

REGRESSIONE LINEARE SEMPLICE

Un modello di regressione lineare semplice spiega una variabile risposta (o dipendente) y in termini di una variabile X detta var. indipend. o regressore o predittore.

La più semplice relazione tra y e X è lineare:

$$y = \alpha + \beta X \quad (1)$$

Ovviamente, noto X , non possiamo prevedere esattamente y , ma osserviamo (1) vero e meno di un errore ε .

Che caratteristiche ha ε ?

Non essendo osservabile, lo consideriamo una v.a. con determinate caratteristiche:

$$\varepsilon \sim N(0, \sigma^2) \quad \text{iid (indipendenti e identicamente distribuiti)}$$

Quindi di fatto, y è v.a.:

$$y = \alpha + \beta X + \varepsilon$$

$$E[y] = \alpha + \beta E[X] + E[\varepsilon] = \alpha + \beta E[X]$$

$$\text{Var}[y] = 0 + \beta^2 \text{Var}[X] + \sigma^2 = \beta^2 \text{Var}[X] + \sigma^2$$

Possiamo calcolare tutto purché X e y sono osservate.
In particolare X è deterministica

α = intercetta = $E[y]$ quando non ho x ($x=0$)

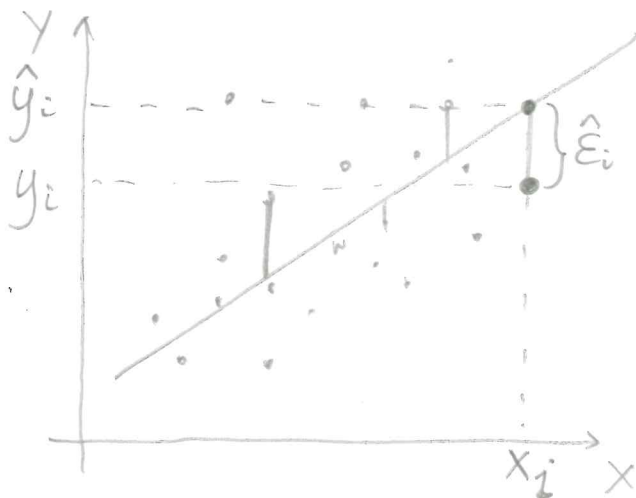
β = coeff. angolare = rappresenta l'incremento di y al crescere di un'unità di x

Nel modello, avendo n osservazioni di y e x accoppiate: $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$

Cerchiamo α e β t.c. minimizziamo ^{la somma di} $|y_i - (\alpha + \beta x_i)| = \varepsilon_i$.
Ma visto che alcuni potrebbero essere negativi, minimizziamo $\sum_{i=1}^n \varepsilon_i^2 \rightarrow$ troviamo $\hat{\alpha}$ e $\hat{\beta}$, le stime dei 2 parametri

$$\Rightarrow \hat{y} = \hat{\alpha} + \hat{\beta}x$$
$$y_i = \hat{\alpha} + \hat{\beta}x_i + \hat{\varepsilon}_i \Rightarrow \hat{y}_i - y_i = \hat{\varepsilon}_i$$

METODO DEI MINIMI QUADRATI ↑



• DIAGNOSTICA DEL MODELLO

$$\varepsilon \sim N(0, \sigma^2)$$

istogramma
qqplot

viene per
costruzione

costante $\forall i$
→ omoschedasticità
grafico dei residui
(standardizzati)

ε_i iid

→ dipende dai dati, che
si assumono indipendenti

$$(X_i, y_i) \perp (X_j, y_j) \quad \forall i \neq j$$

• BONTÀ DEL MODELLO

R^2 , p-values regressori

↓

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

Cosa vuol dire che gli ε_i sono iid?

Def iid: $P(\varepsilon_1 \leq e_1, \varepsilon_2 \leq e_2, \dots, \varepsilon_N \leq e_N) = \prod_{i=1}^N P(\varepsilon_i \leq e_i)$