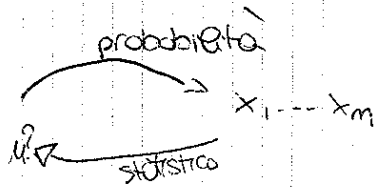
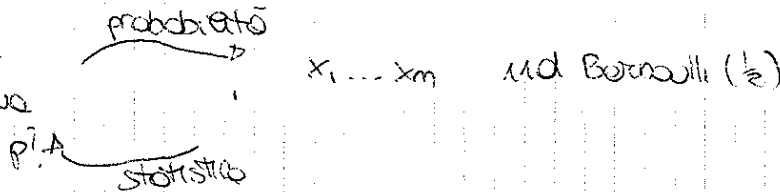


STATISTICA - STIMA

distribuzione di fondo (o popolazione) per esempio $N(202, 14^2)$ nell'esempio contestuale



oppure dato una Bernoulli ($\frac{1}{2}$) cioè numero equo



La statistica inferenziale va dal particolare (campione) al generale (popolazione o distribuzione), dal noto all'ignoto

Se non conosco un parametro, cioè una caratteristica della distribuzione, lo posso **STIMARE** per esempio, \bar{x} stima (è stimatore di) μ .

Similmente, se non conosciamo σ^2 , lo stimiamo tramite la varianza campionaria

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{m-1}$$

E se x_1, \dots, x_m iid Bernoulli (p) e non conosciamo p , lo stimiamo con \hat{p} = frazione di successi su m test-trials

09-11-08

x_1, \dots, x_m iid con densità $f(x)$ si dice campione casuale

Una sua funzione è una STATISTICA.

Nella realtà, spesso i parametri della $f(x)$ sono incogniti. Possiamo usare delle statistiche per **STIMARLI**

Esempio: (METEOROLOGIA)

μ = misurando l'incognito (per esempio la distanza da una stella)

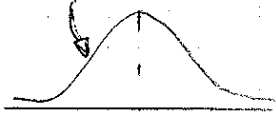
x_1, \dots, x_m m misurazioni indipendenti di μ , ciascuna con densità normale (μ, σ^2)

σ^2 è la varianza e dipende in maniera inversa dalla precisione ($= \frac{1}{\sigma^2}$) dello strumento di misura

in generale $\bar{x} = \frac{\sum x_i}{m}$, la media campionaria, è uno stimatore di μ .

la densità campionaria di $\bar{x} \sim N(\mu, \frac{\sigma^2}{m})$

↑ c'è un fattore $\frac{1}{\sqrt{m}}$ che riduce la deviazione standard.



\bar{x} è centrato su $\mu \rightarrow E(\bar{x}) = \mu \rightarrow$ si dice che \bar{x} è uno stimatore non DISTORTO (NON BIASATO)

la variabilità di \bar{x} intorno a μ è misurata dalla deviazione standard della sua distribuzione campionaria

$SD(\bar{x}) = \sqrt{\text{var}(\bar{x})} = \frac{\sigma}{\sqrt{m}}$ chiamato anche ~~errore standard~~ **ERRORE STANDARD** di \bar{x} .

Esempio 2: sappiamo che da non conoscere σ^2 , lo possiamo stimare con la varianza campionaria

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{m-1}$$

non si divide per m , ma per $m-1$ così il valore atteso di S^2 è centrato $E(S^2) = \sigma^2$ **NON DISTORTO**

$$\text{infatti } E \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{1}{n-1} E \left(\sum (x_i^2 - 2x_i \bar{x} + \bar{x}^2) \right) =$$

$$= \frac{1}{n-1} \left(\sum E(x_i^2) - 2 E \left(\bar{x} \sum_{i=1}^n x_i \right) + n E(\bar{x}^2) \right) =$$

$$E(x_i^2) - (E(x_i))^2 = \text{Var}(x_i) = \sigma^2$$

$$\sqrt{E(x_i)^2 = \sigma^2 + \mu^2}$$

$$= \frac{1}{n-1} \left(\sum (\sigma^2 + \mu^2) - 2n E(\bar{x}^2) + n E(\bar{x}^2) \right) =$$

Analogamente $E(\bar{x}^2) = \text{Var}(\bar{x}) + (E(\bar{x}))^2 = \frac{\sigma^2}{n} + \mu^2$

$$= \frac{1}{n-1} \left(n\sigma^2 + n\mu^2 - n \left(\frac{\sigma^2}{n} + \mu^2 \right) \right) =$$

$$= \frac{1}{n-1} \left(n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2 \right) = \frac{(n-1)\sigma^2}{n-1} = \sigma^2$$

S^2 è stimatore corretto di σ^2

Se x_1, \dots, x_n i.i.d. Bernoulli (p), ~~è~~ e p è incognito, lo stimiamo usando la
proporzione campionaria di successi $\hat{p} = \frac{\text{nr. di successi}}{n} = \frac{\sum x_i}{n}$

(caso particolare di \bar{x} quando le x_i sono binarie)

\hat{p} è centrato su p , e la sua varianza è $\frac{p(1-p)}{n}$

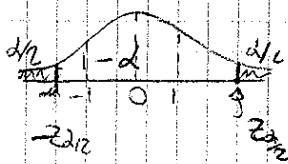
l'errore standard di \hat{p} è $\sqrt{\frac{p(1-p)}{n}} \leq \sqrt{\frac{1}{4n}}$

INTERVALLO DI CONFIDENZA

Vogliamo arrivare ad una STIMA INTERVALLARE di un parametro
INTERVALLO di CONFIDENZA

partiamo dalla distribuzione campionaria di $\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ quando x_1, \dots, x_n i.i.d. $N(\mu, \sigma^2)$.

Possiamo scrivere $\frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$



$$P\left(-z_{\alpha/2} \leq \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

risultato dalla
distribuzione campionaria
di \bar{x}

$$P\left(-z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \bar{x} - \mu \leq z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha \quad (\times \text{ sempre } \sigma/p)$$

$$P\left(-\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq -\mu \leq -\bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

$$P\left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

L'intervallo casuale (destinato)

$\left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$ contiene ("cattura") μ con probabilità $1 - \alpha$

lo chiamiamo INTERVALLO DI CONFIDENZA di quello $1 - \alpha$

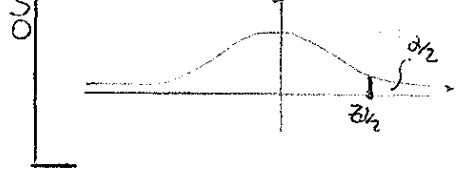
Brevemente possiamo scrivere $\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ intervallo

L'ampiezza dell'intervallo (che è centrato su \bar{x} , stimatore di μ) è $2z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

se n cresce, l'ampiezza decresce. Se σ cresce, anche l'ampiezza cresce

se è quello $1 - \alpha$ a crescere, cresce pure l'ampiezza

Dato un campione casuale x_1, \dots, x_m normale $N(\mu, \sigma^2)$ con σ^2 nota, un i.c. (= intervallo di confidenza) per μ è $(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{m}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{m}})$ abbreviato in $\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{m}}$ dove $z_{\alpha/2}$ è



Per esempio (8.5): il segnale μ è stimato 13.3 ± 1.86 con confidenza del 95%

Usiamo il termine confidenza per indicare il livello di affidabilità della procedura

Non possiamo parlare di "probabilità" perché μ non è aleatorio

Di tutti i possibili campioni che potrei osservare, il 95% (cioè $100(1-\alpha)\%$) contiene μ mentre il 5% (500%) non lo fa.

L'ampiezza dell'intervallo è $2 z_{\alpha/2} \frac{\sigma}{\sqrt{m}}$. Tale ampiezza può servire per determinare m , la numerosità campionaria necessaria (n° di osservazioni).

Supponiamo di desiderare un'ampiezza D fissa

$$2 z_{\alpha/2} \frac{\sigma}{\sqrt{m}} \leq D$$

$$2 z_{\alpha/2} \frac{\sigma}{D} \leq \sqrt{m}$$

$$\Rightarrow \left(2 z_{\alpha/2} \frac{\sigma}{D} \right)^2 \leq m$$

FORMULA PER PROGRAMMARE L'ESPERIMENTO

Esempio 8.3 pag 324

A volte conviene usare INTERVALLI UNILATERALI, per esempio $(b, +\infty)$ per profitti (e altre caratteristiche del tipo "più ce n'è, meglio è") oppure $(-\infty, a)$ per la situazione "meno ce n'è meglio è"

Per esempio, dalla disuguaglianza probabilistica $p\left(\frac{\bar{X}-\mu}{\sigma/\sqrt{m}} \leq z_\alpha\right) = 1-\alpha$ discende l'intervallo di confidenza unilaterale $(\bar{x} - z_\alpha \frac{\sigma}{\sqrt{m}}, \infty)$ di livello $1-\alpha$

↳ si dice anche che $\bar{x} - z_\alpha \frac{\sigma}{\sqrt{m}}$ è un estremo di confidenza inferiore di livello $1-\alpha$

Simmetricamente $(-\infty, \bar{x} + z_\alpha \frac{\sigma}{\sqrt{m}})$ i.c. unilaterale sinistro,

ovvero $\bar{x} + z_\alpha \frac{\sigma}{\sqrt{m}}$ è un estremo di confidenza superiore di livello $1-\alpha$

Cosa succede se σ non è noto? Dobbiamo stimarlo, per esempio usando S ,

la radice della varianza campionaria

Ritroppo $\frac{\bar{X}-\mu}{\sigma/\sqrt{m}} \sim N(0,1)$, ma $\frac{\bar{X}-\mu}{S/\sqrt{m}} \sim N(0,1)$ [vero approssimativamente se n è grande, ma non per n piccolo]
 σ/\sqrt{m} è una v.a. fissa e meno variabile, S/\sqrt{m} è una v.a. e + variabile

Gosset (piumi '800) scoprì in base alla distribuzione di S (vedi 7.6),

$\frac{\bar{X}-\mu}{S/\sqrt{m}}$ ha una distribuzione normale standard, bensì "t" con $m-1$ gradi di libertà

$\bar{x} \pm t_{m-1, \alpha/2} \frac{S}{\sqrt{m}}$ è un intervallo di confidenza per μ quando σ è incognito e m piccolo.

$t_{m-1, \alpha/2}$ è un punto percentile analogo a $z_{\alpha/2}$ tratto dalla tavola di t

Analogamente, $\bar{x} - t_{m-1, \alpha} \frac{s}{\sqrt{m}}$ è un estremo INFERIORE di confidenza di quello $1-\alpha$

$\bar{x} + t_{m-1, \alpha} \frac{s}{\sqrt{m}}$ è un estremo SUPERIORE di confidenza di quello $1-\alpha$

• x_1, \dots, x_m iid Bernoulli(p). Pienzo, voglio stimare p .

Una stima puntuale di p è $\hat{p} = \frac{\sum x_i}{m}$ = proporzione campionaria di successi.

Adesso costruiamo un i.c. per p .

Abbiamo visto che, per il Teorema del Limite Centrale, se m grande,

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{m}\right)$$

Quindi, ripetendo la procedura paradigmatica, otteniamo

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{p(1-p)}{m}}$$

Problema: non conosciamo p , è proprio quello che stiamo stimando

Una prima soluzione è usare $\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{m}}$

che richiede di almeno due approssimazioni: \rightarrow T.L.C. (Teorema del Limite Centrale)
 \rightarrow usare doppiamente \hat{p} al posto di p

Un'alternativa CONSERVATIVA è usare il massimo possibile di $\frac{p(1-p)}{m}$, cioè $\frac{1}{4m}$, ottenendo $\hat{p} \pm \frac{z_{\alpha/2}}{2\sqrt{m}}$ che dà intervalli molto ampi.

Un estremo inferiore di confidenza approssimato di quello $100(1-\alpha)\%$ è

$$\hat{p} - z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{m}}$$

Un estremo superiore di confidenza approssimato di quello $100(1-\alpha)\%$ è

$$\hat{p} + z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{m}}$$

VERIFICA DI IPOTESI

19-11-09

Un'ipotesi è un'affermazione riguardo a un meccanismo probabilistico che non conosciamo completamente.

In particolare, un'ipotesi statistica è un'affermazione riguardo a un parametro incognito che regola la distribuzione di probabilità di un campione osservabile.

Esempio: Una certa moneta è equa??

- ipotesi (nulla): la moneta è equa H_0
- ipotesi (alternativa): la moneta non è equa H_1
(ma favorevole teste)

Noi possiamo osservare le componenti X_1, \dots, X_n iid Bernoulli (p) con p incognita

$$H_0: p = \frac{1}{2}$$

$$H_1: p \neq \frac{1}{2}$$

Un test di ipotesi è una procedura che in base ai dati rifiuta o non rifiuta H_0 .
 L'idea è di rifiutare H_0 se i dati forniscono sufficiente evidenza contro H_0 .

Procediamo con l'osservazione del campione: TCTTT

se $p = \frac{1}{2}$ $P(\text{damo 4 teste su } n=5 \text{ era}) = \binom{5}{4} \left(\frac{1}{2}\right)^4 \left(1 - \frac{1}{2}\right)^{5-4} + \binom{5}{5} \left(\frac{1}{2}\right)^5 \left(1 - \frac{1}{2}\right)^{5-5}$
 \uparrow H_0 è vero \downarrow $= P(X \geq 4 \text{ con } p = \frac{1}{2}) = 1 - P(X < 3 \text{ con } p = \frac{1}{2}) = 0.1875$
 $= 1 - \text{pbinom}(3, 5, 0.5)$

C'è il 19% di probabilità di osservare un numero di teste pari a 4 o più se H_0 è vero.

Questo probabilità è chiamato **p-value** (valore-p)

Quindi per ora non rifiutiamo H_0 .

In generale,

p-value = probabilità di osservare dati così estremi o più estremi di quelli disponibili se H_0 fosse vero.

pag. 386 es 3.8

$$H_0: p = \frac{1}{2} \quad (\text{meglio } H_0: p \leq \frac{1}{2})$$

$$H_1: p > \frac{1}{2}$$

p = proporzione di preoccupati

Spesso:

H_0 = "status quo"

H_1 = quello che il ricercatore vuole dimostrare

$n = 920$ osserviamo $x = 478$ sì $\rightarrow 1 - \text{pbinom}(477, 920, .5)$

$p\text{-value} = P(X \geq 478 \mid p = \frac{1}{2}) \approx 0.1242$ è abbastanza alto \Rightarrow non rifiutiamo H_0 .

Il **p-value** è una misura probabilistica di compatibilità dei dati con l'ipotesi nulla H_0 .
 Nell'esempio 3.8, $p\text{-value} = 0.12$, ancora abbastanza alto.

Idea del test di ipotesi: identifico una regione di rifiuto R che mi porta a rifiutare H_0 se i dati sono in R , e non rifiuto altrimenti.

R sarà una regione di **BASSA COMPATIBILITÀ** DEI DATI con H_0 (basso p-value).

	H_0 vero	H_0 non vero $\rightarrow H_1$
refutare H_0	errore di prima specie	OK
non refutare H_0	OK	errore di seconda specie

In genere si fissa la probabilità di errore di 1^a specie a un livello condizionale

$$(\alpha = 0.05, \alpha = 0.01, \alpha = 0.001)$$

invece azioni basate sul dato

Nell'esempio 3.8 esteso, se scegliamo come soglia 486, cioè

$$R = \{m : m \geq 486\}$$

allora l'errore di 1^a specie è $P(\text{refuto } H_0 \mid H_0 \text{ vero}) = 0.046 \approx 0.05$

Refuterò se $x \geq 486$, ovvero se $p\text{-value} \leq 0.05$

Vediamo come applicare questi concetti al caso del test di ipotesi sul valore atteso di una distribuzione normale

Il sistema di ipotesi che consideriamo è:

$$\begin{aligned} H_0 &: \mu = \mu_0 \\ H_1 &: \mu \neq \mu_0 \end{aligned}$$

dove μ_0 è un certo valore fisso di interesse

Esempio 3.1: $\mu_0 = 10$

$$\begin{aligned} H_0 &: \mu = 10 \\ H_1 &: \mu \neq 10 \end{aligned}$$

La regione di rifiuto corrisponderà a \bar{x} molto lontano da μ_0 , cioè:

(1) $|\bar{x} - \mu_0|$ grande

(2) cioè $\left| \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \right|$ è grande

Quanto grande? Tanto grande che $P\left(\left| \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \right| \geq z_{\alpha/2} \mid H_0\right) = \alpha$
probabilità errore di 1^a specie

Regione con (2) invece che con (1) ci permette di individuare $z_{\alpha/2}$ come la soglia che cerchiamo

Ponendo $\alpha =$ probabilità dell'errore di 1^a specie (o quello di significatività) uguale a un valore basso

23-11-09

Se x_1, \dots, x_n i.i.d. $N(\mu, \sigma^2)$ σ^2 noto
oppure

se n è grande e x_1, \dots, x_n i.i.d. con valore atteso incognito μ , possiamo costruire un test dell'ipotesi

$H_0: \mu = \mu_0$ contro l'alternativa $H_1: \mu \neq \mu_0$ con μ_0 valore di "statistico"

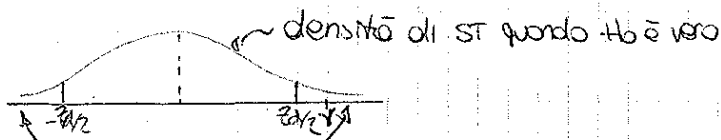
Costruiamo la statistica del test (ST):

$$ST = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

e rifiutiamo quando $|ST| \geq z_{\alpha/2}$

REGOLA DI DECISIONE
(la probabilità di errore di 1^a specie è α)

Il p-value è invece



regioni di rifiuto di R : insieme di valori con poco supporto per H_0

Supponiamo ora di fare il test e osservare il valore γ per lo ST
 ↳ nel caso in figura $\gamma > z_{\alpha/2} \rightarrow$ quindi rifiuto H_0

il valore-p è $P(Z > \gamma)$

\rightarrow p-value = il più piccolo valore di α che porterebbe al rifiuto di H_0

Test di ipotesi unilaterale: $H_0: \mu = \mu_0$ ovvero $H_0: \mu \leq \mu_0$
 $H_1: \mu > \mu_0$

Esempio 9.5

X_1, \dots, X_m sono le variazioni del livello di colesterolo in 40 pazienti.

X_i = colesterolo dopo + colesterolo prima del trattamento

$H_0: \mu = 0$

6.8 = riduzione medio osservata, valore osservato di \bar{x}

12.1 = valore osservato di S

$$T = \frac{m(\bar{x} - \mu_0)}{S} = 3.55 \quad \text{p-value} = 0.0001$$

CAPITOLO 10 (10.2, 10.3)

$X_1, \dots, X_m \text{ i.i.d. } N(\mu_x, \sigma_x^2)$
 $Y_1, \dots, Y_m \text{ i.i.d. } N(\mu_y, \sigma_y^2)$
 \rightarrow indipendenti

Vogliamo testare $H_0: \mu_x = \mu_y$

La statistica test $\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_x^2}{m} + \frac{\sigma_y^2}{m}}}$ se σ_x^2 e σ_y^2 sono note

(1)
$$\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_x^2}{m} + \frac{\sigma_y^2}{m}}}$$

oppure (2)
$$\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{m}}}$$
 se m e m grandi, per il TLC

www.opesoonline.com

Se H_0 è vero, lo ST si comporta come $Z \sim N(0, 1)$

quindi \rightarrow tabella pg. 403 per (1)
 \rightarrow tabella pg. 416 per (2)

Rifiuto H_0 se ST è "estremo", cioè su uno delle due code per il test bilaterale, o su uno delle due per il test unilaterale

Esempio 10.3

In questo esempio poiché: $H_0: \mu_x \leq \mu_y$ rifiuto H_0 per valori \bar{X} molto
 $H_1: \mu_x > \mu_y$

più grandi di \bar{Y} , cioè per valori $\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_x^2}{m} + \frac{\sigma_y^2}{m}}}$ sulla coda destra

poiché osservo $\gamma = 1.34 \rightarrow$ p-value = 0.08

poiché p-value è alto rifiuteremo H_0 solo per valori di α maggiori di 0.08, troppo alti
 quindi non rifiuto H_0 (la prova è non significativa)

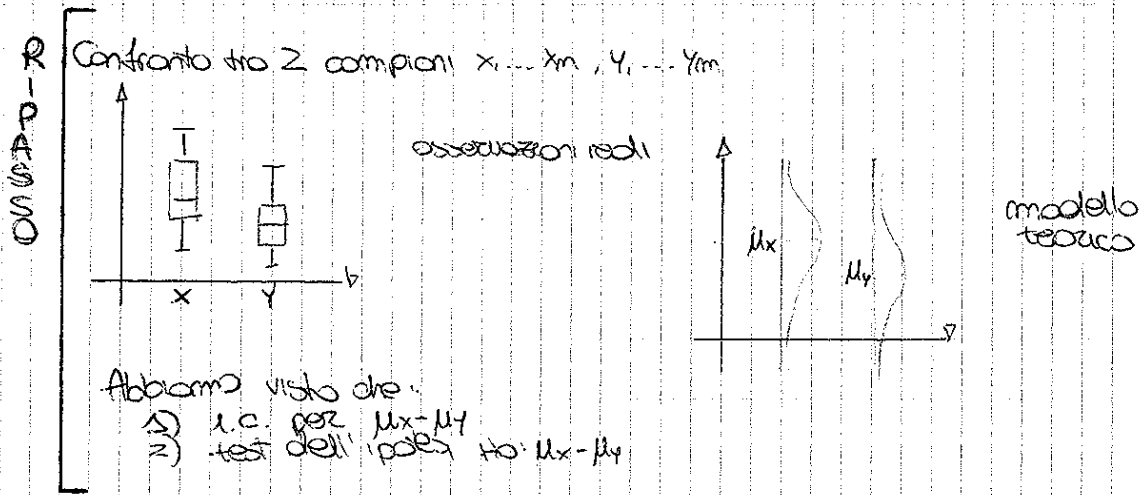
Relazione tra i.c. e test: spesso un test su un parametro θ con ipotesi nulla della forma $H_0: \theta = \theta_0$ rifiuta H_0 se e solo se un appropriato i.c. per θ non contiene θ_0

Esempio: $H_0: \mu_x = \mu_y$ ovvero $H_0: \mu_x - \mu_y$
 $H_1: \mu_x \neq \mu_y$

Test (o quello d) rifiuta H_0 se un i.c. (di quello 1-d) per $\mu_x - \mu_y$ non contiene 0

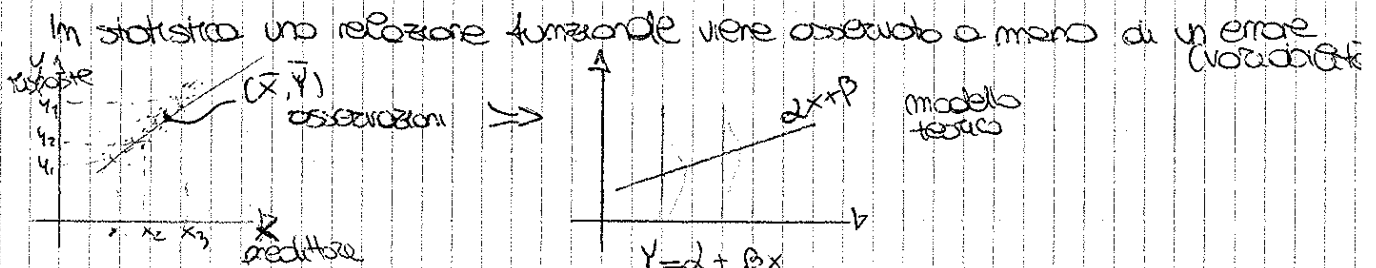
REGRESSIONE - CAPITOLO 12

30-11-08



Oggi parleremo di dipendenza di una variabile risposta Y da una variabile predittore X .

Relazione piú semplice: lineare $\rightarrow Y = a + bx$



Esempi:

X	Y
alimento	consumo
altezza padre	altezza figlio

$E(Y) = a + \beta x$ ovvero $Y = E(Y) + \text{errore} = a + \beta x + e$

REGRESSIONE LINEARE SEMPLICE

dipendenza statistica $E(Y) = f(x)$

$f(x)$ lineare

$f(x)$ funzione della sola x (non stocastica = non deontica)

Il predittore x è una variabile sotto il controllo dello sperimentatore
 la risposta Y è una variabile aleatoria (stocastica) il cui valore atteso è funzione di x , $E(Y) = f(x)$

Esempio 12.4

DATI: osservazioni y_1, \dots, y_m corrispondenti a diversi livelli del predittore x_1, \dots, x_m

OGGETTO DI INFERENZA: i parametri α e β e il varianza dell'errore

METODO DI STIMA di α e β : metodo dei minimi quadrati

↓ trova α e β in modo da minimizzare

$$\sum (y_i - (\alpha + \beta x_i))^2$$

SCARTO DELLA RETTA

La somma dei quadrati degli scarti della retta

$$\min_{\alpha, \beta} \sum (y_i - (\alpha + \beta x_i))^2$$

Si ottiene la seguente soluzione in forma chiusa:

$$\begin{cases} \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \\ \hat{\beta} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{cases}$$

La retta stimata che ne risulta $\hat{Y} = \hat{\alpha} + \hat{\beta}x$ è chiamata **RETTA DI REGRESSIONE STIMATA** o **RETTA DEI MINIMI QUADRATI**.

Il modello di **REGRESSIONE LINEARE SEMPLICE**

$$E(Y) = \alpha + \beta x \quad \text{o} \quad Y = \alpha + \beta x + e$$

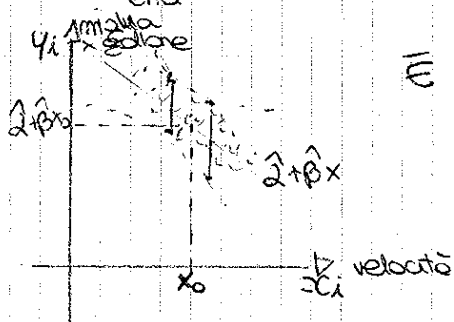
→ stima dei minimi quadrati $\hat{\alpha}$ e $\hat{\beta}$

Se, in aggiunta, supponiamo che e e Y sono normali, con lo stesso varianza incognita σ^2

$$\text{Bisogna stimare } \sigma^2 = \hat{\sigma}^2 = \frac{\sum (y_i - \hat{\alpha} - \hat{\beta} x_i)^2}{m-2} = \frac{SS_R}{m-2}$$

↳ mean square error

parametro incognito	stima
α	$\hat{\alpha}$
β	$\hat{\beta}$
σ^2	$\hat{\sigma}^2$
$\alpha + \beta x_0$	$\hat{\alpha} + \hat{\beta} x_0$



È possibile anche costruire degli i.c. per i parametri incogniti

$$E(Y; x_0) = \alpha + \beta x_0 = \text{valore atteso di } Y \text{ corrispondente al nuovo valore di } x_0$$

? TOD i.c. non sono nel testo ma sono ampiamente disponibili

INTERVALLI DI PREDIZIONE

02-12-08

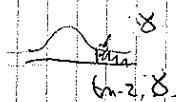
Le differenze $y_i - \hat{\alpha} - \hat{\beta} x_i$ sono chiamate **residui**

Una volta osservati i valori y_i di Y_i i residui diventano i numeri $y_i - \hat{\alpha} - \hat{\beta} x_i$

Test di ipotesi: $H_0: \beta = 0$ (x non ha influenza su y)

Rifiutiamo H_0 se $\hat{\beta}$ è troppo distante da 0, ovvero

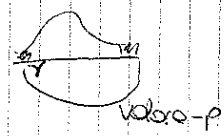
rifiutiamo H_0 se $\sqrt{\frac{(m-2)SS_R}{SS_B}} \hat{\beta} = ST$ è più grande, in valore assoluto,
 statistico test

di t_{n-2}, δ ... 

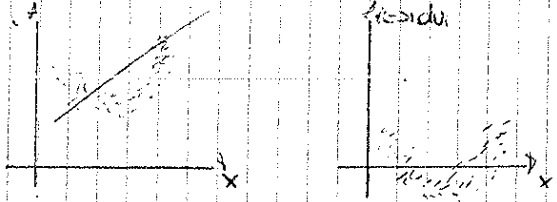
Equivalentemente, il valore-p del test è il seguente. Supponiamo di osservare

$$ST = Y$$

Allora valore-p = $P(|T_{n-2}| \geq |Y|)$



Rifutiamo H_0 se valore-p < δ

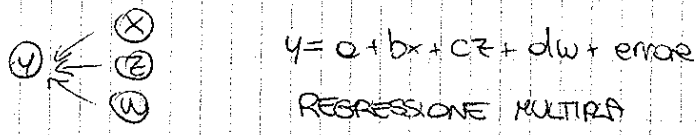


esempio di grafico dei residui che denota una deficienza dell'interpolazione lineare

Sospetto: Y dipende in maniera non lineare da x

Per esempio $y = \log x$ oppure $\log y = x$ \rightarrow trasformazioni dei dati
 oppure $y = a + bx + cx^2$ \rightarrow due predittori = regressione multipla

Se y dipende da x e da altre variabili z, w



REGRESSIONE MULTIPLA:

(*) $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$

Labels: Y (variabile risposta), β_0 (intercetta), β_1, \dots, β_k (coefficienti di regressione), x_1, \dots, x_k (k predittori x_i)

Partendo da una relazione di regressione multipla del tipo (*) possiamo chiederci

- 1) Ci sono alcuni predittori inutili? (e qualche $\beta = 0$?)
- 2) abbiamo dimenticato qualche predittore importante?
- 3) e la relazione lineare sufficientemente buona?

1) \rightarrow test di ipotesi sui singoli β : $H_0: \beta_1 = 0$

Sulla base dei dati calcoliamo le stime dei minimi quadrati

$\beta_0, \beta_1, \dots, \beta_k$ di $\beta_0, \beta_1, \dots, \beta_k$
 minimizzanti $\sum (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_k x_{ik})^2$

con x_{ij} = valore che il predittore j -esimo (j -mo colonna) assume nell'unità statistica i -esima (i -esimo rigo)

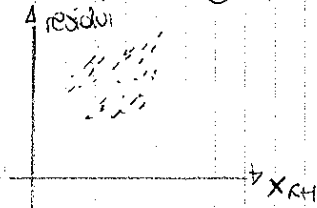
Il test dell'ipotesi $H_0: \beta_1 = 0$ viene fatto usando β_1 e altre informazioni...

Il valore-p per questa ipotesi può essere letto nell'output del programma di regressione

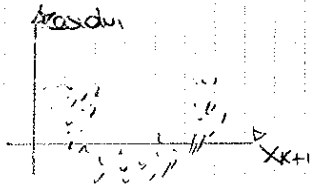
2) è più difficile, usiamo i residui. i residui sono:

$$(Y_i - B_0 - B_1 X_{i1} - \dots - B_k X_{ik} \quad i = 1 \dots n)$$

Potrebbe fare un grafico di questi residui su un eventuale ulteriore predittore X_{k+1}



no pattern = X_{k+1} non spiega nulla in più.



andamento particolare: X_{k+1} "spiega" i residui

COEFFICIENTE DI DETERMINAZIONE

3) i residui servono anche per valutare la bontà di un fit di regressione tramite il coefficiente di determinazione

$$R^2 = \frac{S_{yy} - SSR}{S_{yy}} = \frac{\text{varianza di } Y \text{ spiegata dall'intersezione } \hat{y}}{\text{varianza totale di } Y}$$

dove $S_{yy} = \sum (Y_i - \bar{Y})^2$ varianza totale di Y

$SSR = \sum (Y_i - \hat{Y}_i)^2$ varianza di Y "spiegata" dall'intersezione \hat{Y}

Più R^2 si avvicina a 1, più buona è la regressione

R^2 si usa anche per la regressione multipla. Se c'è un solo predittore X ,

$R^2 =$ quadrato del coefficiente di correlazione lineare tra Y e X

$$-1 \leq R \leq 1 \quad \text{solo univariato}$$

$$0 \leq R^2 \leq 1 \quad \text{anche multivariato}$$