

Lezione 9

La regressione lineare semplice

PROF. ROBERTO COSTA

SCIENZE DELL'EDUCAZIONE - STATISTICA SOCIALE (305SF)

Dove eravamo rimasti

Nella scorsa lezione abbiamo visto come possiamo quantificare le relazioni tra due variabili, quando abbiamo a disposizione i dati in forma tabellare (tabella di distribuzione congiunta).

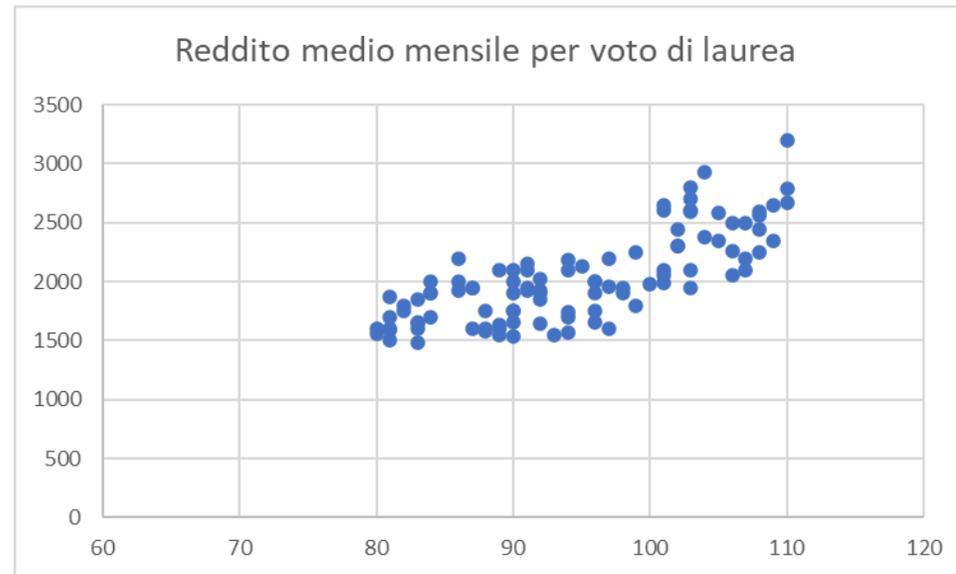
Abbiamo imparato a calcolare diversi indici: chi quadrato e phi quadro di Pearson e V di Cramer.

Ci sono dei dubbi?

La regressione

Partiamo dal comprendere che cosa significa, in statistica, il termine **regressione**.

La regressione è un indicatore statistico che indica l'esistenza o meno di una relazione significativa tra due (analisi bivariata) o più variabili (analisi multivariata) quantitative.



Cosa intendiamo per regressione

In statistica, il termine regressione è stato utilizzato per la prima volta dal biologo inglese Francis Galton nel 1886, quando parlò di «regressione verso la media».

Nell'ambito dei suoi studi sull'ereditarietà dei caratteri, Galton raccolse le stature di 928 figli adulti e dei loro 205 genitori (maschi e femmine). Esaminando le altezze di genitori e figli, notò una relazione tra le due variabili: più alti erano i genitori, più alti erano i figli e viceversa.

Partendo dalla statura media dei genitori ('mid-parent's stature') scoprì che i figli più alti della media avevano genitori ancora più alti di loro e i figli più bassi della media avevano genitori ancora più bassi. A questo fenomeno diede il nome di regressione verso la media.

Ripartiamo da un esempio

Partiamo da un caso concreto con due variabili (dipendente e indipendente) di tipo quantitativo.

Possiamo chiederci che relazione esiste tra il numero di ore dedicate alla preparazione di un esame e il voto conseguito all'esame stesso.

Oppure potremmo chiederci se esiste una relazione tra il voto conseguito all'esame di laurea e il reddito mensile medio.

In entrambi i casi abbiamo definito una variabile dipendente (il voto all'esame e il reddito mensile medio) e una variabile indipendente (il numero di ore dedicate alla preparazione di un esame e il voto conseguito all'esame di laurea).

Il diagramma di dispersione

Il **diagramma di dispersione** è la rappresentazione grafica di una possibile relazione tra due variabili.

Sull'asse X troviamo la variabile indipendente e sull'asse Y la variabile dipendente.

L'insieme dei punti che si crea indica come covariano (variano insieme) le due variabili.

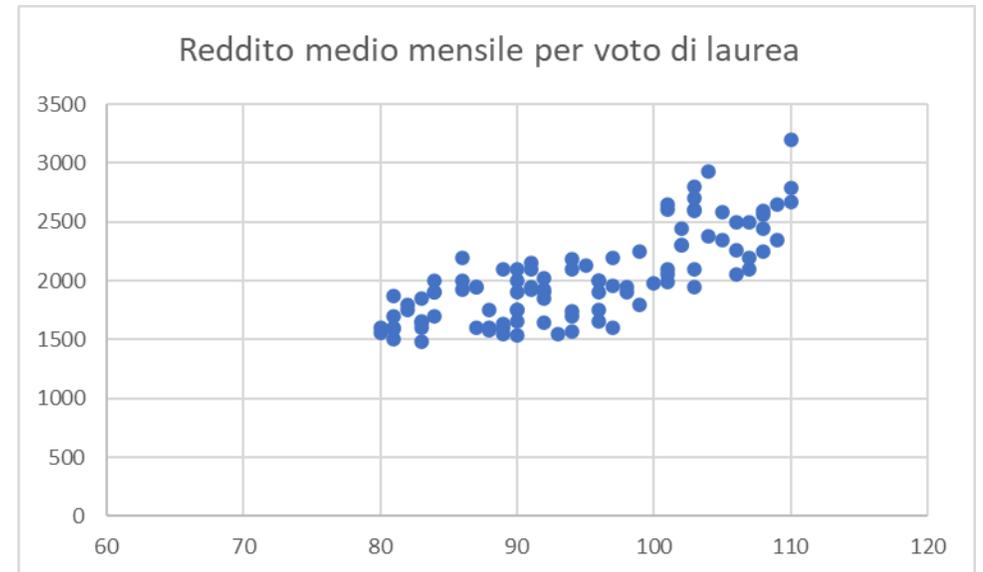
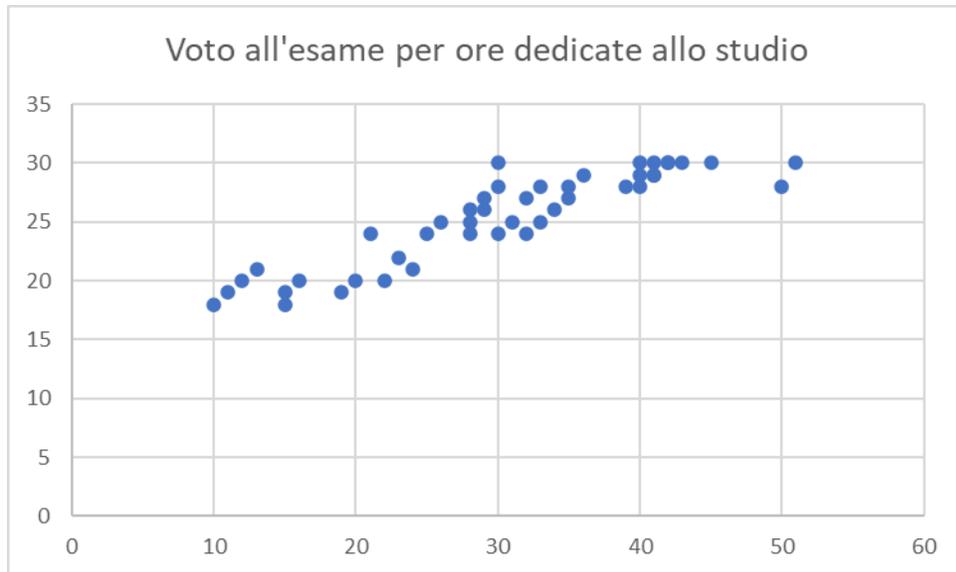
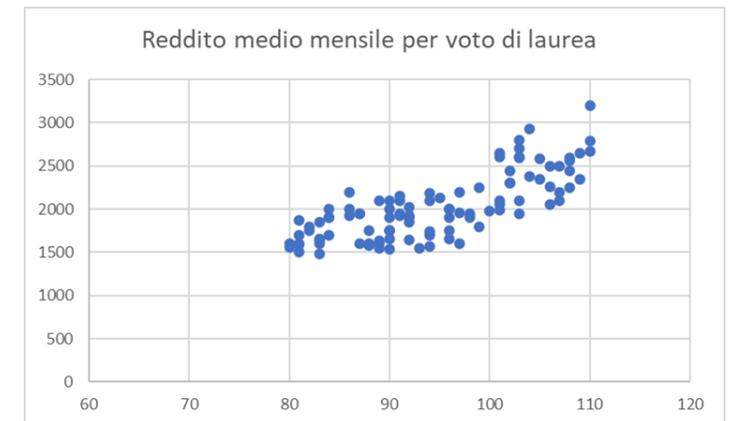


Diagramma di dispersione

L'osservazione del diagramma di dispersione ci consente di trarre alcune conclusioni:

- Le due variabili covariano in modo sistematico, ovvero sono legate da una relazione,
- La distribuzione della «nuvola di punti» dalla sinistra in basso alla destra in alto ci indica che la relazione è positiva, ovvero al crescere della variabile indipendente cresce anche la variabile dipendente.
- La disposizione dei punti ci suggerisce che siamo in presenza di una relazione lineare, ovvero la variabile Y tende a variare sempre nella stessa direzione e nella stessa misura al variare di X.

Non ci aiuta a misurare l'effetto causale, ovvero a quantificare la variazione della variabile dipendente al variare di quella indipendente.



Equazione lineare

Se vogliamo quantificare l'intensità della relazione tra le due variabili, la dobbiamo esprimere attraverso un'equazione matematica.

Ogni equazione è definita dalla sua **forma funzionale** e dai valori che assumono i suoi **parametri**.

Parlando di forma funzionale prenderemo in considerazione solo quella lineare, che possiamo esprimere semplicemente così:

$$Y = \alpha + \beta X$$

Il valore di Y è dato dal parametro α , che è costante, più X moltiplicato per il parametro β .

La seguente è invece la funzione di una parabola:

$$Y = \alpha X^2 + \beta X + \gamma$$

Equazione lineare

Facciamo un esempio concreto:

Ho i due parametri α e β

$$\alpha = 3$$

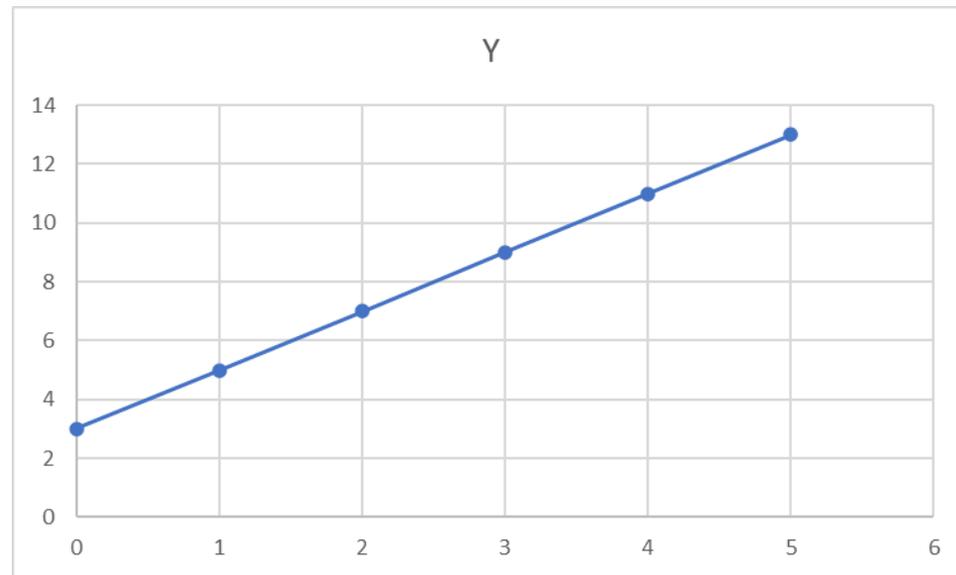
$$\beta = 2$$

La nostra equazione lineare $Y = \alpha + \beta X$

diventa $Y = 3 + 2X$

X	Y
0	3
1	5
2	7
3	9
4	11
5	13

Equazione lineare



X	Y
0	3
1	5
2	7
3	9
4	11
5	13

Equazione lineare

Facciamo un esempio concreto

$$\alpha = 5$$

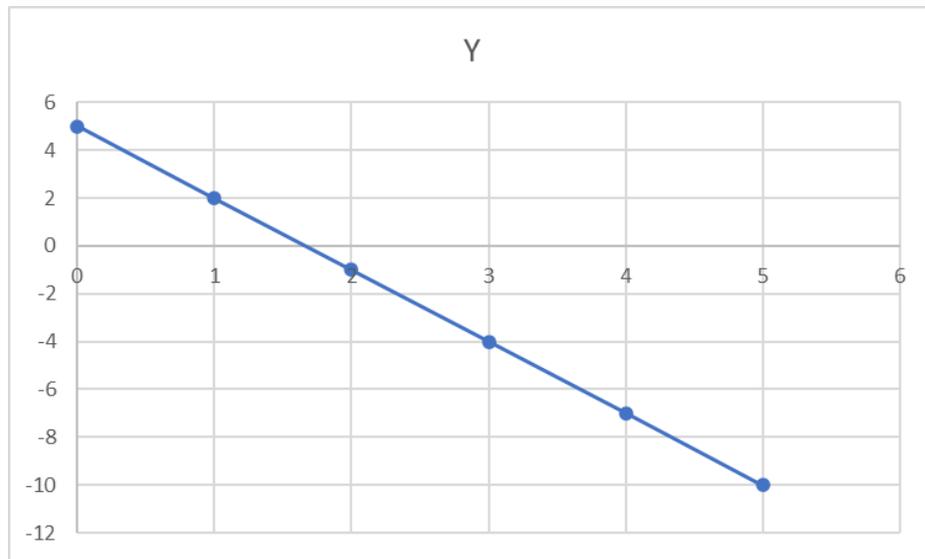
$$\beta = -3$$

La nostra equazione lineare $Y = \alpha + \beta X$

diventa $Y = 5 - 3X$

X	Y
0	5
1	2
2	-1
3	-4
4	-7
5	-10

Equazione lineare



X	Y
0	5
1	2
2	-1
3	-4
4	-7
5	-10

Equazione lineare

Abbiamo espresso la relazione tra X e Y con una linea retta.

Come influiscono i due parametri α e β ?

α stabilisce la distanza dall'asse orizzontale, ovvero il valore di Y in corrispondenza dello 0 della X. Questo parametro viene definito anche come **intercetta o costante**.

β determina **l'inclinazione o coefficiente angolare** e ci dice di quanto varia la Y al variare di X. Questo valore ci spiega l'intensità dell'effetto della variabile indipendente sulla variabile dipendente.

Se β è positivo ci troviamo di fronte a una relazione diretta, mentre se è negativo la relazione è inversa.

Equazione lineare

Abbiamo espresso la relazione tra X e Y con una linea retta $Y = \alpha + \beta X$.

Come abbiamo ripetuto più volte, nelle scienze sociali non è possibile rappresentare esattamente con un'equazione lineare una relazione complessa, come, nel nostro esempio, quella tra reddito e voto di laurea.

Dalle immagini che abbiamo visto in precedenza, al medesimo voto di laurea possono corrispondere diversi livelli di reddito medio.

La relazione tra due variabili non può essere rappresentata esattamente da un'equazione lineare, tuttavia non possiamo negare che la nuvola di punti ci suggerisca una tendenza precisa, ovvero al crescere delle ore di studio aumenta il voto all'esame, oppure al crescere del voto di laurea corrisponda un reddito medio più elevato.

Un'equazione lineare ci aiuta a stimare i due parametri α e β , e ad approssimare la covarianza tra le due variabili.

Modello di regressione lineare semplice

La formula che approssima la relazione tra X e Y con una linea retta è la seguente:

$$\hat{Y}_i = \alpha + \beta X_i$$

\hat{Y}_i rappresenta il **valore atteso** sulla base dei parametri stimati α e β e non quello osservato.

Se vogliamo esprimere i valori osservati di Y dobbiamo aggiungere all'equazione lineare un ulteriore elemento ε_i , che rappresenta la componente erratica (gli errori di predizione).

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

La componente ε_i rappresenta la differenza tra il valore osservato e il valore atteso, derivante dal modello di regressione lineare, ovvero:

$$\varepsilon_i = Y_i - \hat{Y}_i$$

se sostituiamo \hat{Y}_i con $\alpha + \beta X_i$ (visto che $\hat{Y}_i = \alpha + \beta X_i$) otteniamo che:

$$\varepsilon_i = Y_i - \alpha - \beta X_i$$

I residui

Gli errori di predizione vengono chiamati **residui**, dal momento che corrispondono a quella parte di Y che non viene spiegata dall'effetto lineare di X.

ε_i rappresenta:

- l'influenza su Y di tutti i fattori casuali che non sono stati introdotti nel modello di regressione lineare utilizzato.
- Il fatto che la relazione tra X e Y non è detto sia perfettamente lineare.
- Nelle scienze sociali i comportamenti umani sono di consueto caratterizzati da una componente di casualità che nessun modello di regressione, anche il più sofisticato, sarebbe in grado di stimare con precisione il valore di Y.

Scelta della retta di regressione

Abbiamo chiarito che l'obiettivo della regressione lineare è stimare i valori dei parametri α e β che consentono di approssimare nel modo migliore la covarianza tra X e Y.

Questo significa che la retta di regressione migliore è quella che minimizza i valori osservati di Y e quelli predetti attraverso il modello.

Dal momento che la differenza tra valori attesi e valori osservati:

$$\varepsilon_i = Y_i - \hat{Y}_i$$

La migliore retta di regressione è quella che minimizza gli errori di predizione.

Scelta della retta di regressione

Se proviamo a sommare gli errori da una determinata retta scopriremo che si annullano, ovvero

$$\sum_{i=1}^N \varepsilon_i = 0$$

Come abbiamo già visto in altre situazioni, possiamo elevare al quadrato gli scarti e la migliore retta di regressione è quella che minimizza il quadrato della somma dei residui.

Questo significa che rende minima la seguente quantità:

$$\sum_{i=1}^N (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^N \varepsilon_i^2$$

I valori dei parametri α e β che soddisfano questo criterio sono detti **stime dei minimi quadrati**.

Metodo dei minimi quadrati

Tralasciamo la dimostrazione e vediamo qual è la formula che ci consente di stimare i due parametri.

$$\beta = \frac{\sum_{i=1}^N (x_i - M(X))(y_i - M(Y))}{\sum_{i=1}^N (x_i - M(X))^2} = \frac{\text{Codev}(X,Y)}{\text{Dev}(X)}$$

$$\alpha = M(Y) - \beta M(X)$$

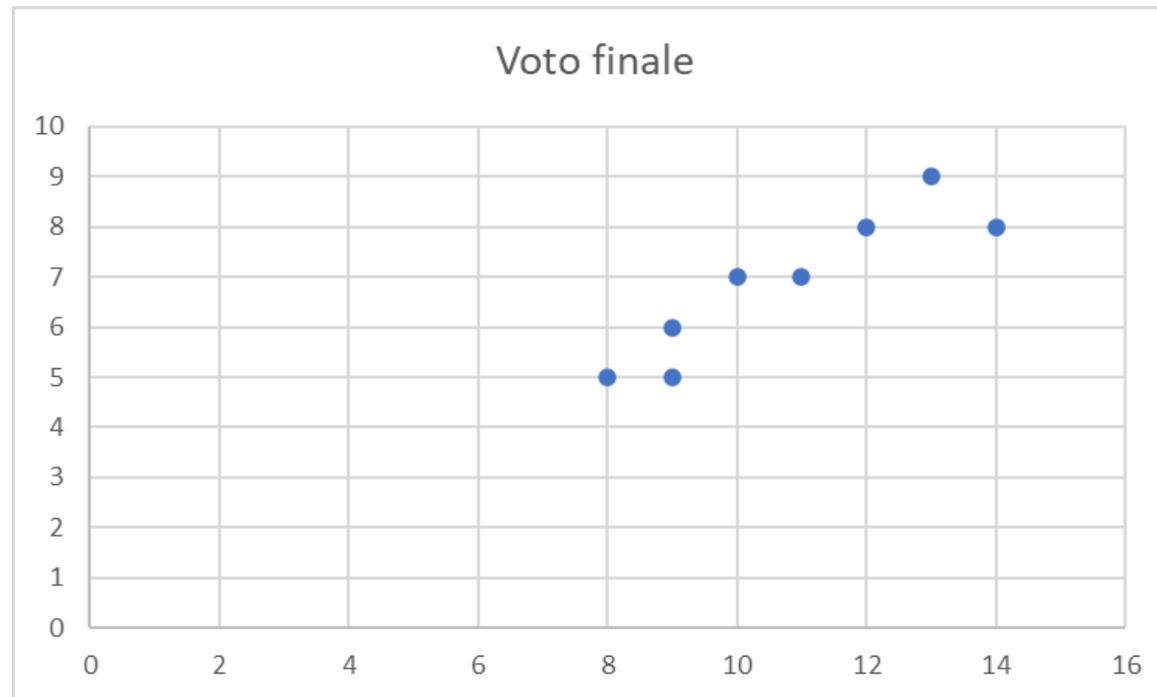
Esercitiamoci

(Esempio di pag. 150 del libro) Abbiamo due variabili misurate su 8 studenti: voto del test di ingresso a inizio anno scolastico e voto finale di matematica.

Studente	Test ingresso	Voto finale
1	12	8
2	10	7
3	14	8
4	9	5
5	9	6
6	13	9
7	11	7
8	8	5

Esercitiamoci

Come prima cosa costruiamoci uno scatterplot (diagramma di dispersione), per vedere come si dispongono sul piano cartesiano i punti individuati dalle coppie di valori.



Esercitiamoci

Dalla rappresentazione grafica possiamo individuare un andamento lineare positivo.

Procediamo quindi con il calcolo dei vari indicatori che ci servono:

$$M(X) = 10,75 \quad M(Y) = 6,88$$

$$\text{Var}(X) = 3,94 \quad \text{Var}(Y) = 1,86$$

$$\text{Dev}(X) = 31,50 \quad \text{Dev}(Y) = 14,88$$

$$\text{Codev}(X;Y) = 19,75$$

$$\text{Cov}(X;Y) = \text{Codev}(X;Y)/N = 2,47$$

Esercitiamoci

Possiamo calcolare il valore di ρ coefficiente di correlazione lineare di Pearson.

$$\rho = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Dev}(X)\text{Dev}(Y)}}$$

$$\rho = \frac{19,75}{\sqrt{31,50*14,88}} = \frac{19,75}{21,65} = 0,91$$

Il coefficiente ρ pari a 0,91 indica che c'è correlazione positiva tra le due variabili.

Esercitiamoci

Possiamo calcolare il parametri della retta di regressione.

$$\beta = \frac{\text{Codev}(X,Y)}{\text{Dev}(X)} = \frac{19,75}{31,50} = 0,63$$

$$\alpha = M(Y) - \beta M(X) = 6,88 - (0,63 * 10,75) = 0,13$$

Esercitiamoci

La nostra retta di regressione è la seguente:

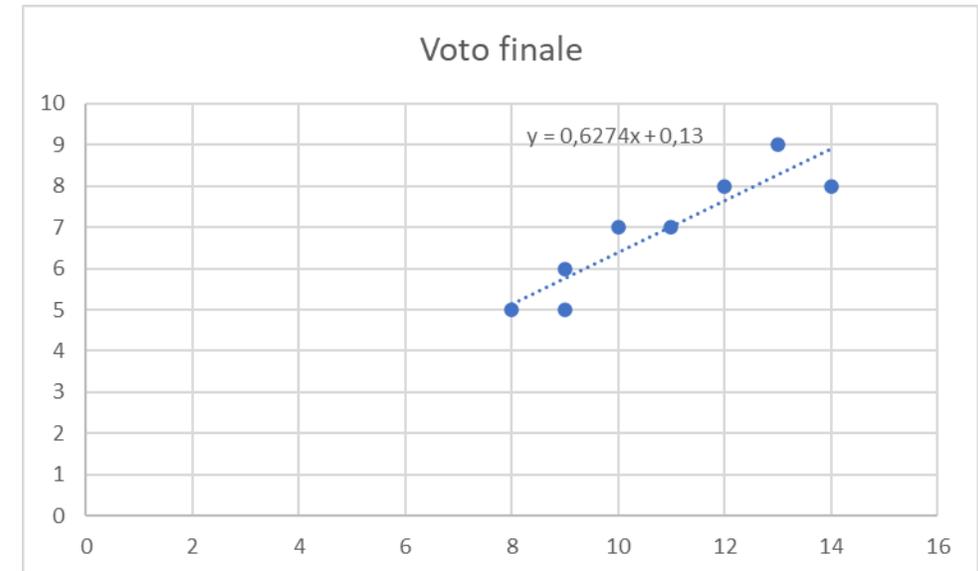
$$\hat{Y}_i = \alpha + \beta X_i$$

$$\hat{Y}_i = 0,13 + 0,63X_i$$

Possiamo costruire la retta di regressione con due punti:

l'intercetta (ovvero il valore di \hat{Y}_i con $\beta = 0$) che corrisponde a (0 e 0,13) e

il valore individuato da $M(X)$ e $M(Y)$, ovvero (10,75 e 6,88).



La bontà del modello

Come abbiamo visto, il metodo dei minimi quadrati ci garantisce che la retta individuata sia la migliore possibile.

Questo però non ci assicura che la retta sia il miglior modello per rappresentare i dati.

Come possiamo capire se la retta di regressione è adatta a rappresentare i dati:

- Calcoliamo un apposito indice;
- Analizziamo graficamente i residui.

L'indice di determinazione R^2

Dal punto di vista teorico possiamo dire che un modello di regressione è tanto migliore quanto i valori della Y e quelli ottenuti con la retta di regressione hanno una correlazione vicina a 1.

L'indice di determinazione, detto anche coefficiente R^2 è il quadrato del coefficiente di correlazione lineare fra Y e \hat{Y} .

$$R^2 = [\rho(Y; \hat{Y})]^2$$

L'indice di determinazione varia tra 0 e 1.

È pari a 1 quando la variabilità totale di Y è totalmente spiegata dalla retta di regressione.

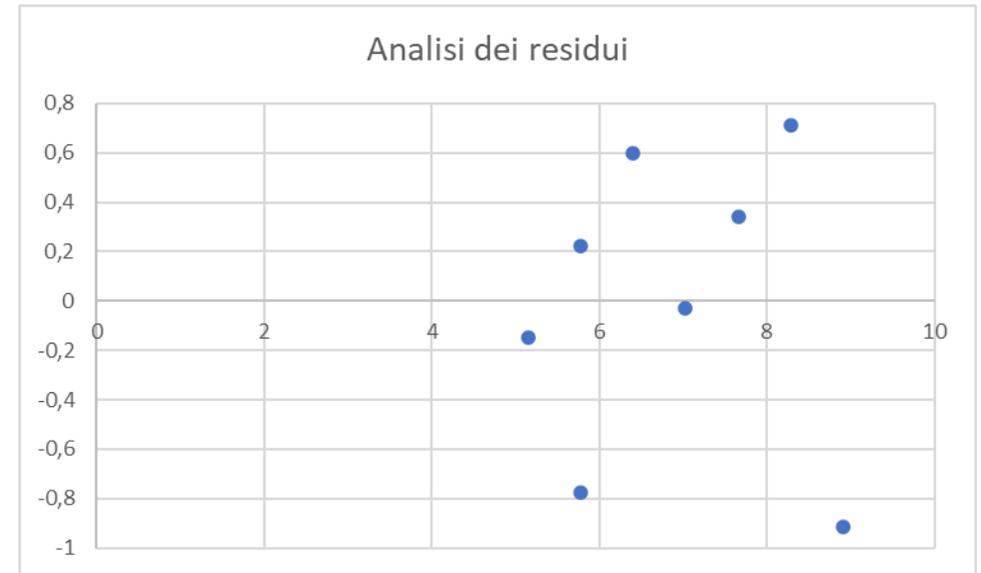
È pari a 0 quando la variabilità totale di Y non è per nulla spiegata dalla retta di regressione.

Analizzare graficamente i residui

Affinché la retta di regressione possa essere considerata una buona approssimazione della relazione tra X e Y , i residui devono avere un andamento casuale rispetto ai valori della X .

Ad esempio se all'aumentare dei valori di X crescessero sistematicamente anche i residui, ci potremmo trovare in presenza di una relazione non lineare e quindi la retta di regressione non è il modello più adeguato.

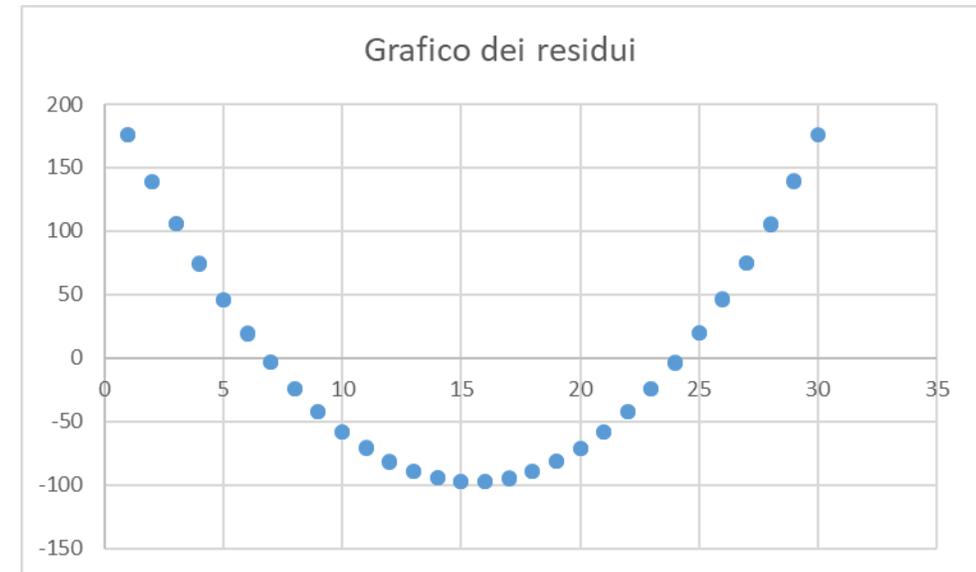
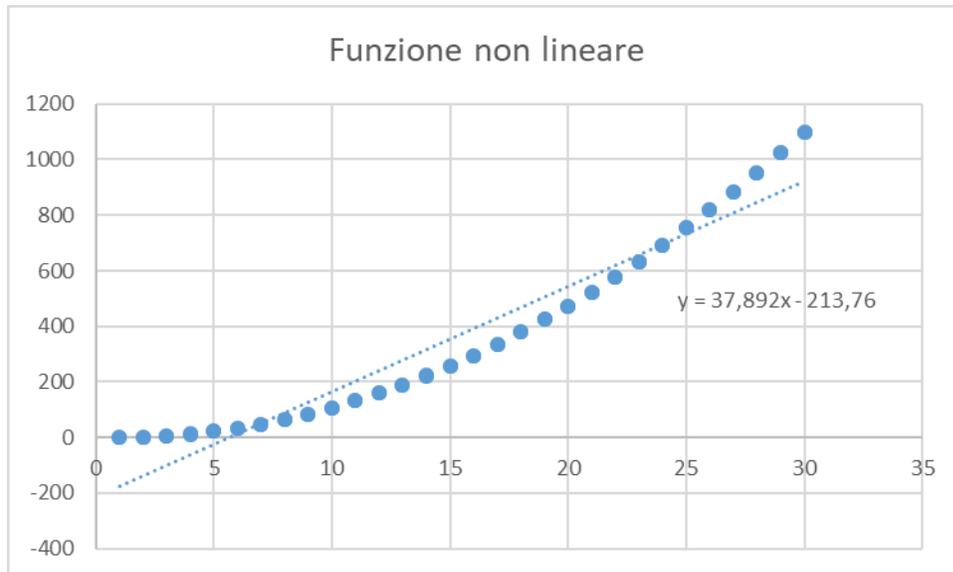
Dovremo quindi costruire un diagramma di dispersione con i valori di x_i e i relativi residui.



Analizzare graficamente i residui

Partiamo da un funzione non lineare e calcoliamo la nostra retta di regressione.

Dal punto di vista grafico la retta di regressione non sembra così inadeguata, ma guardando il grafico dei residui possiamo vedere che questi non si distribuiscono in modo casuale.



Appendice: devianza e codevianza

Nella statistica univariata si usa come misura di dispersione, la **devianza**, ovvero la somma del quadrato degli scarti dalla media:

$$Dev(x) = \sum_{i=1}^N (x_i - M(x))^2$$

La devianza divisa per N ci dà la varianza.

Nella statistica bivariata invece possiamo utilizzare la codevianza ovvero il prodotto degli scarti:

$$Codev(x) = \sum_{i=1}^N (x_i - M(x))(y_i - M(y))$$

La codevianza divisa per N ci dà la covarianza.

Se X e Y sono indipendenti allora $Cov(X;Y) = 0$

Appendice: devianza e codevianza

Il coefficiente di correlazione lineare di Bravais – Pearson si calcola rapportando la covarianza al suo massimo.

Si dimostra che il massimo della covarianza è pari a al prodotto degli scarti quadratici medi per X e Y.

$$\text{Max Cov}(X, Y) = \sqrt{\text{Var}(x)\text{Var}(y)} = \sigma_x \sigma_y$$

$$\rho = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y}$$

Moltiplicando tutto per N otteniamo che:

$$\rho = \frac{\text{Codev}(X,Y)}{\sqrt{\text{Dev}(x)\text{Dev}(y)}}$$

Per concludere

Abbiamo visto come analizzare la relazione tra due fenomeni quantitativi.

In generale, nell'ambito dell'analisi bivariata, abbiamo sempre utilizzato il seguente approccio:

- Abbiamo individuato un modello per definire la correlazione
- Ne abbiamo misurato la bontà.

Giovedì 15 dicembre faremo, attraverso un test, un veloce ripasso di quanto abbiamo visto durante questo corso.

Grazie e in bocca al lupo!

Taming the Statistical Beast

