

Appunti di Statistica

Prof. di Gesu

Appunti del corso
Prof. Maurizio Pratelli

- lezioni su teams - Registra (ma è meglio registrare)

- Usa e-book

- mail: giacomo.digesu@unipi.it

(testo: "Probabilità e statistica per l'ingegneria e le scienze" S. Ross)

Statistica → scienza dei dati (analisi e interpretazione)

↳ al giorno d'oggi ce n'è una quantità enorme.

In questo corso → Rudimenti di "statistica classica"

statistica descrittiva

statistica inferenziale

STATISTICA DESCRITTIVA

Abbiamo un insieme molto grande di dati

Problema: come li descriviamo in modo sintetico, efficiente, suggestivo?

esempio neuroscienziati cercano di comprendere quali parti del cervello si occupano di elaborare l'informazione visuale data dai colori.



Esperimento: presento una persona e danno uno stimolo visuale che consiste nel mostrare delle figure in movimento alternando figure a colori e figure in B/N, in istanti $t=1, \dots, T$ fanno una tomografia a risonanza del cervello.

- Otteniamo una quantità enorme di dati.

- Calcolano "correlazioni" tra l'input dato dallo stimolo visuale e la risposta $x(t)$ al variare del tempo.

- selezionano i cuvetti che sono fortemente correlati con lo stimolo

- generano un'immagine che evidenzia in rosso le parti del cervello selezionate.

Tutta l'analisi descrittiva è stata fatta usando il software R

STATISTICA INFERENZIALE

Partiamo sempre da un insieme di dati. Vogliamo utilizzare i dati per:

- fare delle previsioni

↳ basarci su un modello matematico

- prendere delle decisioni

↳ della realtà che vogliamo descrivere.

I modelli matematici considerati sono probabilistici.

Teoria della probabilità ci fornisce il linguaggio fondamentale per poter fare della statistica inferenziale.

La teoria della probabilità è una teoria astratta che descrive gli "esperimenti aleatori".

Esempi di esperimenti aleatori:

• lancio dado / moneta

• osservazione del meteo nei giorni futuri

• osservazione velocità di propagazione di un virus nel prossimo anno.

Vedremo due teoremi fondamentali:

- legge dei grandi numeri

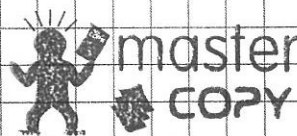
- teorema limite centrale



- esempio ① Problema dell' "Exit Poll": la popolazione vota per una carica politica da due candidati A e B.
 (statistica inferenziale) Viene estratto un campione (a caso) dalla popolazione che ci dice la preferenza espressa nel voto.
 Vogliamo fare una previsione sull'esito del voto basandoci sulle indicazioni del campione estratto.
 (concetto di "stimatori", "intervalli di confidenza")
- ② Raccogliamo i dati ottenuti lanciando 100 volte un dado.
 Dobbiamo decidere in base ai dati se considerare il dado come "regolare" o "truccato". Come faccio a decidere?
 ("test statistico" o "test d'ipotesi")

Nel corso tratteremo:

- Elementi di statistica descrittiva
- Introduzione all'utilizzo del software R / fogli elettronici
- Elementi di teoria della probabilità
- Elementi di Statistica inferenziale



04/03

Software R → sviluppato a partire dagli anni 90 da Ross Ihaka e Robert Gentleman, basandosi sul precedente software S.

→ è un linguaggio di programmazione e un ambiente di sviluppo specifico per l'analisi statistica dei dati.

- ↳ caratteristiche:
- è gratuito
 - cross-platform
 - open-source
 - molto diffuso per l'analisi statistica dei dati (gran parte delle tecniche statistiche classiche e moderne sono già implementate in R)

> "prompt"

frase tra virgolette ""

3x200 → scivo 3*2

nello script # commento, se voglio che sia stampata ""

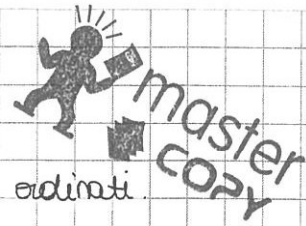
05/03

La media (aritmetica e empirica) dei numeri x_1, \dots, x_n è $\bar{x} = \frac{x_1 + \dots + x_n}{n}$

Il simbolo di sommatoria \sum , (sigma) $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Moda, media e mediana sono indicatori di centralità, ossia riassumono in un solo valore centrale tutti i dati (con criteri diversi).

09/03



3 quartili

Dati x_1, \dots, x_n ordinati ($x_1 \leq \dots \leq x_n$), se non lo sono vanno ordinati.

es: $n=13$ 1, 5, 5, 7, 9, 12, 29, 43, 72, 84, 84, 86, 86

mediana

↳ detta anche secondo quartile

valore minimo \rightarrow quartile zero 1

valore massimo \rightarrow quarto quartile 86

Dividendo ulteriormente: 5, 5, 7, 9, 12, 29, 43, 72, 84, 84, 86, 86
7 \rightarrow primo quartile
84 \rightarrow terzo quartile

3 quartili 1°, 2°, 3° non sono sempre univocamente determinati (quando non lo sono per la sua approssimazione, ad esempio la media)

Significato dei quartili:

- quartile zero \rightarrow non ci sono valori più piccoli
- primo quartile \rightarrow il 25% dei valori sono più piccoli
- secondo quartile \rightarrow il 50% dei valori sono più piccoli
- terzo quartile \rightarrow il 75% dei valori sono più piccoli
- quarto quartile \rightarrow tutti i valori sono più piccoli

Alcuni indicatori di dispersione (o indici di variabilità)

Dati x_1, \dots, x_n

$\delta = \frac{1}{n} (|x_1 - \bar{x}| + \dots + |x_n - \bar{x}|)$ più è grande, più i dati differiscono dalla media. \Rightarrow più c'è variabilità

↳ $\delta = 0 \Leftrightarrow x_1 = \dots = x_n = \bar{x}$

$\delta =$ SCARTO MEDIO ASSOLUTO

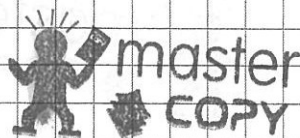
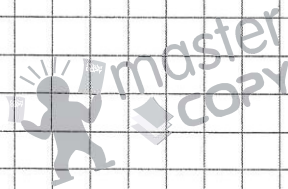
$s = \sqrt{\frac{1}{n} [(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2]}$ $s = \sqrt{\frac{1}{n-1} [(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2]}$

$s =$ DEVIAZIONE STANDARD (campionaria)

↳ è più usata

$s^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2$

$s^2 =$ VARIANZA CAMPIONARIA



11/03

dati numerici

Il coefficiente di correlazione

dati dalle coppie di dati $(x_1, y_1), \dots, (x_n, y_n)$

$$x = (x_1, \dots, x_n)$$

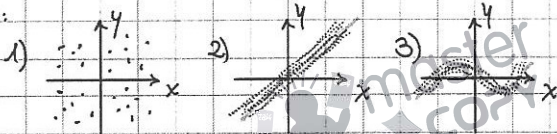
$$y = (y_1, \dots, y_n)$$

es: - Input, Risposte

- sondaggio e chiedere: peso e altezza // età e guadagno // h e €/anno ...

Il coefficiente di correlazione misura il "legame di natura lineare" tra i dati x e i dati y.

es:



in 2) e 3) c'è un legame qual'è la differenza? Tipologia di legame

legame lineare

legame periodico

$\rightarrow \text{cor}(x, y)$ basso.

i dati sono vicini ad una retta

i dati sono vicini ad una curva.

RETTA DI REGRESSIONE $y = mx + c$



In generale il coefficiente di correlazione, $\text{cor}(x, y)$, tra i dati x e i dati y è un numero compreso tra -1 e 1 e si ha:

- $\text{cor}(x, y) = 0 \rightarrow$ nessuna correlazione

- $\text{cor}(x, y) = 1 \rightarrow$ dati su una retta a pendenza positiva ($>x \Rightarrow >y$)

- $\text{cor}(x, y) = -1 \rightarrow$ dati su una retta a pendenza negativa ($>x \Rightarrow <y$)

- $\text{cor}(x, y)$ vicino a 1 \rightarrow i dati si trovano vicino ad una retta a pendenza positiva ⁽²⁾

- $\text{cor}(x, y)$ vicino a -1 \rightarrow i dati si trovano vicino ad una retta a pendenza negativa.

$$\text{cor}(x, y) = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2} \cdot \sqrt{\sum_{k=1}^n (y_k - \bar{y})^2}}$$

Possiamo riscrivere $\text{cor}(x, y)$ in termini della covarianza campionaria $\text{cov}(x, y)$ e di S_x, S_y (deviazioni standard di x e y)

Covarianza campionaria
$$\text{cov}(x, y) = \frac{1}{n-1} [(x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})]$$

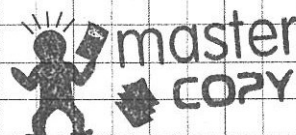
$$= \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

osservazione:

$$S_x^2 = \text{var}(x) = \text{cov}(x, x)$$

$$S_y^2 = \text{var}(y) = \text{cov}(y, y)$$

$$\text{cor}(x, y) = \frac{\text{cov}(x, y)}{S_x \cdot S_y}$$



Per la disuguaglianza di Cauchy-Schwartz: $|\text{cov}(x, y)| \leq S_x S_y$
 $\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}$

dim: 1) consideriamo il caso $\bar{x}, \bar{y} = 0$

Allora $\text{cov}(x, y) = \sum x_k y_k$ e $\text{var}(x) = \sum x_k^2$, $\text{var}(y) = \sum y_k^2$

$$|\sum x_k y_k| \stackrel{?}{\leq} \sqrt{\sum x_k^2} \sqrt{\sum y_k^2} \rightarrow \text{è la disuguaglianza di C.S.!!}$$

(vale sempre)

2) se $\bar{x} \neq 0$ e $\bar{y} \neq 0$



master
COPY

Definisco $u = x - \bar{x}$ e $w = y - \bar{y}$ ($u_1 = x_1 - \bar{x}, \dots; w_1 = y_1 - \bar{y}, \dots$)

allora $\bar{u} = 0$ e $\bar{w} = 0$

quindi mi riconduco a 1) applicato a u e w .

Quindi abbiamo dimostrato che $-1 \leq \text{cor}(x, y) \leq 1$

Determinare m e c della retta di regressione $y = mx + c$
 retta di regressione $m = m_x$, $c = c_x$ con m_x, c_x tali che

$$\min_{m, c} \sum_{k=1}^n (y_k - mx_k - c)^2 = \sum_{k=1}^n (y_k - m_x x_k - c_x)^2$$

Tra tutte le rette, quella di regressione è quella più vicina ai dati

Proposizione:

i) vale $m_x = \frac{\text{cov}(x, y)}{\text{var}(x, y)}$ ($= \text{cor}(x, y) \frac{S_y}{S_x}$) e $c_x = \bar{y} - \frac{\text{cov}(x, y)}{\text{var}(x)} \bar{x}$

ii) vale $\sum_{k=1}^n (y_k - m_x x_k - c_x)^2 = \sum_{k=1}^n (y_k - \bar{y})^2 (1 - [\text{cor}(x, y)]^2)$

la distanza dei punti dalla retta di regressione

Conseguenze

1) da i) segue subito: $\text{cor}(x, y) \leq 0 \Leftrightarrow m_x \leq 0$

2) da ii) segue subito: $|\text{cor}(x, y)| = 1 \Rightarrow \sum |y_k - m_x x_k - c_x|^2 = 0$

ovvero i dati si trovano esattamente sulla retta di regressione

3) se $\text{cor}(x, y) = 0 \Rightarrow m_x = 0$, $c_x = \bar{y}$ e quindi la migliore approssimazione a y_k è $m_x x_k + c_x = \bar{y}$ (indipendente da x_k !)

(sapere x_k non aiuta a indovinare y , il miglior "guess" è sempre \bar{y})

4) da ii) segue che più $|\text{cor}(x, y)|$ è vicino a 1 e più i dati sono vicini alla retta di regressione

dim: (idea: si considera la funzione di due variabili

$$(m, c) \mapsto f(m, c) = \sum (y_k - mx_k - c)^2 \text{ e si trova il minimo assoluto che è proprio } (m_x, c_x))$$

