

Régression Simple.

(1)

→ On suppose qu'il existe une relation linéaire entre Y et x .

$$\text{Modèle : } Y_i = \mu + \alpha x_i + E_i$$

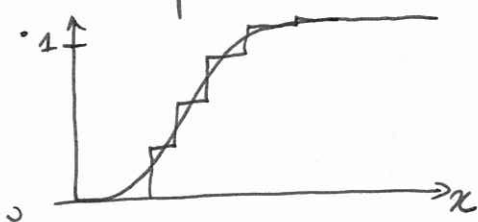
} - E_i iid

} - $V(E_i) = \sigma^2$: variance résiduelle. (à minimiser).

→ 1^{ère} chose: regarder si les hypothèses sont vérifiées
→ Graphique des résidus: pour: homoscedasticité?
tendance?
indép?.

→ Test d'ajustement des résidus à une gaussienne $N(0, \sigma^2)$.

Exemple: Kolmogoroff, Smirnov, fondé sur les fonctions de répartition



$$D_n = \sup_x |F_n(x) - F(x)|$$

= Statistique de test dont la loi ne dépend pas de F sous H_0 .

↳ tendance à rejeter l'hypothèse de normalité.

• Q-Q plot: visualiser.

Conseil: Modèle linéaire robuste à l'hypothèse de normalité.

→ Graphique des résidus

Lecture de la table d'ANOVA.

Source :	df	SS	MS
Model : Variabilité expliquée par x .	1	$\sum (\hat{Y}_i - \bar{Y})^2$	$\frac{SCM}{1}$
Error : Variabilité de Y autour de la droite de reg.	$n-2$	$\sum (Y_i - \hat{Y}_i)^2$	$\frac{SCE}{n-2}$
Total : Variabilité intrinsèque de Y (\perp du modèle)	$n-1$	$\sum (Y_i - \bar{Y})^2$	$\frac{SCT}{n-1}$

• 1^{er} test de Fisher : test du modèle nul contre Complet.

$$H_0 = \left\{ Y_i = \mu + \varepsilon_i \right\} \quad H_1 = \left\{ Y_i = \mu + \alpha x_i + \varepsilon_i \right\}$$

$$\downarrow$$

$$SCR_0, n - k_0 \text{ df}$$

$$\downarrow$$

$$SCR_1, n - k_1 \text{ df}$$

$$F_2 = \frac{(SCR_0 - SCR_1) / (k_1 - k_0)}{SCR_1 / (n - k_1)} \sim \mathcal{F}(k_1 - k_0, n - k_1)$$

↳ revient à faire l'hypothèse sur le paramètre α
 $H_0 = \{ \alpha = 0 \}$.

• Indicateurs : $R^2 = \frac{SCM}{SCT} =$ $\left. \begin{array}{l} \text{interprétation explicative (non} \\ \text{pas prédictif)} \end{array} \right\}$

• $MSE = \hat{\sigma}^2$ On cherche σ petit

Lois des estimateurs et tests sur les paramètres

(2)

$$\text{Lois: } \begin{pmatrix} \hat{\mu} \\ \hat{\alpha} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu \\ \alpha \end{pmatrix}, \begin{pmatrix} V(\hat{\mu}) & \text{cov}(\hat{\mu}, \hat{\alpha}) \\ & V(\hat{\alpha}) \end{pmatrix} \right).$$

Gauss Markov: Estimateurs sans biais de α, μ fonction linéaire des Y_i , ceux de variance minimale

$$V(\hat{\alpha}) = \sigma^2 / \sum (x_i - \bar{x})^2.$$

$$V(\hat{\mu}) = \sigma^2 \left(\frac{1}{n} + \frac{(\bar{x})^2}{\sum (x_i - \bar{x})^2} \right).$$

$$\text{Test: } H_0: \{ \alpha = 0 \} \rightarrow T = \frac{\hat{\alpha}}{S_{\hat{\alpha}}} \quad \text{avec } S_{\hat{\alpha}} = \frac{\hat{\sigma}^2}{\sum (x_i - \bar{x})^2}$$

$$T \sim \mathcal{C}(n-2).$$

ou F_p , si $T \sim \mathcal{C}(p)$ alors $T^2 \sim F(1, p)$

↳ On retrouve les m valeurs de statistiques dans la table d'ANOVA.

Intervalle de Confiance, Intervalle de Prédiction

IC de la droite de régression:

$$IC_{1-\alpha}(\hat{\alpha}) = \left[\hat{\alpha} \pm t_{n-2; 1-\alpha/2} \times S_{\hat{\alpha}} \right]$$

$$IC_{1-\alpha}(\hat{\mu}) = \left[\hat{\mu} \pm t_{n-2; 1-\alpha/2} \times S_{\hat{\mu}} \right].$$

$$\hat{Y}_i \sim \mathcal{N} \left(\mu + \alpha x_i, \sigma \sqrt{\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_j - \bar{x})^2}} \right)$$

$$\Rightarrow IC_{1-\alpha}(\hat{Y}_i) = \left[\hat{\mu} + \hat{\alpha} x_i \pm t_{n-2; 1-\alpha/2} \times \hat{\sigma} \times \sqrt{\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_j - \bar{x})^2}} \right]$$

Intervalle de prediction:

→ On veut predire Y_0 pour un x_0 non observé.

$$Y_0 = \hat{\mu} + \hat{\alpha} x_0 + E_0.$$

$$V(Y_0) = V(\hat{\mu} + \hat{\alpha} x_0) + V(E_0). \rightarrow Y_0 \text{ ne depend pas de } x_0$$

étant donné qu'on n'a pas observé E_0 , on a bien \perp .

$$V(Y_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right) + \underline{\underline{\sigma^2}}.$$

$$IP_{1-\alpha}(Y) = \left[\hat{\mu} + \hat{\alpha} x_0 \pm t_{n-2, 1-\alpha/2} \times \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}} \right].$$