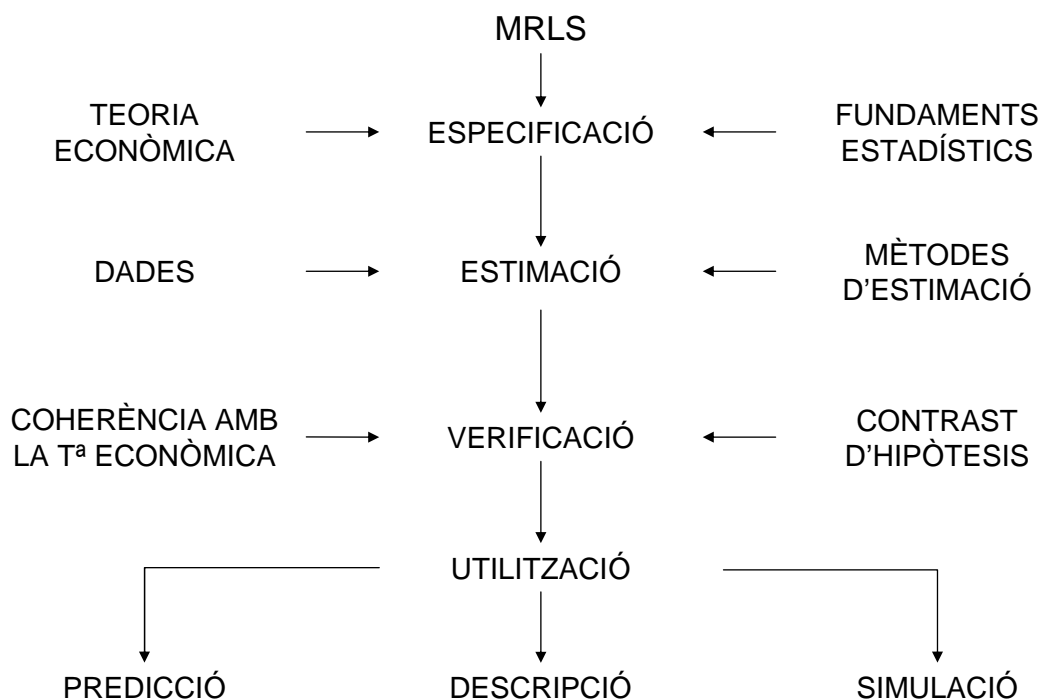


Tema 5. Model de Regressió Lineal Simple

5.1. Objecte i natura del model



López-Tamayo, Jordi

1

Tema 5. Model de Regressió Lineal Simple

5.1. Objecte i natura del model

Exemples:

a) Funció de Consum: $C=a+bR$

Relació lineal entre el consum i la renda disponible. On $0 < b < 1$ representa la proporció marginal a consumir.

b) Funció de demanda: $Q=a+bP$

Relació lineal entre les vendes i el preu de venda. On $b < 0$

Quan observem els consums (vendes) de diferents famílies (empreses) amb la mateixa renda (preus) no totes presenten els mateixos consums (vendes).

És a dir, les relacions econòmiques (empresarials) NO SÓN DETERMINISTES. Existeixen altres aspectes que influeixen sobre les variables observades (gustos, mida familiar, mida empresarial, territori, sector productiu, situació econòmica general, etc..)

Tots aquests aspectes els hem de recollir a algun lloc. Per això, incorporem un TERME DE PERTORBACIÓ (u o ε) amb el que procurem recollir tots aquests aspectes i sobre el que farem alguns supòsits. Així:

$$\begin{cases} C = a + bR + U \\ Q = a + bP + U \end{cases} \quad \text{Amb aquesta operació, convertim les relacions teòriques, DETERMINISTES, en relacions ESTOCÀSTIQUES.}$$

López-Tamayo, Jordi

2

Tema 5. Model de Regressió Lineal Simple

5.1. Hipòtesis per a l'especificació del model

3. No existeixen variables rellevants OMESES i

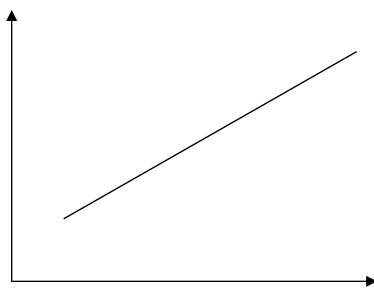
per això les pertorbacions no estan $\longrightarrow \text{cor}(\varepsilon_i, \varepsilon_j) = E(\varepsilon_i * \varepsilon_j) = 0 \quad \forall i \neq j$

AUTOCORRELACIONADES

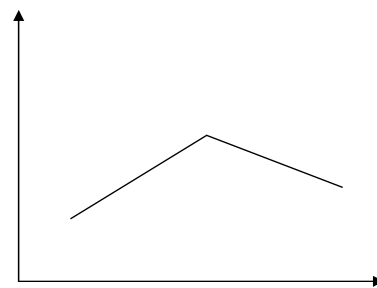
4. Normalitat del terme de pertorbació $\longrightarrow \varepsilon_i \approx N(0, \sigma) \quad \forall i$

5. X és DETERMINISTA (NO ALEATÒRIA) és fixada per l'investigador o de forma EXÒGENA al model.

6. Existeix permanència ESTRUCTURAL, només existeix una recta que representa el COMPORTAMENT DE LA RELACIÓ



SI



NO

López-Tamayo, Jordi

5

Tema 5. Model de Regressió Lineal Simple

5.1. Hipòtesis per a l'especificació del model

Un cop ASUMITS AQUESTS SUPÒSITS, es procedeix a caracteritzar el comportament PROBABILÍSTIC de la VARIABLE ALEATÒRIA ENDÒGENA.

$$E(Y_i) = E(\alpha + \beta X_i + \varepsilon_i) = E(\alpha) + E(\beta X_i) + E(\varepsilon_i) = \underbrace{\alpha + \beta X_i}_{\text{Recta de Regressió}} + \underbrace{0}_{\text{hipòtesi del model}}$$

$$V(Y_i) = E(Y_i - E(Y_i))^2 = E((\alpha + \beta X_i + \varepsilon_i) - (\alpha + \beta X_i))^2 = E(\varepsilon_i)^2 \stackrel{\text{hipòtesi del model}}{=} \sigma^2$$

Quina forma té la distribució d'Y?

És una combinació lineal on l'única variable que presenta aleatorietat és el terme de pertorbació. Atès que aquest és NORMAL, per les propietats de la normal, la variable Y es distribuirà normalment. Així:

$$Y_i \approx N(\alpha + \beta X_i, \sigma)$$

ECET 5.1

López-Tamayo, Jordi

6

Tema 5. Model de Regressió Lineal Simple

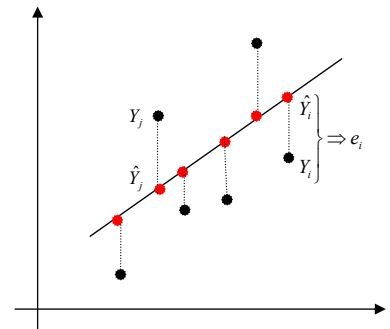
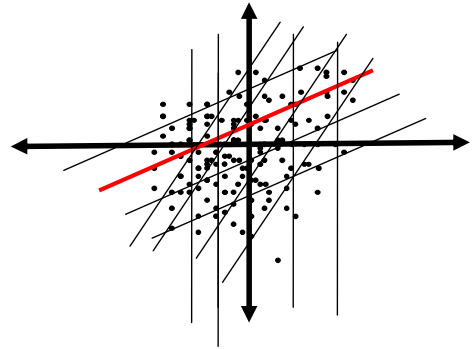
5.3. Estimació dels paràmetres de la recta de regressió per Mínims Quadrats Ordinaris

Per a concretar quina de totes les possibles normals hem d'estimar els paràmetres α , β i σ^2 .

De quina informació disposem? \rightarrow mostra (X i Y) de "n" individus.

Existeixen diferents mètodes d'estimació. Hem vist, durant aquest curs dos: MM, MV. Ara presentem el mètode de MÍNIMS QUADRATS ORDINARIS.

IDEA: La millor recta de regressió serà aquella que minimitzi els errors que cometem.



Tema 5. Model de Regressió Lineal Simple

5.3. Estimació dels paràmetres de la recta de regressió per Mínims Quadrats Ordinaris

CRITERI: Minimització de la suma dels errors al quadrat (no m'importa si m'equivoco per sobre o per baix, m'importa el segment, no el seu signe)

$$\text{Min} \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2$$

$$\begin{cases} \frac{\partial \sum_{i=1}^n e_i^2}{\partial \hat{\alpha}} = -2 \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i) = 0 \Rightarrow \sum_{i=1}^n Y_i - n\hat{\alpha} - \hat{\beta} \sum_{i=1}^n X_i = 0 \\ \frac{\partial \sum_{i=1}^n e_i^2}{\partial \hat{\beta}} = -2 \sum_{i=1}^n X_i (Y_i - \hat{\alpha} - \hat{\beta}X_i) = 0 \Rightarrow \sum_{i=1}^n X_i Y_i - \hat{\alpha} \sum_{i=1}^n X_i - \hat{\beta} \sum_{i=1}^n X_i^2 = 0 \end{cases}$$

$$\begin{cases} \bar{Y} = \hat{\alpha} + \hat{\beta}\bar{X} \\ \sum_{i=1}^n X_i Y_i = (\bar{Y} - \hat{\beta}\bar{X}) \sum_{i=1}^n X_i + \hat{\beta} \sum_{i=1}^n X_i^2 \Rightarrow \hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{S_{XY}}{S_X^2} \end{cases}$$

Tema 5. Model de Regressió Lineal Simple

5.3. Estimació dels paràmetres de la recta de regressió per Mínims Quadrats Ordinaris

CARACTERÍSTIQUES DE LA REGRESSIÓ:

1. Sempre passa pel CENTRE del núvol de punts.

$$2. \sum e_i = 0 \quad i \quad E(X_i e_i) = 0$$

3. Atès que l'estimador de β , depèn de la covariància, el

seu signe depèn d'aquesta. També pot expressar-se com: $\longrightarrow \hat{\beta} = r_{XY} \frac{S_Y}{S_X}$

4. L'ERROR D'ESTIMACIÓ no és el mateix que el TERME DE PERTORBACIÓ

$$\text{NIVELL POBLACIONAL} \longrightarrow Y_i = \alpha + \beta X_i + \varepsilon_i = E(Y_i) + \varepsilon_i$$

$$\text{NIVELL MOSTRAL} \longrightarrow Y_i = \hat{\alpha} + \hat{\beta} X_i + e_i = \hat{Y}_i + e_i$$

ECET 5.2

López-Tamayo, Jordi

9

Tema 5. Model de Regressió Lineal Simple

5.3. Estimació dels paràmetres de la recta de regressió per Mínims Quadrats Ordinaris

PROPIETATS DELS ESTIMADORS MQO:

Si es compleixen les hipòtesis bàsiques del model

1. LINEALITAT: Són funcions lineals de les variables.

2. NO ESBIAXATS $E(\hat{\alpha}) = \alpha$ i $E(\hat{\beta}) = \beta$

3. Variància mínima:

$$V(\hat{\alpha}_*) \geq V(\hat{\alpha}_{MQO})$$

$$V(\hat{\beta}_*) \geq V(\hat{\beta}_{MQO})$$

$$V(\hat{\beta}) = V \left(\frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n ((X_i - \bar{X})^2)^2} V((Y_i - \bar{Y}))$$
$$= \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$V(\hat{\alpha}) = V(\bar{Y} - \hat{\beta}\bar{X}) =$$
$$= V(\bar{Y}) + \bar{X}^2 V(\hat{\beta}) - 2\bar{X} \text{cov}(\hat{\beta}, \bar{X}) =$$
$$= \sigma^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]$$

López-Tamayo, Jordi

10

Tema 5. Model de Regressió Lineal Simple

5.3. Estimació dels paràmetres de la recta de regressió per Mínims Quadrats Ordinaris

4. Distribució Normal.

$$\hat{\alpha} \approx N\left(\alpha, \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2}\right)}\right) \quad i \quad \hat{\beta} \approx N\left(\beta, \sqrt{\frac{\sigma^2}{\sum (X_i - \bar{X})^2}}\right)$$

Ara bé, per a poder fer inferència sobre els dos paràmetres poblacionals α i β , necessitem proposar una estimació de σ^2 :

$$\sigma^2 = E(\varepsilon_i)^2 \Rightarrow \text{es proposa com estimador } \hat{\sigma}^2 = \frac{\sum e_i^2}{n-2}$$

A nivell empíric, el sumatori dels errors es pot descomposar de la següent manera:

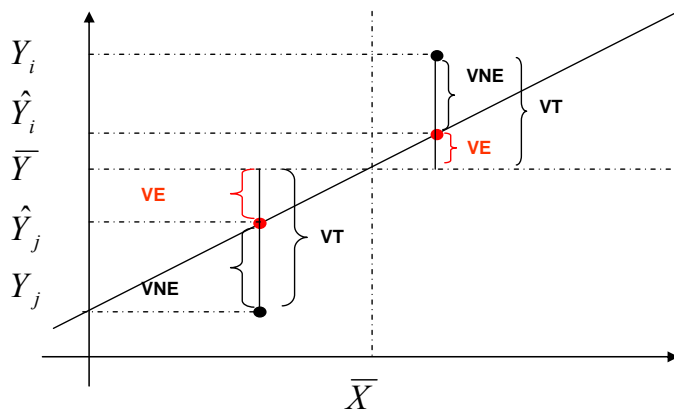
$$\begin{aligned} \sum e_i^2 &= \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2 = \sum (Y_i - (\bar{Y} - \hat{\beta}\bar{X}) - \hat{\beta}X_i)^2 = \\ &= \sum ((Y_i - \bar{Y}) - \hat{\beta}(X_i - \bar{X}))^2 = \sum (Y_i - \bar{Y})^2 + \hat{\beta}^2 \sum (X_i - \bar{X})^2 - 2\hat{\beta} \sum (Y_i - \bar{Y})(X_i - \bar{X}) = \\ &= \sum (Y_i - \bar{Y})^2 + \hat{\beta}^2 \sum (X_i - \bar{X})^2 - 2\hat{\beta}\hat{\beta} \sum (X_i - \bar{X})^2 = \\ &= \sum (Y_i - \bar{Y})^2 - \hat{\beta}^2 \sum (X_i - \bar{X})^2 = nS_Y^2 - \hat{\beta}^2 nS_X^2 = n(S_Y^2 - \hat{\beta}^2 S_X^2) \end{aligned}$$

ECET 5.3

Tema 5. Model de Regressió Lineal Simple

5.4/5.5 Contrastació del model i bondat de l'ajust

La BONDAT DE L'AJUST valora fins a quin punt La recta estimada s'ajusta al núvol de punts. Per Això, valora quina part de la VARIABILITAT TOTAL de la VARIABLE ENDÒGENA és explicada per la RECTA "AJUSTADA" DE REGRESSIÓ



Tenint present el centre del núvol de punts, La distància entre la observació real de la variable endògena i la seva mitjana (Variació Total- VT), es pot descomposar en dos segments:

- Variació Explicada (VE). La distància que existeix entre l'ajust obtingut i la mitjana de la variable endògena.
- Variació NO Explicada (VNE). La distància entre l'observació real i la mitjana de la variable endògena.

$$Y_i - \bar{Y} = (Y_i - \hat{Y}_i) + (\hat{Y}_i - \bar{Y}) = \text{si elevem al quadrat i sumem tots els segments}$$

$$\underbrace{\sum (Y_i - \bar{Y})^2}_{SQT \text{ Suma Quadrats total}} = \underbrace{\sum (Y_i - \hat{Y}_i)^2}_{SQE \text{ Suma Quadrats Errors}} + \underbrace{\sum (\hat{Y}_i - \bar{Y})^2}_{SOR \text{ Suma Quadrats Regressió}} + \underbrace{2 \sum (Y_i - \hat{Y}_i) \sum (\hat{Y}_i - \bar{Y})}_{\text{Atès que } \sum e_i = 0}$$

Tema 5. Model de Regressió Lineal Simple

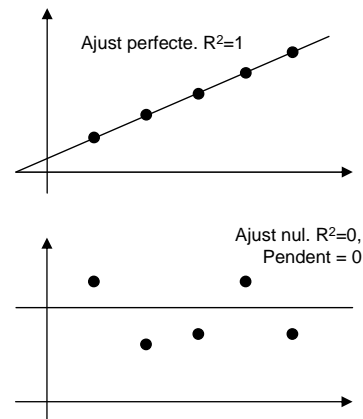
5.4/5.5 Contrastació del model i bondat de l'ajust

El COEFICIENT DE DETERMINACIÓ és la mesura que ens ajudarà a valorar la BONDAT DE L'AJUST

$$R^2 = \frac{SQR}{SQT} = 1 - \frac{SQE}{SQT} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = 1 - \frac{\sum (Y_i - \hat{Y}_i)^2}{\sum (Y_i - \bar{Y})^2}$$

$$R^2 = r_{XY}^2$$

- Si $R^2=1$, els errors són nuls i totes les observacions mostrals es troben a sobre de la recta de regressió. $SQE=0$.
- Si $R^2=0$, tota la variació de la variable endògena se explica pels errors i, aleshores, la recta de regressió no explica res de la variable endògena. $SQR=0$. (El coeficient de correlació entre les variables és 0)
- Per tant $0 < R^2 < 1$. Quan més a prop estigui d'1 millor eera l'ajust i quan més a prop estigui de 0, serà pitjor.



Tema 5. Model de Regressió Lineal Simple

5.4/5.5 Contrastació del model i bondat de l'ajust

CONTRASTACIÓ DEL MODEL: INTERVAL DE CONFIANÇA

Sota les HIPÒTESIS DEL MODEL. L'interval de confiança per a l'estimador del **paràmetre β** és:

$$\hat{\beta} \approx N \left(\beta, \sqrt{\frac{\sigma^2}{\sum (X_i - \bar{X})^2}} \right) \quad \text{Atès que } \Rightarrow \quad \sigma^2 \text{ es desconeguda } \Rightarrow \quad \hat{\sigma}^2 = \frac{\sum e_i^2}{n-2}$$

per tant, $\sigma_{\hat{\beta}} = \sqrt{\frac{\sum e_i^2 / (n-2)}{\sum (X_i - \bar{X})^2}}$ i en aquest cas sabem que la distribució segueix una t-student.

Així, l'interval de confiança es construeix com: $P[\hat{\beta} - t_{\alpha/2, n-2} \hat{\sigma}_{\hat{\beta}} \leq \beta \leq \hat{\beta} + t_{\alpha/2, n-2} \hat{\sigma}_{\hat{\beta}}] = 1 - \alpha$

I en el cas del **paràmetre α** és:

$$\hat{\alpha} \approx N \left(\alpha, \sqrt{\frac{\sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right)}{\hat{\sigma}_{\hat{\alpha}}}} \right) \quad \text{Per tant, } P[\hat{\alpha} - t_{\alpha/2, n-2} \hat{\sigma}_{\hat{\alpha}} \leq \alpha \leq \hat{\alpha} + t_{\alpha/2, n-2} \hat{\sigma}_{\hat{\alpha}}] = 1 - \alpha$$

Tema 5. Model de Regressió Lineal Simple

5.4/5.5 Contrastació del model i bondat de l'ajust

Finalment, l'interval de confiança per a l'estimador del **paràmetre σ^2** és:

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \approx \chi_{n-2}^2 \quad \text{per tant, } P\left[\frac{(n-2)\hat{\sigma}^2}{\chi_{1-\alpha/2, n-2}^2} \leq \sigma^2 \leq \frac{(n-2)\hat{\sigma}^2}{\chi_{\alpha/2, n-2}^2}\right] = 1 - \alpha$$

ECET 5.4

CONTRAST D'HIPÒTESI sobre els **paràmetres β i α**

$$\begin{cases} H_0: \beta = \beta_0 \\ H_1: \beta \neq \beta_0 \end{cases} \Rightarrow \text{Si } \left| \frac{\hat{\beta} - \beta_0}{\hat{\sigma}_{\hat{\beta}}} \right| \geq t_{n-2, \alpha/2} \Rightarrow RH_0$$

$$\begin{cases} H_0: \alpha = \alpha_0 \\ H_1: \alpha \neq \alpha_0 \end{cases} \Rightarrow \text{Si } \left| \frac{\hat{\alpha} - \alpha_0}{\hat{\sigma}_{\hat{\alpha}}} \right| \geq t_{n-2, \alpha/2} \Rightarrow RH_0$$

I, finalment, sobre la BONDAT DE L'AJUST, és a dir, la significació global del model:

$$\begin{cases} H_0: R^2 = 0 \\ H_1: R^2 > 0 \end{cases} \Rightarrow \text{amb l'estadístic } F^* = \frac{(n-1)R^2}{1-R^2} \Rightarrow \text{Si } F^* > F_{1, n-2, 1-\alpha} \Rightarrow RH_0$$

que en un model de REGRESSIÓ LINEAL SIMPLE és el mateix que fer el contrast sobre el pendent de la recta.

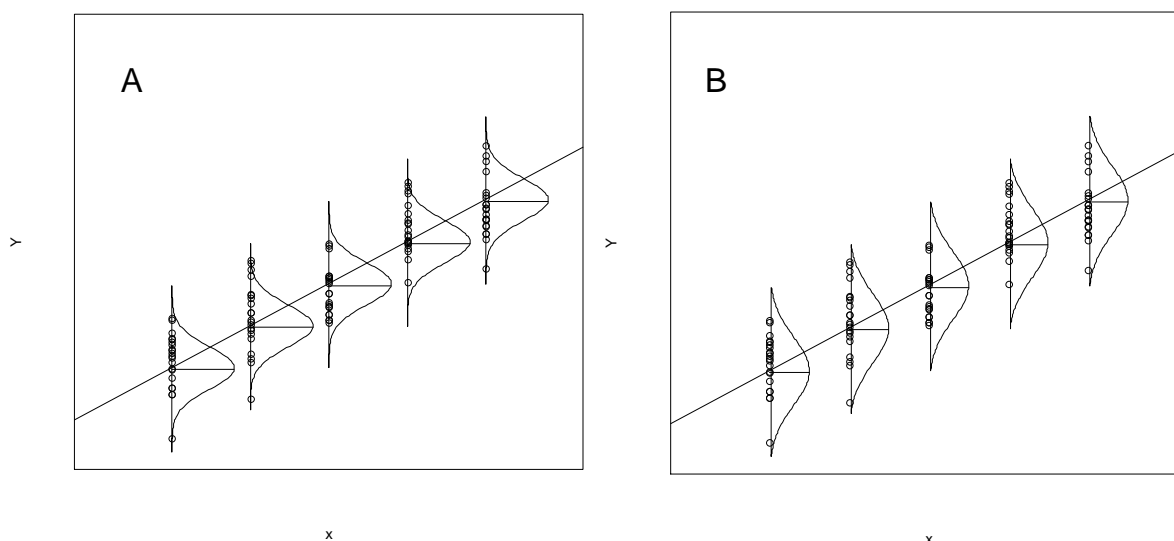
López-Tamayo, Jordi

ECET 5.5

15

Tema 5. Model de Regressió Lineal Simple

ANNEX: Incompliment de les hipòtesis del model: HETEROSCEDASTICITAT



Tant A com B presenten un terme de pertorbació HOMOCEDÀSTIC, malgrat però, a B la variabilitat és major que a A.

HETEROSCEDASTICITAT: Quan aquesta variabilitat no és igual per a cada error "i"

López-Tamayo, Jordi

16

Tema 5. Model de Regressió Lineal Simple

5.6 Predicció puntual i per interval

PREDICCIÓ PUNTUAL:

$$\hat{Y}_{n+1} = \hat{\alpha} + \hat{\beta}X_{n+1}$$

Igual que existeix un ERROR D'ESTIMACIÓ, existeix un ERROR DE PREDICCIÓ.

$$\begin{aligned} e_{n+1} &= Y_{n+1} - \hat{Y}_{n+1} = (\alpha + \beta X_{n+1} + \varepsilon_{n+1}) - (\hat{\alpha} + \hat{\beta}X_{n+1}) = \\ &= \underbrace{(\alpha - \hat{\alpha}) + (\beta - \hat{\beta})X_{n+1}}_{\substack{\text{FONT D'ERROR 1} \\ \text{Per l'error comés a l'hora} \\ \text{d'estimar els paràmetres}}} + \underbrace{\varepsilon_{n+1}}_{\substack{\text{FONT D'ERROR 2} \\ \text{Associat al terme de} \\ \text{pertorbació de l'individu} \\ \text{n+1}}} \end{aligned}$$

La pregunta ara és, com es distribueix aquest error?. Haurem de caracteritzar, la seva esperança, la seva variància i la forma de la distribució.

$$E(e_{n+1}) = E((\alpha - \hat{\alpha}) + (\beta - \hat{\beta})X_{n+1} + \varepsilon_{n+1}) = 0 + 0 + \underbrace{0}_{\substack{\text{Per hipòtesis} \\ \text{del model} \\ E(\varepsilon_{n+1})=0}} = 0$$

Tema 5. Model de Regressió Lineal Simple

5.6 Predicció puntual i per interval

$$V(e_{n+1}) = V((\alpha - \hat{\alpha}) + (\beta - \hat{\beta})X_{n+1} + \varepsilon_{n+1}) = V(\hat{\alpha}) + X_{n+1}^2 V(\hat{\beta}) + V(\varepsilon_{n+1}) + 2X_{n+1} \text{cov}(\hat{\alpha}, \hat{\beta})$$

$$\left\{ \begin{array}{l} V(\hat{\alpha}) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right) \\ V(\hat{\beta}) = \frac{\sigma^2}{\sum (X_i - \bar{X})^2} \\ \text{cov}(\hat{\alpha}, \hat{\beta}) = \frac{\sigma^2 \bar{X}}{\sum (X_i - \bar{X})^2} \end{array} \right\} \Rightarrow \text{Al final: } V(e_{n+1}) = \sigma^2 \left(1 + \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right)$$

Quan més lluny estigui X_{n+1} de \bar{X} més gran serà la variància.

Aquesta és mínima si $X_{n+1} = \bar{X}$.

$$\text{Si } e_{n+1} \approx N \left(0, \underbrace{\sigma^2 \left(1 + \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right)}_{\hat{\sigma}_{e_{n+1}}} \right) \Rightarrow P[\hat{Y}_{n+1} - t_{\alpha/2, n-2} \hat{\sigma}_{e_{n+1}} \leq Y_{n+1} \leq \hat{Y}_{n+1} + t_{\alpha/2, n-2} \hat{\sigma}_{e_{n+1}}] = 1 - \alpha$$