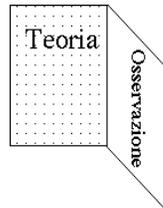


Dei modelli

I problemi statistici incontrati nello studio delle discipline socioeconomiche nascono dal dualismo fra evidenza empirica e analisi teorica.



Tale dualismo è incorporato nella nozione di modello scientifico

Non sono possibili costruzioni teoriche senza ripetizione.

I fenomeni socio-economici non si ripetono se considerati solo in modo superficiale e descrittivo.

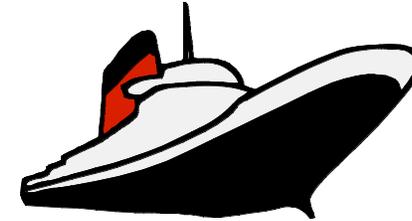
Se però ci si limita ai fattori più rilevanti troviamo delle ricorrenze sulle quali impostare il modello

Semplificazione ed astrazione

Il modello è una rappresentazione semplificata ed astratta di una realtà. Con esso si può lavorare su una realtà più grande, complessa e mutevole.

Esso sta alla realtà come il quadro sta alla fotografia: questa riporta tutto, quello solo ciò che ha colpito l'ispirazione dell'artista.

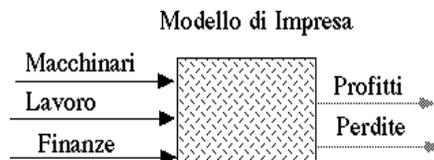
Il modello dà risposte in ragione della sua vicinanza al fenomeno che rappresenta



Per studiare il comportamento della nave non si userà una barchetta di carta, ma una serie di equazioni, disegni e modelli in scala.

Esempio: modello di impresa

Un'impresa può essere rappresentata come una combinazione di inputs per produrre profitti



Il modello descrive e spiega schematicamente una certa situazione imprenditoriale.

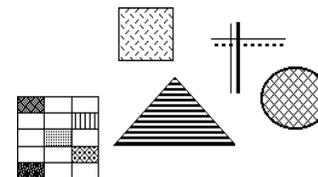
Le conclusioni basate sul modello non sono neutrali: sono legate alle ipotesi in esso inglobate.

Ogni compromesso ragionevole tra semplicità e realtà è un modello

Estensioni del modello

il limite del modello di impresa è che non tiene conto dell'interazione con le altre imprese. Per includerla occorre allargare o cambiare il modello

L'approccio deve essere pluralistico: non c'è un unico, universale modello scientifico.



E' necessario predisporre modelli differenti per affrontare le varie situazioni incontrate nelle scienze economiche e sociali

Gli elementi soggettivi sono tali e tanti che autori diversi, pur lavorando sulla stessa realtà pervengono a modelli diversi e talvolta contrapposti

Tipi di modelli



VERBALI

Descrizione a parole di una situazione verificabile:
La diminuzione del Tasso Ufficiale di Sconto favorisce gli investimenti



MATEMATICI

Gli aspetti essenziali di una situazione sono espressi con delle equazioni
La 1^a legge del moto di Newton

$$f(t) = \beta_0 + \beta_1 t \quad \text{con} \quad \begin{cases} \beta_0 = \text{posizione iniziale} \\ \beta_1 = \text{velocità uniforme} \\ f(t) = \text{distanza percorsa} \end{cases}$$

I modelli studiati dalla statistica -sotto qualsiasi forma- debbono avere precisi riscontri nella realtà

Tipi di modelli/2

FISICI

Si costruisce un apparato che rappresenta IN SCALA la situazione di studio oppure ne rappresenta una parte.
Lo studio del CX per le auto nella galleria del vento



ANALOGICI

Delle relazioni non fisiche sono simulate con dei meccanismi fisici
La diffusione di un dialetto attraverso i cerchi concentrici che si producono sull'acqua



Esempio: il sistema economico

La costruzione della teoria economica consiste nella elaborazione di un sistema di equazioni interconnesse (il modello).

Le equazioni devono comporsi di variabili Indipendenti, dipendenti e Parametri.

Conoscendo i parametri siamo in grado di sapere quali saranno le variazioni nelle dipendenti per ogni combinazione di livelli nelle indipendenti

Esempio di modello
 econometrico
 multiequazionale

$$\begin{cases} Y = C + I \\ C = a_1 + a_2 Y + a_3 r \\ I = a_4 + a_5 r \\ r = 11\% \end{cases} \quad \text{con} \quad \begin{cases} Y = \text{Reddito} \\ C = \text{Consumi} \\ I = \text{Investimenti} \\ r = \text{Tasso di Sconto} \end{cases}$$

$a_1 \ a_2 \ a_3 \ a_4 \ a_5 = \text{parametri}$

Queste relazioni debbono essere coerenti al loro interno (non avere contraddizioni) e rispecchiare la realtà osservata

Il sistema economico/2

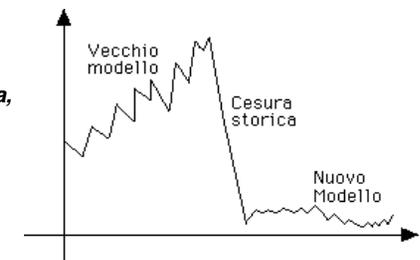
il funzionamento del sistema, PER I SUOI EFFETTI CUMULATIVI, provoca un graduale modificarsi nei parametri

il modello si applica finché i parametri sono invariati o cambiati di poco

Da un certo punto in poi le variazioni nei parametri superano l'elasticità assunta dal modello.

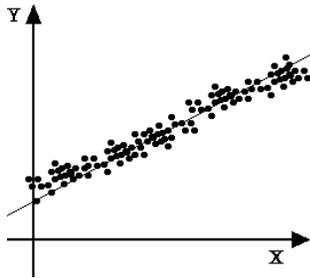
A questo punto il modello crolla perché nel fenomeno in esame c'è un cambiamento strutturale: una vera e propria cesura storica.

Un modello ben costruito dovrebbe spiegare il funzionamento dell'economia, ma contenere anche elementi di autodistruzione.



Relazione tra due variabili

Dopo aver rappresentato graficamente i dati a mezzo dello scatterplot si è interessati a determinare una curva che passi vicino ai punti



- Per sostituire uno schema semplice alla nube dei punti
- Per sintetizzare le tendenze di fondo
- Per ricostruire o determinare il valore della Y noto quello della X o viceversa

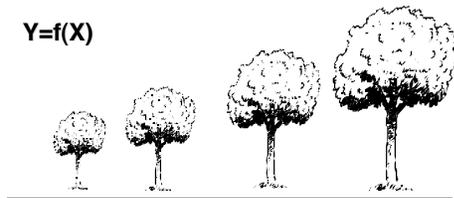
il presupposto è che esiste una variabile (la "X" detta indipendente o esogena) che è causa o comunque agisce sull'altra (la "Y" detta dipendente o endogena).

Esempio di costruzione di un modello

La teoria di un fenomeno può spesso essere sintetizzata da un modello espresso da una equazione.

Sia "Y" l'ampiezza in cm del diametro alla base del tronco di una data specie arborea e sia "X" l'età .

L'idea che il diametro sia più grande secondo l'età può essere espressa dalla relazione funzionale:

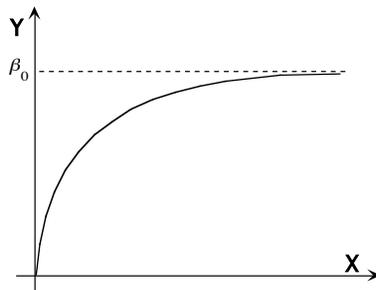


Queste variazioni assicurano all'albero adeguata resistenza e flessibilità.

Esempio di costruzione del modello/2

La funzione "f" è al momento indeterminata: si sa che un certo legame esiste, ma non si riesce a darle una esatta espressione analitica.

E' noto che, a parità di forma, una specie non può superare una dimensione data. Per cui la relazione tra X ed Y è di tipo crescente, ma gli aumenti devono avvenire a ritmo decrescente



il modello potenza è particolarmente adatto per rappresentare tali situazioni.

$$Y = \beta_0 \left[1 - (\beta_1)^x \right], \quad 0 < \beta_1 < 1$$

Esempio di costruzione del modello/3

Nel modello si individuano

Y= variabile ENDOGENA-DIPENDENTE-SPIEGATA-INTERNA-CONSEQUENTE

X= variabile ESOGENA-INDIPENDENTE-ESPLICATIVA-ESTERNA-ANTECEDENTE

Esistono moltissimi fattori che incidono sull'accrescimento: la quota, il tipo di suolo, l'esposizione, l'impianto arboreo, etc.

Tali fattori non solo incidono su "Y" ma si influenzano anche tra di loro determinando una rete complessa di interrelazioni che il modello ignora.

La "X" è un "riassunto" dei fattori determinanti, ovvero si sceglie "X" perché è considerata il risultato del loro comune interagire.

La variabile endogena

Una variabile ha questo ruolo se:

- Rappresenta il fenomeno che si intende spiegare, prevedere, controllare
- E' una risposta ad uno più stimoli in un dato organismo
- E' l'output di un sistema che ha uno o più fattori in input
- Esprime l'obiettivo raggiungibile per uno o più tipi di interventi.

La variabile esogena

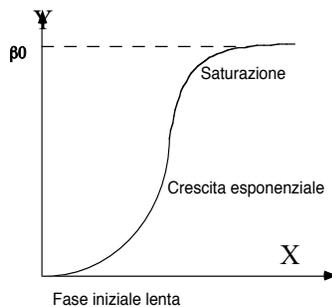
Una variabile ha questo ruolo se:

- E' un riassunto dei fattori determinanti (età dell'albero)
- E' controllabile (spese in pubblicità, aumento delle vendite)
- E' precedente la endogena (quotazione di oggi, quotazione di domani)
- E' ritenuta una causa determinante della endogena (ore di studio e voto d'esame)

Che cos'è il modello?

E l'insieme delle ipotesi e delle equazioni che stabiliscono una certa relazione tra due o più variabili.

La curva logistica $Y = \frac{\beta_0}{1 + e^{-\beta_1 X}}$ ingloba le seguenti ipotesi:



Avvio difficile, crescita lenta e faticosa;

Accelerazione dello sviluppo che appare quasi incontrollabile.

Rallentamento del fenomeno che finisce con l'attestarsi sull'asintoto β_0

Fase iniziale lenta

Tale modello descrive bene l'aumento delle popolazioni di persone, di imprese, di batteri, di prodotti, etc.

Esempio: modello Keynesiano semplice

Intende fornire una spiegazione del funzionamento del sistema economico

La variabile esplicativa è la domanda globale ovvero gli investimenti e segnatamente quelli di natura pubblica "F".

La variabile "F" si configura come "esogena" in quanto soggetta -almeno in parte- a controllo governativo.

La variabile endogena Y è il livello di produzione e del reddito della nazione in condizione di sottoccupazione

$$Y = f(F)$$

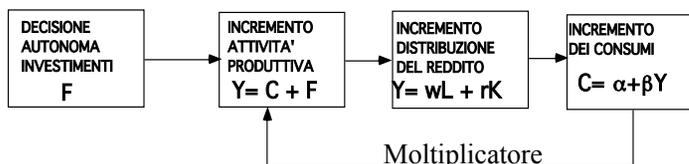


il modello Keynesiano semplice/2

F variabile autonoma
 $Y = C + F$
 $wL + rK = Y$
 $C = \alpha + \beta Y$

$$Y = \alpha + \beta Y + F \Rightarrow (1 - \beta)Y = \alpha + F$$

$$Y = \left[\frac{1}{1 - \beta} \right] F + \frac{\alpha}{1 - \beta}$$



La realizzazione degli investimenti attiva la produzione che a sua volta genera un aumento di reddito ai fattori lavoro e capitale.

Tale aumento induce maggiori consumi che si concretizza in una maggiore produzione di beni

il modello Keynesiano semplice/3

La quantità

$$\left[\frac{1}{1 - \beta} \right]$$

rappresenta il “moltiplicatore” e misura l’impatto sul reddito generato da un incremento unitario di investimenti “F”.

“ β ” è la propensione marginale al consumo ovvero l’incremento indotto nei consumi da un aumento unitario di reddito:

$$0 < \beta < 1$$

Se la propensione marginale al consumo vale 0.60 il moltiplicatore varrà 2.5 per cui un investimento pari a 1’000 produrrà -attraverso il circuito economico- un incremento di reddito/produzione pari a 2’500.

Finalità del modello



DESCRITTIVE

al modello si chiede solo di rappresentare bene la realtà osservata.

La capacità dei Chip di Memoria si quadruplica ogni 3 anni



INTERPRETATIVE

il modello deve mettere in evidenza i legami tra i fenomeni coinvolti in forme e modi riconducibili a precise teorizzazioni.

La produzione è una funzione lineare omogenea di capitale e lavoro (equazione Cobb-Douglas)

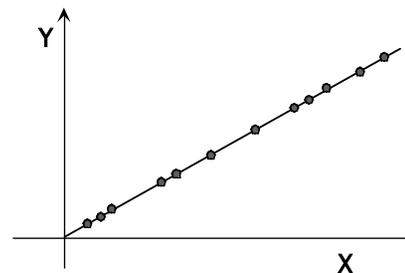


PREVISIONALI

il modello deve fornire previsioni sull’andamento futuro del fenomeno

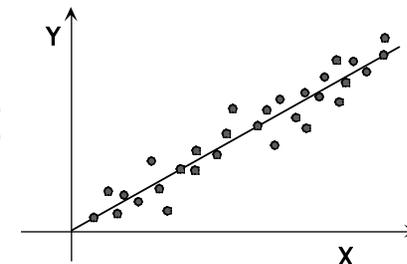
Le esportazioni di beni durevoli aumentano linearmente nel tempo

Relazioni stocastiche e deterministiche



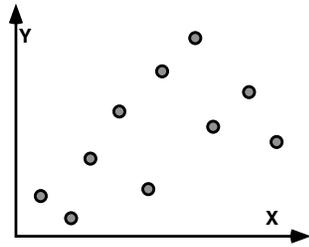
DETERMINISTICA: ad ogni età corrisponde un determinato diametro del tronco

STOCASTICA: il diametro del tronco aumenta con l’età, ma l’incremento non è UNIVOCO: talvolta aumenta di più, altre volte di meno



Relazioni esatte ed approssimate

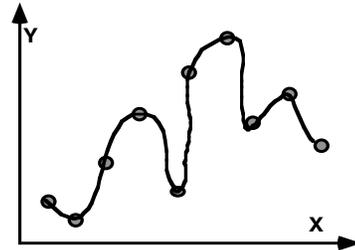
Consideriamo il seguente scatterplot:



Si conoscono i valori Y_i corrispondenti ai valori X_i ;

Siamo alla ricerca di una funzione $f(X)$ i cui valori $f(X_i)$ siano "vicini" alle Y_i .

E' possibile determinare un adattamento perfetto: se (X_i, Y_i) $i=1, \dots, n$ sono "n" coppie di valori distinte allora un polinomio di grado "n-1" o inferiore si adatta perfettamente ai punti.



Relazioni esatte ed approssimate/2

Ad esempio, il polinomio di Lagrange

passa esattamente per gli "n" punti

$$f(X) = \sum_{i=1}^n y_i \left[\frac{\prod_{\substack{j=1 \\ j \neq i}}^n (X - X_j)}{\prod_{\substack{j=1 \\ j \neq i}}^n (X_i - X_j)} \right]$$

il calcolo dei valori può essere facilmente programmato al computer

Di solito i polinomi di Lagrange hanno un grado troppo elevato (comunque non superiore a n-1) per essere usabili in modo rapido e semplice.

Se si aggiunge un nuovo punto il calcolo deve essere ripetuto

In statistica si rinuncia alla perfetta interpolazione matematica per un adattamento approssimato, ma più essenziale e stabile

Relazioni esatte ed approssimate/3

il compromesso tra bontà di adattamento e semplicità del modello deriva da forti dosi di convenzionalismo

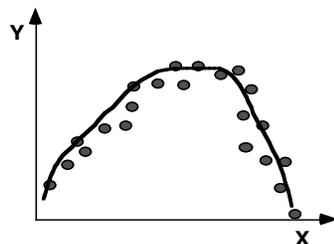


SEMPLICITA' Si definisce un sistema di curve "semplici" e flessibile

$$f(X; \theta_1, \theta_2, \dots, \theta_m)$$



ADATTABILITA' Si cerca di riconoscere la struttura del modello "f" nello scatterplot per poi adattarvi quella più idonea.



L'andamento degli n=22 punti ricorda in modo indiscutibile la parabola.

Parabole ne esistono infinite, quale scegliere?

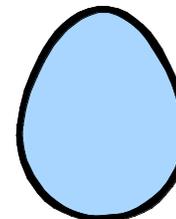
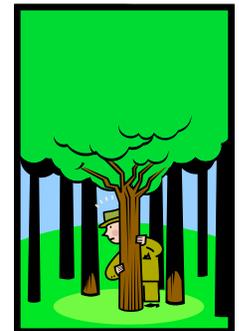
N.B. l'identificazione del modello è ostacolata dalla presenza di errori

Proposta del modello lineare:

il rasoio di Occam

Se è necessario dare una soluzione ad un problema di cui si sa poco, la risposta più semplice comporta meno rischi in caso di errore ed è spesso quella giusta.

Smarriti in una foresta se ne esce spesso procedendo in linea retta.



L'uovo di Colombo

Principio di semplicità di Galilei

La natura procede per vie semplici ed offre così la sicura scelta tra le varie spiegazioni possibili dei suoi fenomeni

Proposta del modello lineare/2

il legame più semplice tra due variabili è quello lineare

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

Ipotizziamo che l'ordinata "Y" sia dovuta alla combinazione ADDITIVA di due valori: la parte deterministica (lineare) ed un errore

Si cammina nel piano fino a che non si guarda l'orizzonte per ricordarsi che la terra è una sfera.

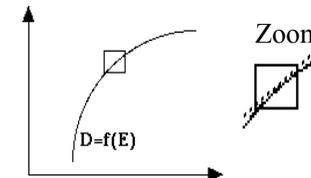
il termine "u" è il risultato di:

- Errori e carenze nella misurazione e nella rilevazione di "Y" e di "X"
- Insufficienza del solo fattore X a "spiegare" da solo la Y
- Inadeguatezza della "semplice" relazione lineare

Una ragione di più

Teorema di Taylor.

Se la funzione "f" che lega "D" ad "E" ha derivate prime e seconde continue in un intorno del punto E0, in tale intorno la "f" è ben approssimata dalla retta



In generale si può dire che la scelta del modello lineare è motivata da

- Ragioni di semplicità
- Esigenze di sintesi
- Approssimazione funzionale

Limiti del modello lineare

I sistemi dinamici hanno la proprietà di non poter essere compresi se non in modo globale.

Questa regola ammette una sola eccezione: i sistemi integrabili, fra i quali si collocano in prima fila i sistemi lineari".

Si perde però di vista l'instabilità potenziale.

Se un fenomeno ha effetti cumulativi tali che:

$$Y_{n+1} = 100 * Y_n \Rightarrow Y_n = Y_0 100^n$$

una causa infinitesima può avere effetti catastrofici:



il battito d'ali di una farfalla nei Caraibi imprime al vento una forza pari a 0.000'000'000'000'000'1 nodi, ma dopo solo 9 passaggi il vento ha una forza di 100 nodi che travolgerà New York



Modello di regressione lineare semplice

Supponiamo di disporre di "n" coppie di osservazioni

| y | x |
|----------------|----------------|
| y ₁ | x ₁ |
| y ₂ | x ₂ |
| ⋮ | ⋮ |
| y _i | x _i |
| ⋮ | ⋮ |
| y _n | x _n |

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

↑
COMPONENTE DETERMINISTICA

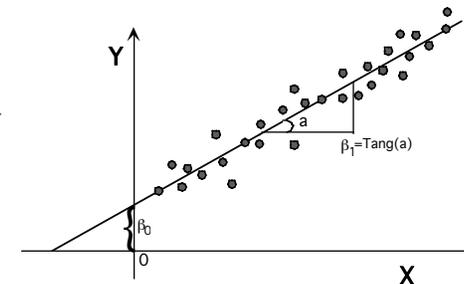
←
COMPONENTE STOCASTICA (Non osservabile)

β_0 = intercetta.

valore a cui tende la Y quando la X è zero.

β_1 = Coefficiente angolare.

variazione in Y per un aumento unitario in X



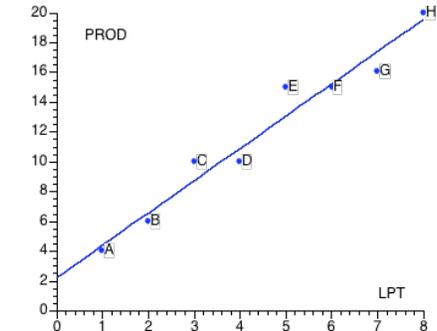
La terminologia

- **Modello.** Perché è un insieme di ipotesi rispetto al legame esistente tra la variabile esogena ed endogena. Le ipotesi in genere danno luogo ad una equazione, lineare nel nostro caso
- **Regressione.** E' una etichetta storica dovuta agli studi di Francis Galton (1889) sull'effetto di regressione: la tendenza a prevalere dei valori medi.
- **Lineare.** Perché i parametri incogniti vi compaiono con potenza 1
- **Semplice.** Perché c'è una sola variabile esplicativa (REGRESSORE) in contrapposizione a *multipla* termine usato quando vi sono più variabili esplicative.

Esempio di modello di regressione

L'ing. Consolata Mirabelli è responsabile della produzione di semilavorati. Nel breve periodo controlla solo il lavoro part-time. L'ing. intende conoscere che relazione (se c'è) tra questo fattore e la produzione

| Prova | L.P.T. | PROD |
|-------|--------|------|
| A | 1 | 4 |
| B | 2 | 6 |
| C | 3 | 10 |
| D | 4 | 10 |
| E | 5 | 15 |
| F | 6 | 15 |
| G | 7 | 16 |
| H | 8 | 20 |

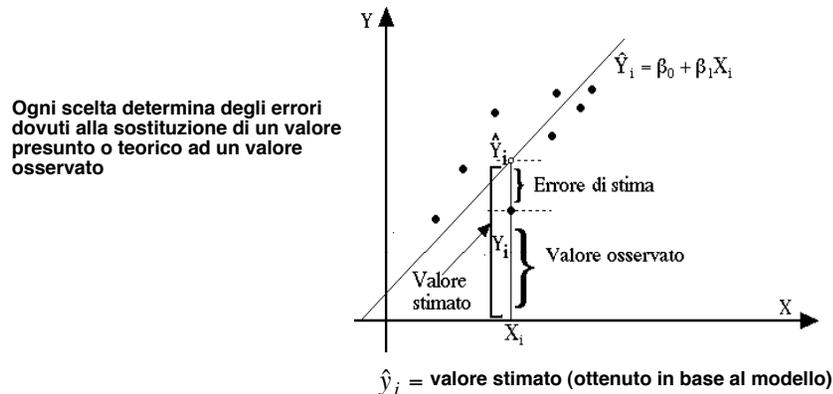


Cosa indica lo scatterPLOT?

- 1) Una sicura tendenza all'aumento della produzione dovuta all' aumento di LPT
- 2) Un'altrettanto sicura dispersione intorno alla tendenza (espressa dalla retta)

Calcolo dei parametri

Se per due punti passa una sola retta fra più di due punti non allineati ne passano infinite.



Occorre stabilire un criterio che ci permetta di scegliere quella che passa più vicino ai punti ovvero si adatta bene allo scatterplot

Criteri di calcolo

La vicinanza della retta è espressa con una sintesi degli scarti relativi tra valori osservati e valori stimati.

Poiché la somma dei valori è fissa il denominatore è ignorato

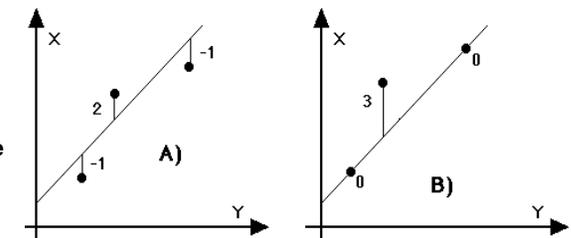
$$Q_r(\beta_0, \beta_1) = \frac{\sum_{i=1}^n |Y_i - \hat{Y}_i|^r}{n} = \frac{\sum_{i=1}^n |Y_i - \hat{Y}_i|^r}{\sum_{i=1}^n |Y_i|^r}$$

(di solito $r=1$ oppure $r=2$)

La scelta del criterio determina la scelta della retta.

Secondo Q_1 l'adattamento è migliore con la retta B.

E' il contrario secondo Q_2 .



$$Q_1(A) = |-1| + |2| + |-1| = 4$$

$$Q_2(A) = 1^2 + 2^2 + 1^2 = 6$$

$$Q_1(B) = |0| + |3| + |0| = 3$$

$$Q_2(B) = 0^2 + 3^2 + 0^2 = 9$$

Scelta tra scarti assoluti e al quadrato

il criterio dei minimi assoluti -proposto dall'astronomo Boscovich e ripreso da Laplace- risale almeno al 1755.

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_i) = 0 \quad \text{Somma degli scarti negativi e somma degli scarti positivi uguali in valore assoluto}$$

$$\sum_{i=1}^n |y_i - \beta_0 - \beta_1 X_i| \quad \text{Minima rispetto ai parametri incogniti } \beta_0, \beta_1$$

Pro Gli scarti hanno lo stesso ordine di grandezza dei valori per cui il criterio risulta semplice e naturale

Contro La soluzione è ottenuta con algoritmi di calcolo numerico.
La soluzione non è necessariamente univoca (si pensi alla mediana per "n" pari se la retta migliore è parallela all'asse della "X")
Le trattazioni delle proprietà statistiche è difficile e poco generale.

Scelta tra scarti assoluti e al quadrati/2

il criterio dei minimi quadrati risale a Legendre e Gauss.

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_i) = 0 \quad \text{Somma degli scarti negativi e somma degli scarti positivi uguali in valore assoluto}$$

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_i)^2 \quad \text{Minima rispetto ai parametri incogniti } \beta_0, \beta_1$$

Pro **Espressione univoca e semplice della soluzione. Trattazione chiara e rigorosa delle proprietà statistiche**

Contro **Peso eccessivo agli scarti più grandi. Dati i due valori $Y_1=12$ e $Y_2=8$ ed ipotizziamo uno scarto del 20% in entrambi, si ottiene:**

$$\frac{(y_1 - \hat{y}_1)^2}{y_1} = \frac{(12 - 9)^2}{12} = 0.75; \quad \frac{|12 - 9|}{12} = 0.25$$

$$\frac{(y_2 - \hat{y}_2)^2}{y_2} = \frac{(8 - 6)^2}{8} = 0.50; \quad \frac{|8 - 6|}{8} = 0.25$$

Soluzione dei minimi quadrati

Partiamo dall'errore i-esimo:

$$(y_i - \hat{y}_i) = y_i - \beta_0 - \beta_1 X_i = y_i - \beta_0 - \beta_1 X_i \pm \bar{y} \pm \beta_1 \bar{x} \\ = (y_i - \bar{y}) + (\bar{y} - \beta_0 - \beta_1 \bar{x}) - \beta_1 (X_i - \bar{x})$$

che evidenzia il ruolo del punto (\bar{x}, \bar{y}) baricentro dello scatterplot

Elevando al quadrato e sviluppando si ottiene:

$$(y_i - \hat{y}_i)^2 = [(y_i - \bar{y}) + (\bar{y} - \beta_0 - \beta_1 \bar{x}) - \beta_1 (x_i - \bar{x})]^2 = (y_i - \bar{y})^2 + (\bar{y} - \beta_0 - \beta_1 \bar{x})^2 + \beta_1^2 (x_i - \bar{x})^2 + 2(y_i - \bar{y})(\bar{y} - \beta_0 - \beta_1 \bar{x}) - 2\beta_1 (y_i - \bar{y})(x_i - \bar{x}) - 2\beta_1 (\bar{y} - \beta_0 - \beta_1 \bar{x})(x_i - \bar{x})$$

Considerando la somma di tutti gli "n" termini e ricordando che la somma degli scarti dalla media aritmetica e nulla si arriva a:

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\bar{y} - \beta_0 - \beta_1 \bar{x})^2 + \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 - 2\beta_1 \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})$$

Soluzione dei minimi quadrati/2

Definiamo: $S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$; $S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2$; $S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$;

Tali quantità sono note come devianze e covarianze

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = S_{yy} + n(\bar{y} - \beta_0 - \beta_1 \bar{x})^2 + \beta_1^2 S_{xx} - 2\beta_1 S_{xy}$$

Ne consegue:

$$= S_{yy} + n(\bar{y} - \beta_0 - \beta_1 \bar{x})^2 + \beta_1^2 S_{xx} - 2\beta_1 S_{xy} + \frac{S_{xy}^2}{S_{xx}} - \frac{S_{xy}^2}{S_{xx}}$$

$$= S_{yy} + n(\bar{y} - \beta_0 - \beta_1 \bar{x})^2 + S_{xx} \left[\beta_1^2 - 2\beta_1 \frac{S_{xy}}{S_{xx}} + \frac{S_{xy}^2}{S_{xx}^2} \right] - \frac{S_{xy}^2}{S_{xx}}$$

$$= S_{yy} + n(\bar{y} - \beta_0 - \beta_1 \bar{x})^2 + S_{xx} \left[\beta_1 - \frac{S_{xy}}{S_{xx}} \right]^2 - \frac{S_{xy}^2}{S_{xx}}$$

La somma degli errori dipende dalle incognite solo attraverso dei termini al quadrato per cui il minimo si ottiene azzerando quei termini e cioè

$$\hat{\beta}_1 = \frac{S_{yx}}{S_{xx}}; \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

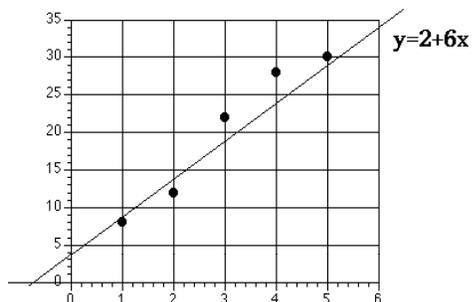
Esempio

| X | Y | (x- \bar{x}) | (y- \bar{y}) | (x- \bar{x})(y- \bar{y}) | (x- \bar{x}) ² |
|----|-----|-----------------|-----------------|--------------------------------|------------------------------|
| 1 | 8 | -2 | -12 | 24 | 4 |
| 2 | 12 | -1 | -8 | 8 | 1 |
| 3 | 22 | 0 | 2 | 0 | 0 |
| 4 | 28 | 1 | 8 | 8 | 1 |
| 5 | 30 | 2 | 10 | 20 | 4 |
| 15 | 100 | | | 60 | 10 |

$$\bar{y} = \frac{100}{5} = 20; \quad \bar{x} = \frac{15}{5} = 3$$

$$\hat{\beta}_1 = \frac{60}{10} = \frac{60}{10} = 6$$

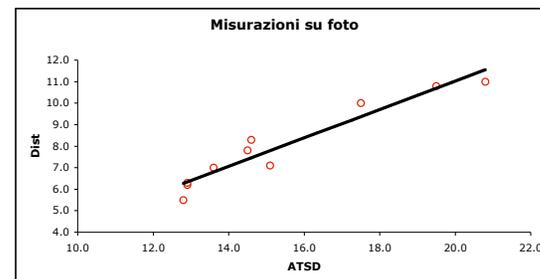
$$\hat{\beta}_0 = 20 - 6 \cdot 3 = 20 - 18 = 2$$



Esempio con Excel

Di seguito si riportano dei dati relativi ad X=ampiezza totale della sede stradale e Y= distanza tra un ciclista e un'auto (ottenuta con misurazioni su foto)

| ATSD | DIST |
|------|------|
| 12.8 | 5.5 |
| 12.9 | 6.2 |
| 12.9 | 6.3 |
| 13.6 | 7.0 |
| 14.5 | 7.8 |
| 14.6 | 8.3 |
| 15.1 | 7.1 |
| 17.5 | 10.0 |
| 19.5 | 10.8 |
| 20.8 | 11.0 |



SUMMARY OUTPUT

| Regression Statistics | |
|-----------------------|--------|
| Multiple R | 0.9607 |
| R Square | 0.9229 |
| Adjusted R Square | 0.9133 |
| Standard Error | 0.5821 |
| Observations | 10 |

Strumenti-->Analisi dei dati-->Regressione

| | Coefficients | Standard Error | t Stat | P-value |
|--------------|--------------|----------------|---------|---------|
| Intercept | -2.1825 | 1.0567 | -2.0654 | 0.0727 |
| X Variable 1 | 0.6603 | 0.0675 | 9.7858 | 0.0000 |

Chiariti altrove

Esercizio

Il prezzo e l'epoca dell'usato per una particolare auto è stato rilevato presso alcuni rivenditori

| X | Y |
|---|------|
| 2 | 10.5 |
| 5 | 3.2 |
| 1 | 11.7 |
| 6 | 4.8 |
| 4 | 7.3 |
| 5 | 3.6 |
| 3 | 8.8 |
| 2 | 9.9 |



- 1) Calcolare le stime dei parametri
- 2) Disegnare lo scatterplot e la retta teorica

Uso della retta di regressione

- INTERPOLAZIONE**
 Lo scopo è trovare i valori della dipendente o di sostituirla con i valori particolarmente anomali, per valori noti della indipendente.
- ESTRAPOLAZIONE**
 Determinazione del valore della dipendente che corrisponde ad un valore della indipendente non necessariamente osservato.
- CONTROLLO**
 Determinazione del valore della indipendente idoneo a determinare un fissato livello della dipendente

In ogni caso si ottengono dei VALORI TEORICI costituenti la stima dei VALORI VERI che rimangono comunque sconosciuti

Esempio

Unità di lavoro part-time e aumento di produzione

| Prova | L.P.T. | PROD |
|-------|--------|------|
| A | 1 | 4 |
| B | 2 | 6 |
| C | 3 | 10 |
| D | 4 | 10 |
| E | 5 | 15 |
| F | 6 | 15 |
| G | 7 | 16 |
| H | 8 | 20 |

$$PROD = 2.25 * 2.1667 * LPT$$

Ogni unità di lavoro part-time addizionale è responsabile di 2.1667 tonn. di produzione.

Se il lavoro part-time non fosse impiegato la produzione media sarebbe a 2.25 tonn.

Supponiamo che si decida di impiegare 10 LPT quale sarà l'incremento di produzione?

$$\hat{y}_{10} = 2.25 + 2.1667 * 10 = 2.25 + 21.667 = 23.917$$

Se invece si volesse stabilire quante LPT impiegare per ottenere 16 semilavorati allora

$$16 = 2.25 + 2.1667 * \hat{X} \Rightarrow \hat{X} = \frac{(16 - 2.25)}{2.1667} = 6.346$$

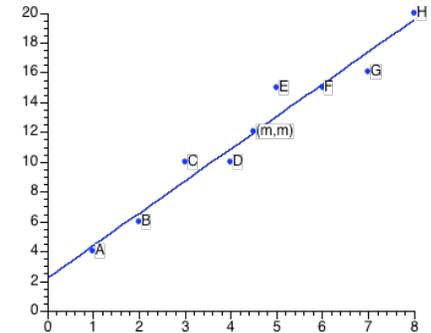
Proprietà della retta di regressione

La retta di regressione passa per il punto di coordinate (\bar{x}, \bar{y})

La retta stimata può essere scritta come:

$$y = \bar{y} + \hat{\beta}_1(x - \bar{x}) \Rightarrow \bar{y} + \hat{\beta}_1(\bar{x} - \bar{x}) = \bar{y}$$

Non si tratta di un vincolo aggiuntivo, ma una caratteristica intrinseca al metodo dei minimi quadrati



Proprietà della retta di regressione/2

La somma degli scarti tra osservate e teoriche è nulla:

$$\sum_{i=1}^n y_i - \hat{y}_i = \sum_{i=1}^n y_i - \bar{y} - \hat{\beta}_1(x - \bar{x}) \Rightarrow \sum_{i=1}^n (y_i - \bar{y}) - \hat{\beta}_1 \sum_{i=1}^n (x - \bar{x}) = 0 - \hat{\beta}_1 * 0 = 0$$

Ciò implica che Media osservate = Media teoriche

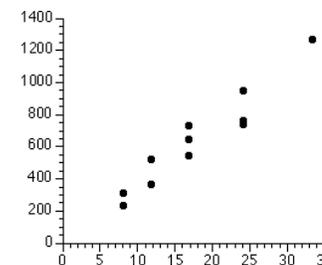
$$\frac{\sum_{i=1}^n \hat{y}_i}{n} = \frac{\sum_{i=1}^n \bar{y} + \hat{\beta}_1(x - \bar{x})}{n} = \frac{n\bar{y} + \hat{\beta}_1 \sum_{i=1}^n (x - \bar{x})}{n} = \frac{n\bar{y} + 0}{n} = \bar{y}$$

Esempio

L'urbanista Palmira Morrone investiga la relazione tra flusso di traffico X (mgl di auto per 24 ore) ed il contenuto di piombo Y nella corteggia degli alberi che fiancheggiano una superstrada (peso a secco in $\mu\text{g/g}$)

| i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|-------|-----|-----|------|------|------|------|------|------|------|------|------|
| x_i | 8.3 | 8.3 | 12.1 | 12.1 | 17.0 | 17.0 | 17.0 | 24.3 | 24.3 | 24.3 | 33.6 |
| y_i | 227 | 312 | 362 | 521 | 640 | 539 | 728 | 945 | 738 | 759 | 1263 |

- Disegnare lo scatterplot;
- Stimare i parametri;
- Calcolare i valori teorici
- Verificare che le proprietà indicate



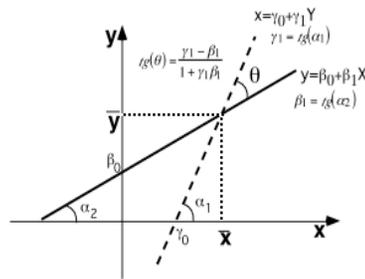
| | y_i | \hat{y}_i | \hat{e}_i |
|-------|----------|-------------|-------------|
| | 227 | 287.48 | -60.48 |
| | 312 | 287.48 | 24.52 |
| | 362 | 424.98 | -62.98 |
| | 521 | 424.98 | 96.02 |
| | 640 | 602.28 | 37.72 |
| | 539 | 602.28 | -63.28 |
| | 728 | 602.28 | 125.72 |
| | 945 | 866.43 | 78.57 |
| | 738 | 866.43 | -128.43 |
| | 759 | 866.43 | -107.43 |
| | 1263 | 1202.94 | 60.06 |
| MEDIE | 639.4545 | 639.4545 | 0.0000 |

Proprietà della retta di regressione/3

il ruolo di esogena ed endogena può essere scambiato:

$$y_i = \beta_0 + \beta_1 x_i + e_i \Rightarrow \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \text{ dove } \begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \end{cases}$$

$$x_i = \gamma_0 + \gamma_1 y_i + e'_i \Rightarrow \hat{x}_i = \hat{\gamma}_0 + \hat{\gamma}_1 y_i \text{ dove } \begin{cases} \hat{\gamma}_0 = \bar{x} - \hat{\gamma}_1 \bar{y} \\ \hat{\gamma}_1 = \frac{S_{xy}}{S_{yy}} \end{cases}$$



Le due rette interpolanti sono legate:

$$\hat{\gamma}_1 = \frac{S_{xy}}{S_{yy}} = \frac{\sqrt{S_{xx}}}{\sqrt{S_{yy}}} * \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}} = \frac{\sqrt{S_{xx}}}{\sqrt{S_{yy}}} * r$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sqrt{S_{yy}}}{\sqrt{S_{xx}}} * \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}} = \frac{\sqrt{S_{yy}}}{\sqrt{S_{xx}}} * r$$

$$\hat{\gamma}_1 * \hat{\beta}_1 = r^2$$

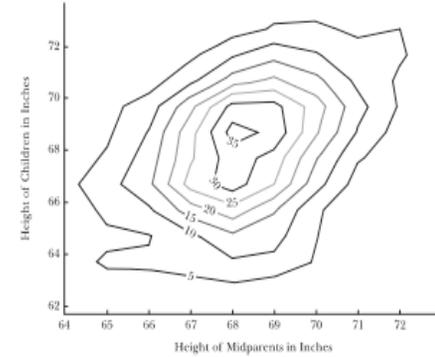
$$tg(\theta) = \frac{r^2 - 1}{r}$$

i coefficienti angolari hanno sempre lo stesso segno per cui le due rette non sono mai perpendicolari

Le due rette sono parallele (e coincidenti) se e solo se Y=X dato che ora tg(theta)=0

Alle origine del concetto di regressione

Figure 4
Galton's Original Smoother+Contour Plots: Contours after Galton Smoother



Sir Francis Galton notò che i figli di padri alti erano più alti della media, ma meno di quanto non eccedessero dalla media i loro padri. I figli di padri bassi erano in media bassi, ma meno bassi della media generale di quanto non lo fossero i padri.

Ipotizzò quindi una generale tendenza al livellamento delle altezze.

L'effetto di regressione

Il principio del ritorno alla media lo si ritrova in varie occasioni

Un docente che loda gli studenti per il buon risultato raggiunto in una prova vedrà un esito peggiore nella prova successiva (Metodo Fata turchina)

il docente che sgrida gli studenti per la pessima riuscita di un test otterrà risultati molto migliori nella seguente prova (Metodo sergente Harman)

Un buon governo sarà seguito da una amministrazione inefficace e ad un premier inadeguato succederà un brillante primo ministro.

Nelle competizioni articolate su due fasi è frequente notare il ribaltamento degli esiti tra la prima e seconda prova: i migliori che peggiorano ed i peggiori che migliorano.

L'effetto di regressione/2

Per ottenere un buon risultato in un'impresa difficile concorrono due fattori:

Talento/Genio

Sorte



il successo in una prova ardua implica che entrambi i fattori hanno agito a favore.

Nella seconda prova il talento/genio magari migliorano o agiscono con la stessa intensità

La Sorte è capricciosa e imprevedibile e non si ripete.

Ed ecco l'effetto di regressione alla media in cui gli scarti si annullano tutti.

Misura dell'adattamento

I minimi quadrati ci garantiscono il miglior adattamento possibile, ma questo potrebbe non essere abbastanza.

Dobbiamo trovare misure standardizzate e normalizzate che siano in grado di quantificare il grado di scostamento tra valori stimati e valori osservati.

Come protagonisti principali avremo

I valori osservati

I valori teorici

Il numero delle osservazioni

Varianza o SQM degli errori

$$s_e = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}} \quad s_e^2 = \frac{S_{yy}}{n-2} [1 - r(x, y)^2]$$

E' nullo solo in caso di perfetta relazione lineare ($r=1$).

Non varia però entro limiti predefiniti. Possiamo solo dire che un adattamento è peggiore di un altro, ma non se un dato adattamento è buono o no

Risente anche delle unità di misura della dipendente. Non è quindi neanche standardizzato.

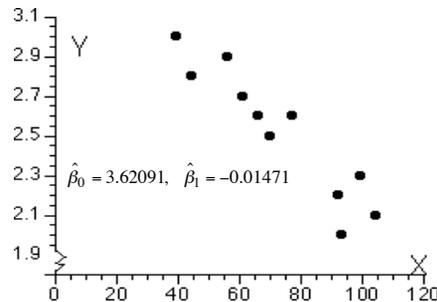
Per il calcolo usiamo quantità già pronte

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{y}_i)^2 &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \sum_{i=1}^n (y_i - \mu_y + \hat{\beta}_1 \mu_x - \hat{\beta}_1 x_i)^2 \\ &= \sum_{i=1}^n (y_i - \mu_y)^2 - \hat{\beta}_1 \sum_{i=1}^n (y_i - \mu_y)(x_i - \mu_x) \\ &= S_{yy} - (\hat{\beta}_1)^2 S_{xx} \end{aligned}$$

Esempio

L'azienda Costantina Tenuta fa parte di una commissione chiamata a valutare una serie di progetti per l'idoneità al finanziamento. Per controllarne la congruità pone in relazione il numero X dei progetti per area e il tempo medio di completamento Y.

| Settori destinatari | Progetti | Tempi medi di compl. |
|----------------------------|----------|----------------------|
| Edilizia demaniale | 105 | 2.1 |
| Opera stradali extraurbane | 94 | 2.0 |
| Disinquinamento | 93 | 2.2 |
| Ferrovie | 78 | 2.6 |
| Edilizia Sanitaria | 71 | 2.5 |
| Edilizia scolastica | 67 | 2.6 |
| Porti commerciali | 62 | 2.7 |
| Infrastrutture urbane | 57 | 2.9 |
| Energia | 45 | 2.8 |
| Smaltimento RSU | 40 | 3.0 |
| Ferrovie Metropolitane | 36 | 3.4 |
| Archivi, Biblioteche | 30 | 3.2 |
| Ferrovie in concessione | 12 | 3.3 |
| Altri | 100 | 2.3 |



$$\sum y = 37.6, \quad \sum y^2 = 103.54, \quad \sum x = 890, \quad \sum x^2 = 67182$$

$$\begin{aligned} s_e &= \sqrt{\frac{S_{yy} - (\hat{\beta}_1)^2 S_{xx}}{n-2}} = \sqrt{\frac{[\sum (\sum y)^2 / n] - (\hat{\beta}_1)^2 [\sum x^2 - (\sum x)^2 / n]}{n-2}} \\ &= \sqrt{\frac{103.54 - (37.6)^2 / 14}{12} - \frac{(-0.01471)^2 [67182 - (890)^2 / 14]}{12}} = 0.1479 \end{aligned}$$

Correlazione teoriche-osservate

Una possibilità di valutare l'adattamento potrebbe basarsi su:

$$\begin{aligned} \frac{Cov(y_i, \hat{y}_i)}{\sigma(y_i)\sigma(\hat{y}_i)} &= \frac{\sum_{i=1}^n (y_i - \bar{y})(\bar{y} + \hat{\beta}_1(x_i - \bar{x}) - \bar{y})}{\sqrt{S_{yy} \sum_{i=1}^n (\bar{y} + \hat{\beta}_1(x_i - \bar{x}) - \bar{y})^2}} = \frac{\hat{\beta}_1 S_{xy}}{\sqrt{S_{yy} \hat{\beta}_1^2 S_{xx}}} \\ &= \frac{\hat{\beta}_1}{|\hat{\beta}_1|} r \end{aligned}$$

Ne consegue che l'adattamento è anche misurabile da:

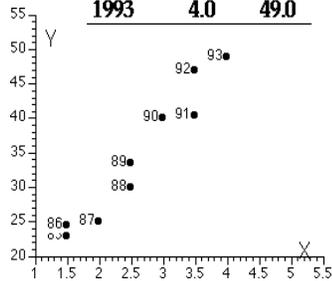
$$\left| \frac{Cov(y_i, \hat{y}_i)}{\sigma(y_i)\sigma(\hat{y}_i)} \right| = \left| \frac{\hat{\beta}_1}{|\hat{\beta}_1|} r \right| = |r|$$

cioè dal valore assoluto del coefficiente di correlazione tra osservate e stimate che coincide con il valore assoluto del coefficiente di correlazione "r" tra X ed Y.

Esempio

La dott.ssa Sarina Bonfiglio, analista finanziario, sta studiando la relazione tra X= Tasso medio sui prestiti nel sistema interbancario e Y=importo della cedola semestrale di un titolo obbligazionario.

| Anni | TMPSI | Ce.Se. |
|------|-------|--------|
| 1985 | 1.5 | 23.0 |
| 1986 | 1.5 | 24.5 |
| 1987 | 2.0 | 25.0 |
| 1988 | 2.5 | 30.0 |
| 1989 | 2.5 | 33.5 |
| 1990 | 3.0 | 40.0 |
| 1991 | 3.5 | 40.5 |
| 1992 | 3.5 | 47.0 |
| 1993 | 4.0 | 49.0 |



- Disegnare lo scatterplot
- Calcolare i parametri
- Misurare l'adattamento con $r(x,y)$
- Supponendo che il dato del 1993 sia non affidabile perché affetto dalla crisi nello SME calcolare il valore interpolato.
- Quale sarà la cedola semestrale se nel 1994 il TMPSI arriva a 5.5?

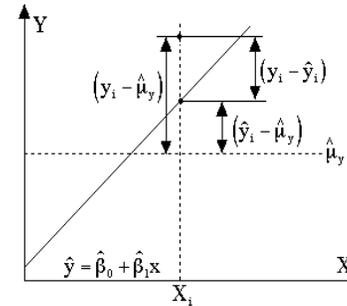
$$\hat{\beta}_0 = 6.448718; \hat{\beta}_1 = 10.602564; s_e^2 = 6.4803$$

$$r(x,y) = 0.9703$$

$$\hat{y}_{93} = 6.448718 + 10.602564 * 4.0 = 48.89$$

$$\hat{y}_{94} = 6.448718 + 10.602564 * 5.5 = 64.76$$

Coefficiente di determinazione (R²)



La variabilità di "Y" può essere scomposta in due parti distinte. Infatti, l'identità

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

rimane anche quando si considerano i quadrati (se la retta è quella dei minimi quadrati)

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n [(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})]^2 = \sum_{i=1}^n [(y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{y})^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y})] \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i - 2\bar{y} \sum_{i=1}^n (y_i - \hat{y}_i) \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2\hat{\beta}_0 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) + 2\hat{\beta}_1 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \end{aligned}$$

Ancora sull'R²

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

Dividendo per "n" si ha la seguente relazione:

$$\text{Varianza totale} = \text{Varianza errori} + \text{Varianza stime}$$

La varianza delle stime è la parte di variabilità (attitudine a presentare modalità diverse) che il nostro modello riesce a spiegare, quella degli errori è la parte che rimane ignota.

$$\text{Varianza totale} = \text{Varianza NON spiegata} + \text{Varianza spiegata}$$

Formula dell' R²

Dividendo i membri per la devianza totale si ha

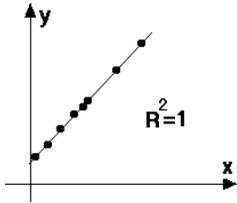
$$1 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} + \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

il 1° addendo è il rapporto tra varianza non spiegata e varianza totale, il 2° è il rapporto tra varianza spiegata e varianza totale.

Questo rapporto è usato come indice della bontà di adattamento ed è noto come il COEFFICIENTE DI DETERMINAZIONE

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

Casi estremi

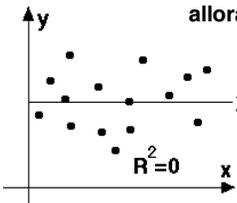


Se tutte le osservate sono allineate su di una retta, teoriche ed osservate coincidono e quindi

$$\text{Se } y_i = \hat{y}_i \text{ per ogni "i"} \Rightarrow R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1;$$

$$\hat{y}_i = \bar{y} + \hat{\beta}_1(x_i - \bar{x})$$

Se la retta di regressione è piatta (coefficiente angolare nullo) allora le teoriche sono tutte pari alla media e quindi



$$\text{Se } \hat{y}_i = \bar{y} \text{ per ogni "i"} \Rightarrow R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - 1 = 0$$

Esempio

Studio della relazione tra il massimo del battito cardiaco sotto stress ed età

| X | Y | (x-M _x) | (y-M _y) | (x-M _x)(y-M _y) | (x-M _x) ² | (y-M _y) ² | y _i | (y - y _i) ² |
|----|-----|---------------------|---------------------|--|----------------------------------|----------------------------------|----------------|------------------------------------|
| 10 | 210 | -25 | 25 | -625 | 625 | 625 | 212.2058 | 4.8656 |
| 20 | 200 | -15 | 15 | -225 | 225 | 225 | 201.3234 | 1.7514 |
| 25 | 195 | -10 | 10 | -100 | 100 | 100 | 195.8822 | 0.7783 |
| 35 | 190 | 0 | 5 | 0 | 0 | 25 | 184.9998 | 25.0020 |
| 40 | 185 | 5 | 0 | 0 | 25 | 0 | 179.5586 | 29.6088 |
| 45 | 175 | 10 | -10 | -100 | 100 | 100 | 174.1174 | 0.7790 |
| 50 | 165 | 15 | -20 | -300 | 225 | 400 | 168.6762 | 13.5144 |
| 55 | 160 | 20 | -25 | -500 | 400 | 625 | 163.2350 | 10.4652 |
| 35 | 185 | | | -1850 | 1700 | 2100 | | 86.7647 |

$$s_e = \sqrt{\frac{86.7647}{6}} = 3.8027$$

$$\beta_0 = 223.0812; \beta_1 = -1.0882;$$

$$R^2 = 1 - \frac{86.7647}{2100} = 0.9572$$

$$r(y, \hat{y}) = \frac{-1850}{\sqrt{1700 * 2100}} = 0.9791$$

$$R^2 = (-0.9791)^2 = 0.9587$$

Esempio

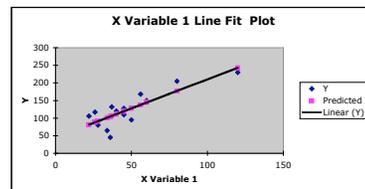
Reddito Superficie SUMMARY OUTPUT

| 22 | 106 |
|-----|-----|
| 26 | 117 |
| 45 | 128 |
| 37 | 132 |
| 28 | 80 |
| 50 | 95 |
| 56 | 168 |
| 34 | 65 |
| 60 | 150 |
| 40 | 120 |
| 45 | 110 |
| 36 | 45 |
| 80 | 205 |
| 120 | 230 |

Si supponga che la proprietaria di un'agenzia immobiliare voglia stabilire la relazione tra reddito familiare e superficie di un appartamento.

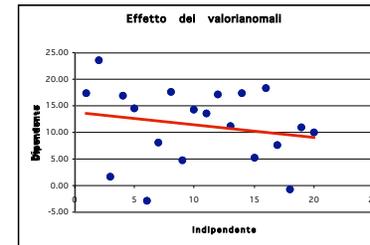
RESIDUAL OUTPUT

| Observation | Predicted Y | Residuals |
|-------------|-------------|--------------|
| 1 | 81.2988 | 24.7012 |
| 2 | 87.9060 | 29.0940 |
| 3 | 119.2901 | 8.7099 |
| 4 | 106.0757 | 25.9243 |
| 5 | 91.2096 | -11.2096 |
| 6 | 127.5491 | -32.5491 |
| 7 | 137.4599 | 30.5401 |
| 8 | 101.1203 | -36.1203 |
| 9 | 144.0671 | 5.9329 |
| 10 | 111.0311 | 8.9689 |
| 11 | 119.2901 | -9.2901 |
| 12 | 104.4239 | -59.4239 |
| 13 | 177.1031 | 27.8969 |
| 14 | 243.1750438 | -13.17504381 |



Effetto dei valori anomali

| X | Y |
|----|--------|
| 1 | 17.27 |
| 2 | 23.49 |
| 3 | 1.72 |
| 4 | 16.87 |
| 5 | 14.58 |
| 6 | -2.76 |
| 7 | 8.04 |
| 8 | 17.65 |
| 9 | 4.75 |
| 10 | 14.39 |
| 11 | 13.61 |
| 12 | 17.22 |
| 13 | 11.22 |
| 14 | 17.39 |
| 15 | 5.31 |
| 16 | 18.39 |
| 17 | 7.64 |
| 18 | -0.73 |
| 19 | 10.86 |
| 20 | 9.89 |
| 27 | 100.00 |
| 34 | 150.00 |
| 41 | 190.00 |
| 48 | 260.00 |

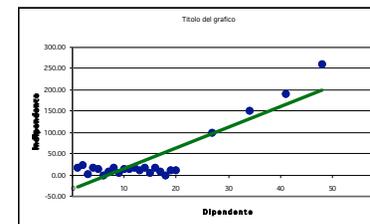


L'impatto della "X" può essere fuorviato dalla presenza di alcuni valori remoti.

| | |
|---------------|---------|
| R multiplo | 0.2064 |
| R al quadrato | 0.0426 |
| R al quadrato | -0.0106 |
| Errore standa | 7.0803 |
| Osservazioni | 20 |

| | Beta | Stat t | p-value |
|---------------|---------|---------|---------|
| Intercetta | 13.9197 | 4.2322 | 0.0005 |
| Variabile X 1 | -0.2457 | -0.8950 | 0.3826 |

I minimi quadrati trascurano il blocco di 20 dati tra i quali non c'è relazione significativa (o è negativa). Invece, pone attenzione ai quattro punti tra i quali c'è una relazione positiva



| | |
|---------------|---------|
| R multiplo | 0.8795 |
| R al quadrato | 0.7735 |
| R al quadrato | 0.7632 |
| Errore standa | 32.7075 |
| Osservazioni | 24 |

| | beta | Stat t | p-value |
|---------------|----------|---------|---------|
| Intercetta | -34.9737 | -3.2383 | 0.0038 |
| Variabile X 1 | 4.9060 | 8.6687 | 0.0000 |

Funzioni lineari

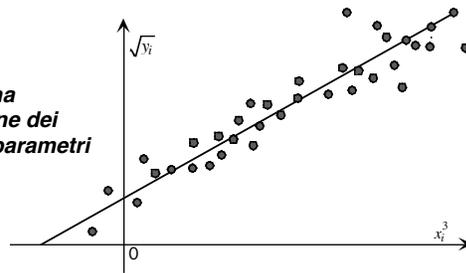
La linearità del modello di regressione è legata solo al modo in cui compaiono i parametri e non alle variabili.

In questo senso i modelli

$$\sqrt{y_i} = \beta_0 + \beta_1 x_i^3 + u_i \quad e^{y_i} = \beta_0 + \beta_1 \ln(|x_i - 7|) + e_i$$

sono lineari dato che i parametri vi compaiono direttamente e con potenza uno.

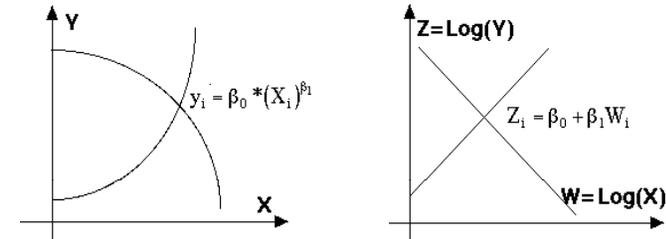
Quello che si riporta sugli assi ha importanza per l'interpretazione dei risultati, non per la stima dei parametri



Funzioni linearizzabili

Il modello di regressione si estende alle forme analitiche LINEARIZZABILI

Si tratta di espressioni che diventano lineari con opportune trasformazioni.



Il modello è linearizzabile se modificando opportunamente gli assi, la relazione appare lineare

Modelli intrinsecamente lineari

Il modello è tale se modificando opportunamente gli assi, la relazione appare lineare, ma non nei parametri originali

ERRORI ADDITIVI

$$y_i = e^a + b^2 x_i + u_i \Rightarrow y_i = \beta_0 + \beta_1 x_i + u_i \quad \beta_0 = e^a, \quad \beta_1 = b^2$$

ERRORI MOLTIPLICATIVI

$$a(z_i)^b = c(w_i)^d e_i \Rightarrow \ln(a) + b \ln(z_i) = \ln(c) + d \ln(w_i) + \ln(e_i)$$

$$\Rightarrow y_i = \beta_0 + \beta_1 x_i + u_i \quad \beta_0 = \frac{\ln(c/a)}{b}, \quad \beta_1 = d/b$$

Da notare che per errori moltiplicativi si deve in genere anche ipotizzare che

$$e_i > 0; \quad E(e_i) = 1$$

Altrimenti sarebbe impossibile la linearizzazione

Esempio

La relazione tra percentuali cumulate di redditi Q_i e percentuali cumulate di redditori P_i può essere rappresentata dalla curva di Lorenz

$$Q_i = P_i^a (2 - P_i)^b e_i$$

Determinare la stima dei parametri

La forma analitica è linearizzabile con la trasformazione seguente:

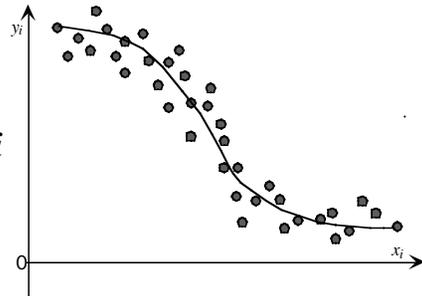
$$\frac{\ln(Q_i)}{\ln(P_i)} = a + b \frac{\ln(2 - P_i)}{\ln(P_i)} + U_i \quad Y_i = \frac{\ln(Q_i)}{\ln(P_i)}; \quad X_i = \frac{\ln(2 - P_i)}{\ln(P_i)}; \quad U_i = \frac{\ln(E_i)}{\ln(P_i)}$$

| P_i | Q_i | Y_i | X_i | $(Y_i - M_y)$ | $(X_i - M_x)$ | $C(x,y)$ | $(X_i - M_x)^2$ | |
|-------|--------|---------|---------|---------------|---------------|----------|-----------------|-----------------|
| 0.1 | 0.0270 | 1.5686 | -0.2788 | -0.3170 | 0.3099 | -0.0982 | 0.0777 | $M_x = 1.8856$ |
| 0.2 | 0.0762 | 1.5996 | -0.3652 | -0.2860 | 0.2235 | -0.0639 | 0.1334 | $M_y = -0.5887$ |
| 0.3 | 0.1354 | 1.6608 | -0.4407 | -0.2248 | 0.1480 | -0.0333 | 0.1942 | $a' = 1.8001$ |
| 0.4 | 0.2059 | 1.7247 | -0.5129 | -0.1609 | 0.0758 | -0.0122 | 0.2631 | $b' = -0.1442$ |
| 0.5 | 0.2874 | 1.7989 | -0.5850 | -0.0867 | 0.0037 | -0.0003 | 0.3422 | |
| 0.6 | 0.3809 | 1.8895 | -0.6587 | 0.0039 | -0.0700 | -0.0003 | 0.4339 | |
| 0.7 | 0.4891 | 2.0052 | -0.7356 | 0.1196 | -0.1469 | -0.0176 | 0.5411 | |
| 0.8 | 0.6147 | 2.1808 | -0.8171 | 0.2952 | -0.2284 | -0.0674 | 0.6676 | |
| 0.9 | 0.7650 | 2.5425 | -0.9046 | 0.6569 | -0.3159 | -0.2075 | 0.8183 | |
| | | 16.9705 | -5.2985 | 0.0000 | 0.0000 | -0.5007 | 3.4715 | |

Modelli non lineari

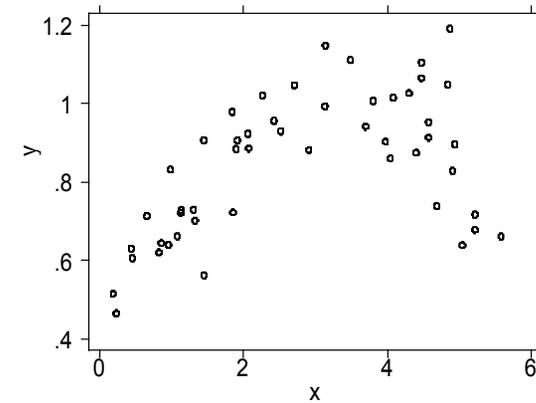
I modelli si dicono NON LINEARI se in nessun modo è possibile ricondurli ad una forma lineare diretta o intrinseca nei parametri

$$y_i = \beta_0 (\beta_1)^{x_i} + e_i$$



In questo caso la stima dei parametri avviene con procedure di ottimizzazione
Non sono semplici da utilizzare, ma diventa sempre più facile utilizzarle

Relazione non lineare



La legge Yerkes-Dodson descrive il legame ad "U rovesciata" tra l'intensità dello stimolo e la qualità attesa della performance.

Regressione per serie evolutive

La situazione è quella di un fenomeno che segua un ordinamento unidimensionale, il cui valore attuale dipende essenzialmente da quelli accaduti in precedenza.

- ☒ L'indice MIB
- ☒ Prospezione verticale di un terreno
- ☒ Spese alimentari

Se $t=1,2,\dots,n$ è l'indice che individua i vari punti nei quali il fenomeno viene rilevato, la regressione per serie evolutive avrà espressione:

$$y_t = \beta_0 + \beta_1 y_{t-1} + u_t$$

dove Y_{t-1} è la cosiddetta variabile "ritardata di lag 1".

Regressione per serie evolutive/2

Si distinguono due situazioni:



STATICA COMPARATA

Le osservazioni sulla endogena e sulla esogena, sono relative allo stesso fenomeno ma rilevato su diverse unità in epoche diverse



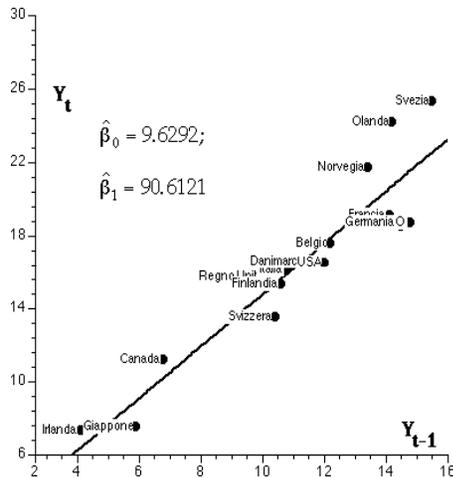
AUTOREGRESSIONE

La variabile esogena è data, per ogni osservazione (cioè in relazione dinamica), dalla endogena ritardata

Esempio: statica comparata

A partire dai dati sulla pressione fiscale in due epoche diverse e per vari paesi "occidentali" determinate i parametri della retta di regressione

| Paesi | Indip. Y_{t-1} | Dip. Y_t |
|-------------|---------------------|---------------|
| Irlanda | 4.10 | 7.30 |
| Olanda | 14.20 | 24.20 |
| Canada | 6.80 | 11.20 |
| Svezia | 15.50 | 25.30 |
| Norvegia | 13.40 | 21.70 |
| Regno Unito | 10.10 | 15.70 |
| Italia | 10.80 | 16.10 |
| Danimarca | 11.50 | 16.60 |
| Finlandia | 10.60 | 15.30 |
| Belgio | 12.20 | 17.60 |
| USA | 12.00 | 16.50 |
| Francia | 14.10 | 19.10 |
| Svizzera | 10.40 | 13.50 |
| Austria | 14.50 | 18.50 |
| Giappone | 5.90 | 7.50 |
| Germania O | 14.80 | 18.70 |



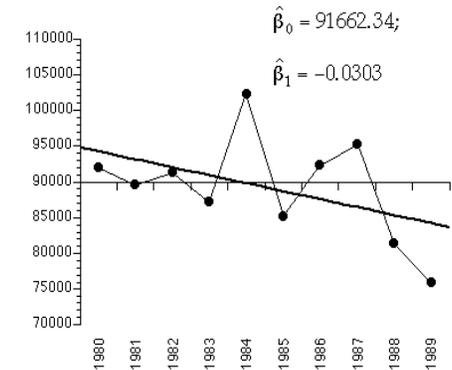
Esempio: autoregressione

Nell'esempio precedente una stessa variabile è osservata in due tempi diversi per le medesime unità. Lo stesso modello può essere applicato in situazioni in cui il valore della dipendente al tempo "t" è legato linearmente al valore della stessa dipendente al tempo "t-1"

$$Y_t = \beta_0 + \beta_1 Y_{t-1} + e_t$$

Produzione di frumento

| Anni | Indip. Y_{t-1} | Dip. Y_t |
|------|---------------------|---------------|
| 1980 | 91952 | 89590 |
| 1981 | 89590 | 91241 |
| 1982 | 91241 | 87174 |
| 1983 | 87174 | 102269 |
| 1984 | 102269 | 85171 |
| 1985 | 85171 | 92287 |
| 1986 | 92287 | 95205 |
| 1987 | 95205 | 81412 |
| 1988 | 81412 | 75882 |



Analisi del trend

il TREND è il sentiero predefinito che si immagina il fenomeno tenda a seguire a meno di piccoli ed incontrollabili errori. Inoltre, se spostato dal trend, tende a ritornarci

Un dato fenomeno è osservato periodicamente e si ipotizza che l'intensità rilevata dipenda proprio dal momento di osservazione

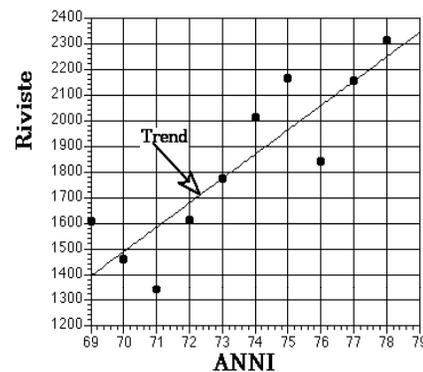
$$Y_t = \beta_0 + \beta_1 t + e_t$$

Pubblicazioni di riviste

| Anni t | Riviste Y_t |
|-----------|------------------|
| 69 | 1603 |
| 70 | 1455 |
| 71 | 1338 |
| 72 | 1605 |
| 73 | 1768 |
| 74 | 2005 |
| 75 | 2160 |
| 76 | 1836 |
| 77 | 2152 |
| 78 | 2309 |

$$\hat{\beta}_0 = -5137.2$$

$$\hat{\beta}_1 = 94.697$$



Analisi del trend/2

In questa formulazione la variabile esogena "t" può variare in un qualsiasi insieme di valori equispaziati:

$$t \in \{1, 2, \dots, n\};$$

$$t \in \{10, 20, \dots, 10 * n\};$$

$$t \in \{1990, 1991, \dots, 1989 + n\};$$

$$t \in \{1.5, 3.0, \dots, 1.5 * n\};$$

$$t \in \{1, 3, 5, 7, \dots, 2n + 1\};$$

Ne consegue che la variabile esogena debba essere interpretata come un insieme fisso di costanti.

Perché fallisce un modello

ERRATA TEORIZZAZIONE

Le relazioni ipotizzate non reggono alla prova dei fatti per cui il modello non si conforma alla realtà osservata (ad esempio manca una variabile o è inserita un'altra non pertinente ovvero si è forzata la linearità)

E' difficile accertare l'influenza di tale eventualità.

Il modello è una visione semplificata della realtà, ma potrebbe esserne una visione semplicistica

Voto d' esame = f(Simpatia ispirata agli esaminatori)

Produzione agraria = f(Entità delle piogge)

Scorte magazzino = f(Vendite)

Si tratta di una limitazione intrinseca alla modellistica che si controlla solo presupponendo la validità del modello.

Perché fallisce un modello/2

ERRATA FORMULAZIONE

Le variabili sono state correttamente individuate, ma usate in modo sbagliato.

Ad esempio, la curva di Gompertz, spesso usata dagli attuari, per la costruzione delle tavole di mortalità ha equazione:

$$y = \beta_0 * e^{-\beta_1 e^{-\beta_2 X}}$$

Se però allo scatterplot viene adattato il modello

$$y = \beta_0 + \beta_1 X + \beta_2 X^2$$

le sue capacità esplicative saranno limitate ed occorre riformulare il modello.

Perché fallisce un modello/3

SCARSA QUALITÀ DEI DATI

Se i dati acquisiti sui fenomeni coinvolti nel modello sono inattendibili sarà scadente anche il modello

GARBAGE IN -----> GARBAGE OUT



I risultati di una elaborazione statistica non possono essere più attendibili dei dati da essa utilizzati

Invece di utilizzare un numero indice sintetico dei prezzi per l'intera collettività nazionale si utilizza un indice per la scala mobile dei salari.

Gli strumenti di misurazione contengono errori sistematici o sono stati volontariamente alterati

Perché fallisce un modello/4

ERRORE NELLA PROCEDURA STATISTICA

Se i dati sono contaminati da errori di rilevazione e/o di misurazione in misura moderata è ancora possibile ottenere buoni risultati.

Si devono però utilizzare tecniche statistiche robuste rispetto a questo tipo di difetti.

Ovvero procedure che filtrino gli errori lasciando la buona sostanza delle informazioni acquisite.