. Si vous le souhaitez, on pourra accompagner la démonstration d'un dessin.

2 points : 0.5 pour l'énoncé de la décomposition puis 1.5 pour la preuve bien faite. Si l'énoncé et la preuve ne sont faites que pour le cas réel, mettre 1 point en tout. Pas de bonus particulier pour les copies ayant fourni un dessin.

Solution. Pour $\hat{\theta} \in \mathbb{R}^k$ un estimateur de $\theta \in \mathbb{R}^k$, de matrice de covariance finie, la décomposition du risque quadratique $R_\theta(\hat{\theta})$ s'écrit :

$$R_\theta(\hat{\theta}) = \|b_\theta(\hat{\theta})\|^2 + \text{Tr}(\text{Var}_\theta(\hat{\theta})),$$

où $\text{Var}_\theta(\hat{\theta})$ est la matrice de covariance de $\hat{\theta}$, et $b_\theta(\hat{\theta})$ est son biais défini par

$$b_\theta(\hat{\varphi}) := \mathbb{E}_\theta[\hat{\varphi} - \varphi(\theta)] = \mathbb{E}_\theta[\hat{\varphi}] - \varphi(\theta) \in \mathbb{R}^d.$$

Pour la preuve, Il suffit d'écrire, pour tout $\theta \in \Theta$,

$$\begin{aligned} R_\theta(\hat{\varphi}) &= \mathbb{E}_\theta[\|\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}] + \mathbb{E}_\theta[\hat{\varphi}] - \varphi(\theta)\|^2] \\ &= \mathbb{E}_\theta[\|\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}]\|^2] + \mathbb{E}_\theta[\|\mathbb{E}_\theta[\hat{\varphi}] - \varphi(\theta)\|^2] + 2(\mathbb{E}_\theta[\hat{\varphi}] - \varphi(\theta))^T \mathbb{E}_\theta[\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}]] \\ &= \mathbb{E}_\theta[\|\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}]\|^2] + \|b_\theta(\hat{\varphi})\|^2, \end{aligned}$$

et par cyclicité et linéarité de la trace, et linéarité de l'espérance,

$$\begin{aligned} \mathbb{E}_\theta[\|\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}]\|^2] &= \mathbb{E}_\theta[\text{Tr}((\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}])^T (\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}]))] \\ &= \mathbb{E}_\theta[\text{Tr}((\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}])(\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}])^T)] \\ &= \text{Tr}[\mathbb{E}_\theta[(\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}])(\hat{\varphi} - \mathbb{E}_\theta[\hat{\varphi}])^T]] \\ &= \text{Tr}(\text{Var}_\theta(\hat{\varphi})). \end{aligned}$$

2. Soient X_1, \dots, X_n des variables aléatoires i.i.d., de moyenne μ , de variance finie $\sigma^2 > 0$, toutes deux inconnues. On note \bar{X} la moyenne empirique et $\hat{\sigma}^2$ l'estimateur non biaisé de la variance. Démontrer que $\frac{\sqrt{n}(\bar{X} - \mu)}{\sqrt{\hat{\sigma}^2}}$ converge en loi lorsque $n \rightarrow \infty$ vers une loi à préciser.

2 points : 1 pour la loi des grands nombres bien justifiée et appliquée à $\hat{\sigma}^2$ (0.5 si non justifiés), puis 0.5 pour le TCL appliqué à $\frac{\sqrt{n}(\bar{X} - \mu)}{\sqrt{\hat{\sigma}^2}}$, même non justifié. Enfin, 0.5 pour Slutsky.

3. Dans le modèle linéaire $Y = X\theta + \varepsilon$, supposé identifiable, rappeler sans justification la forme explicite de l'estimateur des moindres carrés $\hat{\theta}$ en fonction de X et Y . Puis, redémontrer qu'il est sans biais et calculer sa matrice de covariance en fonction de X et σ^2 .

2 points. Écriture de $\hat{\theta}$: 0.5pt, calcul biais 0.5 pt, Variance 1 pt via par exemple $\text{Var}(AX) = A\text{Var}(X)A^T$.

Exercice 2 – Modèle de Cobb-Douglas en économétrie (9 points)

Dans tout cet exercice, on pourra se référer au tableau de quantiles donnée à la fin de l'énoncé.

En 1928, les économistes Paul Douglas et Charles Cobb proposent un modèle de fonction de production, aujourd'hui connu sous le nom de modèle de Cobb-Douglas. Ce modèle décrit la relation entre la valeur produite V_i d'une entreprise (par exemple, en euros), son capital K_i (en euros) et sa force de travail L_i (nombre d'heures travaillées). Il suppose que la fonction de production de ces entreprises suit la forme suivante :

$$V_i = \lambda L_i^\beta K_i^\gamma,$$

où λ, β, γ sont des paramètres réels avec $\lambda > 0$. En passant au logarithme, tout cela se réécrit :

$$\log V_i = \alpha + \beta \log L_i + \gamma \log K_i$$

avec $\alpha = \log \lambda$. Le modèle linéaire associé au modèle de Cobb-Douglas est :

$$\log V_i = \alpha + \beta \log L_i + \gamma \log K_i + \varepsilon_i, \quad (1)$$

où les ε_i sont supposées i.i.d. de loi $\mathcal{N}(0, \sigma^2)$ où $\sigma^2 > 0$ est inconnu.

Dans cet exercice, on suppose qu'on observe les triplets (V_i, L_i, K_i) pour n entreprises et qu'ils suivent le modèle (1).

1. Écrire le modèle (1) sous sa forme matricielle $Y = X\theta + \varepsilon$ avec $\theta = \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix}$ en précisant Y et X en fonction des V_i, L_i et K_i .

1 point.

Solution. Avec les notations

$$Y = \begin{bmatrix} \log V_1 \\ \vdots \\ \log V_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & \log L_1 & \log K_1 \\ \vdots & \vdots & \vdots \\ 1 & \log L_n & \log K_n \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix},$$

le modèle s'écrit sous la forme matricielle $Y = X\theta + \varepsilon$.

Pour les données récoltées auprès de 1658 entreprises, on a appliqué la méthode des moindres carrés et on a obtenu les résultats suivants :

$$\hat{\theta} = \begin{pmatrix} 3.136 \\ 0.738 \\ 0.282 \end{pmatrix} \quad \text{et} \quad \text{SCR}(\hat{\theta}) = \|Y - X\hat{\theta}\|^2 = 148.27.$$

On donne aussi :

$$(X^T X)^{-1} = \begin{bmatrix} 0.0288 & 0.0012 & -0.0034 \\ 0.0012 & 0.0016 & -0.0010 \\ -0.0034 & -0.0010 & 0.0009 \end{bmatrix}$$

2. Rappeler la forme de $\hat{\sigma}^2$, estimateur sans biais de σ^2 . Calculer ici sa valeur approchée.

1 point : 0.5 pour la formule, 0.5 pour l'application numérique.

Solution. un estimateur sans biais de σ^2 est $\hat{\sigma}^2 = \frac{\text{SCR}(\hat{\theta})}{n-3}$, l'application numérique donne $\hat{\sigma}^2 = \frac{148.27}{1655} \approx 0.09$.

3. Donner un intervalle de confiance au niveau 95% pour le coefficient β . On donnera d'abord sa forme générale avant l'application numérique.

1.5 point : 1 pour la formule (même si non justifiée), puis 0.5 pour l'application numérique.

Solution. Puisqu'on sait que $\frac{\hat{\beta} - \beta}{\hat{\sigma} \sqrt{[(X^T X)^{-1}]_{2,2}}} \sim \mathcal{T}(1655)$, un intervalle de confiance à 95% pour β est donné par

$$\left[\hat{\beta} \pm t_{1-\alpha/2}^{\mathcal{T}(n-p)} \hat{\sigma} \sqrt{[(X^T X)^{-1}]_{2,2}} \right],$$

avec $t_{\beta}^{\mathcal{T}(n-p)}$ le quantile de la loi $\mathcal{T}(n-p)$ d'ordre β . L'application numérique donne

$$I(\beta) = [0.738 - 1.96 \times \sqrt{0.09 \times 0.0016}, 0.738 + 1.96 \times \sqrt{0.09 \times 0.0016}] \approx [0.71; 0.76].$$

On dit que les rendements sont :

- *constants* si, lorsqu'une entreprise double sa force de travail et double son capital, toutes choses égales par ailleurs, sa production double exactement ;
 - *croissants* si, lorsqu'une entreprise double sa force de travail et double son capital, toutes choses égales par ailleurs, sa production fait plus que doubler.
4. Traduire sur les coefficients α, β, γ les hypothèses de rendements constants ou croissants dans le cadre du modèle.

0.5 point : si la réponse pour les rendements croissants est $\beta + \gamma \geq 1$, mettre les points quand même.

Solution. Les rendements sont constants si $\beta + \gamma = 1$, croissants si $\beta + \gamma > 1$.

5. On veut tester H_0 : les rendements sont constants, contre H_1 : les rendements sont croissants. Construire un test de niveau d'erreur 5% répondant à la question. Faire l'application numérique. Conclure.

2 points : 0.5 pour la bonne stat de test, 0.5 si la zone de rejet est de la forme optimale du corrigé, 0.5 pour l'appli numérique, et 0.5 pour la conclusion.

Solution. On teste donc $H_0 : \beta + \gamma = 1$ contre $H_1 : \beta + \gamma > 1$. Il s'agit d'un test de Student de type $a^T \theta = c$ avec $a = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$ et $c = 1$. La statistique de test associée est donc

$$T = \frac{a^T \hat{\theta} - c}{\hat{\sigma} \sqrt{a^T (X^T X)^{-1} a}}$$

qui suit une loi $\mathcal{T}(n-p)$ sous H_0 . Au vu de H_1 (T est typiquement plus grande sous H_1 que sous H_0), on va prendre une région de rejet unilatère de la forme

$$\left\{ T > q_{0.95}^{\mathcal{T}(n-p)} \right\}.$$

Ici T se réécrit

$$T = \frac{\hat{\beta} + \hat{\gamma} - 1}{\hat{\sigma} \sqrt{H_{1,1} + H_{2,2} + 2H_{1,2}}},$$

avec $H = (X^T X)^{-1}$. L'application numérique donne :

$$T = \frac{0.738 + 0.282 - 1}{\sqrt{0.09 \times (0.0016 + 0.0009 - 2 \times 0.0010)}} \approx 2.98 > 1.645.$$

On rejette H_0 : les rendements ne sont pas constants, on accepte qu'ils sont croissants au niveau 5%.

6. On a obtenu $\sum_{i=1}^n Y_i^2 = 16650.41$ et $\sum_{i=1}^n Y_i = 5199.49$. Calculer le coefficient du R^2 du modèle. Commenter.

2 points : 0.5 pour une formule du R^2 , 1 point pour le calcul du SCT et l'appli numérique, 0.5 pour le commentaire (même bref).

Solution. On rappelle que le coefficient de détermination est défini par :

$$R^2 = 1 - \frac{SCR}{SCT},$$

où

$$SCT = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n\bar{Y}^2, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

On a les valeurs suivantes :

$$\sum_{i=1}^n Y_i^2 = 16650.41, \quad \sum_{i=1}^n Y_i = 5199.49, \quad n = 1658, \quad \text{et} \quad SCR = 148.27.$$

On calcule d'abord la moyenne :

$$\bar{Y} = \frac{5199.49}{1658} \approx 3.136.$$

Ainsi,

$$SCT = 16650.41 - 1658 \times 3.136^2 \approx 16650.41 - 16305.60 = 344.81.$$

On en déduit :

$$R^2 = 1 - \frac{148.27}{344.81} \approx 1 - 0.43 = 0.57.$$

Le modèle explique environ 57% de la variabilité des $\log(V_i)$. Cela traduit un ajustement modéré : la force de travail et le capital expliquent une part importante, mais non totale, de la production. Il subsiste environ 43% de la variance non expliquée, probablement due à d'autres facteurs économiques ou à des effets spécifiques aux entreprises. Ce n'est pas si satisfaisant, quand on pense qu'on est à l'échelle log!

7. Le modèle étudié ici permet-il de mettre en lumière le lien de causalité entre la force de travail, le capital et la production ? Commenter.

1 point : on n'attend pas une rédaction aussi longue que dans le corrigé. Juste dire que corrélation et causalité ne sont pas la même chose, c'est ça qui était attendu.

Solution. Le modèle de Cobb–Douglas estimé ici met en évidence une corrélation entre la production, la force de travail et le capital, mais il ne permet pas d'établir une relation causale. En effet, il s'agit d'un modèle linéaire estimé par moindres carrés ordinaires, fondé sur des données observationnelles. Sans dispositif expérimental ni hypothèses fortes d'exogénéité des variables L_i et K_i , on ne peut exclure la présence de variables omises (comme la technologie, la qualité de la main-d'œuvre ou l'efficacité managériale) ou d'une causalité inverse (par exemple, les entreprises les plus productives investissent davantage en capital et en travail). De plus, le modèle est très simplifié, puisqu'il ne considère que deux facteurs explicatifs, alors que la production dépend en réalité de nombreux autres déterminants. Ainsi, le modèle met en lumière des relations statistiques intéressantes, mais ne permet pas à lui seul de conclure à un effet causal de la force de travail et du capital sur la production.

$1 - \alpha$ \ loi	$\mathcal{T}(k)$ avec $k \geq 1000$	$\mathcal{N}(0, 1)$
0.90	1.282	1.282
0.95	1.645	1.645
0.975	1.960	1.960

Tableau de quantiles d'ordre $1 - \alpha$ pour différentes lois.

Exercice 3 – QCM (5 points)

Ces questions n'ont aucun lien entre elles. Indiquer sur la copie, sans justification, l'unique bonne réponse à chaque question. Chaque bonne réponse rapporte 1 point, une mauvaise réponse ne rapporte ni n'enlève aucun point.

Comme l'énoncé l'indique : 1 point par bonne réponse, pas de notation négative.

Solution. D/B/D/D/A

- On veut tester H_0 contre H_1 . Supposons qu'on dispose d'un test de niveau d'erreur α , dont la zone de rejet s'écrit $\{T > -1.5\}$ pour T une statistique de test. Avec les données dont je dispose, T prend la valeur -0.8 .
 - Je conserve H_0 , avec une probabilité d'erreur d'au moins $1 - \alpha$.
 - Je rejette H_0 , avec une probabilité d'erreur que je ne contrôle pas.
 - La statistique de test T est toujours plus petite sous H_0 que sous H_1 .
 - Je m'expose à un risque de première espèce.

Solution. Je rejette H_0 d'après la forme de la région de rejet. Lorsque je fais cela, le cas où je me trompe, c'est celui où j'ai rejeté à tort H_0 . C'est donc l'erreur de première espèce. Réponse D.

2. On se place dans un modèle où X_1, \dots, X_n sont des variables aléatoires i.i.d. réelles, de loi \mathbb{P}_θ de densité $x \mapsto f(x - \theta)$ par rapport à la mesure de Lebesgue sur \mathbb{R} , où f est une densité connue sur \mathbb{R} , et $\theta \in \mathbb{R}$ est un paramètre inconnu. On suppose que ce modèle est régulier et on note \mathbb{E}_θ l'espérance sous \mathbb{P}_θ . Dans ce modèle, l'information de Fisher :

A) est une variable aléatoire qui dépend de θ .

B) ne dépend pas de θ .

C) vaut $n \times \mathbb{E}_\theta[(\log f(X_1 - \theta))^2]$.

D) vaut $n \times \mathbb{E}_\theta \left[-\frac{f'(X_1 - \theta)}{f(X_1 - \theta)} \right]$.

Solution. Sous les hypothèses de l'énoncé, on a que la log-vraisemblance du modèle s'écrit

$$\ell(\theta; X) = \sum_{i=1}^n \log f(X_i - \theta),$$

puis en dérivant,

$$\dot{\ell}(\theta; X) = - \sum_{i=1}^n \frac{f'(X_i - \theta)}{f(X_i - \theta)},$$

qui est toujours centrée. On a donc

$$\begin{aligned} I(\theta) &= \text{Var}[\dot{\ell}(\theta; X)] = n \text{Var} \left[-\frac{f'(X_1 - \theta)}{f(X_1 - \theta)} \right] \\ &= n \mathbb{E} \left[\left(-\frac{f'(X_1 - \theta)}{f(X_1 - \theta)} \right)^2 \right] \\ &= n \int_{\mathbb{R}} \frac{f'(x - \theta)^2}{f(x - \theta)^2} f(x - \theta) dx \\ &\stackrel{u=x-\theta}{=} n \int_{\mathbb{R}} \frac{f'(u)^2}{f(u)} du, \end{aligned}$$

qui ne dépend plus de θ . Réponse B.

3. Soit $(Z_n)_{n \geq 0}$ une suite de variables aléatoires réelles strictement positives telle que $\sqrt{n}(Z_n - e) \rightarrow \mathcal{N}(0, 4)$, en loi, lorsque n tend vers ∞ . Alors, lorsque n tend vers ∞ , $\sqrt{n}(\log Z_n - 1)$ tend en loi vers :

- A) $\mathcal{N}(0, 4e)$ B) $\mathcal{N}(0, 4)$ C) $\mathcal{N}(0, 4/e)$ D) $\mathcal{N}(0, 4/e^2)$

Solution. On applique la méthode Delta pour une transformation $f : x \mapsto \log x$, dérivable en $x = e$ et de dérivée égale à $f'(e) = 1/e$. Donc, d'après la méthode Delta, $\sqrt{n}(\log Z_n - 1)$ tend en loi vers $f'(e)\mathcal{N}(0, 4) = \mathcal{N}(0, 4/e^2)$ (en loi). Réponse D.

4. On se place dans un modèle où X_1, \dots, X_n sont des variables aléatoires i.i.d. de loi $\mathcal{N}(\mu, 1)$ où $\mu \in \mathbb{R}$ est un paramètre inconnu. On note \bar{X} la moyenne empirique des X_i et q_β le quantile d'ordre β de la loi $\mathcal{N}(0, 1)$. On veut tester : $H_0 : \mu = 0$ contre $H_1 : \mu = 1$. Soit $\alpha \in]0, 1[$. Le test de région de rejet R est uniformément plus puissant de niveau d'erreur α :

A) lorsque $R = \left\{ \exp \left(-\frac{1}{2} \sum_{i=1}^n (X_i - 1)^2 + \frac{1}{2} \sum_{i=1}^n X_i^2 \right) > nq_{1-\alpha}^2 \right\}$.

B) lorsque $R = \left\{ \sum_{i=1}^n (X_i - 1)^2 - \sum_{i=1}^n X_i^2 < q_{1-\alpha}^2/n \right\}$.

C) lorsque $R = \left\{ \bar{X} > q_{1-\alpha} \right\}$.

D) lorsque $R = \left\{ \bar{X} > q_{1-\alpha}/\sqrt{n} \right\}$.

Solution. On applique le théorème de Neyman-Pearson, qui nous dit qu'on peut essayer de chercher une telle région \mathcal{R} test sous la forme

$$\mathcal{R} = \left\{ \frac{L_1(X)}{L_0(X)} > k_\alpha \right\},$$

où k_α est à fixer pour avoir une taille α et $L_{0/1}$ désigne la vraisemblance sous $H_{0/1}$. En regardant les calculs, on a

$$\begin{aligned} \frac{L_1(X)}{L_0(X)} > k_\alpha &\iff \frac{\frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (X_i - 1)^2\right)}{\frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n X_i^2\right)} > k_\alpha \\ &\iff \exp\left(-\frac{1}{2} \sum_{i=1}^n (X_i - 1)^2 + \frac{1}{2} \sum_{i=1}^n X_i^2\right) > k_\alpha \\ &\iff \exp\left(\sum_{i=1}^n X_i - n/2\right) > k_\alpha \\ &\iff \frac{1}{n} \sum_{i=1}^n X_i > 1/2 + \log(k_\alpha)/n. \end{aligned}$$

Or, on connaît la loi de $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ sous H_0 : c'est celle d'une variable de loi $1/\sqrt{n} \times \mathcal{N}(0, 1)$. Le test ci-dessus est de taille α si et seulement si $1/2 + \log(k_\alpha)/n = q_{1-\alpha}/\sqrt{n}$. Réponse D.

5. Soit $Z \sim \mathcal{N}(0_d, I_d)$, et $A \in \mathbb{R}^{d \times d}$ une matrice de rang r . On définit la variable aléatoire Y de la façon suivante :

$$Y := Z - \arg \min_{u \in \mathbb{R}^d : Au = AZ} \|u\|^2.$$

- A) Y et $Z - Y$ sont indépendants et $\mathbb{E}[\|Y\|^2] = d - r$.
- B) Y et $Z - Y$ sont indépendants et $\mathbb{E}[\|Y\|^2] = r$.
- C) Z et $Z - Y$ sont indépendants et $\mathbb{E}[\|Y\|^2] = d - r$.
- D) Z et $Z - Y$ sont indépendants et $\mathbb{E}[\|Y\|^2] = r$.

Solution. L'espace affine sur lequel porte la minimisation s'écrit

$$\{u \in \mathbb{R}^d : Au = AZ\} = Z + \text{Ker } A.$$

Cela entraîne la réécriture

$$Y = Z - (Z + \arg \min_{v \in \text{Ker } A} \|Z + v\|^2) = Z - (Z - \Pi_{\text{Ker } A}(Z)) = \Pi_{\text{Ker } A}(Z),$$

où $\Pi_{\text{Ker } A}$ est le projecteur orthogonal sur $\text{Ker } A$. D'après le théorème de Cochran, $Y = \Pi_{\text{Ker } A}(Z)$ et $Z - Y = (\text{id} - \Pi_{\text{Ker } A})(Z) = \Pi_{(\text{Ker } A)^\perp}(Z)$ sont indépendantes, et $\|Y\|^2$ suit une loi $\chi^2(\dim \text{Ker } A) = \chi^2(d - r)$ par le théorème du rang. On conclut en utilisant $\mathbb{E}[\chi^2(\ell)] = \ell$. Réponse A.

Fin du sujet.