

# Bandes de confiance simultanées par la méthode de Scheffé dans le cas de la régression linéaire simple

J.M. Azaïs \*

21 mai 2012

Considérons un échantillon de femmes enceintes qui étaient parfaitement réglées. Par échographie on constitue un échantillon des deux variables suivantes

X l'âge du foetus mesuré en semaine d'aménorrhée Y la taille du fémur du foetus

Y	X	Y	X
178	48.2	152	39.9
185	47.4	168	44
157	37.7	106	17.1
226	58.8	119	21.6
227	58.3	186	46.7
153	38	175	45.4
184	52.6	125	26.2
227	59	139	33.5
140	30	148	35.6
139	31.6	155	38
230	62.4	150	36.8
122	24.1	154	38.4
180	45.2		

données disponibles à

<http://dl.dropbox.com/u/17864747/femur.m>

Il est clair que c'est la variable  $X$  qui a une influence sur  $Y$  et il est donc logique de faire la régression de  $Y$  sur  $X$  (on peut vérifier qu'une régression linéaire simple suffit).

Maintenant supposons qu'une femme qui ne connaît pas l'âge de son foetus subit une échographie et la longueur du foetus est mesurée à 120 mm. Quel est l'âge estimé du foetus et comment construire un intervalle de confiance pour cet âge? Comme nous allons le voir la réponse à cette question demande de

---

\*Université de Toulouse, Statistique et Probabilités, IMT, Mél : jean-marc.azais@math.univ-toulouse.fr

construire un intervalle de confiance pour la droite de régression, un tel intervalle est dit simultané (voir définition plus bas).

## 1 Le modèle

On considère donc le modèle de régression linéaire simple.

$$Y_i = \mu + \beta X_i + \epsilon_i \quad i = 1, \dots, n. \quad (1)$$

avec les hypothèses classiques : les  $\epsilon_i$  sont i.i.d. de loi  $N(0, \sigma^2)$ .

Nous supposons pour simplifier que le nombre d'observations  $n$  est grand ( $n > 100$ ) (ce qui n'est pas le cas dans l'exemple numérique) Sans perte de généralité on peut également supposer que la variable  $X$  est centrée et d'inertie 1 :  $\sum_{i=1}^n X_i = 0$  ;  $\sum_{i=1}^n X_i^2 = 1$ .

On sait bien qu'alors

$$\begin{aligned} \hat{\mu} &= \bar{Y} && \text{est gaussien de variance } \sigma^2/n \\ \hat{\beta} &= \sum_{i=1}^n Y(t)X(t) && \text{est gaussien de variance } \sigma^2, \end{aligned}$$

et ces deux estimateurs sont indépendants.

La prédiction de la réponse au point  $x$  est donnée par

$$\hat{Y}(x) = \hat{\mu} + \hat{\beta}x \quad \text{de variance } \sigma^2\left(\frac{1}{n} + x^2\right).$$

En centrant et réduisant on définit

$$\bar{Z}(x) := \frac{\hat{Y}(x) - (\mu + \beta x)}{\sqrt{\sigma^2(x^2 + 1/n)}},$$

$\sigma^2$  qui est inconnu. Mais comme le nombre d'observations est grand, l'estimateur

$$\hat{\sigma}^2 := 1/(n-2) \sum_{i=1}^n (Y_i - \hat{\mu} - \hat{\beta}X_i)^2.$$

est presque exact. et on peut remplacer sans conséquence  $\sigma^2$  par  $\hat{\sigma}^2$  pour obtenir la quantité

$$Z(x) := \frac{\hat{Y}(x) - (\mu + \beta x)}{\sqrt{\hat{\sigma}^2(x^2 + 1/n)}},$$

qui est maintenant calculable.

Ce qui donne l'intervalle de confiance

$$Y(x) := (\mu + \beta x) \in IC(x) := [\hat{Y}(x) \pm Z_\alpha \sqrt{\hat{\sigma}^2(x^2 + 1/n)}], \quad (2)$$

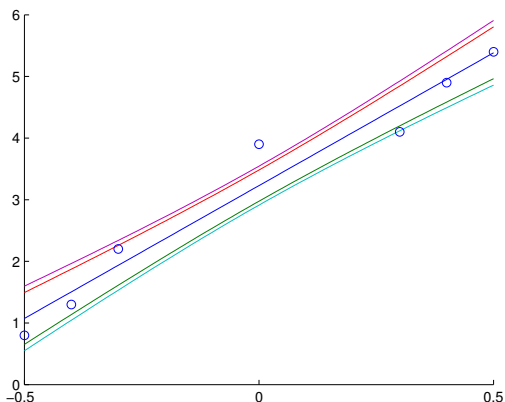


FIGURE 1 – Exemple de données avec la droite de régression, la bande de confiance non-simultanée, la bande de confiance simultanée (à l'extérieur)

où  $Z_\alpha$  est la valeur critique au niveau  $\alpha$  d'une loi normale centrée réduite, c'est à dire le fractile  $1 - \alpha/2$ .

Si on fait varier  $x$  on obtient une bande autour de la droite de régression délimitée par une hyperbole.

Montrons qu'il s'agit bien d'une hyperbole : l'équation ( 2) peut se réécrire

$$y = A + Bx \pm C\sqrt{1 + Dx^2}$$

ce qui se réécrit encore

$$(y - A - Bx)^2 = C^2 + C^2 Dx^2.$$

C'est donc bien une conique et comme les coefficients de  $y$  et  $x$  dans le carré sont opposés, on voit que  $x$  et  $Y$  vont varier de  $-\infty$  à  $+\infty$  c'est donc une hyperbole et non une ellipse ou une parabole.

La bande ainsi construite a la propriété **non simultanée** suivante. Pour  $x$  fixé

$$P\{Y(x) \in IC(x)\} = 1 - \alpha.$$

Mais si  $x$  varie dans un ensemble  $\mathcal{X}$ , la bande ci-dessus n'est pas simultanée dans le sens où en général

$$P\{\forall x \in \mathcal{X} : Y(x) \in IC(x)\} \text{ est en général beaucoup plus petite que } 1 - \alpha.$$

Pour obtenir pour une nouvelle bande  $B(x)$  **simultanée** (plus grande) vérifiant

$$P\{\forall x \in \mathcal{X} : Y(x) \in B(x)\} = 1 - \alpha, \quad (3)$$

on a recours à

## 2 La méthode de Scheffé

On sait que

$$\begin{bmatrix} \frac{\hat{\mu} - \mu}{\sqrt{\hat{\sigma}^2/n}} \\ \frac{\hat{\beta} - \beta}{\hat{\sigma}} \end{bmatrix} \sim N(0, I_2).$$

Donc en définissant, pour bien normaliser,  $\nu := \sqrt{n}\mu$

$$\frac{(\hat{\nu} - \nu)^2}{\hat{\sigma}^2} + \frac{(\hat{\beta} - \beta)^2}{\hat{\sigma}^2} \sim \chi^2(2),$$

approximativement. On obtient donc une région de confiance pour  $(\nu, \beta)$  qui est un disque en écrivant

$$(\hat{\nu} - \nu)^2 + (\hat{\beta} - \beta)^2 \leq \hat{\sigma}^2 \chi_{1-\alpha}^2(2), \quad (4)$$

où  $\chi_{1-\alpha}^2(2)$  est le fractile  $1 - \alpha$  de la loi  $\chi^2(2)$ .

Maintenant écrivons

$$\mu + \beta x = \sqrt{1/n + x^2} \langle (\nu, \beta), \left( \frac{n^{-1/2}}{\sqrt{1/n + x^2}}, \frac{x}{\sqrt{1/n + x^2}} \right) \rangle.$$

Notez que le second vecteur dans le produit scalaire est de norme 1 et que dans (4) le rayon du disque vaut  $\sqrt{\hat{\sigma}^2 \chi_{1-\alpha}^2(2)}$ . Il est facile d'en déduire que

$$|\mu + \beta x - (\hat{\mu} + \hat{\beta}x)| \leq \sqrt{(1/n + x^2)(\hat{\sigma}^2 \chi_{1-\alpha}^2(2))},$$

avec proba  $1 - \alpha$ , uniformément en  $x$ . La région est donc simultanée au sens de (3).

## 3 Suggestions

- Détailler toutes les étapes qui conduisent à la construction de (2).
- Détailler soigneusement la preuve de la méthode de Sheffé.
- Étendre la méthode de Sheffé au cas où  $n$  n'est pas grand.
- Construire les deux régions sur l'exemple et en déduire l'intervalle de confiance pour l'âge quand la longueur du fémur est 120 mm.
- Étendre la méthode à un modèle de régression multiple en utilisant un argument d'orthogonalisation de Gram-Schmidt. Plus précisément si

$$Y_i = \beta_0 X_i^0 + \beta_1 X_i^1 + \dots + X_i^{p-1} + \epsilon_i \quad i = 1, \dots, n, \quad (5)$$

est un modèle de régression où  $X^0$  est le vecteur composé de 1, et  $X^1 \dots X^{p-1}$  sont les autres régresseurs (l'exposant ne signifie pas une puissance!). On peut se ramener au cas où tous ces régresseurs sont orthogonaux et d'inertie 1. Dans ce cas, la méthode se généralise facilement.

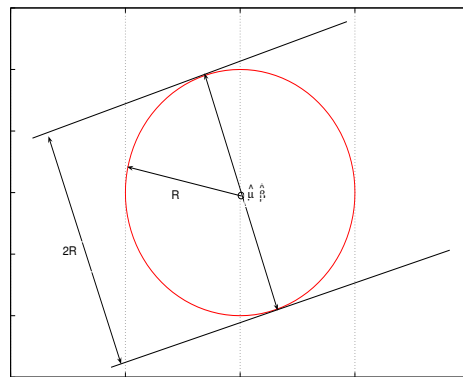


FIGURE 2 – Représentation du disque de confiance autour de la valeur  $\hat{\nu}, \hat{\beta}$  et de sa projection sur une droite quelconque. On constate bien que cette projection a la même taille qu'un diamètre.