

UGA L3 Miash:
Économétrie 1
Modèle de régression linéaire et moindres carrés

Michal W. Urdanivia*

* Université de Grenoble Alpes, Faculté d'Économie, GAEL,
e-mail: michal.wong-urdanivia@univ-grenoble-alpes.fr

18 septembre 2017

Contenu

1. Modèle de régression linéaire et moindres carrés

Contenu

1. Modèle de régression linéaire et moindres carrés
2. Géométrie des moindres carrés

Contenu

1. Modèle de régression linéaire et moindres carrés
2. Géométrie des moindres carrés
3. Intervalles de confiance

Contenu

1. Modèle de régression linéaire et moindres carrés
2. Géométrie des moindres carrés
3. Intervalles de confiance
4. Tests d'hypothèses

Section 1

Modèle de régression linéaire et moindres carrés

Définitions

- On s'intéresse à l'effet d'un groupe de variables $X \in \mathbb{R}^K$, traditionnellement appelées *régresseurs*, sur une autre variable $Y \in \mathbb{R}$ traditionnellement appelée *variable dépendante*.
- On dispose de données $\{(Y_i, X_i)\}_{i=1}^n$ où Y_i est une variable aléatoire et X_i est un vecteur $K \times 1$,

$$X_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ \cdot \\ \cdot \\ \cdot \\ X_{iK} \end{pmatrix}$$

- (Y_i, X_i) est appelée observation (sous entendu de (Y, X)), et l'ensemble des n observations est un *échantillon*.
- Le vecteur X_i contient les valeurs des K variables pour l'observation i .

Définitions

- Pour des *données en coupe* il est souvent supposé que toutes les observations sont tirées indépendamment les unes des autres à partir d'une même distribution (i.e., loi de probabilité).
- On dit dans ce cas que l'échantillon d'observations $\{(Y_i, X_i)\}_{i=1}^n$ est un échantillon aléatoire ou de manière équivalente que les observations sont identiquement et indépendamment distribuées (i.i.d.).
- L'hypothèse d'observations i.i.d. ne signifie pas que Y_i et X_i soient indépendants, mais plutôt que l'observation (Y_i, X_i) est indépendante de toute autre observation (Y_j, X_j) pour $i \neq j$, n'excluant donc pas que Y_i et X_i puissent être liés.
- Un outil pour étudier la relation entre la variable dépendante et les régresseurs est l'espérance conditionnelle de Y_i sachant X_i , $\mathbb{E}(Y_i|X_i)$, laquelle vue comme une fonction de X_i est appelée *fonction de régression*.
- La différence entre Y_i et son espérance conditionnelle (i.e., sa fonction de régression) est appelée *terme d'erreur* (ou plus succinctement *erreur*),

$$U_i = Y_i - \mathbb{E}(Y_i|X_i) \quad (1)$$

Définitions

- U_i n'est pas une variable observable par l'analyste étant donné que l'espérance conditionnelle lui est inconnue.
- Dans un cadre *paramétrique* ou *semi-paramétrique*, il est souvent supposé que l'espérance conditionnelle est connue à un ensemble de *paramètres* près.
- Ainsi dans le modèle de régression linéaire on suppose que $\mathbb{E}(Y_i|X_i)$ est linéaire par rapport à un vecteur de paramètres inconnus,

$$\mathbb{E}(Y_i|X_i) = X_{i1}\beta_1 + X_{i2}\beta_2 + \dots + X_{iK}\beta_K = X_i^\top \beta \quad (2)$$

où,

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{pmatrix}$$

Définitions

est un vecteur de K paramètres constants.

- La linéarité de $\mathbb{E}(Y_i|X_i)$ peut être justifiée, si par exemple, la distribution des observations $\{(Y_i, X_i)\}_{i=1}^n$ est une loi normale multivariée.
- Rappelons néanmoins que lorsque $\mathbb{E}(Y_i|X_i)$ n'est pas linéaire il est possible de caractériser β de manière à ce que (2) constitue la *meilleure prédiction linéaire* de la variable dépendante par les régresseurs.
- Notons aussi que comme

$$\beta_k = \frac{\delta \mathbb{E}(Y_i|X_i)}{\delta X_{ik}}, \quad k = 1, 2, \dots, K.$$

le vecteur β est le vecteur des *effets marginaux* des régresseurs, i.e., β_k donne la variation dans l'espérance conditionnelle de Y_i lorsque le régresseur X_{ik} varie, pour des valeurs fixes des autres régresseurs X_{il} , $l = 1, 2, \dots, K$, $l \neq k$.

- Ceci est une des raisons pour lesquelles un des principaux objectifs est l'estimation du vecteur inconnu β à partir des données.

Définitions

- Les équations (1) et (2) permettent d'écrire,

$$Y_i = X_i^\top \beta + U_i \quad (3)$$

où par définition de (1)

$$\mathbb{E}(U_i|X_i) = 0 \quad (4)$$

- Ceci implique que les régresseurs ne contiennent aucune information quant à l'écart entre Y_i et son espérance conditionnelle.
- En outre, *la loi des conditionnements successifs* implique que les erreurs ont une espérance nulle : $\mathbb{E}(U_i) = 0$. Notons aussi qu'avec des observations i.i.d. les erreurs sont aussi i.i.d.
- Une hypothèse fréquente sur les erreurs consiste à supposer qu'ils sont *homoscédastiques* (on parle d'hypothèse d'homoscédasticité), par quoi on entend que leur variance est indépendante des régresseurs, et la même pour toutes les observations,

$$\text{Var}(U_i|X_i) = \sigma^2$$

pour une constante $\sigma^2 > 0$.

Hypothèses

- Introduisons les notations vectorielles et matricielles suivantes,

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} X_1^\top \\ X_2^\top \\ \vdots \\ X_n^\top \end{pmatrix} = \begin{pmatrix} X_{11} & X_{12} & \cdot & \cdot & \cdot & X_{1K} \\ X_{21} & X_{22} & \cdot & \cdot & \cdot & X_{2K} \\ \vdots & \vdots & & & & \\ X_{n1} & X_{n2} & \cdot & \cdot & \cdot & X_{nK} \end{pmatrix}$$

Hypothèses

$$\mathbf{U} = \begin{pmatrix} U_1 \\ U_2 \\ \cdot \\ \cdot \\ U_n \end{pmatrix}$$

- Le modèle de régression linéaire consiste dans les hypothèses suivantes :

Hypothèse H1

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{U}$$

Hypothèse H2

$$\mathbb{E}(\mathbf{U}|\mathbf{X}) = 0 \text{ p.s.}$$

Hypothèse H3

$$\text{Var}(\mathbf{U}|\mathbf{X}) = \sigma^2 \mathbf{I}_n \text{ p.s.}$$

Hypothèse H4

$$\text{Rang}(\mathbf{X}) = K \text{ p.s.}$$

Hypothèses

- Pour l'inférence nous supposons parfois que,

Hypothèse H5

$$\mathbf{U}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$$

- Les hypothèses H1-H5 définissent alors le *modèle de régression linéaire normal* avec,

$$\mathbf{Y}|\mathbf{X} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2\mathbf{I}_n)$$

- Remarquons qu'étant donné que les covariances dans H5 sont toutes nulles, H5 implique l'indépendance des erreurs.
- Les hypothèses H1-H4 seules, n'impliquent pas l'indépendance entre les observations.
- En fait, plusieurs résultats importants n'exigent pas d'observations indépendantes.

Hypothèses

- Néanmoins, nous supposons parfois l'indépendance sans la normalité.

Hypothèse H6

Les observations $\{(Y_i, X_i)\}_{i=1}^n$ sont i.i.d..

Dans le cas de régresseurs fixes cette hypothèse peut être remplacé par celle d'erreurs i.i.d.

- L'hypothèse H2 dit que \mathbf{U} est indépendant de \mathbf{X} en espérance, ce qui est une hypothèse forte.
- Cependant, plusieurs résultats importants peuvent être obtenus avec une hypothèse plus faible de non corrélation.

Hypothèse H7

Pour $i = 1, 2, \dots, n$, $\mathbb{E}(U_i X_i) = 0$, et $\mathbb{E}(U_i) = 0$.

- Toutefois sous cette condition $X_i^\top \beta$ ne peut pas s'interpréter comme une espérance conditionnelle, auquel cas (3) doit être vu comme un *processus générateur des données*.

Hypothèses

- L'hypothèse H3 implique que les erreurs U_i ont la même variance pour tout i , et ne sont pas corrélés entre eux, i.e., $\mathbb{E}(U_i U_j | \mathbf{X}) = 0$ pour $i \neq j$.
- Notons que l'indépendance entre les erreurs peut aussi être obtenue avec la condition H5 ou sous les conditions H1 et H6.
- L'hypothèse H4 exige que les colonnes de \mathbf{X} soient linéairement indépendantes.
- Que cette hypothèse ne soit pas vérifiée signifie qu'un ou plus de régresseurs duplique l'information contenue dans les autres, et ce faisant doit être écarté.
- Souvent, une des colonnes de \mathbf{X} est le vecteur unitaire et le paramètre qui lui est associé est appelé *constante*.
- La constante du modèle donne la valeur moyenne de la variable dépendante lorsque tous les régresseurs sont égaux à zéro.

Estimation par la méthode des moments

- Nous allons construire des estimateurs des paramètres β et σ^2 .
- Une des méthodes les plus ancienne pour construire des estimateurs est la *méthode des moments*(MM).
- La MM consiste à construire des estimateurs pour des paramètres définis par des moments théoriques en considérant les contreparties empiriques de ces moments appelées alors moments empiriques.
- Par exemple si un paramètre est défini au travers d'une espérance(moment théorique), son estimateur sera construit à partir d'une moyenne(moment empirique) calculée sur les observations.
- Les hypothèses **H1**, et **H2** ou **H7** impliquent que la vraie valeur de β doit satisfaire,

$$\mathbb{E}(U_i X_i) = \mathbb{E}\left((Y_i - X_i^\top \beta) X_i\right) = 0 \quad (5)$$

Estimation par la méthode des moments

- Un *estimateur des moments* (i.e., obtenu selon la MM) de β , $\hat{\beta}$, est obtenu en remplaçant l'espérance dans (5) par la moyenne empirique,

$$n^{-1} \sum_{i=1}^N (Y_i - X_i^\top \hat{\beta}) X_i = n^{-1} \sum_{i=1}^N X_i Y_i - n^{-1} \sum_{i=1}^N X_i X_i^\top \hat{\beta} = 0 \quad (6)$$

- En résolvant par rapport à $\hat{\beta}$ on obtient,

$$\hat{\beta} = \left(n^{-1} \sum_{i=1}^N X_i X_i^\top \right)^{-1} n^{-1} \sum_{i=1}^N X_i Y_i \quad (7)$$

qui peut s'écrire alternativement,

$$\hat{\beta} = \left(\sum_{i=1}^N X_i X_i^\top \right)^{-1} \sum_{i=1}^N X_i Y_i \quad (8)$$

Estimation par la méthode des moments

ou,

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (9)$$

où l'on note que la matrice $\sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^\top = \mathbf{X}^\top \mathbf{X}$ est inversible sous l'hypothèse H4.

- On définit les *valeurs ajustées* ou *prédictions*, ainsi qu'un vecteur $n \times 1$ des valeurs ajustées ou des prédictions, par respectivement,

$$\hat{Y}_i = \mathbf{X}_i^\top \hat{\beta}, \quad \hat{\mathbf{Y}} = (\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_n)^\top$$

- On définit les *résidus*, et le vecteur $n \times 1$ des résidus, par respectivement,

$$\hat{U}_i = Y_i - \mathbf{X}_i^\top \hat{\beta}, \quad \hat{\mathbf{U}} = (\hat{U}_1, \hat{U}_2, \dots, \hat{U}_n)^\top$$

Estimation par la méthode des moments

- Du fait de (6) le vecteur des résidus vérifie les K *équations normales*,

$$\sum_{i=1}^N \hat{U}_i X_i = \begin{pmatrix} \sum_{i=1}^N \hat{U}_i X_{i1} \\ \sum_{i=1}^N \hat{U}_i X_{i2} \\ . \\ . \\ . \\ \sum_{i=1}^N \hat{U}_i X_{iK} \end{pmatrix} = 0 \quad (10)$$

ou en notation matricielle,

$$\mathbf{X}^\top \hat{\mathbf{U}} = 0 \quad (11)$$

Estimation par la méthode des moments

- Si le modèle contient une constante alors il résulte des équations normales que $\sum_{i=1}^N \hat{U}_i = 0$ (il suffit en effet de considérer que, par exemple, le premier régresseur est constant et égal à 1).
- Afin d'estimer σ^2 considérons,

$$\sigma^2 = \mathbb{E}(U_i^2) = \mathbb{E}\left((Y_i - X_i^\top \beta)^2\right)$$

- Dans la mesure où β , est inconnu un estimateur sera obtenu en remplaçant β par son estimateur des moments,

$$\hat{\sigma}^2 = n^{-1} \sum_{i=1}^N (Y_i - X_i^\top \hat{\beta})^2 \quad (12)$$

Moindres carrés

- Soit le problème consistant à minimiser l'erreur de prédiction quand on cherche à prédire Y_i par son espérance conditionnelle, $\mathbb{E}(Y_i|X_i)$, supposée être une fonction linéaire telle que (2).
- Plus précisément, $Y_i - \mathbb{E}(Y_i|X_i)$ étant l'erreur de prédiction on cherche β qui minimise un critère de perte quadratique,

$$\beta \in \arg \min_{b \in \mathbb{R}^K} S(b)$$

où $S(b) = \mathbb{E}((Y_i - X_i^\top b)^2)$.

- La contrepartie empirique de ce problème permet de définir un estimateur de β par,

$$\hat{\beta} \in \arg \min_{b \in \mathbb{R}^K} S_n(b)$$

où $S_n(b) = n^{-1} \sum_{i=1}^N ((Y_i - X_i^\top b)^2)$, est la contrepartie empirique de la fonction objectif $S(b)$.

Moindres carrés

- Nous pouvons montrer que l'estimateur des moments de la section précédente est aussi l'estimateur des moindres carrés.
- La fonction objectif précédente peut s'écrire (voir notes),

$$S_n(b) = (\mathbf{Y} - \mathbf{X}\hat{\beta})^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}) + (\hat{\beta} - b)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - b)$$

- La minimisation de $S_n(b)$ équivaut à minimiser $(\hat{\beta} - b)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - b)$ car $(\mathbf{Y} - \mathbf{X}\hat{\beta})^\top (\mathbf{Y} - \mathbf{X}\hat{\beta})$ ne fait pas intervenir b .
- Sous l'hypothèse H4 la matrice \mathbf{X} est de plein rang, et dans ce cas $\mathbf{X}^\top \mathbf{X}$ est définie positive,

$$(\hat{\beta} - b)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - b) \geq 0$$

et $(\hat{\beta} - b)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - b) = 0$ ssi $\hat{\beta} = b$.

- Alternativement, nous pouvons montrer que $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ est l'estimateur des moindres carrés de β (i.e., il minimise $S_n(b)$).

Moindres carrés

- Pour cela, écrivons,

$$S(b) = \mathbf{Y}^\top \mathbf{Y} - 2b^\top \mathbf{X}^\top \mathbf{Y} + b^\top \mathbf{X}^\top \mathbf{X} b$$

- En utilisant le fait que pour une matrice symétrique \mathbf{A} ,

$$\frac{\delta(\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\delta \mathbf{x}} = 2\mathbf{A} \mathbf{x}$$

la condition du premier ordre est,

$$\frac{\delta S_n(\hat{\beta})}{\delta b} = -2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X} \hat{\beta} = 0$$

ce qui permet d'obtenir,

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- Remarquons aussi que les conditions du premier ordre peuvent s'écrire $\mathbf{X}^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}) = 0$, ce qui correspond aux équations normales vue précédemment.

Propriétés de l'estimateur des moindres carrés

- Nous allons présenter un certain nombre de propriétés de l'estimateur des moindres carrés.

Proposition 1

$\hat{\beta}$ est un estimateur linéaire.

- Un estimateur b est linéaire s'il peut s'écrire comme $b = \mathbf{A}\mathbf{Y}$, où \mathbf{A} est une matrice quelconque qui dépend de \mathbf{X} uniquement, et ne dépend pas de \mathbf{Y} .
- Pour l'estimateur des moindres carrés nous avons, $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

Proposition 2

Sous les hypothèses *H1*, *H2*, et *H4*, $\hat{\beta}$ est sans biais, i.e.,

$$\mathbb{E}(\hat{\beta}) = \beta$$

Propriétés de l'estimateur des moindres carrés

- Pour montrer cette propriété écrivons, en utilisant l'hypothèse H1,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \mathbf{U}) = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U}$$

Calculons l'espérance conditionnelle de $\hat{\beta}$,

$$\mathbb{E}(\hat{\beta} | \mathbf{X}) = \mathbb{E}(\beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U} | \mathbf{X}) = \beta + \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U} | \mathbf{X})$$

Notons que,

$$\mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U} | \mathbf{X}) = (\mathbf{X} \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{U} | \mathbf{X}) = 0$$

car sous l'hypothèse H2, $\mathbb{E}(\mathbf{U} | \mathbf{X}) = 0$. Nous avons donc,

$$\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta \tag{13}$$

et par la loi des espérances itérées,

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(\mathbb{E}(\hat{\beta} | \mathbf{X})) = \beta$$

Propriétés de l'estimateur des moindres carrés

- L'équation (13) montre que $\hat{\beta}$ est conditionnellement sans biais sachant \mathbf{X} .
- On remarque aussi que pour que $\hat{\beta}$ soit sans biais l'hypothèse H7 n'est pas suffisante.

Proposition 3

Sous les hypothèses H1, H2, et H4,

$$\text{Var}(\hat{\beta}|\mathbf{X}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{U}\mathbf{U}^\top | \mathbf{X}) \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}$$

et avec des erreurs homoscédastiques(i.e., sous l'hypothèse H3),

$$\text{Var}(\hat{\beta}|\mathbf{X}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

Propriétés de l'estimateur des moindres carrés

- Pour montrer ces résultats, partons de la définition de la variance conditionnelle de $\hat{\beta}$,

$$\begin{aligned}
 \text{Var}(\hat{\beta}|\mathbf{X}) &= \mathbb{E} \left(\left(\hat{\beta} - \mathbb{E}(\hat{\beta}|\mathbf{X}) \right) \left(\hat{\beta} - \mathbb{E}(\hat{\beta}|\mathbf{X}) \right)^\top | \mathbf{X} \right) \\
 &= \mathbb{E} \left(\left(\hat{\beta} - \beta \right) \left(\hat{\beta} - \beta \right)^\top | \mathbf{X} \right) \\
 &= \mathbb{E} \left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} | \mathbf{X} \right) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{U} \mathbf{U}^\top | \mathbf{X}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}
 \end{aligned}$$

- Et avec des erreurs homoscédastiques, $\mathbb{E}(\mathbf{U} \mathbf{U}^\top | \mathbf{X}) = \sigma^2 \mathbf{I}_n$, de sorte que,

$$\begin{aligned}
 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{U} \mathbf{U}^\top | \mathbf{X}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}
 \end{aligned}$$

Propriétés de l'estimateur des moindres carrés

- Notons qu'avec des régresseurs fixes $\text{Var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$.

Proposition 4

Sous les hypothèses *H1* - *H5*,

$$\hat{\beta}|\mathbf{X} \sim \mathcal{N}\left(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\right)$$

- Il est suffisant ici de montrer ici que conditionnellement à \mathbf{X} la distribution de $\hat{\beta}$ est normale.
- On aura alors que, $\hat{\beta}|\mathbf{X} \sim \mathcal{N}\left(\mathbb{E}(\hat{\beta}|\mathbf{X}), \text{Var}(\hat{\beta}|\mathbf{X})\right)$.
- Néanmoins la normalité de $\hat{\beta}|\mathbf{X}$ résulte ici de ce que $\hat{\beta}$ est une fonction de linéaire de \mathbf{Y} , et que sous l'hypothèse *H5* $\mathbf{Y}|\mathbf{X}$ est normale.

Propriétés de l'estimateur des moindres carrés

- Dans le cas de régresseur fixes, il suffit d'omettre le conditionnement par rapport à \mathbf{X} et,

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$$

Proposition 5

(Éfficacité ou théorème de Gauss-Markov.) Sous les hypothèses H1-H4, l'estimateur des moindres carrés est le meilleur estimateur linéaire sans biais de β , dans le sens où il s'agit de l'estimateur, dans la classe des estimateurs linéaires et sans biais, qui présente la plus petite variance. i.e., pour tout estimateur linéaire sans biais, b , la matrice $\text{Var}(b|\mathbf{X}) - \text{Var}(\hat{\beta}|\mathbf{X})$ doit être semi-définie positive :

$$\text{Var}(b|\mathbf{X}) - \text{Var}(\hat{\beta}|\mathbf{X}) \geq 0$$

En outre, si $\tilde{\beta}$ est un estimateur linéaire et sans biais et $\text{Var}(\tilde{\beta}|\mathbf{X}) = \text{Var}(\hat{\beta}|\mathbf{X})$, alors $\tilde{\beta} = \hat{\beta}$ p.s.

Propriétés de l'estimateur des moindres carrés

- Avant de démontrer ce résultat notons qu'il discute la variance conditionnelle de l'estimateur des moindres carrés, et ce faisant il se réfère à des estimateurs conditionnellement sans biais.
- Soit b un estimateur linéaire sans biais de β . Il doit ainsi vérifier,

$$b = \mathbf{A}\mathbf{Y}, \quad \mathbb{E}(b|\mathbf{X}) = \beta$$

- Ces deux conditions impliquent que $\mathbf{A}\mathbf{X} = \mathbf{I}_K$ p.s. En effet,

$$\begin{aligned}\mathbb{E}(b|\mathbf{X}) &= \mathbb{E}(\mathbf{A}(\mathbf{X}\beta + \mathbf{U})) \\ &= \mathbf{A}\mathbf{X}\beta + \mathbf{A}\mathbb{E}(\mathbf{U}|\mathbf{X})\end{aligned}$$

- Par l'hypothèse H2, $\mathbb{E}(\mathbf{U}|\mathbf{X}) = 0$, et par conséquent, pour que b soit sans biais nous avons besoin de $\mathbf{A}\mathbf{X} = \mathbf{I}_K$.

Propriétés de l'estimateur des moindres carrés

- Montrons maintenant que $\text{Cov}(\hat{\beta}, b|\mathbf{X}) = \text{Var}(\hat{\beta}|\mathbf{X})$,

$$\begin{aligned}
 \text{Cov}(\hat{\beta}, b|\mathbf{X}) &= \mathbb{E}((\hat{\beta} - \beta)(b - \beta)^\top) \\
 &= \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{A}^\top | \mathbf{X}) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{U} \mathbf{U}^\top | \mathbf{X}) \mathbf{A}^\top \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{A}^\top \text{ (car sous H3, } \mathbb{E}(\mathbf{U} \mathbf{U}^\top | \mathbf{X}) = \sigma^2 \mathbf{I}_n) \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \text{ (car, } \mathbf{X}^\top \mathbf{A}^\top = \mathbf{I}_K) \\
 &= \text{Var}(\hat{\beta}|\mathbf{X})
 \end{aligned}$$

Finalement,

$$\begin{aligned}
 \text{Var}(\hat{\beta} - b|\mathbf{X}) &= \text{Var}(\hat{\beta}|\mathbf{X}) - \text{Cov}(\hat{\beta}, b|\mathbf{X}) - \text{Cov}(b, \hat{\beta}|\mathbf{X}) + \text{Var}(b|\mathbf{X}) \\
 &= \text{Var}(b|\mathbf{X}) - \text{Var}(\hat{\beta}|\mathbf{X})
 \end{aligned} \tag{14}$$

Propriétés de l'estimateur des moindres carrés

et notons que dans la mesure où toute matrice de variance-covariances est semi-définie positive, nous avons,

$$\text{Var}(b|\mathbf{X}) - \text{Var}(\hat{\beta}|\mathbf{X}) \geq 0$$

- Pour démontrer l'unicité, considérons un estimateur linéaire sans biais $\tilde{\beta}$ tel que $\text{Var}(\tilde{\beta}|\mathbf{X}) = \text{Var}(\hat{\beta}|\mathbf{X})$.
- Alors, par (14), $\text{Var}(\hat{\beta} - b|\mathbf{X}) = 0$, et par conséquent, $\tilde{\beta} = \hat{\beta} + c(\mathbf{X})$ pour une fonction $c(\mathbf{X})$ à valeurs dans \mathbb{R}^K qui dépend uniquement de \mathbf{X} .
- Cependant, comme $\hat{\beta}$ et $\tilde{\beta}$ sont conditionnellement sans biais sachant \mathbf{X} , il s'en suit que $c(\mathbf{X}) = 0$ p.s.
- Notons que l'hypothèse H3, $\mathbb{E}(\mathbf{U}\mathbf{U}^\top|\mathbf{X}) = \sigma^2\mathbf{I}_n$, joue un rôle crucial dans la démonstration du résultat précédent.
- Sans elle, il ne serait pas possible de tirer des conclusions quant à l'efficacité de l'estimateur des moindres carrés.

Section 2

Géométrie des moindres carrés

Matrices de projection

- Nous pouvons penser à \mathbf{Y} et aux colonnes \mathbf{X} comme des éléments de l'espace euclidien à n dimensions, \mathbb{R}^n .
- Considérons le sous-espace de \mathbb{R}^n appelé l'*espace des colonnes* de la matrice $n \times K$, \mathbf{X} .
- C'est la collection de tous les vecteurs dans \mathbb{R}^n qui peuvent s'écrire comme des combinaisons linéaires des colonnes de \mathbf{X} ,

$$\mathcal{S}(\mathbf{X}) = \left\{ z \in \mathbb{R}^n : z = \mathbf{X}b, b = (b_1, b_2, \dots, b_K) \in \mathbb{R}^K \right\}$$

- Rappelons maintenant qu'étant donné deux vecteurs a, b , dans \mathbb{R}^n , la distance entre a et b est donné par la norme euclidienne de leur différence $\|a - b\| = \sqrt{(a - b)^\top (a - b)}$.
- Par conséquent, le problème de la minimisation de la somme des carrés des erreurs, $(\mathbf{Y} - \mathbf{X}b)^\top (\mathbf{Y} - \mathbf{X}b)$, consiste à trouver, parmi tous les éléments de $\mathcal{S}(\mathbf{X})$, celui dont la distance par rapport à \mathbf{Y} est la plus petite,

$$\min_{\tilde{\mathbf{Y}} \in \mathcal{S}(\mathbf{X})} \|\mathbf{Y} - \tilde{\mathbf{Y}}\|^2$$

Matrices de projection

- Une solution au problème des moindres carrés, $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ doit être choisie de sorte que le vecteur des résidus, $\hat{\mathbf{U}} = \mathbf{Y} - \hat{\mathbf{Y}}$ soit orthogonal(perpendiculaire) à chaque colonne de \mathbf{X} ,

$$\hat{\mathbf{U}}^\top \mathbf{X} = 0$$

- Un résultat de cela est que $\hat{\mathbf{U}}$ est orthogonal à chaque élément de $\mathcal{S}(\mathbf{X})$: si $\mathbf{z} \in \mathcal{S}(\mathbf{X})$, alors il existe $\mathbf{b} \in \mathbb{R}^K$ tel que $\mathbf{z} = \mathbf{X}\mathbf{b}$, et,

$$\begin{aligned}\hat{\mathbf{U}}^\top \mathbf{z} &= \hat{\mathbf{U}}^\top \mathbf{X}\mathbf{b} \\ &= 0\end{aligned}$$

- La collection des éléments de \mathbb{R}^n orthogonaux à $\mathcal{S}(\mathbf{X})$ est appelée *complément orthogonal* de $\mathcal{S}(\mathbf{X})$,

$$\mathcal{S}^\perp(\mathbf{X}) = \left\{ \mathbf{z} \in \mathbb{R}^n : \mathbf{z}^\top \mathbf{X} = 0 \right\}$$

- Tout élément de $\mathcal{S}^\perp(\mathbf{X})$ est orthogonal à chaque élément de $\mathcal{S}(\mathbf{X})$.

Matrices de projection

- La solution au problème des moindres carrés est donnée par,

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\ &= \mathbf{P}_\mathbf{X} \mathbf{Y}\end{aligned}$$

où

$$\mathbf{P}_\mathbf{X} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

est appelée *matrice de projection orthogonale*.

- Pour tout vecteur $\mathbf{Y} \in \mathbb{R}^n$,

$$\mathbf{P}_\mathbf{X} \mathbf{Y} \in \mathcal{S}(\mathbf{X})$$

- En outre, le vecteur des résidus est dans $\mathcal{S}^\perp(\mathbf{X})$,

$$\mathbf{Y} - \mathbf{P}_\mathbf{X} \mathbf{Y} \in \mathcal{S}^\perp(\mathbf{X}) \tag{15}$$

Matrices de projection

- Pour montrer (15), notons d'abord, qu'étant donné que les colonnes de \mathbf{X} sont dans $\mathcal{S}(\mathbf{X})$,

$$\begin{aligned}\mathbf{P}_\mathbf{X}\mathbf{X} &= \mathbf{X}(\mathbf{X}\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X} \\ &= \mathbf{X}\end{aligned}$$

et comme $\mathbf{P}_\mathbf{X}$ est une matrice symétrique,

$$\mathbf{X}^\top\mathbf{P}_\mathbf{X} = \mathbf{X}^\top$$

Maintenant,

$$\begin{aligned}\mathbf{X}^\top(\mathbf{Y} - \mathbf{P}_\mathbf{X}\mathbf{Y}) &= \mathbf{X}^\top\mathbf{Y} - \mathbf{X}^\top\mathbf{P}_\mathbf{X}\mathbf{Y} \\ &= \mathbf{X}^\top\mathbf{Y} - \mathbf{X}^\top\mathbf{Y} \\ &= 0\end{aligned}$$

Matrices de projection

- Ainsi, par définition, les résidus $\mathbf{Y} - \mathbf{P}_X \mathbf{Y} \in \mathcal{S}^\perp(\mathbf{X})$. Les résidus peuvent s'écrire,

$$\begin{aligned}\hat{\mathbf{U}} &= \mathbf{Y} - \mathbf{P}_X \mathbf{Y} \\ &= (\mathbf{I}_n - \mathbf{P}_X) \mathbf{Y} \\ &= \mathbf{M}_X \mathbf{Y}\end{aligned}$$

où,

$$\begin{aligned}\mathbf{M}_X &= \mathbf{I}_n - \mathbf{P}_X \\ &= \mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top\end{aligned}$$

est une matrice de projection dans $\mathcal{S}^\perp(\mathbf{X})$.

- Les matrices \mathbf{P}_X et \mathbf{M}_X présentent les propriétés suivantes.
 1. $\mathbf{P}_X + \mathbf{M}_X = \mathbf{I}_n$. Ceci implique, que pour tout $\mathbf{Y} \in \mathbb{R}^n$,

$$\mathbf{Y} = \mathbf{P}_X \mathbf{Y} + \mathbf{M}_X \mathbf{Y}$$

Matrices de projection

2. P_X et M_X sont symétriques,

$$P_X^T = P_X, \quad M_X^T = M_X$$

3. P_X et M_X sont idempotentes,

$$P_X P_X = P_X, \quad M_X M_X = M_X$$

En effet,

$$\begin{aligned} P_X P_X &= \left(X(X^T X)^{-1} X^T \right) \left(X(X^T X)^{-1} X^T \right) \\ &= X(X^T X)^{-1} X^T \\ &= P_X \end{aligned}$$

de même,

$$\begin{aligned} M_X M_X &= (I_n - P_X)(I_n - P_X) \\ &= I_n - 2P_X + P_X P_X \\ &= I_n - P_X \\ &= M_X \end{aligned}$$

Matrices de projection

4. P_X et M_X sont orthogonales,

$$\begin{aligned}P_X M_X &= P_X (I_n - P_X) \\&= P_X - P_X P_X \\&= P_X - P_X \\&= 0\end{aligned}$$

Cette propriété implique que $M_X X = 0$. En effet,

$$\begin{aligned}M_X X &= (I_n - P_X) X \\&= X - P_X X \\&= X - X \\&= 0\end{aligned}$$

- Dans la discussion ci-dessus, aucune des hypothèses quant au modèle de régression n'ont été utilisées.
- Étant donné des données, Y et X , nous pouvons toujours calculer l'estimateur des moindres carrés, indépendamment du processus générateur des données derrière les données.

Matrices de projection

- Néanmoins, nous avons besoin d'un modèle(i.e., d'hypothèses) pour pouvoir discuter des propriétés d'un estimateur(e.g., le fait qu'il soit ou non sans biais, etc).

Propriétés de $\hat{\sigma}^2$

- σ^2 est estimé par,

$$\begin{aligned}\hat{\sigma}^2 &= n^{-1} \sum_{i=1}^N \hat{U}_i^2 \\ &= n^{-1} \hat{\mathbf{U}}^\top \hat{\mathbf{U}}\end{aligned}$$

- Mais sous H1 - H4 $\hat{\sigma}^2$ est biaisé.
- En effet,

$$\begin{aligned}\hat{\mathbf{U}} &= \mathbf{M}_X \mathbf{Y} \\ &= \mathbf{M}_X (\mathbf{X}\beta + \mathbf{U}) \\ &= \mathbf{M}_X \mathbf{U}\end{aligned}$$

où la dernière égalité résulte de ce que $\mathbf{M}_X \mathbf{X} = 0$.

Propriétés de $\hat{\sigma}^2$

- En outre,

$$\begin{aligned} n\hat{\sigma}^2 &= \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \\ &= \mathbf{U}^\top \mathbf{M}_X \mathbf{M}_X \mathbf{U} \\ &= \mathbf{U}^\top \mathbf{M}_X \mathbf{U} \end{aligned}$$

- Étant donné que $\mathbf{U}^\top \mathbf{M}_X \mathbf{U}$ est un scalaire,

$$\mathbf{U}^\top \mathbf{M}_X \mathbf{U} = \text{Tr} \left(\mathbf{U}^\top \mathbf{M}_X \mathbf{U} \right)$$

où $\text{Tr}(A)$ désigne la trace de la matrice A .

Propriétés de $\hat{\sigma}^2$

- Nous avons,

$$\begin{aligned}\mathbb{E} \left(\mathbf{U}^\top \mathbf{M}_X \mathbf{U} | \mathbf{X} \right) &= \mathbb{E} \left(\text{Tr}(\mathbf{U}^\top \mathbf{M}_X \mathbf{U}) | \mathbf{X} \right) \\ &= \mathbb{E} \left(\text{Tr}(\mathbf{M}_X \mathbf{U} \mathbf{U}^\top) | \mathbf{X} \right) \text{ (car } \text{Tr}(ABC) = \text{Tr}(BCA)) \\ &= \text{Tr} \left(\mathbf{M}_X \mathbb{E} \left(\mathbf{U} \mathbf{U}^\top \right) | \mathbf{X} \right) \\ &\quad \text{(car l'opérateur trace et l'espérance sont linéaires)} \\ &= \sigma^2 \text{Tr}(\mathbf{M}_X)\end{aligned}$$

- La dernière égalité résulte de ce que par l'hypothèse H3, $\mathbb{E}(\mathbf{U}^\top \mathbf{U}) = \sigma^2 \mathbf{I}_n$.

Propriétés de $\hat{\sigma}^2$

- Maintenant,

$$\begin{aligned}\text{Tr}(\mathbf{M}_X) &= \text{Tr} \left(\mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \\ &= \text{Tr}(\mathbf{I}_n) - \text{Tr} \left(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \\ &= \text{Tr}(\mathbf{I}_n) - \text{Tr} \left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \right) \\ &= \text{Tr}(\mathbf{I}_n) - \text{Tr}(\mathbf{I}_K) \\ &= n - K\end{aligned}$$

- Il s'en suit que,

$$\mathbb{E}(\hat{\sigma}^2) = \frac{n-K}{n} \sigma^2 \quad (16)$$

- L'estimateur $\hat{\sigma}^2$ est biaisé, mais le résultat précédent suggère qu'il est aisé de le modifier afin d'obtenir un estimateur sans biais.

Propriétés de $\hat{\sigma}^2$

- Pour cela, définissons,

$$\begin{aligned}s^2 &= \hat{\sigma}^2 \frac{n}{n-K} \\ &= (n-K)^{-1} \sum_{i=1}^N \hat{U}_i^2\end{aligned}$$

- Et il résulte de (16) que,

$$\mathbb{E}(s^2) = \sigma^2$$

Régression partitionnée

- Considérons,

$$\mathbf{X} = (\mathbf{X}_1 \ \mathbf{X}_2)$$

et écrivons le modèle comme suit,

$$\mathbf{Y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{U}$$

où \mathbf{X}_1 est une matrice $(n \times K_1)$, \mathbf{X}_2 est une matrice $(n \times K_2)$,
 $K_1 + K_2 = K$, et,

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

β_1 et β_2 étant des vecteurs de paramètres, $(K_1 \times 1)$ et $(K_2 \times 1)$.

- Concentrons nous sur \mathbf{X}_1 et β_1 .
- Soit l'estimateur des moindres carrés de β ,

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

Régression partitionnée

- Nous pouvons écrire la version suivante des équations normales,

$$(\mathbf{X}^\top \mathbf{X}) \hat{\beta} = \mathbf{X}^\top \mathbf{Y}$$

comme suit,

$$\begin{pmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{X}_2 \\ \mathbf{X}_2^\top \mathbf{X}_1 & \mathbf{X}_2^\top \mathbf{X}_2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1^\top \mathbf{Y} \\ \mathbf{X}_2^\top \mathbf{Y} \end{pmatrix}$$

- On peut obtenir des expressions pour $\hat{\beta}_1$ et $\hat{\beta}_2$ par inversion de la matrice partitionnée à gauche de l'équation ci-dessus.
- Alternativement, définissons \mathbf{M}_2 comme la matrice de projection sur l'espace orthogonal à l'espace $\mathcal{S}(\mathbf{X}_2)$,

$$\mathbf{M}_2 = \mathbf{I}_n - \mathbf{X}_2(\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top$$

alors,

$$\hat{\beta}_1 = (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y} \quad (17)$$

Régression partitionnée

- Pour montrer cela, commençons par écrire,

$$\mathbf{Y} = \mathbf{X}_1 \hat{\beta}_1 + \mathbf{X}_2 \hat{\beta}_2 + \hat{\mathbf{U}} \quad (18)$$

- Notons que par construction,

$$\mathbf{M}_2 \hat{\mathbf{U}} = \hat{\mathbf{U}} (\hat{\mathbf{U}} \text{ est orthogonal à } \mathbf{X}_2)$$

$$\mathbf{M}_2 \mathbf{X}_2 = 0$$

$$\mathbf{X}_1^\top \hat{\mathbf{U}} = 0$$

$$\mathbf{X}_2^\top \hat{\mathbf{U}} = 0$$

Régression partitionnée

- Substituons l'équation (18) dans la partie droite de l'équation (17),

$$\begin{aligned} \left(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1\right)^{-1} \mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y} &= \left(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1\right)^{-1} \mathbf{X}_1^\top \mathbf{M}_2 \left(\mathbf{X}_1 \hat{\beta}_1 + \mathbf{X}_2 \hat{\beta}_2 + \hat{\mathbf{U}}\right) \\ &= \left(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1\right)^{-1} \mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1 \hat{\beta}_1 \\ &\quad + \left(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1\right)^{-1} \mathbf{X}_1^\top \hat{\mathbf{U}} \quad (\text{car } \mathbf{M}_2 \mathbf{X}_2 = 0 \text{ et } \mathbf{M}_2 \hat{\mathbf{U}} = \hat{\mathbf{U}}) \\ &= \hat{\beta}_1 \end{aligned}$$

- Étant donné que \mathbf{M}_2 est symétrique et idempotente, on peut écrire,

$$\begin{aligned} \hat{\beta}_1 &= \left((\mathbf{M}_2 \mathbf{X}_1)^\top (\mathbf{M}_2 \mathbf{X}_1)\right)^{-1} (\mathbf{M}_2 \mathbf{X}_1)^\top (\mathbf{M}_2 \mathbf{Y}) \\ &= \left(\tilde{\mathbf{X}}_1^\top \tilde{\mathbf{X}}_1\right)^{-1} \tilde{\mathbf{X}}_1 \tilde{\mathbf{Y}} \end{aligned}$$

Régression partitionnée

où,

$$\begin{aligned}\tilde{\mathbf{X}}_1 &= \mathbf{M}_2 \mathbf{X}_1 \\ &= \mathbf{X}_1 - \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{X}_1\end{aligned}$$

à savoir les résidus de la régression de \mathbf{X}_1 sur \mathbf{X}_2 . Et où,

$$\begin{aligned}\tilde{\mathbf{Y}} &= \mathbf{M}_2 \mathbf{Y} \\ &= \mathbf{Y} - \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{Y}\end{aligned}$$

à savoir les résidus de la régression de \mathbf{Y} sur \mathbf{X}_2 .

- Ainsi, pour obtenir les coefficients de K_1 premiers régresseurs, plutôt que de réaliser la régression avec les $K_1 + K_2 = K$ régresseurs,
 - on peut régresser \mathbf{Y} sur \mathbf{X}_2 pour obtenir les résidus $\tilde{\mathbf{Y}}$,
 - régresser \mathbf{X}_1 sur \mathbf{X}_2 pour obtenir les résidus $\tilde{\mathbf{X}}_1$,
 - et alors régresser $\tilde{\mathbf{Y}}$ sur $\tilde{\mathbf{X}}_1$ pour obtenir $\hat{\beta}_1$.
- Autrement dit, $\hat{\beta}_1$ décrit l'effet de \mathbf{X}_1 une fois que ceux de \mathbf{X}_2 ont été contrôlés.

Régression partitionnée

- De manière similaire que pour $\hat{\beta}_1$, nous avons pour $\hat{\beta}_2$,

$$\hat{\beta}_2 = (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{Y}$$

où,

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$$

- Prenons comme exemple le modèle suivant,

$$Y_i = \beta_1 + \beta_2 X_i + U_i, \quad i = 1, 2, \dots, n$$

Soit $\mathbf{1}_n$ le vecteur unitaire ($n \times 1$), i.e.,

$$\mathbf{1}_n = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix}$$

Régression partitionnée

- La matrice des régresseurs est alors,

$$(\mathbf{1}_n \ X) = \begin{pmatrix} 1 & X_1 \\ 1 & X_2 \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ 1 & X_n \end{pmatrix}$$

- Considérons,

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{1}_n(\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top$$

et,

$$\hat{\beta}_2 = \frac{X^\top \mathbf{M}_1 \mathbf{Y}}{X^\top \mathbf{M}_1 X}$$

Régression partitionnée

- Nous avons, $\mathbf{1}_n^\top \mathbf{1}_n = n$, par conséquent,

$$\mathbf{M}_1 = \mathbf{I}_n - \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n}$$

et,

$$\begin{aligned}\mathbf{M}_1 \mathbf{X} &= \mathbf{X} - \mathbf{1}_n \frac{\mathbf{1}_n^\top \mathbf{X}}{n} \\ &= \mathbf{X} - \bar{X} \mathbf{1}_n \\ &= \begin{pmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \end{pmatrix}\end{aligned}$$

Régression partitionnée

où,

$$\begin{aligned}\bar{X} &= \frac{\mathbf{1}_n^\top X}{n} \\ &= n^{-1} \sum_{i=1}^N X_i\end{aligned}$$

- Ainsi la matrice \mathbf{M}_1 transforme le vecteur X en un vecteur dont les éléments sont les écarts des observations X_i à leur moyenne.

Régression partitionnée

- Et nous pouvons écrire,

$$\begin{aligned}\hat{\beta}_2 &= \frac{\sum_{i=1}^N (X_i - \bar{X}) Y_i}{\sum_{i=1}^N (X_i - \bar{X})^2} \\ &= \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^N (X_i - \bar{X})^2}\end{aligned}$$

Qualité de l'ajustement et coefficient de détermination ou R^2

- Écrivons,

$$\begin{aligned} \mathbf{Y} &= \mathbf{P}_X \mathbf{Y} + \mathbf{M}_X \mathbf{Y} \\ &= \hat{\mathbf{Y}} + \hat{\mathbf{U}} \end{aligned}$$

où par construction,

$$\begin{aligned} \hat{\mathbf{Y}}^\top \hat{\mathbf{U}} &= (\mathbf{P}_X \mathbf{Y})^\top (\mathbf{M}_X \mathbf{Y}) \\ &= \mathbf{Y}^\top \mathbf{P}_X \mathbf{M}_X \mathbf{Y} \\ &= 0 \end{aligned}$$

- Supposons que le modèle contienne une constante, par exemple la première colonne de la matrice des régresseurs \mathbf{X} est le vecteur unitaire $\mathbf{1}_n$.

Qualité de l'ajustement et coefficient de détermination ou R^2

- La *variation totale* dans \mathbf{Y} est,

$$\begin{aligned}\sum_{i=1}^N (Y_i - \bar{Y})^2 &= \mathbf{Y}^\top \mathbf{M}_1 \mathbf{Y} \\ &= (\hat{\mathbf{Y}} + \hat{\mathbf{U}})^\top \mathbf{M}_1 (\hat{\mathbf{Y}} + \hat{\mathbf{U}}) \\ &= \hat{\mathbf{Y}}^\top \mathbf{M}_1 \hat{\mathbf{Y}} + \hat{\mathbf{U}}^\top \mathbf{M}_1 \hat{\mathbf{U}} + 2\hat{\mathbf{Y}}^\top \mathbf{M}_1 \hat{\mathbf{U}}\end{aligned}$$

où $\bar{Y} = n^{-1} \sum_{i=1}^N Y_i$.

- Comme le modèle contient une constante,

$$\mathbf{1}_n^\top \hat{\mathbf{U}} = 0$$

et,

$$\mathbf{M}_1 \hat{\mathbf{U}} = \hat{\mathbf{U}}$$

Qualité de l'ajustement et coefficient de détermination ou R^2

- Cependant, $\hat{\mathbf{Y}}^\top \hat{\mathbf{U}} = 0$, et par conséquent,

$$\mathbf{Y}^\top \mathbf{M}_1 \mathbf{Y} = \hat{\mathbf{Y}}^\top \mathbf{M}_1 \hat{\mathbf{Y}} + \hat{\mathbf{U}}^\top \hat{\mathbf{U}}$$

ou,

$$\sum_{i=1}^N (Y_i - \bar{Y})^2 = \sum_{i=1}^N (\hat{Y}_i - \bar{\hat{Y}})^2 + \sum_{i=1}^N \hat{U}_i^2$$

$$\text{où } \bar{\hat{Y}} = n^{-1} \sum_{i=1}^N \hat{Y}_i.$$

Qualité de l'ajustement et coefficient de détermination ou R^2

- Notons que,

$$\begin{aligned}\bar{Y} &= \frac{\mathbf{1}_n^\top \mathbf{Y}}{n} \\ &= \frac{\mathbf{1}_n^\top \hat{\mathbf{Y}}}{n} + \frac{\mathbf{1}_n^\top \hat{\mathbf{U}}}{n} \\ &= \frac{\mathbf{1}_n^\top \hat{\mathbf{Y}}}{n} \\ &= \overline{\hat{Y}}\end{aligned}$$

- Ainsi, la moyenne des Y_i et celle de leurs valeurs ajustées \hat{Y}_i étant égales, nous pouvons écrire,

$$\sum_{i=1}^N (Y_i - \bar{Y})^2 = \sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^N \hat{U}_i^2$$

Qualité de l'ajustement et coefficient de détermination ou R^2

ou,

$$SCT = SCE + SCR$$

où, $SCT := \sum_{i=1}^N (Y_i - \bar{Y})^2$ est la *somme des carrés totale*,

$SCE := \sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2$ est la *somme des carrés expliqués*, et $SCR := \sum_{i=1}^N \hat{U}_i^2$ est la *somme des carrés des résidus*.

Qualité de l'ajustement et coefficient de détermination ou R^2

- Le rapport de la SCE à la SCT est appelé coefficient de détermination (On l'appelle/prononce généralement "R deux") ou R^2 ,

$$\begin{aligned} R^2 &= \frac{SCE}{SCT} \\ &= \frac{\sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^N (Y_i - \bar{Y})^2} \\ &= 1 - \frac{\sum_{i=1}^N \hat{U}_i^2}{\sum_{i=1}^N (Y_i - \bar{Y})^2} \\ &= 1 - \frac{\hat{\mathbf{U}}^\top \hat{\mathbf{U}}}{\mathbf{Y}^\top \mathbf{M}_1 \mathbf{Y}} \end{aligned}$$

Propriétés du R^2

1. Le R^2 est borné entre 0 et 1 ainsi que cela est indiqué par sa décomposition. Remarquez néanmoins que ceci n'est plus vrai dans un modèle sans constante, et dans ce cas il est indiqué de ne pas utiliser la définition précédente du R^2 . Remarquez aussi que si $R^2 = 1$ alors $\hat{\mathbf{U}}^\top \hat{\mathbf{U}} = 0$, ce qui sera vrai seulement si $\mathbf{Y} \in \mathcal{S}(\mathbf{X})$, i.e., \mathbf{Y} est *exactement* une combinaison linéaire des colonnes de \mathbf{X} .
2. Le R^2 augmente avec le nombre de régresseurs.
3. Le R^2 indique la part de la variation de \mathbf{Y} dans l'échantillon qui est expliquée par \mathbf{X} . Cependant notre objectif n'est pas d'expliquer des variations dans l'échantillon mais celle de la population (dont est tiré l'échantillon). Il en résulte qu'un R^2 élevé n'est pas nécessairement un indicateur d'un bon modèle de régression et un R^2 faible n'est pas non plus un argument en défaveur du modèle considéré.
4. Il est toujours possible de trouver une matrice de régresseurs \mathbf{X} pour laquelle $R^2 = 1$, il suffit de prendre n vecteurs linéairement indépendants. En effet, un tel ensemble de vecteurs génère tout l'espace \mathbb{R}^n de sorte que tout vecteur $\mathbf{Y} \in \mathbb{R}^n$ peut s'écrire comme une combinaison linéaire exacte des colonnes de \mathbf{X} .

R^2 ajusté

- Étant donné que le R^2 augmente avec le nombre de régresseurs, une mesure alternative pour juger de la qualité de la régression est le R^2 *ajusté*,

$$\begin{aligned}\bar{R}^2 &= 1 - \frac{n-1}{n-K}(1 - R^2) \\ &= 1 - \frac{\hat{\mathbf{U}}^\top \hat{\mathbf{U}} / (n-K)}{\mathbf{Y}^\top \mathbf{M}_1 \mathbf{Y} / (n-1)}\end{aligned}$$

- Le R^2 ajusté diminue la qualité de ajustement lorsque le nombre de régresseurs augmente relativement au nombre d'observations de sorte que \bar{R}^2 peut diminuer avec le nombre de régresseurs.
- Cependant il n'y a pas vraiment d'argument fort pour utiliser une telle mesure de l'ajustement.

Section 3

Intervalles de confiance

Introduction

- On considère le modèle de régression normal défini par les hypothèses H1-H5.
- L'estimateur ponctuel $\hat{\beta}$, n'est pas très informatif dans la mesure où $\mathbb{P}(\hat{\beta} = \beta) = 0$.
- C'est pour cela qu'on s'intéresse ici à des intervalles(régions) aléatoires qui présentent la propriété d'inclure la vraie valeur du paramètre avec une certaine probabilité spécifiée $(1 - \alpha)$, où α est un nombre "petit" appelé *niveau de confiance*(e.g., 0.01, 0.05, 0.10).
- Un intervalle de confiance avec une probabilité $(1 - \alpha)$ de couvrir β est noté $CI_{1-\alpha}$.

Cas scalaire

On cherche à construire un intervalle de confiance pour le paramètre β_1 dans la régression partitionnée,

$$\mathbf{Y} = \beta_1 \mathbf{X}_1 + \mathbf{X}_2 \beta_2 + \mathbf{U}$$

où \mathbf{X}_1 est un vecteur $(n \times 1)$ contenant les valeurs observées du premier régresseur.

L'estimateur des moindres carrés de β_1 est,

$$\hat{\beta}_1 = \frac{\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{Y}}{\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1}$$

où $\mathbf{M}_2 = \mathbf{I}_n - \mathbf{X}_2(\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2$.

Une méthode pour construire un intervalle de confiance consiste à considérer des intervalles symétriques autour de l'estimateur ponctuel,

$$\text{CI}_{1-\alpha} = \left[\hat{\beta}_1 - c, \hat{\beta}_1 + c \right] \quad (19)$$

Cas scalaire

Comme $\hat{\beta}_1$ est une fonction de l'échantillon aléatoire, l'intervalle de confiance donné dans (19) l'est aussi.

Le problème maintenant est de choisir c tel que,

$$\mathbb{P}(\beta_1 \in \text{CI}_{1-\alpha} | \mathbf{X}) = 1 - \alpha$$

où $\mathbf{X} = (\mathbf{X}_1 \ \mathbf{X}_2)$.

Pour choisir c , nous avons besoin de connaître la distribution de $\hat{\beta}_1 | \mathbf{X}$.

Cas scalaire

Sous les hypothèses H1-H5,

$$\hat{\beta}_1 | \mathbf{X} \sim \mathcal{N} \left(\beta_1, \sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1) \right)$$

et par conséquent,

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2 (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} | \mathbf{X} \sim \mathcal{N}(0, 1) \quad (20)$$

Pour montrer ce résultat, notons que $\hat{\beta}_1$ est un estimateur linéaire, et écrivons $\hat{\beta}_1 = \beta_1 + (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}) / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)$.

Soit z_τ le quantile τ de la distribution normale standard ; autrement dit, si $Z \sim \mathcal{N}(0, 1)$,

$$\mathbb{P}(Z \leq z_\tau) = \tau$$

Notons qu'étant donné que la distribution normale standard est symétrique autour de zéro, nous avons,

$$Z_\alpha = -Z_{(1-\alpha)}$$

Cas scalaire

et par conséquent,

$$\mathbb{P}(-z_{1-\alpha/2} \leq Z \leq z_{1-\alpha/2}) = 1 - \alpha$$

Par exemple, pour $\alpha = 0.05$, $z_{1-0.05/2} = z_{0.975} = 1.96$, et $z_{0.025} = -1.96$.

σ^2 est connu

- Supposons pour le moment que σ^2 soit connu et que ce faisant nous puissions calculer la variance de $\hat{\beta}_1$ (et non pas un estimateur).
- Posons,

$$c = z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\beta}_1 | \mathbf{X})} = z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}$$

- Montrons maintenant que,

$$\mathbb{P} \left(\beta_1 \in \left[\hat{\beta}_1 - z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}, \hat{\beta}_1 + z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \right] \mid \mathbf{X} \right) = 1 - \alpha$$

σ^2 est connu

■ En effet,

$$\begin{aligned}
 & \mathbb{P} \left(\hat{\beta}_1 - z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \leq \beta_1 \leq \hat{\beta}_1 + z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \mid \mathbf{X} \right) \\
 &= \mathbb{P} \left(-z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \leq \beta_1 - \hat{\beta}_1 \leq z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \mid \mathbf{X} \right) \\
 &= \mathbb{P} \left(-z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \leq \hat{\beta}_1 - \beta_1 \leq z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \mid \mathbf{X} \right) \\
 &= \mathbb{P} \left(-z_{1-\alpha/2} \leq \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \leq z_{1-\alpha/2} \mid \mathbf{X} \right)
 \end{aligned} \tag{21}$$

■ Le résultat découle de (20), (21), et de la définition de $z_{1-\alpha/2}$.

σ^2 est inconnu

- On peut ici suivre une approche similaire à la précédente mais en remplaçant dans un premier temps σ^2 par son estimateur,

$$s^2 = \hat{\mathbf{U}}^\top \hat{\mathbf{U}} / (n - K)$$

- Cependant $(\hat{\beta}_1 - \beta_1) / \sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}$ n'est pas normalement distribué car c'est une fonction non-linéaire des termes aléatoires $\hat{\beta}_1$ et s^2 .
- Il s'en suit que nous ne pouvons pas utiliser les quantiles de la distribution normale pour la construction des intervalles de confiance.
- En fait, il s'avère que,

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} | \mathbf{X} \sim t_{n-K} \quad (22)$$

- Rappelons que la distribution t_{n-K} est définie comme suit,

$$Z / \sqrt{V / (n - K)}$$

où $Z \sim \mathcal{N}(0, 1)$, $V \sim \chi_{n-K}^2$, et Z et V sont indépendantes.

σ^2 est inconnu

- Écrivons,

$$\begin{aligned}\frac{\hat{\beta}_1 - \beta_1}{\sqrt{s^2/(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} &= \left(\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2/(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \right) / \frac{s^2}{\sigma^2} \\ &= \left(\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2/(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \right) / \sqrt{\frac{\hat{\mathbf{U}}^\top \hat{\mathbf{U}}}{\sigma^2} / (n - K)}\end{aligned}\quad (23)$$

- Nous savons déjà que, dans l'expression précédente,
 $(\hat{\beta}_1 - \beta_1) / \sqrt{\sigma^2/(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} | \mathbf{X} \sim \mathcal{N}(0, 1)$.
- Nous allons montrer maintenant que conditionnellement à \mathbf{X} ,

$$\frac{\hat{\mathbf{U}}^\top \hat{\mathbf{U}}}{\sigma^2} | \mathbf{X} \sim \chi_{n-K}^2 \quad (24)$$

σ^2 est inconnu

- Pour cela nous avons besoin du résultat suivant,

Lemme 3.1

- Supposons que le vecteur $(n \times 1)$ $U \sim \mathcal{N}(0, \mathbf{I}_n)$. Soit A une matrice $(n \times n)$ symétrique et idempotente avec $\text{Rang}(A) = r \leq n$.
- Alors $U^\top A U \sim \chi_r^2$.

Démonstration.

(voir notes de cours)



- A présent pour montrer (24), écrivons,

$$\frac{\hat{\mathbf{U}}^\top \hat{\mathbf{U}}}{\sigma^2} = \left(\frac{\mathbf{U}}{\sigma^2} \right)^\top \mathbf{M}_{\mathbf{X}} \left(\frac{\mathbf{U}}{\sigma^2} \right) \quad (25)$$

où,

$$\mathbf{M}_{\mathbf{X}} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

σ^2 est inconnu

- Par l'hypothèse H5,

$$\frac{\mathbf{U}}{\sigma} | \mathbf{X} \sim \mathcal{N}(0, \mathbf{I}_n) \quad (26)$$

- Étant donné que \mathbf{M}_X est symétrique et idempotente, ses valeurs propres sont soit zéro ou un.
- Par conséquent,

$$\begin{aligned} \text{Rang}(\mathbf{M}_X) &= \text{Tr}(\mathbf{M}_X) \\ &= n - K \end{aligned} \quad (27)$$

- Le résultat dans (24) découle de (25), (26), (27), et du lemme 3.1.
- Finalement, montrons que $\hat{\beta}_1 - \beta_1$ et $\hat{\mathbf{U}}^\top \hat{\mathbf{U}}$ dans (23) sont indépendants conditionnellement à \mathbf{X} .
- Écrivons,

$$\begin{aligned} \hat{\beta}_1 - \beta_1 &= (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}) / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1) \\ \hat{\mathbf{U}}^\top \hat{\mathbf{U}} &= \mathbf{U}^\top \mathbf{M}_X \mathbf{U} \end{aligned}$$

σ^2 est inconnu

- Il suffit de montrer l'indépendance de $\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}$ et $\mathbf{M}_X \mathbf{U}$.
- Comme $\hat{\beta}_1$ est une fonction de $\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}$, et $\hat{\mathbf{U}}^\top \hat{\mathbf{U}}$ est une fonction de $\mathbf{M}_X \mathbf{U}$, l'indépendance de $\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}$ et $\mathbf{M}_X \mathbf{U}$ implique l'indépendance de $\hat{\beta}_1$ et $\hat{\mathbf{U}}^\top \hat{\mathbf{U}}$.
- Premièrement, montrons que les termes ne sont pas corrélés,

$$\begin{aligned}\text{Cov}(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}, \mathbf{M}_X \mathbf{U} | \mathbf{X}) &= \mathbb{E}(\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U} \mathbf{U}^\top \mathbf{M}_X | \mathbf{X}) \\ &= \mathbf{X}_1^\top \mathbf{M}_2 \mathbb{E}(\mathbf{U} \mathbf{U}^\top | \mathbf{X}) \mathbf{M}_X \\ &= \mathbf{X}_1^\top \mathbf{M}_2 (\sigma^2 \mathbf{I}_n) \mathbf{M}_X \\ &= \sigma^2 \mathbf{X}_1^\top \mathbf{M}_2 \mathbf{M}_X \\ &= \sigma^2 \mathbf{X}_1^\top \mathbf{M}_X \text{ (voir section précédente)} \\ &= 0\end{aligned}$$

- Dans la mesure où $\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}$ et $\mathbf{M}_X \mathbf{U}$ sont des fonctions linéaires de \mathbf{U} , elles sont normalement distribuées conditionnellement à \mathbf{X} .

σ^2 est inconnu

- Étant donné qu'elles ne sont pas corrélées, la normalité implique qu'elles sont indépendantes.
- En conséquence, $\hat{\beta}_1 - \beta_1$, fonction de $\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{U}$, et $\hat{\mathbf{U}}^\top \hat{\mathbf{U}}$ sont aussi indépendants.
- Nous avons montré (22).
- En conséquence, en construisant des intervalles de confiance, si l'on remplace l'inconnue σ^2 par s^2 , on doit remplacer $z_{1-\alpha/2}$ par les quantiles de la t distribution, $t_{n-K, 1-\alpha/2}$,

$$CI_{1-\alpha} = \left[\hat{\beta}_1 - t_{n-K, 1-\alpha/2} \sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}, \hat{\beta}_1 + t_{n-K, 1-\alpha/2} \sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \right]$$

- L'expression $s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)$ qui apparaît dans l'équation ci-dessus est la variance estimée de $\hat{\beta}_1$,

$$\widehat{\text{Var}}(\hat{\beta}_1 | \mathbf{X}) = s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)$$

σ^2 est inconnu

- Ainsi, on construit un intervalle de confiance de niveau α pour β_k , $k = 1, 2, \dots, K$, comme suit,

$$\text{CI}_{1-\alpha} = \left[\hat{\beta}_1 - t_{n-K, 1-\alpha/2} \sqrt{\widehat{\text{Var}}(\hat{\beta}_1|\mathbf{X})}, \hat{\beta}_1 + t_{n-K, 1-\alpha/2} \sqrt{\widehat{\text{Var}}(\hat{\beta}_1|\mathbf{X})} \right] \quad (28)$$

Cas vectoriel

- Supposons que l'on s'intéresse au vecteur des paramètres $\beta = (\beta_1, \beta_2, \dots, \beta_K)^\top$.
- L'équation (28) décrit comment construire des intervalles de confiance "individuels" pour les éléments de β .
- Ces intervalles concernent les distributions marginales des éléments de β , et leur simple combinaison ne produit pas un ensemble qui inclue tout le vecteur β avec une probabilité souhaitée.
- Dans cette partie, nous considérons la construction de régions aléatoires qui incluent β avec une certaine probabilité pré-spécifiée $1 - \alpha$.
- Nous conservons la notation $CI_{1-\alpha}$, malgré le fait que $CI_{1-\alpha}$ est maintenant un sous-ensemble de \mathbb{R}^K .
- Ce qui suit est une approche simple et conventionnelle pour construire des régions de confiance.
- Nous cherchons une région de confiance $CI_{1-\alpha} = \{b \in \mathbb{R}^K\}$ tel que $\mathbb{P}(\beta \in CI_{1-\alpha} | \mathbf{X}) = 1 - \alpha$.

Cas vectoriel

- Considérons une forme quadratique par rapport à $(\hat{\beta} - \beta)$,

$$\begin{aligned} (\hat{\beta} - \beta)^\top \left(\widehat{\text{Var}}(\hat{\beta}|\mathbf{X}) \right)^{-1} (\hat{\beta} - \beta) / K &= (\hat{\beta} - \beta)^\top \left(s^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right)^{-1} (\hat{\beta} - \beta) / K \\ &= \frac{(\hat{\beta} - \beta)^\top (\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})^{-1} (\hat{\beta} - \beta) / K}{s^2 / \sigma^2} \\ &= \frac{(\hat{\beta} - \beta)^\top (\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})^{-1} (\hat{\beta} - \beta) / K}{\left(\frac{\hat{\mathbf{U}}^\top \hat{\mathbf{U}}}{\sigma^2} \right) / (n - K)} \end{aligned} \quad (29)$$

- Montrons maintenant que l'expression dans (29) possède une distribution $F_{K, n-K}$ conditionnellement à \mathbf{X} .
- La distribution $F_{K, n-K}$ est définie comme la distribution de,

$$\frac{V/K}{W/(n-K)}$$

où $V \sim \chi_K^2$, et $W \sim \chi_{n-K}^2$ sont indépendantes.

Cas vectoriel

- De la discussion dans la partie précédente nous savons que, $\hat{\mathbf{U}}^\top \hat{\mathbf{U}} / \sigma^2 | \mathbf{X} \sim \chi^2_{n-K}$ qui est indépendant du numérateur dans (29).
- Il résulte de cela, que nous devons montrer que

$$(\hat{\beta} - \beta)^\top \left(\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right)^{-1} (\hat{\beta} - \beta) | \mathbf{X} \sim \chi^2_K \quad (30)$$

- Pour cela nous avons besoin du lemme suivant,

Lemme 3.2

Supposons que le vecteur $(K \times 1)$, $U \sim \mathcal{N}(0, \Sigma)$, où Σ est une matrice définie positive de variances-covariances. Alors, $U^\top \Sigma^{-1} U \sim \chi^2_K$.

Démonstration.

(voir notes de cours)



- Le résultat dans (30) découle du lemme 3.2.
- En conséquence,

$$\frac{(\hat{\beta} - \beta)^\top (s^2 (\mathbf{X}^\top \mathbf{X})^{-1})^{-1} (\hat{\beta} - \beta)}{K} | \mathbf{X} \sim F_{K, n-K}$$

Cas vectoriel

- Soit, $F_{K,n-K,\tau}$ le quantile τ de la distribution F . La région de confiance de niveau α se construit comme suit,

$$CI_{1-\alpha} = \left\{ b \in \mathbb{R}^K : (\hat{\beta} - b)^\top \left(s^2(\mathbf{X}^\top \mathbf{X})^{-1} \right)^{-1} (\hat{\beta} - b) / K \leq F_{K,n-K,1-\alpha} \right\}$$

- La discussion précédente implique que,

$$\begin{aligned} \mathbb{P}(\beta \in CI_{1-\alpha} | \mathbf{X}) &= \mathbb{P} \left((\hat{\beta} - \beta)^\top \left(s^2(\mathbf{X}^\top \mathbf{X})^{-1} \right)^{-1} (\hat{\beta} - \beta) / K \leq F_{K,n-K,1-\alpha} | \mathbf{X} \right) \\ &= 1 - \alpha \end{aligned}$$

Remarque 1

- *La région/intervalle de confiance $CI_{1-\alpha}$ est une fonction de l'échantillon $\{(Y_i, X_i)\}_{i=1}^n$, et il est ce faisant aléatoire, ce qui nous permet de parler de la probabilité que $CI_{1-\alpha}$ contienne la vraie valeur de β .*
- *D'un autre côté, la réalisation de $CI_{1-\alpha}$ n'est pas aléatoire. Une fois que l'intervalle de confiance est calculé pour des observations données, il n'y a plus de sens à parler de la probabilité qu'il inclue β . C'est soit zéro, soit un.*

Section 4

Tests d'hypothèses

Test d'une hypothèse par rapport à un seul coefficient

Nous poursuivons notre discussion sur le modèle de régression linéaire normal, i.e., le modèle défini par les hypothèses [H1-H5](#). On considère,

$$\mathbf{Y} = \beta_1 \mathbf{X}_1 + \mathbf{X}_2 \beta_2 + \mathbf{U}$$

où \mathbf{X}_1 est le vecteur $(n \times 1)$ d'observations du premier régresseur.

Supposons que la variance des erreurs σ^2 soit connue.

Soit $\hat{\beta}_1$ l'estimateur des moindres carrés de β_1 .

Cherchons à tester,

$$\begin{aligned} H_0 &: \beta_1 = \beta_{1,0} \\ H_1 &: \beta_1 \neq \beta_{1,0} \end{aligned} \tag{31}$$

Une règle de décision pour un test de niveau α peut reposer sur l'intervalle de confiance $CI_{1-\alpha}$.

Test d'une hypothèse par rapport à un seul coefficient

Nous avons,

$$CI_{1-\alpha} = \left[\hat{\beta}_1 - z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}, \hat{\beta}_1 + z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} \right]$$

Considérons le test suivant,

Rejeter H_0 si $\beta_{1,0} \notin CI_{1-\alpha}$

Dans ce cas la région critique est donnée par le complément de $CI_{1-\alpha}$.

On rejette ainsi H_0 si,

$$\beta_{1,0} \leq \hat{\beta}_1 - z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}$$

ou

$$\beta_{1,0} \geq \hat{\beta}_1 + z_{1-\alpha/2} \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}$$

Test d'une hypothèse par rapport à un seul coefficient

De manière équivalente, on rejette si,

$$\left| \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \right| > z_{1-\alpha/2} \quad (32)$$

Un tel test est appelé *bilatéral*, car sous l'hypothèse alternative, la vraie valeur β_1 peut être plus petite ou plus grande que $\beta_{1,0}$.

L'expression à gauche de l'inégalité est une statistique de test.

Pour calculer la probabilité de rejet de l'hypothèse nulle supposons que la vraie valeur soit donnée par β_1 . Écrivons,

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} + \frac{\beta_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \quad (33)$$

Nous avons que,

Test d'une hypothèse par rapport à un seul coefficient

$$\frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} | \mathbf{X} \sim \mathcal{N} \left(\frac{\beta_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}}, 1 \right)$$

Si l'hypothèse nulle est vraie alors $\beta_1 - \beta_{1,0} = 0$ et la statistique de test présente une distribution normale standard. Dans ce cas par définition de $z_{1-\alpha/2}$.

$$\begin{aligned} \mathbb{P}(\text{rejeter } H_0 | \mathbf{X}, H_0 \text{ est vraie}) &= \mathbb{P} \left(\left| \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \right| > z_{1-\alpha/2} \right) \\ &= \alpha \end{aligned}$$

Ainsi, le test suggéré a la taille correcte α .

Si l'hypothèse nulle est fausse, la distribution de la statistique de test n'est pas centrée autour de zéro, et l'on verra des taux de rejet supérieurs à α .

Test d'une hypothèse par rapport à un seul coefficient

La probabilité de rejet est une fonction de la vraie valeur β_1 et dépend de la magnitude du deuxième terme dans (33), $|\beta_1 - \beta_{1,0}| / \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}$.

Supposons par exemple que,

$$\begin{aligned}\beta_{1,0} &= 0 \\ \sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)} &= 1 \\ \alpha &= 0.05 \text{ (et } z_{1-\alpha/2} = 1.96)\end{aligned}$$

Soit $Z \sim \mathcal{N}(0, 1)$. Dans ce cas, la *fonction puissance* du test est,

$$\begin{aligned}\pi(\beta_1) &= \mathbb{P} \left(\left| \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \right| > z_{1-\alpha/2} \right) \\ &= \mathbb{P} \left(\left| \frac{\hat{\beta}_1 - \beta_1 + \beta_1 - \beta_{1,0}}{\sqrt{\sigma^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \right| > 1.96 \mid \mathbf{X} \right) \\ &= \mathbb{P} (|Z + \beta_1| > 1.96) \\ &= \mathbb{P} (Z < -1.96 - \beta_1) + \mathbb{P} (Z > 1.96 - \beta_1)\end{aligned}$$

Test d'une hypothèse par rapport à un seul coefficient

Par exemple,

$$\pi(\beta_1) = \begin{cases} 0.52 & \text{pour } \beta_1 = -2 \\ 0.17 & \text{pour } \beta_1 = -1 \\ 0.05 & \text{pour } \beta_1 = 0 \\ 0.17 & \text{pour } \beta_1 = 1 \\ 0.52 & \text{pour } \beta_1 = 2 \end{cases}$$

Dans ce cas la fonction puissance est minimisée en $\beta_1 = \beta_{1,0}$ où $\pi(\beta_1) = \alpha$.

Pour le calcul des *p-values* considérons l'exemple suivant.

Supposons, qu'étant donné des données la statistique de test dans (32) soit égale à 1.88.

Pour la distribution normale standard $\mathbb{P}(Z > 1.88) = 0.03$.

Par conséquent la *p-value* du test est 0.06. On rejeterait l'hypothèse nulle pour tous les tests avec un niveau de significativité supérieur à 0.06.

Test d'une hypothèse par rapport à un seul coefficient

Dans le cas où σ_2 est inconnu, on peut tester (31) en considérant la t -statistique,

$$\begin{aligned} T &= \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} \\ &= \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_1 | \mathbf{X})}} \end{aligned} \quad (34)$$

Le test est donné par la règle de décision suivante,

$$\text{Rejeter } H_0 \text{ si } |T| > t_{n-K, 1-\alpha/2}$$

Dans ce cas (voir section précédente) sous H_0 ,
 $\mathbb{P}(|T| > t_{n-K, 1-\alpha/2} | \mathbf{X}, H_0 \text{ est vraie}) = \alpha$.

On peut aussi considérer des tests *unilatéraux*.

Test d'une hypothèse par rapport à un seul coefficient

Dans le cas de ces tests l'hypothèse nulle et l'hypothèse alternative peuvent être spécifiées comme suit,

$$H_0 : \beta_1 \leq \beta_{1,0}$$

$$H_1 : \beta_1 > \beta_{1,0}$$

Notons que dans ce cas, et H_0 et H_1 sont composées, et la probabilité de rejet varie non seulement selon les valeurs de β_1 spécifiées sous H_1 mais aussi selon H_0 .

Dans ce cas un test valide devra satisfaire la condition,

$$\sup_{\beta_1 \leq \beta_{1,0}} \mathbb{P}(\text{rejeter } H_0 | \mathbf{X}, \beta_1) \leq \alpha \quad (35)$$

i.e., la probabilité maximale de rejeter H_0 quand elle est vraie ne doit pas dépasser α . Soit T telle que définie dans (34) et considérons le test suivant (règle de décision) :

$$\text{Rejeter } H_0 \text{ quand } T > t_{n-K, 1-\alpha}$$

Test d'une hypothèse par rapport à un seul coefficient

Sous H_0 , nous avons,

$$\begin{aligned}\mathbb{P}(\text{rejeter } H_0 | \beta_1 \leq \beta_{1,0}) &= \mathbb{P}(T > t_{n-K, 1-\alpha} | \mathbf{X}, \beta_1 \leq \beta_{1,0}) \\ &= \mathbb{P}\left(\frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} > t_{n-K, 1-\alpha} | \mathbf{X}, \beta_1 \leq \beta_{1,0}\right) \\ &\leq \mathbb{P}\left(\frac{\hat{\beta}_1 - \beta_1}{\sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} > t_{n-K, 1-\alpha} | \mathbf{X}, \beta_1 \leq \beta_{1,0}\right) \text{ (car } \beta_1 \leq \beta_{1,0}) \\ &= \alpha \text{ (étant donné que } \frac{\hat{\beta}_1 - \beta_1}{\sqrt{s^2 / (\mathbf{X}_1^\top \mathbf{M}_2 \mathbf{X}_1)}} | \mathbf{X} \sim t_{n-K})\end{aligned}$$

Ainsi, la condition sur la taille (35) est satisfaite. Notons qu'étant donné qu'il s'agit d'un test unilatéral, la probabilité d'erreur de type 1 est portée uniquement par la queue droite de la distribution.

Test d'une contrainte linéaire simple

Considérons le modèle de régression linéaire normal défini par les hypothèse $H1-H5$,

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{U}$$

Supposons que l'on souhaite tester,

$$H_0 : c^\top \beta = r$$

$$H_1 : c^\top \beta \neq r$$

Dans ce cas c est un vecteur ($K \times 1$), r est un scalaire, et sous l'hypothèse nulle,

$$c_1\beta_1 + c_2\beta_2 + \dots + c_K\beta_K - r = 0$$

Par exemple, en posant $c_1 = 1$, $c_2 = -1$, $c_3 = \dots = c_K = 0$, et $r = 0$, nous pouvons tester l'hypothèse que $\beta_1 = \beta_2$.

Test d'une contrainte linéaire simple

Pour l'estimateur des moindres carrés nous avons,

$$\widehat{\beta}|\mathbf{X} \sim \mathcal{N}\left(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\right) \quad (36)$$

Alors,

$$\frac{\mathbf{c}^\top \widehat{\beta} - \mathbf{c}^\top \beta}{\sqrt{\sigma^2 \mathbf{c}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{c}}} | \mathbf{X} \sim \mathcal{N}(0, 1) \quad (37)$$

Par conséquent, sous H_0 ,

$$\frac{\mathbf{c}^\top \widehat{\beta} - r}{\sqrt{\sigma^2 \mathbf{c}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{c}}} | \mathbf{X} \sim \mathcal{N}(0, 1) \quad (38)$$

Test d'une contrainte linéaire simple

Considérons la t statistique,

$$\begin{aligned} T &= \frac{c^\top \hat{\beta} - r}{\sqrt{s^2 c^\top (\mathbf{X}^\top \mathbf{X})^{-1} c}} \\ &= \left(\frac{c^\top \hat{\beta} - r}{\sqrt{\sigma^2 c^\top (\mathbf{X}^\top \mathbf{X})^{-1} c}} \right) / \sqrt{\frac{\mathbf{U}^\top \mathbf{M}_X \mathbf{U}}{\sigma^2} / (n - K)} \end{aligned}$$

Sous H_0 , le résultat dans (38) est vérifié.

En outre, conditionnellement à \mathbf{X} ,

$$\mathbf{U}^\top \mathbf{M}_X \mathbf{U} / \sigma^2 | \mathbf{X} \sim \chi_{n-K}^2 \text{ indépendant de } \hat{\beta} \quad (39)$$

Par conséquent sous H_0 ,

$$T | \mathbf{X} \sim t_{n-K}$$

Test d'une contrainte linéaire simple

Ainsi, le niveau de significativité α du test bilatéral de $H_0 : c^\top \beta = r$ est donné par,

$$\text{Rejeter } H_0 \text{ si } |T| > t_{n-K, 1-\alpha/2}$$

En posant l'élément j de c , $c_j = 1$ et le restant des éléments de c égaux à zéro on obtient le test discuté dans la sous-section précédente,

$$H_0 : \beta_j = r$$

$$H_1 : \beta_j \neq r$$

On rejette H_0 si,

$$|T| = \left| \frac{\hat{\beta}_j - r}{\sqrt{s^2 [(\mathbf{X}^\top \mathbf{X})^{-1}]_{jj}}} \right| > t_{n-K, 1-\alpha/2}$$

où $[(\mathbf{X}^\top \mathbf{X})^{-1}]_{jj}$ est l'élément (j, j) de la matrice $(\mathbf{X}^\top \mathbf{X})^{-1}$.

Tests de contraintes linéaires multiples

Supposons que l'on souhaite tester,

$$H_0 : \mathbf{R}\beta = r$$

$$H_1 : \mathbf{R}\beta \neq r$$

où \mathbf{R} est une matrice ($q \times K$) est r est un vecteur ($r \times 1$). Par exemple,

■ $\mathbf{R} = \mathbf{I}_K$, $r = 0$. Dans ce cas on teste que $\beta_1 = \dots = \beta_K = 0$.

■ $\mathbf{R} = \begin{pmatrix} 1 & 1 & 0 & 0 & \cdot & \cdot & \cdot & 0 \\ 0 & 0 & 1 & 0 & \cdot & \cdot & \cdot & 0 \end{pmatrix}$, $r = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$. Dans ce cas,
 $H_0 : \beta_1 + \beta_2 = 1, \beta_3 = 0$.

Tests de contraintes linéaires multiples

Considérons la F statistique,

$$F = \left(\mathbf{R}\hat{\beta} - r \right)^{\top} \left(s^2 \mathbf{R}(\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{R}^{\top} \right)^{-1} \left(\mathbf{R}\hat{\beta} - r \right) / q$$

On peut alors montrer (voir notes de cours) que sous H_0 ,

$$F | \mathbf{X} \sim F_{q, n-K} \quad (40)$$

Par conséquent, le test est donné par,

$$\begin{aligned} \text{Rejeter } H_0 \text{ si } F &= \left(\mathbf{R}\hat{\beta} - r \right)^{\top} \left(s^2 \mathbf{R}(\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{R}^{\top} \right)^{-1} \left(\mathbf{R}\hat{\beta} - r \right) / q \\ &> F_{q, n-K, 1-\alpha} \end{aligned}$$