

LEZIONE N.8 (a cura di Teresa Fanelli)

- Forma matriciale del Modello di Regressione Semplice

L'assunzione di base del modello è: $Y_i = \beta_0 + x_i \beta_1 + \varepsilon_i$ $i=1,2,\dots,n$. Lo stesso modello può essere scritto attraverso vettori e matrici, si ha: $Y = X\beta + \varepsilon$, dove:

- Y è un vettore di dimensione $n \times 1$ $\rightarrow Y_{n \times 1} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$

- X è una matrice di dimensione $n \times 2$ $\rightarrow X_{n \times 2} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$

- β è il vettore dei PARAMENTRI di dimensione 2×1 $\rightarrow \beta_{2 \times 1} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$

- ε è un vettore di dimensioni $n \times 1$ $\rightarrow \varepsilon_{n \times 1} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$

Questa forma risulta importante nel modello di regressione con più variabili.

Risolvendo tutti i passaggi ,prima il prodotto:

$$X_{n \times 2} \times \beta_{2 \times 1} = \begin{pmatrix} \beta_0 + x_1 \beta_1 \\ \beta_0 + x_2 \beta_1 \\ \vdots \\ \beta_0 + x_n \beta_1 \end{pmatrix} \left. \vphantom{\begin{pmatrix} \beta_0 + x_1 \beta_1 \\ \beta_0 + x_2 \beta_1 \\ \vdots \\ \beta_0 + x_n \beta_1 \end{pmatrix}} \right\} \begin{array}{l} \text{il vettore della COMPONENTE} \\ \text{DETERMINISTICA, di} \\ \text{dimensione } n \times 1 \end{array}$$

E aggiungendo il termine di errore al vettore ottenuto ,si ha:

$$Y = X\beta + \varepsilon = \begin{pmatrix} \beta_0 + x_1 \beta_1 + \varepsilon_1 \\ \beta_0 + x_2 \beta_1 + \varepsilon_2 \\ \vdots \\ \beta_0 + x_n \beta_1 + \varepsilon_n \end{pmatrix}$$

Riprendendo le IPOTESI sul termine di errore,quali:

1. $E(\varepsilon_i) = 0 \quad \forall i$;
2. $V(\varepsilon_i) = \sigma^2 \quad \forall i$;
3. $\varepsilon_i \perp \varepsilon_j \quad \forall i \neq j$ dove \perp è il simbolo di indipendenza:

- in alcune considerazioni si ha che $\varepsilon_i \sim N(0, \sigma^2) \quad \forall i$;

risulta importante vedere come queste ipotesi vengono applicate al vettore ε . Si ha:

1. $E(\varepsilon) = 0$, ciò vuol dire che ε è un vettore con tutti elementi nulli, di dimensione $n \times 1$;
2. $V(\varepsilon) = \sigma^2 I$ dove σ^2 è la varianza comune e I è la matrice identità = $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$, da cui si ottiene

$$\rightarrow \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & \\ \dots & \dots & \sigma^2 & \dots \\ 0 & \dots & \dots & \sigma^2 \end{pmatrix}$$

Questa matrice corrisponde alla matrice varianza-covarianza, dove la diagonale principale è data dalle varianze e gli 0 corrispondono alle covarianze.

3. $\varepsilon \sim N(0, \sigma^2 I)$ questo implica che gli errori sono INDIPENDENTI. Tale ipotesi sintetizza le altre ipotesi viste per ε . Nel modello di regressione con più variabili basta utilizzare solo quest'ipotesi.

- Stimatori

Se si conosce β , vuol dire che si dispone del vettore $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$, si ha:

$\hat{y} = x\beta$ è un vettore di dimensione $n \times 1$, dove \hat{y} è il vettore delle previsioni che si possono calcolare per ogni soggetto. La differenza tra il vettore dei dati osservati e il vettore delle previsioni, dà il vettore degli errori

$$\text{previsti di dimensioni } n \times 1 \rightarrow y - \hat{y}_{n \times 1} = \begin{pmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{pmatrix}$$

sulla base di questo vettore si può definire la somma dei quadrati: $S(\beta) = (y - \hat{y})'_{1 \times n} (y - \hat{y})_{n \times 1}$.

Il risultato di questa operazione è uno scalare, si ha:

$$(y_1 - \hat{y}_1, y_2 - \hat{y}_2, \dots, y_n - \hat{y}_n) \begin{pmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{pmatrix} = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2 = \sum (y_i - \hat{y}_i)^2$$

In forma esplicita: $S(\beta) = (y - \beta)'(y - x\beta)$, da cui si ottiene $\hat{\beta}_{2 \times 1} = (x'x)^{-1}_{(2 \times 2)} x'y_{(2 \times 1)}$ → è lo stimatore dei minimi

quadrati, ed è equivalente alle formule viste precedentemente = $\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \bar{y} - \hat{\beta}_1 \bar{x} \\ \frac{Cov(x, y)}{D(x)} \end{pmatrix}$

→ quindi, concludendo, ci sono due modi per poter stimare β_0 e β_1 :

1. Con le formule più lunghe;
2. Attraverso $(x'x)^{-1}(x'y)$

- Diagnostica

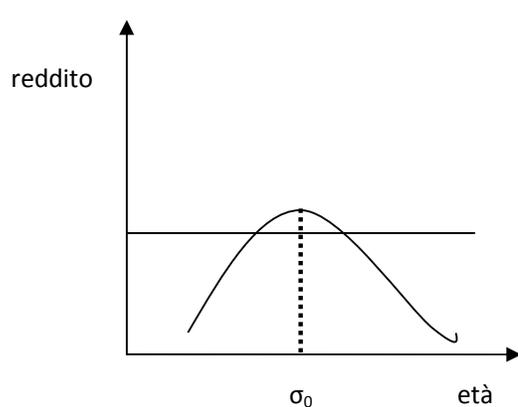
il modello di regressione semplice si basa su due assunzioni:

1. Parte deterministica $\rightarrow E(y_i) = \beta_0 + \beta_1 x_i$, ciò vuol dire che il valore atteso di y è una funzione lineare di x .

In alcuni casi risulta essere più complesso, del tipo: $E(y_i) = \beta_0 + \beta_1 x_i^2$

o anche: $E(y_i) = \beta_0 + \beta_1 x_i^2 + \beta_2 x_i^3$.

Queste situazioni sono presenti soprattutto nell'ambito economico, ad esempio nella relazione tra reddito-età, graficamente:

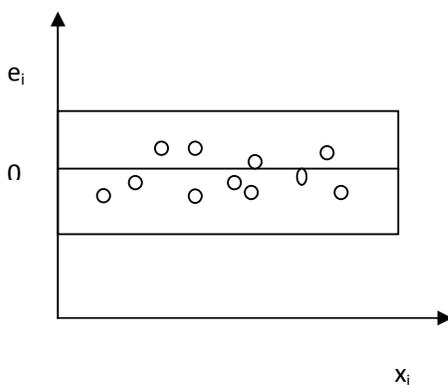


In questo caso il modello di regressione non è adeguato, perché la relazione è di tipo quadratica: $y \uparrow$ e poi \downarrow ; ciò porta a dire che l'età non influenza il reddito il che è del tutto sbagliato.

Se le assunzioni non sono vere, attraverso l'inferenza, le conclusioni possono essere sbagliate, questo perché il modello è del tipo lineare. Per capire se il modello è adeguato vengono utilizzati due indici:

1. $0 \leq r^2 \leq 1$;
2. Statistica F ($H_0: \beta_i = 0$)

Per avere un'analisi più completa si utilizzano i residui di regressione: cioè quantità calcolate in corrispondenza di ogni stima $\rightarrow e_i = y_i - \hat{y}_i$ (differenza tra valore osservato e valore stimato). Il metodo più semplice è rappresentato dal grafico dei residui: per ogni soggetto si ha (x_i, e_i) :

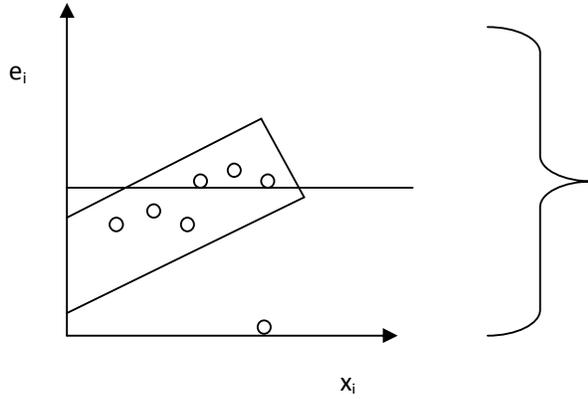


-Se le assunzioni sono vere i punti si concentrano intorno allo zero formando così un rettangolo; questo perché $\frac{1}{n} \sum e_i = 0$, ciò vuol dire che ci sono valori positivi negativi che si compensano tra di loro.

-Se ci sono problemi sulle assunzioni i punti portano

Risulta,così importante riconoscere alcuni grafici:

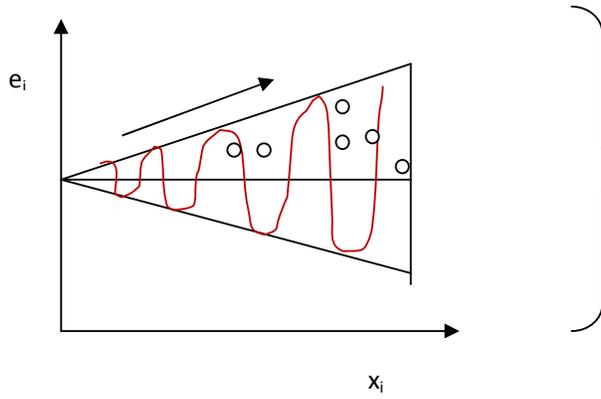
1. Problema:componente deterministica non specificata:



-I residui sono intorno allo zero,ma primo sotto e poi diventano positivi. Ciò vuol dire che il modello non è stato specificato nella forma lineare.

-Esiste,quindi una relazione sistematica tra x e e_i .

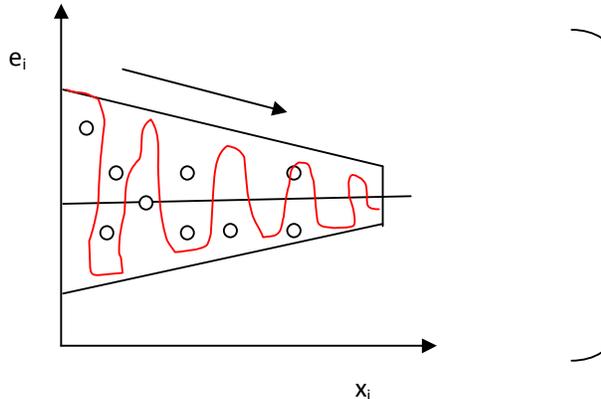
2. Problema:varianza costante(eteroschedasticità):



La varianza degli errori aumenta con x . L'andamento della varianza si nota dalle oscillazioni:

-oscillazione piccola,varianza bassa,viceversa nel caso contrario.

In questo caso non sono valide le conclusioni relative alla verifica di ipotesi e all'intervallo di confidenza.



Caso contrario: σ^2 diminuisce con la x .

-assunzioni non rispettate;

-è valida sempre l'eteroschedasticità.

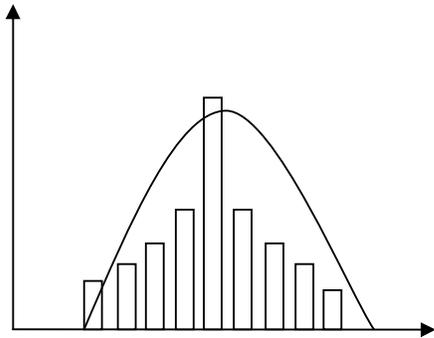
A volte può accadere di avere grafici in cui al posto di x_i vi è y_i ; ciò succede perché x_i è una trasformazione lineare di y_i e viceversa. Lo stesso accade per il residuo standardizzato $e^* = \frac{e_i}{\sqrt{s^2}}$.

Un'altra possibilità è quella di verificare la normalità; questo è possibile attraverso diversi modi:

1.

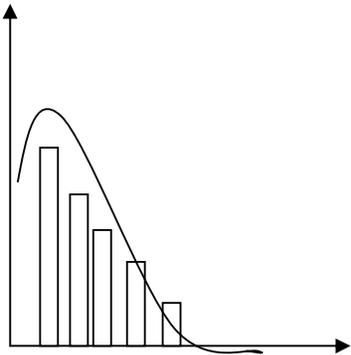
Modo: utilizz

o dell'istogramma:

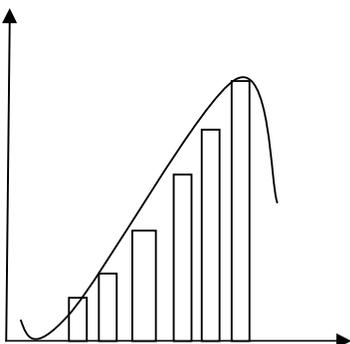


se le assunzioni sono rispettate l'istogramma avrà la forma di una normale; questa può essere verificata sovrapponendo la curva della normale.

Tuttavia può accadere che queste assunzioni non vengano rispettate, si avranno così grafici differenti:



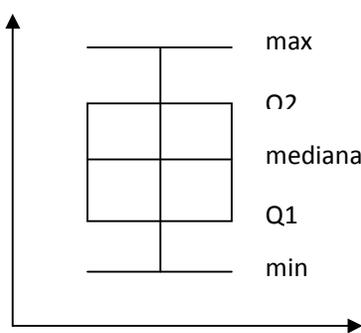
- Distribuzione asimmetrica verso destra;
- si hanno delle conseguenze sulla verifica d'ipotesi e sugli intervalli di confidenza;
- gli errori sono pochi, ma molto forti



- Caso opposto a quello precedente:
- asimmetria negativa;
 - pochi errori e molto piccoli.

2.

Modo: box plot(ϵ_i):

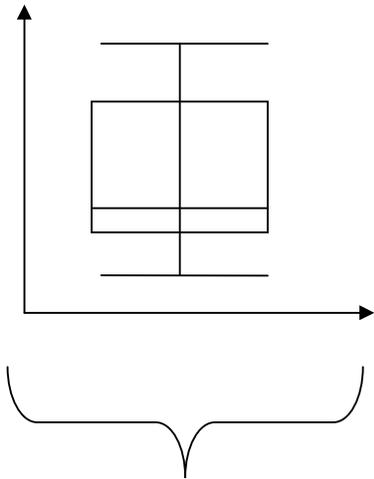


Se è normale, si ha:

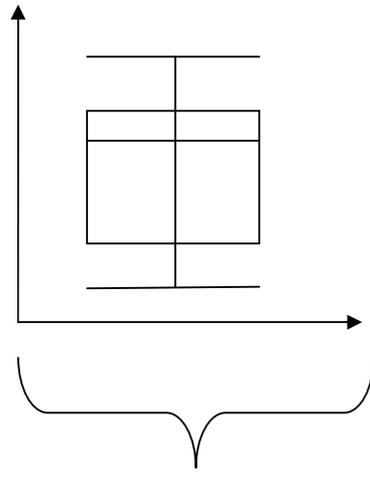
-mediana al centro;

-scatola tra il max e il min.

Possono verificarsi delle situazioni in cui non si ha la normale:

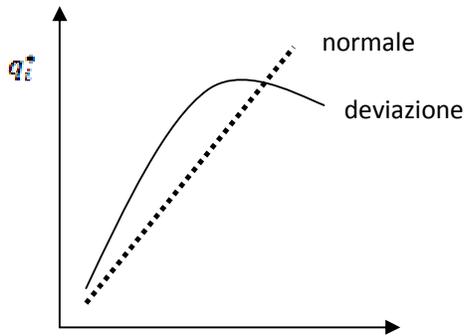


Mediana vicino al 1° quantile



Mediana vicino al 2° quantile

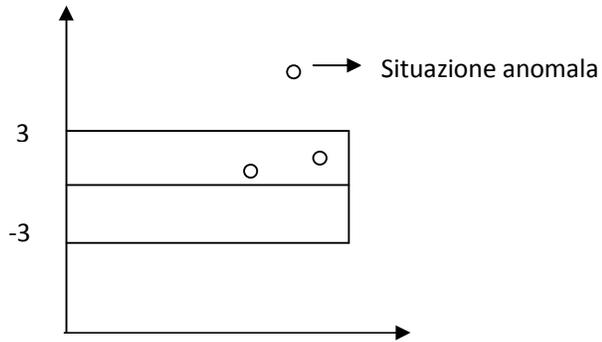
3. Modo: Q-Qplot: rappresentazione dei quantili della distribuzione di ϵ_i , con una percentuale tra 1% e 99%, e dei quantili teorici della normale (q_i, q_i^*):



Se c'è normalità, i punti si collocano su una retta.

Ogni volta che i punti deviano dalla retta, la condizione di normalità non è rispettata.

È possibile verificare che nel grafico ci siano delle osservazioni anomale, dette outlier; non sono altro che osservazioni che hanno una relazione diversa dagli altri punti. Questo è possibile attraverso un grafico:



Si ha un outlier quando un punto si trova fuori dalla fascia (-3,+3).

È importante standardizzare. Residui standardizzati perdono l'unità di misura; il vantaggio è che per ogni data set si ha la stessa soglia.

È importante individuare outlier, poiché possono esserci delle omissioni.