

Exemple: Soit les $n=5$ observations suivantes du couple (X_t, Y_t) :

t	X_t	Y_t	calculs: X_t^2	$X_t Y_t$
1	2	2	4	4
2	3	5	9	15
3	6	3	36	18
4	9	6	81	54
5	12	5	144	60
Totaux:	32	21	274	151

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{t=1}^n X_t Y_t - n \bar{X} \bar{Y}}{\sum_{t=1}^n X_t^2 - n \bar{X}^2} = \frac{151 - (5)(21/5)(32/5)}{274 - (5)(32/5)^2} = \frac{83}{346} \approx 0.2399$$

$$\Rightarrow \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = \frac{21}{5} - \left(\frac{83}{346} \right) \times \left(\frac{32}{5} \right) \approx 2.6647$$

$$\Rightarrow \text{Modèle de régression: } Y_t = 2.6647 + 0.2399 X_t + \varepsilon_t$$

\Rightarrow Prévisions et résidus:

t	$\hat{Y}_t = \hat{\beta}_0 + \hat{\beta}_1 X_t$	$\hat{\varepsilon}_t$
1	3.1445	-1.1445
2	3.3844	1.6156
3	4.1041	-1.1041
4	4.8238	1.1762
5	5.5435	-0.5435

$$\sum_{t=1}^5 \hat{\varepsilon}_t \approx -0.0003 \approx 0!$$

Logiciel R:

```

x ← c(2, 3, 6, 9, 12)
y ← c(2, 5, 3, 6, 5)
reg ← lm(y ~ x)
summary(reg)
fitted(reg)
residuals(reg)

```

... entrer les données

... estimation des paramètres

... résumé de l'estimation

... valeurs de \hat{y}_t

... résidus $\hat{\varepsilon}_t$

Logiciel SAS:

```

data donnees_exemple;
input x y;
cards;
2 2
3 5
6 3
9 6
12 5
run;
proc reg data=donnees_exemple;
model y=x;
run;

```

... entrer les données

... estimation des paramètres

2.2.2) Caractéristiques du terme d'erreur:

Rappelons que :

$$Y_t = \beta_0 + \beta_1 X_t + \varepsilon_t$$

$$\Rightarrow Y_t = E(Y_t | X_t) + \varepsilon_t$$

On a que:

$$1) E(\varepsilon_t) = 0 \quad \dots \text{par définition par que } E(Y_t) = E(Y_t | X_t)$$

$$2) \text{Var}(\varepsilon_t) = \sigma^2 \quad \dots \text{constante par hypothèse}$$

$$3) \text{Cov}(\varepsilon_t, \varepsilon_s) = 0, t \neq s \quad \dots \text{pas de corrélation par hypothèse}$$

2.3) Propriétés de l'estimateur des moindres carrés (EMC):

2.3.1) Propriété 1: Estimateur sans biais:

$$\bullet E(\hat{\beta}_1) = E\left(\frac{\sum_{t=1}^n (X_t - \bar{X})(Y_t - \bar{Y})}{\sum_{t=1}^n (X_t - \bar{X})^2}\right)$$

$$= \frac{\sum_{t=1}^n (X_t - \bar{X}) E(Y_t - \bar{Y})}{\sum_{t=1}^n (X_t - \bar{X})^2}$$

$$= \frac{\sum_{t=1}^n (X_t - \bar{X}) (E(Y_t) - E(\bar{Y}))}{\sum_{t=1}^n (X_t - \bar{X})^2}$$

$$= \frac{\sum_{t=1}^n (X_t - \bar{X}) (E(Y_t) - E(\bar{Y}))}{\sum_{t=1}^n (X_t - \bar{X})^2}$$

Puisque le point (\bar{X}, \bar{Y}) est toujours sur la droite de régression:

$$E(\hat{\beta}_1) = \frac{\sum_{t=1}^n (X_t - \bar{X}) \times (\beta_0 + \beta_1 X_t - \beta_0 - \beta_1 \bar{X})}{\sum_{t=1}^n (X_t - \bar{X})^2}$$

$$= \frac{\sum_{t=1}^n (X_t - \bar{X}) \beta_1 (X_t - \bar{X})}{\sum_{t=1}^n (X_t - \bar{X})^2}$$

$$= \beta_1 \frac{\sum_{t=1}^n (X_t - \bar{X})^2}{\sum_{t=1}^n (X_t - \bar{X})^2}$$

$$\Rightarrow \boxed{E(\hat{\beta}_1) = \beta_1}$$

Par conséquent,

$$E(\hat{\beta}_0) = E(\bar{Y} - \hat{\beta}_1 \bar{X}) = E(\bar{Y}) - \bar{X} E(\hat{\beta}_1)$$

$$= \beta_0 + \beta_1 \bar{X} - \bar{X} \beta_1$$

$$= \beta_0$$

$$\Rightarrow \boxed{E(\hat{\beta}_0) = \beta_0}$$

2.3.2) Propriété 2: Variances et covariances:

$$\bullet \text{Var}(\hat{\beta}_1) = \text{Var}\left(\frac{\sum_{t=1}^n (x_t - \bar{x})(y_t - \bar{y})}{\sum_{t=1}^n (x_t - \bar{x})^2} \right)$$

$$= \frac{\text{Var}\left(\sum_{t=1}^n (x_t - \bar{x})y_t - \sum_{t=1}^n (x_t - \bar{x}) \cdot \bar{y} \right)}{\left[\sum_{t=1}^n (x_t - \bar{x})^2 \right]^2}$$

$$= \frac{\text{Var}\left(\sum_{t=1}^n (x_t - \bar{x})y_t \right) + \cancel{\text{Var}\left(\bar{y} \sum_{t=1}^n (x_t - \bar{x}) \right)}}{\left[\sum_{t=1}^n (x_t - \bar{x})^2 \right]^2}$$

$$= \frac{\sum_{t=1}^n (x_t - \bar{x})^2 \text{Var}(y_t) + \text{Var}\left(\bar{y} \times [n\bar{x} - n\bar{x}] \right)}{\left[\sum_{t=1}^n (x_t - \bar{x})^2 \right]^2}$$

$$= \frac{\sum_{t=1}^n (x_t - \bar{x})^2 \text{Var}(\beta_0 + \beta_1 x_t + \varepsilon_t) + 0}{\left[\sum_{t=1}^n (x_t - \bar{x})^2 \right]^2}$$

$$= \frac{\sum_{t=1}^n (x_t - \bar{x})^2 \text{Var}(\varepsilon_t)}{\left[\sum_{t=1}^n (x_t - \bar{x})^2 \right]^2}$$

$$= \frac{\sum_{t=1}^n (x_t - \bar{x})^2 \sigma^2}{\left[\sum_{t=1}^n (x_t - \bar{x})^2 \right]^2}$$

⊛ IMP

$$\Rightarrow \boxed{\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{t=1}^n (x_t - \bar{x})^2}}$$

$$\begin{aligned} \bullet \text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) \\ &= \text{Var}(\bar{y}) + \text{Var}(\hat{\beta}_1 \bar{x}) - 2 \text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x}) \\ &= \text{Var}\left(\frac{\sum_{t=1}^n y_t}{n}\right) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2 \bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1) \\ &= \frac{n \times \text{Var}(y)}{n^2} + \bar{x}^2 \times \left(\frac{\sigma^2}{\sum_{t=1}^n (x_t - \bar{x})^2} \right) - 2 \bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1) \quad \text{⊛} \end{aligned}$$

$$\text{⊛ } \text{Cov}(\bar{y}, \hat{\beta}_1) = \text{Cov}\left(\frac{\sum_{t=1}^n y_t}{n}, \frac{\sum_{s=1}^n (x_s - \bar{x})(y_s - \bar{y})}{\sum_{s=1}^n (x_s - \bar{x})^2}\right)$$

$$\begin{aligned} &= \frac{1}{n} \times \frac{1}{\sum_{s=1}^n (x_s - \bar{x})^2} \times \text{Cov}\left(\sum_{t=1}^n y_t, \sum_{s=1}^n (x_s - \bar{x}) y_s - \bar{y} \sum_{s=1}^n (x_s - \bar{x})\right) \\ &= \frac{1}{n} \times \frac{1}{\sum_{s=1}^n (x_s - \bar{x})^2} \sum_{t=1}^n \sum_{s=1}^n \text{Cov}(y_t, (x_s - \bar{x}) y_s) \end{aligned}$$

↑ = 0 !!!

$$= \frac{1}{n} \times \frac{1}{\sum_{s=1}^n (x_s - \bar{x})^2} \sum_{t=1}^n \sum_{s=1}^n (x_s - \bar{x}) \text{Cov}(e_t, e_s)$$

$$= \frac{1}{n} \times \frac{1}{\sum_{s=1}^n (x_s - \bar{x})^2} \left(\sum_{\substack{t=1 \\ t \neq s}}^n \sum_{s=1}^n (x_s - \bar{x}) \times 0 + \sum_{\substack{t=1 \\ t=s}}^n \sum_{s=1}^n (x_s - \bar{x}) \times \sigma^2 \right)$$

$$= \frac{1}{n} \times \frac{1}{\sum_{s=1}^n (x_s - \bar{x})^2} \times \left(\sigma^2 \sum_{t=1}^n (x_s - \bar{x}) \right)$$

$$= \frac{1}{n} \times \frac{1}{\sum_{s=1}^n (x_s - \bar{x})^2} \times \sigma^2 (n\bar{x} - n\bar{x})$$

$$= 0 !!!$$

Donc;

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{n} + \frac{\bar{x}^2 \sigma^2}{\sum_{t=1}^n (x_t - \bar{x})^2}$$

IMP

Finalement,

$$\begin{aligned} \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) &= \text{Cov}(\bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1) \\ &= \text{Cov}(\bar{y}, \hat{\beta}_1) - \bar{x} \text{Var}(\hat{\beta}_1) \\ &= 0 - \bar{x} \times \frac{\sigma^2}{\sum_{t=1}^n (x_t - \bar{x})^2} \end{aligned}$$

↑
voir ⊗

IMP

2.3.3) Propriété 3: Optimalité:

Le théorème de Gauss-Markov établit que l'estimateur des moindres carrés est l'estimateur non biaisé à variance minimale.

Idée de la preuve:

(1) Considérer l'estimateur $\hat{\beta}^* = \sum_{t=1}^n C_t \cdot y_t$

(2) Minimiser $\text{Var}(\hat{\beta}^*)$ sous la contrainte que $E(\hat{\beta}^*) = \beta$; où

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

2.4) Régression passant par l'origine:

Dans certaines situations, il est possible que l'on souhaite forcer la droite de régression à passer par l'origine.

ex: y_t : consommation d'essence en L d'une voiture
 x_t : # de km parcouru

$$\dots x_t = 0 \Rightarrow y_t = 0 !!!$$

Dans ce cas, on peut postuler le modèle suivant:

$$y_t = \beta x_t + \varepsilon_t$$

On peut montrer que

$$\hat{\beta} = \frac{\sum_{t=1}^n x_t y_t}{\sum_{t=1}^n x_t^2}$$

... même procédure qu'en (2.2.1)

2.5) Analyse de la variance:

Un tableau d'analyse de la variance permet d'évaluer la qualité de l'ajustement du modèle aux observations

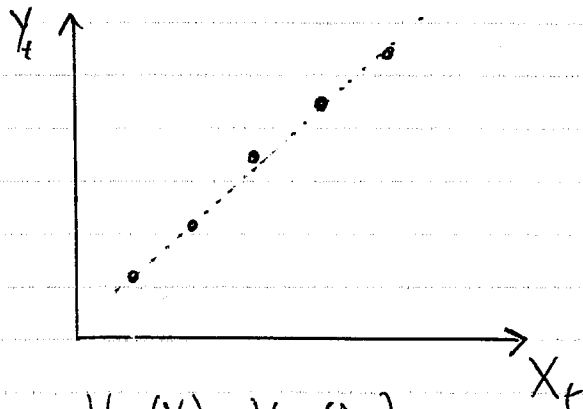
Idee:

(I) Si on décide de modéliser Y_t sans la régression (= analyse statistique des risques actuels...), alors Y est vu comme une v.a. avec une certaine variance = $\text{Var}(Y)$

(II) En utilisant la régression pour modéliser Y_t en fonction de X_t une partie de $\text{Var}(Y_t)$ est "expliquée" par $\text{Var}(X_t)$, alors que l'autre partie reste "inexpliquée"

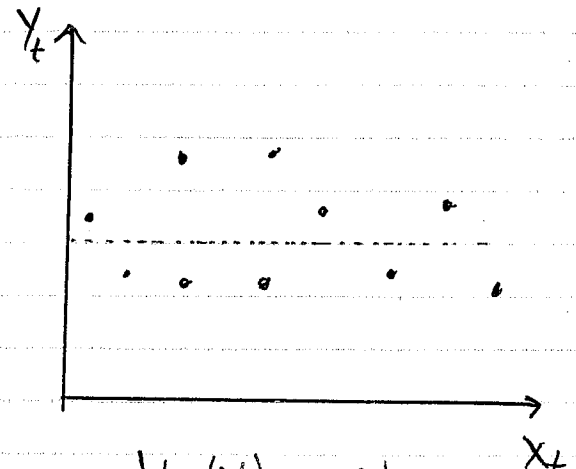
(III) L'utilité de la régression = Proportion de $\text{Var}(Y_t)$ expliquée par $\text{Var}(X_t)$

exemple:



$$\text{Var}(Y_t) = \text{Var}(X_t)$$

⇒ $\text{Var}(Y_t) = 100\% \times \text{Var}(X_t)$: rég. utile!



$$\text{Var}(Y_t) = 0\% \underbrace{\text{Var}(X_t)}_{\text{expliqué}} + 100\% \underbrace{\text{Var}(\varepsilon_t)}_{\text{inexpliqué}}$$

⇒ rég. inutile ⇒ utiliser stat. de base

2.5.1) Background 1: Somme des carrés:

La variance totale dans les Y_t est proportionnelle à:

$$SST = \sum_{t=1}^n (Y_t - \bar{Y})^2$$

Décomposition:

$$(Y_t - \bar{Y}) = Y_t - \hat{Y}_t + \hat{Y}_t - \bar{Y}$$

$$\Rightarrow (Y_t - \bar{Y}) = (Y_t - \hat{Y}_t) + (\hat{Y}_t - \bar{Y})$$

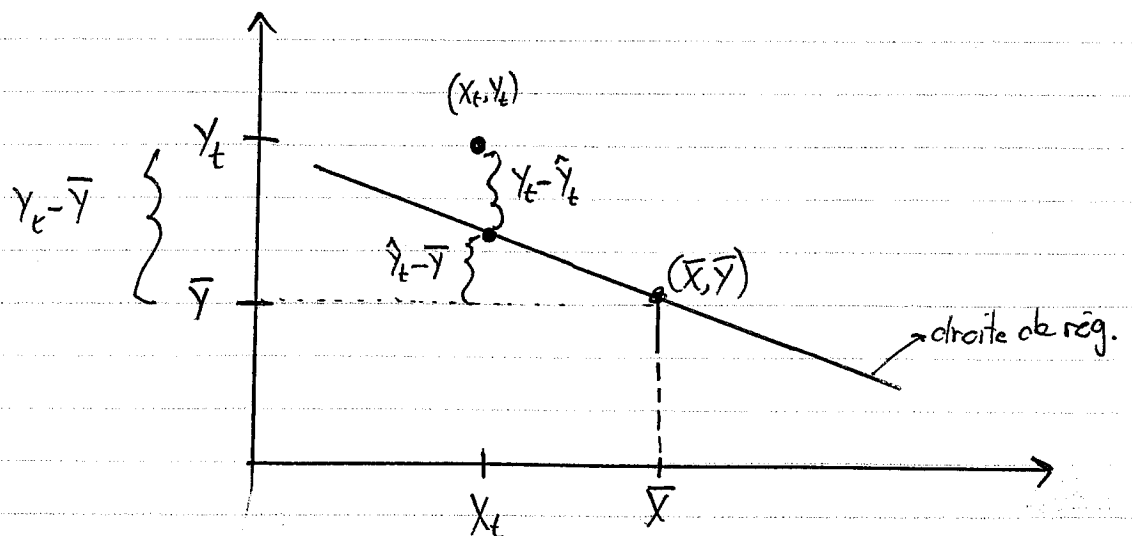
$$\Rightarrow (Y_t - \bar{Y}) = (\hat{Y}_t - \bar{Y}) + (Y_t - \hat{Y}_t)$$

Variation
totale
de Y_t

Variation
de \hat{Y}_t
=
Variation
expliquée
par la
régression

$\hat{\epsilon}_t$
=
résidu
=
Variation inexpliquée
par la régression

Illustration:



Pour conséquent, on a que:

$$\begin{aligned}
 SST &= \sum_{t=1}^n \left[(\hat{y}_t - \bar{y}) + (y_t - \hat{y}_t) \right]^2 \\
 &= \underbrace{\sum_{t=1}^n (\hat{y}_t - \bar{y})^2}_{SSR} + \underbrace{\sum_{t=1}^n (y_t - \hat{y}_t)^2}_{SSE} + 2 \sum_{t=1}^n (\hat{y}_t - \bar{y})(y_t - \hat{y}_t)
 \end{aligned}$$

↑
régression
↑
erreur
★

$$\begin{aligned}
 \text{★} &= 2 \sum_{t=1}^n (\hat{y}_t - \bar{y})(y_t - \hat{y}_t) \\
 &= 2 \sum_{t=1}^n (\hat{\beta}_0 + \hat{\beta}_1 X_t - \hat{\beta}_0 - \hat{\beta}_1 \bar{X})(y_t - \bar{y} + \bar{y} - \hat{y}_t) \\
 &= 2 \sum_{t=1}^n \hat{\beta}_1 (X_t - \bar{X})(y_t - \bar{y} + \hat{\beta}_0 + \hat{\beta}_1 \bar{X} - \hat{\beta}_0 - \hat{\beta}_1 X_t) \\
 &= 2 \sum_{t=1}^n \hat{\beta}_1 (X_t - \bar{X})((y_t - \bar{y}) - \hat{\beta}_1 (X_t - \bar{X})) \\
 &= 2 \hat{\beta}_1 \sum_{t=1}^n (X_t - \bar{X})(y_t - \bar{y}) - 2 \hat{\beta}_1^2 \sum_{t=1}^n (X_t - \bar{X})^2 \\
 &= 2 \hat{\beta}_1 (S_{xy} - \hat{\beta}_1 S_{xx}) \\
 &= 2 \hat{\beta}_1 \left(S_{xy} - \left(\frac{S_{xy}}{S_{xx}} \right) S_{xx} \right) \\
 &= 2 \hat{\beta}_1 (S_{xy} - S_{xy}) = 0!
 \end{aligned}$$

Section
(2.2.1)
remarque
(4)

Ainsi:

$$SST = SSR + SSE$$

Variation
expliquée
par la
régression

Variation
inexpliquée,
ou résiduelle

Intuitivement:

- Dans un bon modèle (régression utile), on aimerait que

$$\Rightarrow SST \approx SSR \quad \dots \text{Var}(Y_t) \approx \text{Var}(X_t)$$

$$\text{ou} \Rightarrow SSE \approx 0 \quad \dots \text{variation résiduelle faible!}$$

- On définit le coefficient de détermination:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

: % de la variance
dans Y_t expliquée
par la régression

$$= 1 - \left(\begin{array}{l} \% \text{ de la variance} \\ \text{dans } Y_t \text{ inexpliquée} \\ \text{par la régression} \end{array} \right)$$

$$\Rightarrow R^2 \in [0, 1]$$

$$\Rightarrow R^2 = 100\% \Rightarrow \text{régression utile} ; R^2 = 0\% \Rightarrow \text{rég. inutile!}$$

rendu
ici
14-09-2010