

6 Modello di regressione lineare semplice

↓
 considera una variabile dipendente Y (= **VARIABILE RISPOSTA**)
 in funzione lineare di una sola variabile
 indipendente x (= **VARIABILE ESPLICATIVA**)

$$f(x) = \underbrace{\beta_0}_{\text{INTERCETTA}} + \underbrace{\beta_1 x}_{\text{COEFFICIENTE ANGOLARE}}$$

coefficienti di regressione

RELAZIONE STATISTICA

$Y = f(x) + \varepsilon$ → **TERMINE D'ERRORE** = rappresenta il contributo di tutti gli altri fattori ($\neq x$) non misurabili che influenzano Y

↓
 è una **VARIABILE CASUALE**

↓
 ed è conseguenza anche y
 è una **VARIABILE CASUALE**

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

non conosciamo questi parametri
 quindi dobbiamo **STIMARLI**

↓
 per farlo dobbiamo fare delle **ASSUNZIONI**

RESIDUI

$$e_i = y_i - \hat{y}_i$$

valore reale

valore stimato dal modello

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot x_i$$

PROPRIETA'

$$\sum e_i = 0$$

la somma di tutti i residui deve essere zero!!

ASSUNZIONI

1. La funzione di regressione deve essere lineare

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad i=1, 2, \dots, n$$

2. Le ε_i sono v.c. INDIPENDENTI TRA LORO e indipend. dal valore di x_i

$$E(\varepsilon_i) = 0 \quad \text{varianza costante} \rightarrow \text{OMOSCEDASTICITA'}$$

in media gli errori commessi si compensano

3. I valori x_i della variabile esplicativa X sono noti senza errori.

METODO DEI MINIMI QUADRATI

→ metodo di stima dei coeff. di regressione β_0 e β_1 che consiste nel cercare le stime che MINIMIZZANO la FUNZIONE DI PERDITA

escludi dagli errori della retta

$$= \frac{\sum xy - n \cdot \bar{x} \cdot \bar{y}}{\sum (x_i - \bar{x})^2} \rightarrow \text{DEVIANZA}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum xy}{\sum x^2}$$

$$= \frac{\sum xy / \sum x}{\sum y / \sum x}$$

sono stimatori CORRETTI

$$E(\hat{\beta}_0) = \beta_0 \quad E(\hat{\beta}_1) = \beta_1$$

ed EFFICIENTI (varianza minima)

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{var}(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right] \quad \text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

PROPRIETA'

1. La retta di regressione stimata passa sempre per il BARICENTRO $(\bar{x}; \bar{y})$

$$\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n y_i$$

$$E(y_i | x = x_i) = \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$4. \quad SQT = SQR + SQE$$

DECOMPOSIZIONE DELLA DEVIANZA TOTALE

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2$$

somma totale dei quadrati

$$\text{DEV}(Y) = \sum y^2 - n \cdot \bar{y}^2$$

somma dei quadrati della regressione (DEVIANZA SPIEGATA) (VARIABILITA' DEL MODELLO)

somma dei quadrati degli errori (DEVIANZA RESIDUA) (VARIABILITA' DEGLI ERRORI!!!)

COEFFICIENTE DI DETERMINAZIONE → MISURA LA BONTÀ DI ADATTAMENTO DI REGRESSIONI AI DATI

significa che del totale delle variabili $R^2_{xy} = \frac{SQR}{SQT} = 1 - \frac{SQE}{SQT}$

$R^2_{xy} = (r_{xy})^2$ visto che ci stiamo riferendo ad un modello semplice R^2_{xy} può essere misurato anche come quadrato del coeff di correlazione lineare

campo di variazione $[0; 1]$
 rel. perf. (SQE=0)
 rel. lineare tra x e y (SQT=0)

$r_{xy} = \frac{COV(X,Y)}{\sigma_X \cdot \sigma_Y}$
 $R^2 = \left[\frac{n \cdot COVAR(X,Y)}{\sqrt{n \cdot \sigma_X} \cdot \sqrt{n \cdot \sigma_Y}} \right]^2 = \left[\frac{CODEV(X,Y)}{\sqrt{DEV(X)} \cdot \sqrt{DEV(Y)}} \right]^2$

PROPRIETÀ DELLO STIMATORE \hat{Y}_i

$E(\hat{Y}_i) = \beta_0 + \beta_1 x_i$ → STIMATORE CORRETTO

$V(\hat{Y}_i) = \sigma^2_{\epsilon} \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_j - \bar{x})^2} \right]$

ERRORE STANDARD DI REGRESSIONE → SCARTO QUADRATICO MEDIO DELLA STIMA

$S = \sqrt{S^2} = \sqrt{\frac{\sum \hat{e}_i^2}{n-2}}$ → stima corretta della varianza dei residui σ^2_{ϵ}

ASSUNZIONE DI NORMALITÀ DEGLI ERRORI

$\epsilon_i \sim N(0; \sigma^2)$

↓
 I SINGOLI ERRORI SONO INDIPENDENTI TRA LORO

$E(\epsilon_i \epsilon_j) = 0 \quad \forall i \neq j$

di conseguenza l'intero modello di regressione si distribuisce normalmente

$Y_i \sim N(\beta_0 + \beta_1 x_i; \sigma^2)$

www.unidocs.it - Appunti e dispense per superare i tuoi esami universitari

www.unidocs.it - Appunti e dispense per superare i tuoi esami universitari

DISTRIBUZIONE DEGLI STIMATORI $\hat{\beta}_0$ e $\hat{\beta}_1$

1. $\hat{\beta}_0 \sim N(\beta_0; \sigma_{\hat{\beta}_0}^2)$ $\hat{\beta}_1 \sim N(\beta_1; \sigma_{\hat{\beta}_1}^2)$

insieme assumono le forme di DISTRIBUZIONE NORMALE BIVARIATA

2. $z_0 = \frac{\hat{\beta}_0 - \beta_0}{S(\hat{\beta}_0)} \sim \mathcal{N}_{n-2}$ $z_1 = \frac{\hat{\beta}_1 - \beta_1}{S(\hat{\beta}_1)} \sim \mathcal{N}_{n-2}$

stima della SD

quando $n \rightarrow \infty$ le distri. convergono ad una $N(0, 1)$

INTERVALLI DI CONFIDENZA PER I PARAMETRI

$$P\left(\hat{\beta}_0 - t_{\frac{\alpha}{2}} \cdot S(\hat{\beta}_0) \leq \beta_0 \leq \hat{\beta}_0 + t_{\frac{\alpha}{2}} \cdot S(\hat{\beta}_0)\right) = 1 - \alpha$$

$$P\left(\hat{\beta}_1 - t_{\frac{\alpha}{2}} \cdot S(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + t_{\frac{\alpha}{2}} \cdot S(\hat{\beta}_1)\right) = 1 - \alpha$$

VERIFICA DI IPOTESI PER I PARAMETRI

se vogliamo verificare $H_0: \beta_1 = b_1$

sotto H_0 : $\frac{\hat{\beta}_1 - b_1}{S(\hat{\beta}_1)} \sim \mathcal{N}_{n-2}$

TAVOLA DELL'ANALISI DELLA VARIANZA (ANOVA)

la nitrova uguale in JMP!

SORGENTE DI VARIAZIONE	SOMMA DEI QUADRATI	GDL	MEDIA DEI QUADRATI	F
REGRESSIONE	SQR	1	$MQR = \frac{SQR}{1}$	$F = \frac{MQR}{MGE}$
RESIDUO	SQE	$n-2$	$MGE = \frac{SQE}{n-2}$	

Statistiche test dell'ANOVA
↓
si distribuisce come una F di Fisher

il suo P-value è uguale a quello per β_1

SQR
 $n-1$ gdl.

nel modello di regressione lineare semplice
il test statistico per la verifica di ipotesi su β_1