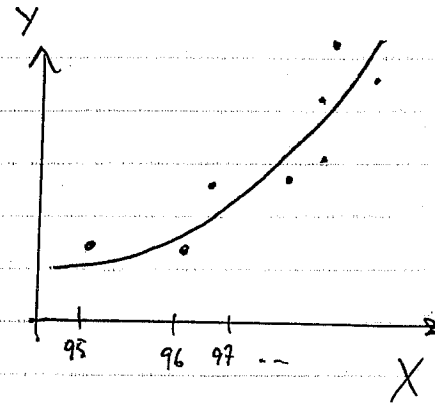


(4)

3) Régression exponentielle:

$$Y = \beta_0 \times e^{\beta_1 X} * \varepsilon$$

\downarrow \uparrow temps
 Sévérité des sinistres Erreur aléatoire



NOTE: $\ln(Y) = \ln(\beta_0) + \beta_1 X + \ln(\varepsilon)$

$$Y^* = \beta_0^* + \beta_1 X + \varepsilon^*$$

... peut se voir comme une régression linéaire simple

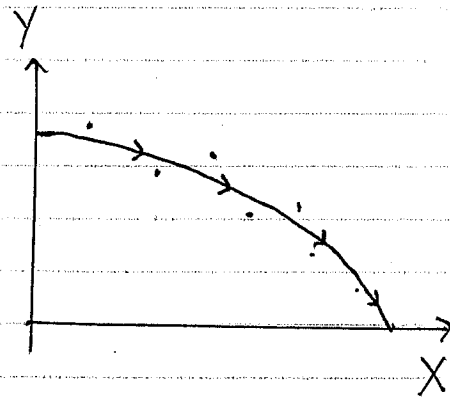
NOTE: Ce modèle est aussi appelé:

- Régression multiplicative
- Régression log-linéaire

4) Régression quadratique:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$$

\downarrow \downarrow \downarrow \downarrow
 Hauteur d'un objet lancé Temps Erreur aléatoire
||
Vent etc...



Note: En posant $X_1 = X$ et $X_2 = X^2$, ce modèle peut se voir comme :

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon \quad \dots \text{une régression multiple !!!}$$

(5)

* Dans ce cours, nous allons nous concentrer sur les modèles linéaires

- Plus simples
- Plusieurs modèles peuvent se ramener à un modèle linéaire (ex: reg. multiplicative, reg. quadratique)
- Constituent souvent une très bonne approximation de la réalité qui peut être très complexe (ex: assurance!!!)
- Se généralisent facilement (ex: GLM: generalized linear models)

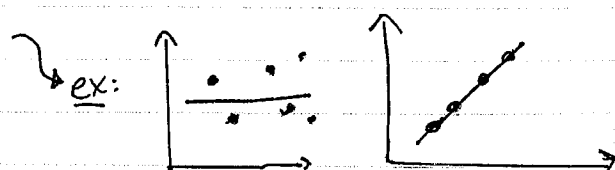
Le problème: Trouver les paramètres $\beta_0, \beta_1, \dots, \beta_p$ tels que $\varepsilon = Y - f(X_1, \dots, X_p; \beta_0, \beta_1, \dots, \beta_p)$ est "petite"

Quelle mesure d'erreur utiliser?

1) Erreur totale? $\sum_{t=1}^n \varepsilon_t = \sum_{t=1}^n (Y_t - (\beta_0 + \beta_1 X_t))$

→ Facile à mettre à 0

→ Pas fiable!



2) Erreur absolue? $\sum_{t=1}^n |\varepsilon_t| = \sum_{t=1}^n |Y_t - (\beta_0 + \beta_1 X_t)|$

→ Très robuste

→ Complicé mathématiquement (minimiser $\sum |\varepsilon_t|$ implique de dériver cette fonction!!!)

3) Erreur quadratique?:
$$\sum_{t=1}^n \varepsilon_t^2 = \sum_{t=1}^n \left[Y_t - (\beta_0 + \beta_1 X_t) \right]^2 \quad (6)$$

→ Plus simple mathématiquement

→ Donne beaucoup de poids aux grandes erreurs

* Cette mesure est celle que nous allons étudier *

2.2) Le modèle de régression linéaire simple:

On dispose de n paires d'observations $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$

Exemple: X_t : # d'années de scolarité de l'actuaire t

Y_t : Salaire de l'actuaire t

Idee:

- Par $X_t = 0$; on a $Y_t = \beta_0$ (ex: β_0 = Salaire d'un stagiaire!)
«en moyenne»

- Par chaque année additionnelle de scolarité, le salaire augmente de β_1 unités «en moyenne»

- Ainsi «en moyenne» on a:

$$E(Y_t | X_t) = \beta_0 + \beta_1 X_t$$

Habituellement, la relation n'est pas exacte dans la réalité.
La «différence» (ou l'erreur) est notée ε_t et est assumée aléatoire:

$$\varepsilon_t = Y_t - E(Y_t | X_t) = Y_t - (\beta_0 + \beta_1 X_t)$$

En réorganisant ; on trouve le modèle suivant :

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

Diagram illustrating the components of the regression model:

- Y_t : Salaire
- β_0 : Ordonnée à l'origine
|| Salaire si $X_t = 0$
- β_1 : Pente
|| Taux de croissance du salaire avec X_t
- X_t : # années de scolarité
- ϵ_t : Erreur aléatoire
|| résidu

But: Trouver les paramètres β_0 et β_1 de façon à minimiser l'erreur ϵ_t

Note: Si ϵ_t est minimal, cela veut dire que $Y_t \approx \beta_0 + \beta_1 X_t$

\Rightarrow Donc la droite de régression $(\beta_0 + \beta_1 X_t)$ est une bonne approximation de Y_t .

2.2.1) Coefficients de régression:

Les paramètres (ou coefficients) β_0 et β_1 sont déterminés en minimisant l'erreur quadratique (méthode des moindres carrés).

$$\begin{aligned} S(\beta_0, \beta_1) &= \sum_{t=1}^n \epsilon_t^2 \\ &= \sum_{t=1}^n (Y_t - (\beta_0 + \beta_1 X_t))^2 \\ &= \sum_{t=1}^n (Y_t - \beta_0 - \beta_1 X_t)^2 \end{aligned}$$

(8)

Minimisation:Note: on met un "hat" pour représenter l'estimation
... les paramètres "optimaux"

$$① \frac{\partial S(\hat{\beta}_0, \hat{\beta}_1)}{\partial \hat{\beta}_0} = 0 \Rightarrow -2 \sum_{t=1}^n (Y_t - \hat{\beta}_0 - \hat{\beta}_1 X_t) = 0$$

$$\Rightarrow \boxed{\sum_{t=1}^n Y_t - n \times \hat{\beta}_0 - \hat{\beta}_1 \sum_{t=1}^n X_t = 0}$$

$$② \frac{\partial S(\hat{\beta}_0, \hat{\beta}_1)}{\partial \hat{\beta}_1} = 0 \Rightarrow -2 \sum_{t=1}^n (Y_t - \hat{\beta}_0 - \hat{\beta}_1 X_t) \times X_t = 0$$

$$\Rightarrow \boxed{\sum_{t=1}^n X_t Y_t - \hat{\beta}_0 \sum_{t=1}^n X_t - \hat{\beta}_1 \sum_{t=1}^n X_t^2 = 0}$$

Deux équations (① et ②) et deux inconnus ($\hat{\beta}_0$ et $\hat{\beta}_1$):

$$\text{De ①: } \hat{\beta}_0 = \left(\frac{\sum_{t=1}^n Y_t}{n} \right) - \hat{\beta}_1 \left(\frac{\sum_{t=1}^n X_t}{n} \right)$$

IMP:

$$\Rightarrow \boxed{\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}}$$

$$\text{De ②: } \hat{\beta}_1 = \frac{\sum_{t=1}^n X_t Y_t - \hat{\beta}_0 \sum_{t=1}^n X_t}{\sum_{t=1}^n X_t^2}$$

$$\text{On met ① dans ②} \quad \hat{\beta}_1 = \frac{\sum_{t=1}^n X_t Y_t - (\bar{Y} - \hat{\beta}_1 \bar{X}) \sum_{t=1}^n X_t}{\sum_{t=1}^n X_t^2}$$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{t=1}^n X_t Y_t - (\bar{Y} - \hat{\beta}_1 \bar{X}) \times n \times \bar{X}}{\sum_{t=1}^n X_t^2}$$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{t=1}^n X_t Y_t - n \bar{X} \bar{Y} + \hat{\beta}_1 \times n \times \bar{X}^2}{\sum_{t=1}^n X_t^2}$$

On isole $\hat{\beta}_1$:

*IMP:

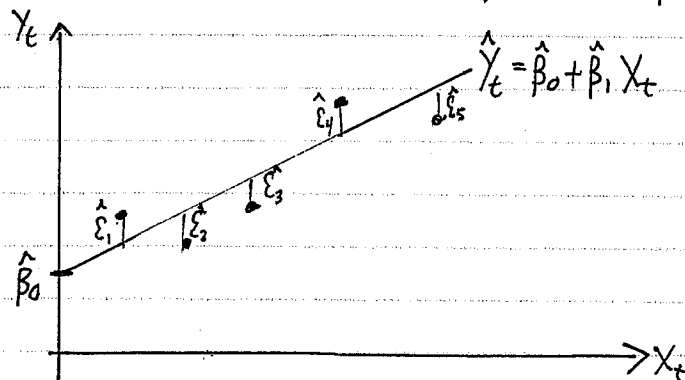
$$\hat{\beta}_1 = \frac{\sum_{t=1}^n X_t Y_t - n \bar{X} \bar{Y}}{\sum_{t=1}^n X_t^2 - n \bar{X}^2}$$

Remarques:

1) On note $\hat{\varepsilon}_t$ les «résidus» générés par le modèle estimé:

$$\hat{\varepsilon}_t = Y_t - \hat{Y}_t$$

$$\Rightarrow \varepsilon_t = Y_t - (\hat{\beta}_0 + \hat{\beta}_1 X_t) \quad ; \text{ pour } t=1, \dots, n$$



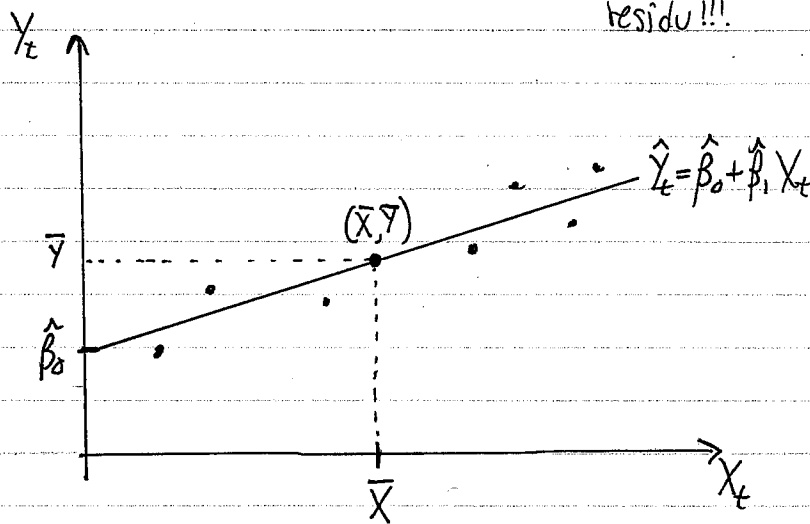
2) Le "centre de gravité" (ou centre de masse) des données (\bar{X}, \bar{Y}) se trouve sur la droite de régression

Preuve:

$$\text{De ①: } \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\Rightarrow \bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X} + 0$$

\equiv
↑
aucun
résidu!!!



$$\begin{aligned} 3) \text{ On note : } \rightarrow S_{XX} &= \sum_{t=1}^n (X_t - \bar{X})^2 = \sum_{t=1}^n (X_t^2 - 2X_t\bar{X} + \bar{X}^2) \\ &= \sum_{t=1}^n X_t^2 - 2\bar{X} \sum_{t=1}^n X_t + n\bar{X}^2 \\ &= \sum_{t=1}^n X_t^2 - 2\bar{X} \times n\bar{X} + n\bar{X}^2 \\ &= \left(\sum_{t=1}^n X_t^2 \right) - n\bar{X}^2 \end{aligned}$$

$$\rightarrow S_{XY} = \sum_{t=1}^n (X_t - \bar{X})(Y_t - \bar{Y}) \stackrel{\text{démon.}}{=} \left(\sum_{t=1}^n X_t Y_t \right) - n\bar{X}\bar{Y}$$

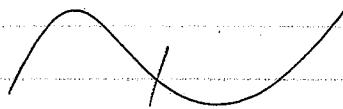
4) Sans la remarque (3); on a que:

$$\boxed{\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}} \quad \text{et} \quad \boxed{\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}}$$

5) La somme des résidus de tout modèle de régression linéaire est nulle:

Preuve:

$$\begin{aligned} \sum_{t=1}^n \hat{\varepsilon}_t &= \sum_{t=1}^n (Y_t - (\hat{\beta}_0 + \hat{\beta}_1 X_t)) \\ &\stackrel{\text{de } 4)}{=} \sum_{t=1}^n (Y_t - (\bar{Y} - \hat{\beta}_1 \bar{X}) - \hat{\beta}_1 X_t) \\ &= \sum_{t=1}^n Y_t - \sum_{t=1}^n \bar{Y} + \hat{\beta}_1 \sum_{t=1}^n \bar{X} - \hat{\beta}_1 \sum_{t=1}^n X_t \\ &= n \times \bar{Y} - n \times \bar{Y} + \hat{\beta}_1 \times n \times \bar{X} - \hat{\beta}_1 \times n \times \bar{X} \\ &= 0 !!! \end{aligned}$$



↑↑↑↑

Rendu ici
07-09-2010