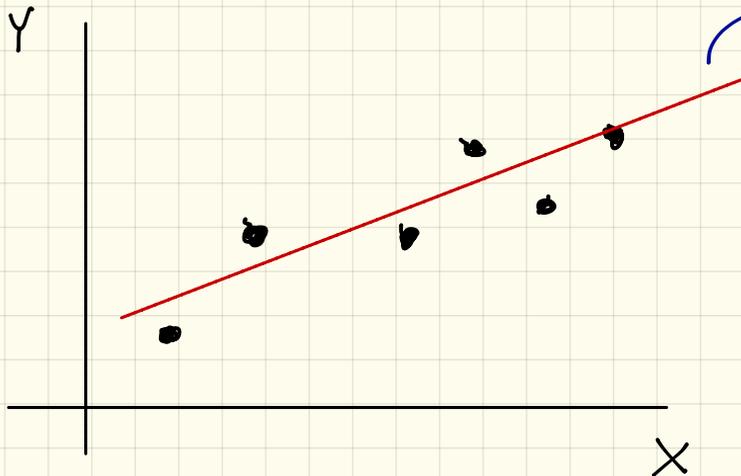


Lezione n. 34

(martedì, 1 marzo 2016)

	Task Description	Due Date	Status
<input type="checkbox"/>	Richiami sulla regressione e differenza nel caso di indagini campionarie		
<input type="checkbox"/>	Le ipotesi alla base del modello classico di regressione lineare		
<input type="checkbox"/>	Teorema di Gauss - Markov		
<input type="checkbox"/>	Le ipotesi alla base del modello classico normale di regressione		
<input type="checkbox"/>	Stima della varianza dell'errore		
<input type="checkbox"/>	Stima intervallare e verifiche di ipotesi per i parametri di regressione		
<input type="checkbox"/>	ARGOMENTO AGGIUNTIVO R^2 e formule inverse per calcolare $Der(R)$ e $Der(E)$		
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			

RICHIAMI REGRESSIONE (1)



obiettivo:
"migliore" retta che
interpoli il sistema
di punti

critero minimi quadrati

$$\min_{b_0, b_1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$b_0 + b_1 x$$

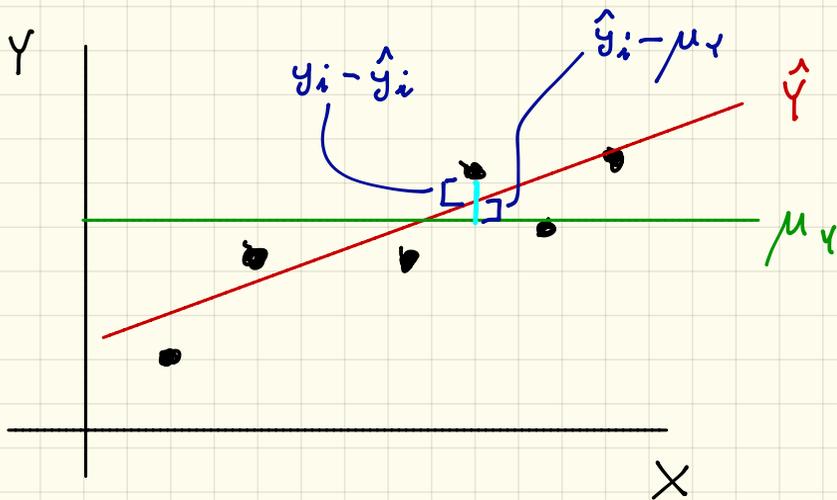
soluzioni dei minimi
quadrati

$$b_1 = \frac{\text{covar}(X, Y)}{\text{Dev}(X)} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$$

$$b_0 = \mu_Y - b_1 \mu_X$$

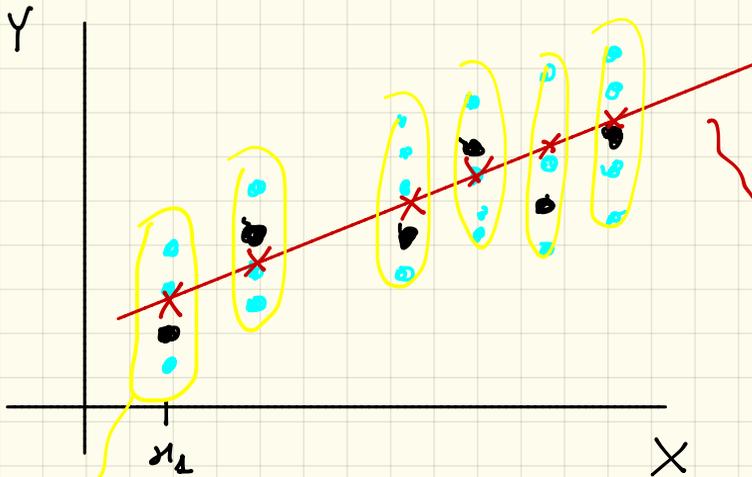
RICHIAMI REGRESSIONE (2)

$$R^e = \frac{\text{Der. spiegata}}{\text{Der. totale}} = \frac{\text{Der}(R)}{\text{Der}(Y)} = 1 - \frac{\text{Der}(E)}{\text{Der}(Y)}$$
$$= \rho_{xy}^e = \hat{\beta}_1^e \frac{\text{Der}(X)}{\text{Der}(Y)}$$



MEMO:

$$\sum_i (y_i - \mu_Y)^e = \sum_i (y_i - \hat{y}_i)^e + \sum_i (\hat{y}_i - \mu_Y)^e$$



● valori osservati
 (campione)

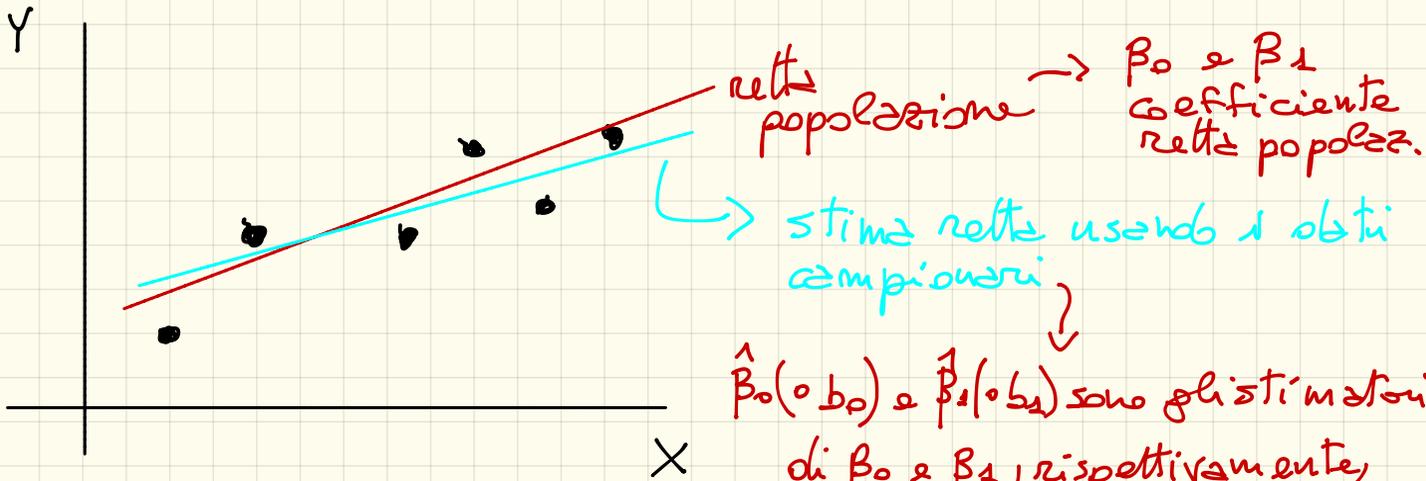
x $\mu_{Y|X=x_i}$
 medie condizionate

$Y | X = x_1$

Retta di regressione:
 luogo geometrico su cui
 giacciono le medie
 condizionate

RETTA REGRESSIONE POPOLAZIONE

le medie condizionate $Y | X = x_i$ cadono su una retta
 ↳ legame lineare tra X e Y



$\hat{\beta}_0$ ($\circ b_0$) e $\hat{\beta}_1$ ($\circ b_1$) sono gli stimatori di β_0 e β_1 , rispettivamente, usando il criterio dei minimi quadrati

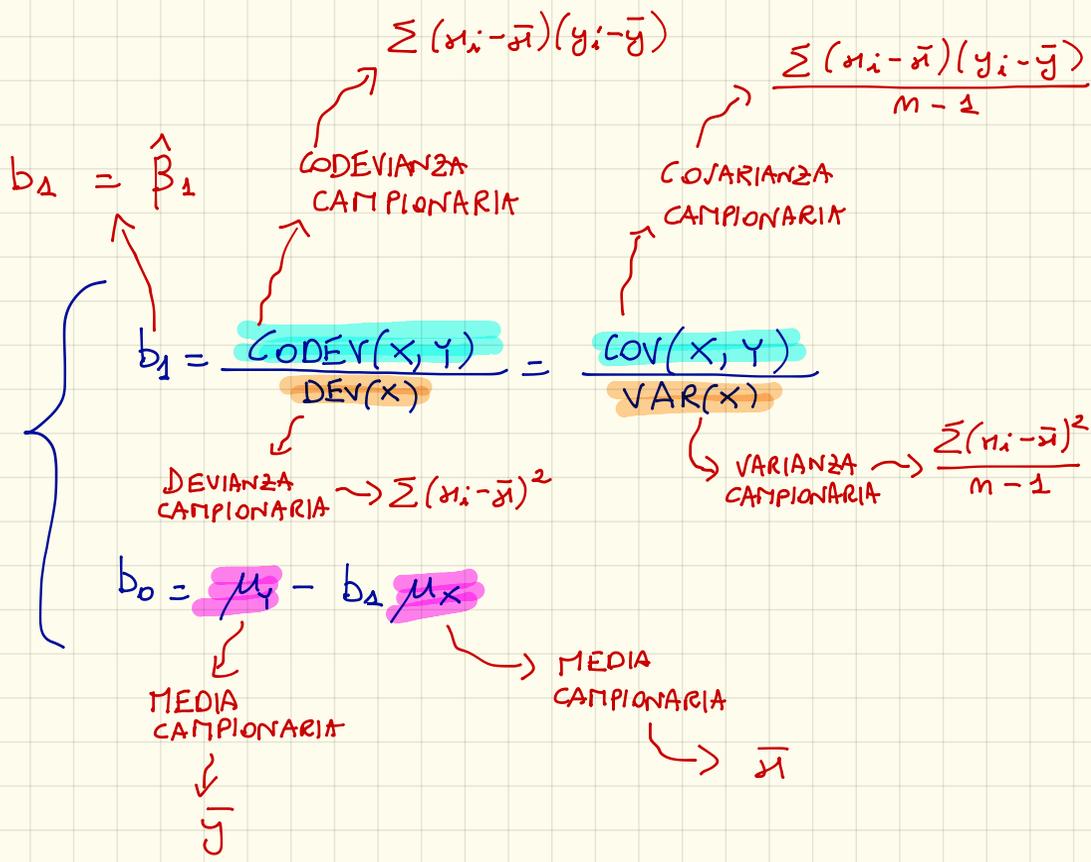
POPOLAZIONE : due parametri

$$\begin{array}{c}
 \beta_0 \\
 \downarrow \\
 b_0 \equiv \hat{\beta}_0
 \end{array}$$

$$\begin{array}{c}
 \beta_1 \\
 \downarrow \\
 b_1 \equiv \hat{\beta}_1
 \end{array}$$

stimatori usando il criterio di minimi quadrati

Come cambiano le formule nel passaggio ai dati campionari:



Soluzioni minimi quadrati (prima parte del corso)

IPOTESI (ASSUNZIONI) MODELLO CLASSICO DI REGRESSIONE

① Esiste un legame lineare tra X e Y

↳ le medie condizionate delle distribuzioni $Y|X=x_i$ giacciono su una retta

↳ $E(Y|X=x_i) =$ funzione lineare $= \beta_0 + \beta_1 x_i$

funzione di regressione

intercetta o termine noto

pendenza o coeff. angolare

se estendo il punto di vista alle intere distribuzioni condizionate $Y|X=x_i$

$$Y = \underbrace{\beta_0 + \beta_1 x_i}_{\text{parte lineare}} + \underbrace{\varepsilon}_{\text{parte casuale}}$$

deterministica

v.c. errore

↳ sintesi di tutto quello che spiega Y ma che non ho inserito nel modello

Su ε (r.c. errore) si fanno una serie di ipotesi

① $E(\varepsilon) = 0$

↳ dovrei scrivere $E(\varepsilon | X = x_i)$

non ci sono fattori sistematici che non ho considerato nel modello

↳ tecnicamente

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

$$E(Y | X = x) = E(\underbrace{\beta_0 + \beta_1 X}_{\text{deterministica}} + \underbrace{\varepsilon}_{\text{casuale}} | X = x) =$$

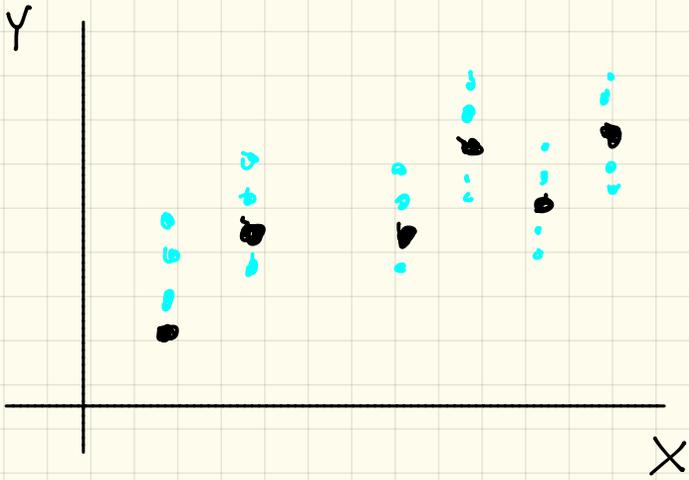
$$= E(\beta_0 + \beta_1 X | X = x) + E(\varepsilon | X = x) =$$

$$= \beta_0 + \beta_1 X \quad \underbrace{= 0}_{\text{sotto l'ipotesi ①}}$$

② $\text{Var}(\varepsilon) = \sigma_{\varepsilon}^2 = \sigma^2$ \rightarrow costante

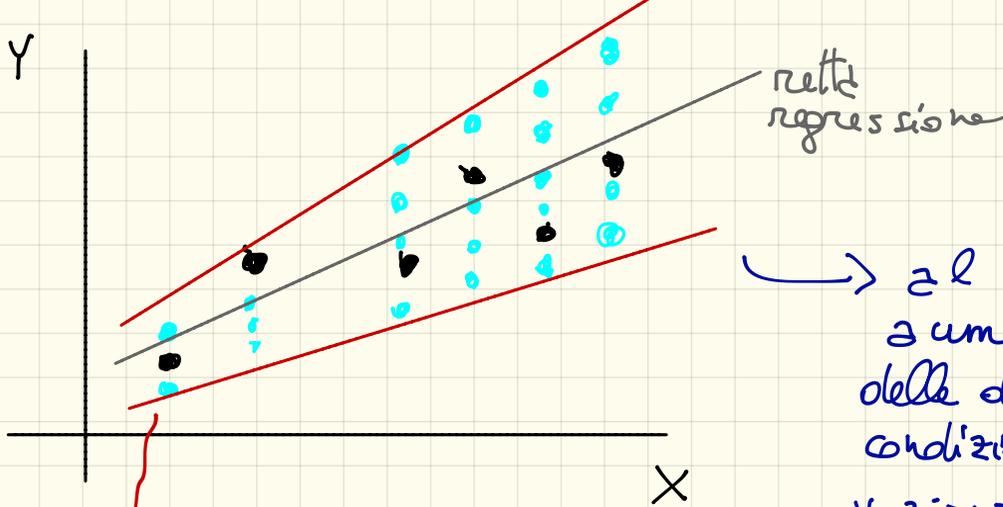
$\hookrightarrow \text{Var}(\varepsilon | X = x) = \sigma^2$ non dipende da X

ipotesi di omoschedasticità



\rightarrow le distribuz. condizionate $Y | X = x_i$ sono tutte caratterizzate dalla stessa variabilità (σ_{ε}^2)

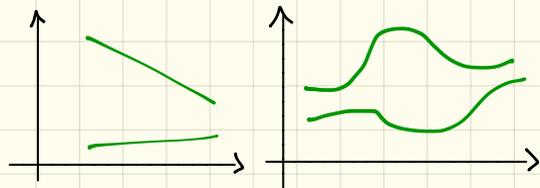
\hookrightarrow da un punto di vista pratico questo richiede di stimare un unico parametro (σ^2) per "governare" la v.c. ε



→ al crescere di X
 aumenta la varianza
 delle distribuzioni
 condizionate, ovvero la
 varianza di $\varepsilon|X=x$



esempio ETEROSCHEDASTICITA'



per questa
 distribuzione
 condizionate c'è
 minore probabilità
 di sbagliare usando il
 modello di regressione

NOTA: $\text{Var}(\varepsilon|X=x) = \sigma^2 \implies \text{Var}(Y|X=x) = \sigma^2$

③ le distribuzioni condizionate $Y|X=x$ sono incorrelate



ovvero le ε sono incorrelate

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$$

↳ più precisamente

$$\text{Cov}(\varepsilon | X=x_i, \varepsilon | X=x_j) = 0$$

Sotto le tre ipotesi del modello classico vale il teorema di

Gauss-Markov:

$$\hat{\beta}_0(b_0) \text{ e } \hat{\beta}_1(b_1)$$

sono BLUE estimators

best (efficienti in senso assoluto)
↳ unbiased (non distorto)
↳ linear

Se si è interessati solo alla stima puntuale di β_0 e β_1 , il teorema di Gauss-Markov, se valgono le ipotesi:

$$E(\varepsilon|X) = 0 \implies E(Y|X) = \beta_0 + \beta_1 X$$

$$\text{Var}(\varepsilon|X) = \sigma^2 \implies \text{Var}(Y|X) = \sigma^2$$

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$$

ci assicura che:

gli stimatori sono non distorti $\left\{ \begin{array}{l} E(\hat{\beta}_0) = \beta_0 \\ E(\hat{\beta}_1) = \beta_1 \end{array} \right.$ $\left. \begin{array}{l} \text{Var}(\hat{\beta}_0) = \text{minima} \\ \text{Var}(\hat{\beta}_1) = \text{minima} \end{array} \right\}$ gli stimatori sono efficienti in senso assoluto

Il teorema di Gauss-Markov non fornisce però alcuna informazione circa la distribuzione di $\hat{\beta}_0$ e $\hat{\beta}_1$:

$$\hat{\beta}_0 \sim (?)$$

$$\hat{\beta}_1 \sim (?)$$

Senza informazioni su queste distribuzioni campionarie, non è possibile costruire stime intervallari o effettuare test di ipotesi

Per procedere, è necessario modificare le ipotesi di base

↳ modello classico normale



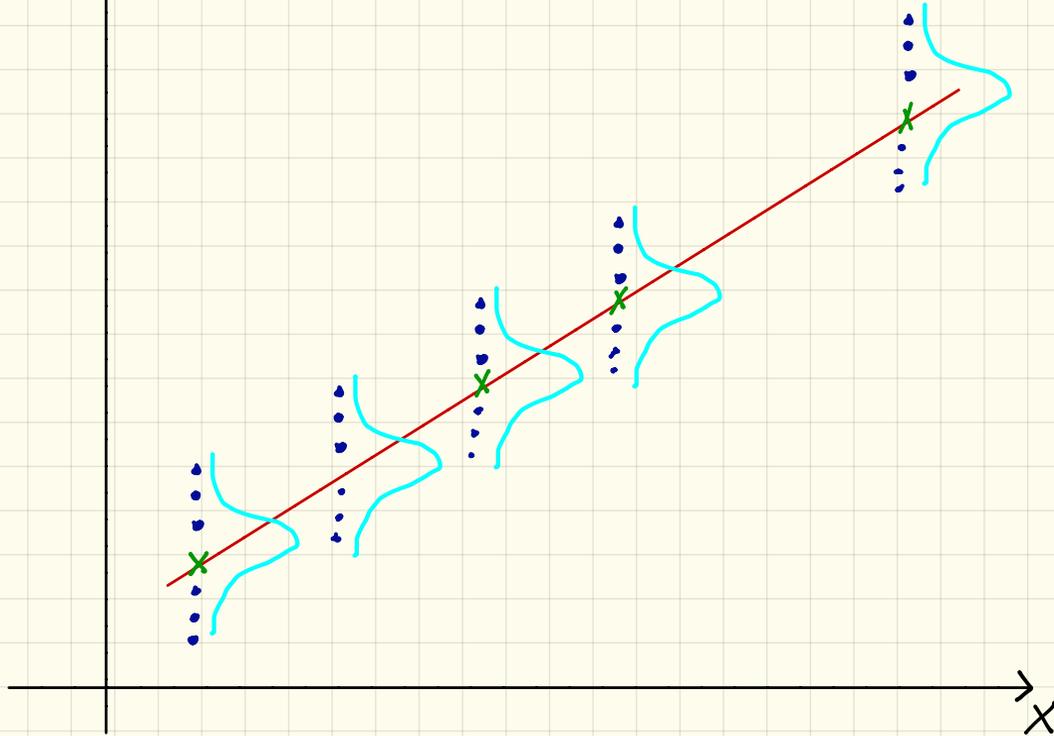
- si introduce l'ipotesi di normalità degli errori
- l'ipotesi di incorrelazione viene rafforzata in ipotesi di indipendenza



sotto questo quadro aggiornato si ha che:

- $Y \sim N$
- $\hat{\beta}_0$ e $\hat{\beta}_1$, essendo combinazioni lineari delle Y_i , sono a loro volta normali

introdurre l'ipotesi di normalità degli errori
equivale a dire che le distribuzioni condizionate
della Y sono normali



POPOLAZIONE $\leadsto Y = \beta_0 + \beta_1 X + \epsilon$

componenti deterministiche

modello classico

parametri

stimatori

$\hat{\beta}_0$ (oppure b_0)
 $\hat{\beta}_1$ (oppure b_1)

$E(Y|X) = \beta_0 + \beta_1 X$

$\text{Var}(Y|X) = \sigma^2$

- $E(\epsilon | X) = 0$
- $\text{Var}(\epsilon | X) = \sigma^2 = \sigma_\epsilon^2$
- ϵ_i e ϵ_j sono incorrel. $i \neq j$

TEOREMA DI GAUSS MARKOV
 gli stimatori dei minimi quadrati sono BLUE

non permette di ottenere gli intervalli di confidenza e i test per β_0 e β_1

modello classico normale

- $E(\epsilon | X) = 0$
- $\text{Var}(\epsilon | X) = \sigma^2 = \sigma_\epsilon^2$
- incorrelazione viene rafforzata (ϵ_i e ϵ_j sono indipendenti)

NUOVA $\leadsto \epsilon \sim N$

$\epsilon_i \sim N(0, \sigma^2)$
 con ϵ_i ed ϵ_j indipend. quando $i \neq j$

Sotto le ipotesi del modello classico normale si ha che $Y \sim N$



è possibile mostrare che $\hat{\beta}_0$ e $\hat{\beta}_1$ sono combinazioni lineari delle Y_i

↳ essendo le $Y_i \sim N$ allora $\hat{\beta}_0$ e $\hat{\beta}_1$ saranno distribuite normalmente

Se consideriamo $\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} = \frac{\sum (X_i - \bar{X}) Y_i}{\sum (X_i - \bar{X})^2}$

$\sum (X_i - \bar{X}) Y_i - \sum (X_i - \bar{X}) \bar{Y} = \sum (X_i - \bar{X}) Y_i - \bar{Y} \underbrace{\sum (X_i - \bar{X})}_{=0}$

essendo le X_i deterministiche nel modello classico si vede come $\hat{\beta}_1$ è una combinazione lineare delle Y_i

Lo stesso avviene per $\hat{\beta}_0$:

$\hat{\beta}_0 = \bar{y}_n - \hat{\beta}_1 \bar{x}_n = \sum \left[1 - \frac{(x_i - \bar{x}) \bar{y}}{\sum (x_i - \bar{x})^2} \right] y_i$ → anche $\hat{\beta}_0$ è esprimibile come combinazione lineare delle Y

Per definire completamente $\hat{\beta}_0$ e $\hat{\beta}_1$ ci servono le corrispondenti varianze:

$$\sigma_{\hat{\beta}_1}^2 = \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

varianza v.c. errore

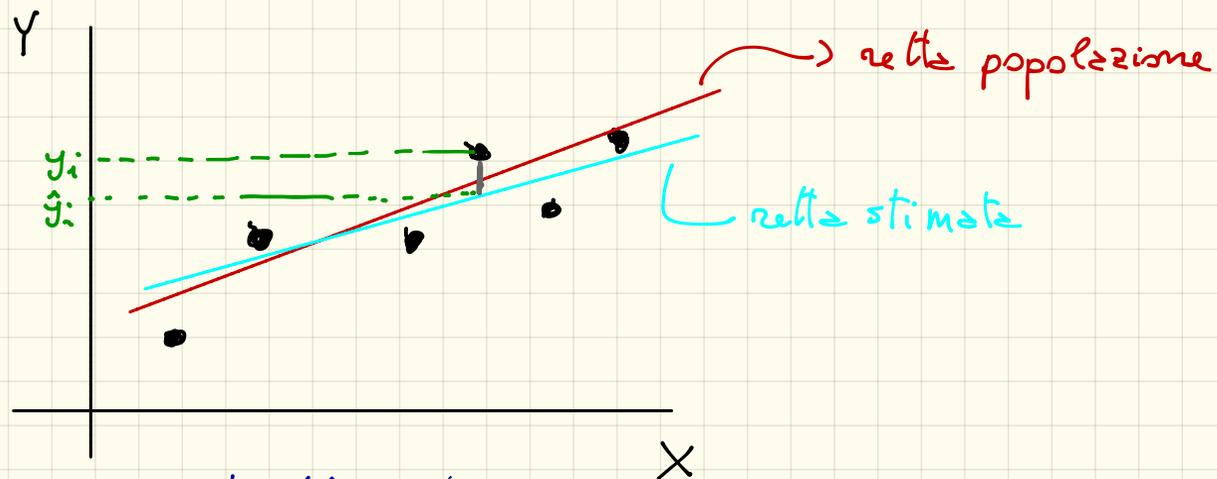
$(n-1) s_x^2$

$$\sigma_{\hat{\beta}_0}^2 = \text{Var}(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right] = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{(n-1) s_x^2} \right]$$

$$\left. \begin{aligned} \hat{\beta}_0 &\sim N\left(\beta_0, \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right]\right) \\ \hat{\beta}_1 &\sim N\left(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right) \end{aligned} \right\} \begin{aligned} \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right]}} &\sim N(0, 1) \\ \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}}} &\sim N(0, 1) \end{aligned}$$

Since σ^2 is not known we have to estimate it starting from data

samples: we use the residuals or rather $y_i - \hat{y}_i$ to estimate σ^2



NOTE: properties of residuals

The $e_i = y_i - \hat{y}_i$ are such that $\sum e_i = \sum (y_i - \hat{y}_i) = 0$

To estimate σ^2 consider the distances to the regression line:

$$\sum (y_i - \hat{y}_i)^2 = \text{Dev}(E)$$

$$\hat{\sigma}^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n-2} = \frac{\text{Dev}(\text{E})}{n-2}$$

stimatore per σ^2

si perde un grado di libert  per
ciascun parametro

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right]}} \sim t_{n-2}$$

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum (x_i - \bar{x})^2}}} \sim t_{n-2}$$

direttamente le statistiche pivotali
e le statistiche test

Invertendo le statistiche pivotali per β_0 e β_1 si hanno:

$$IC_{1-\alpha}(\beta_0) = \left[\hat{\beta}_0 \pm t_{n-2, \frac{\alpha}{2}} \underbrace{\sqrt{\text{Var}(\hat{\beta}_0)}}_{\text{errore standard}} \right] = \left[\hat{\beta}_0 \pm t_{n-2, \frac{\alpha}{2}} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{y}^2}{\sum (n_i - \bar{n})^2} \right]} \right]$$

$$IC_{1-\alpha}(\beta_1) = \left[\hat{\beta}_1 \pm t_{n-2, \frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\beta}_1)} \right] = \left[\hat{\beta}_1 \pm t_{n-2, \frac{\alpha}{2}} \sqrt{\frac{\hat{\sigma}^2}{\sum (n_i - \bar{n})^2}} \right]$$

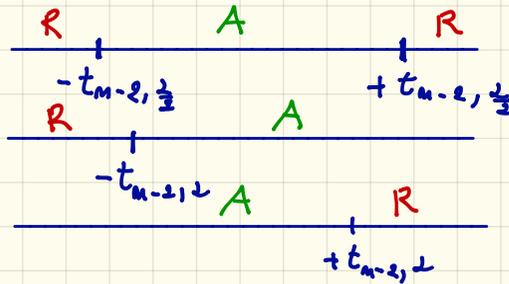
Le due statistiche test sono:

$$H_0: \beta_0 = \text{valore}$$

$$H_1: \beta_0 \neq \text{valore}$$

<

>



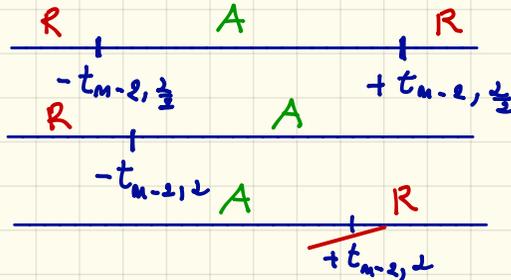
$$T_m = \frac{\hat{\beta}_0 - \beta_0 \text{ sotto } H_0}{\sqrt{\text{Var}(\hat{\beta}_0)}}$$

$$H_0: \beta_1 = \text{valore}$$

$$H_1: \beta_1 \neq \text{valore}$$

<

>



$$T_m = \frac{\hat{\beta}_1 - \beta_1 \text{ sotto } H_0}{\sqrt{\text{Var}(\hat{\beta}_1)}}$$

Di solito si utilizza il valore 0 come valore del parametro sotto H_0 :

$$H_0: \beta_0 = 0$$

$$H_1: \beta_0 \neq 0$$

$\hat{>}$

$$T_m = \frac{\hat{\beta}_0 - 0}{\sqrt{\text{var}(\hat{\beta}_0)}} = \frac{\hat{\beta}_0}{\sqrt{\text{var}(\hat{\beta}_0)}}$$

↳ test significatività dell'intercetta

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

$\hat{>}$

$$T_m = \frac{\hat{\beta}_1}{\sqrt{\text{var}(\hat{\beta}_1)}}$$

↳ test significatività intercetta

↳ nel modello semplice equivale a testare la significatività del modello

FORMULE ALTERNATIVE CALCOLO R^2

$$R^2 = \frac{\text{Dev}(R)}{\text{Dev}(Y)} = 1 - \frac{\text{Dev}(E)}{\text{Dev}(Y)} = \hat{\rho}_{xy}^2$$

coefficiente di correlazione campionario

Si ha anche che:

$$\begin{aligned} R^2 &= \hat{\rho}_{xy}^2 = \left[\frac{\text{CODEV}(X, Y)}{\sqrt{\text{Dev}(X)} \sqrt{\text{Dev}(Y)}} \right]^2 = \\ &= \frac{\text{Cov}(X, Y)^2}{\text{Dev}(X) \text{Dev}(Y)} \frac{\text{Dev}(X)}{\text{Dev}(X)} = \left[\frac{\text{CODEV}(X, Y)}{\text{Dev}(X)} \right]^2 \frac{\text{Dev}(X)}{\text{Dev}(Y)} = \\ &= \hat{\beta}_1^2 \frac{\text{Dev}(X)}{\text{Dev}(Y)} \end{aligned}$$

Dalle relazioni sopra esposte si ha quindi:

$$R^e = \frac{\text{Dev}(R)}{\text{Dev}(Y)} = \hat{\beta}_{xy}^e \implies \text{Dev}(R) = \hat{\beta}_{xy}^e \text{Dev}(Y)$$

$$R^e = \frac{\text{Dev}(R)}{\text{Dev}(Y)} = \hat{\beta}_1^e \frac{\text{Dev}(X)}{\text{Dev}(Y)} \implies \text{Dev}(R) = \hat{\beta}_1^e \text{Dev}(X)$$

E' possibile cioè calcolare la $\text{Dev}(R)$, e quindi anche la $\text{Dev}(E) = \text{Dev}(Y) - \text{Dev}(R)$, a partire da $\hat{\beta}_{xy}^e$ e $\text{Dev}(Y)$ oppure a partire da $\hat{\beta}_1^e$ e $\text{Dev}(Y)$

↳ questo permette di evitare il calcolo di $\sum (y_i - \hat{y}_i)^2$ e $\sum (\hat{y}_i - \mu_Y)^2$ che può essere laborioso

