



MODERN STATISTICS

INTUITION, MATH, PYTHON, R

MIKE X COHEN

 SINCXPRESS

7- Analyze and Improve Data Quality

7.1 - Data quality matters

mancaenti

ben organizzati // poco rumore // pochi artefatti // pochi dati

- numero tecniche di filtraggio, trasformazioni
- pulire i dati richiede meticolosità, critico valore, basso spreco e imperfezioni, conoscenze statistiche, esperienza pratica, solida comprensione dell'argomento, familiarità con le caratteristiche uniche dei propri dati
- La qualità è fondamentale, su queste nostre analisi sono prese decisioni operative

Garbage in, garbage out (GIGO)

Dati scadenti → risultati scadenti
 Dati di qualità scadenti → risultati sorprendenti non previsti
 ↳ necessari, ma non sufficienti

- i dati sono sempre imperfetti - la domanda è: ho bisogno di ulteriori dati?

7.2 - Data cleaning phases

4 finestre temporali x stato qualità

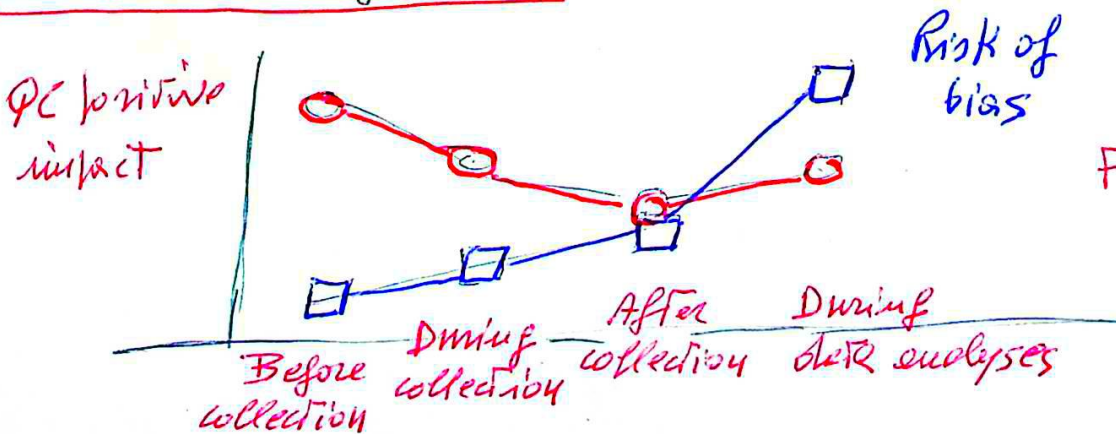


Fig. 7.1

Before getting data

- familiarizzare con pubblicazioni scientifiche, relazioni tecniche, blog, youtube
- condurre uno studio pilota, x intercettare problemi imprevisti

During data collection

- ispezionare i dati in tempo reale x verificare presenza di artefatti, errori di calibrazione dei sensori -
- trovare & risolvere un problema

After data collection

- pulizia dei dati, 2 strategie
- trasformazione dei dati
- rimozione dei dati cattivi
- questo è il momento peggiore x aumentare la qualità

During data analysis

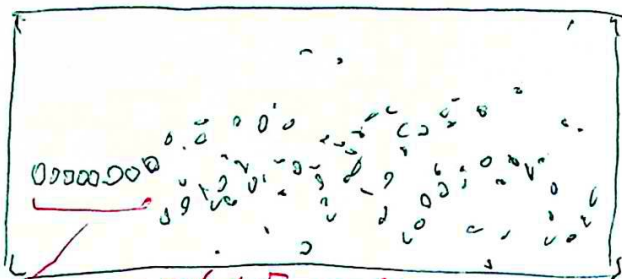
- possono essere introdotti bias, che cambiano il risultato -
- meglio evitare la trasformazione dei dati - se possibile

7.3 Addressing data quality

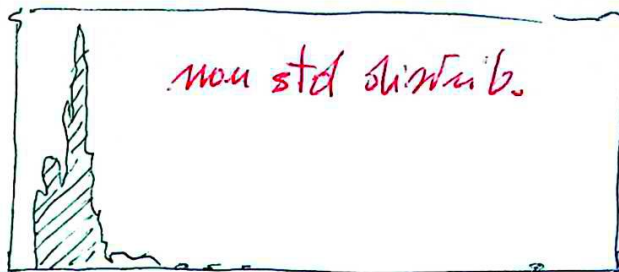
- ci sono 2 famiglie di approcci x fare questo (NON SONO mutuamente esclusive)

Qualitative quality assessment

- visualizzare i dati in modalità differenti, come cercare in queste visualizzazioni?
 - l'intervallo dei dati, tenere d'occhio i valori non quelli attesi
- la forma della distribuzione
 - facciamo qualche esempio con l'età di profici



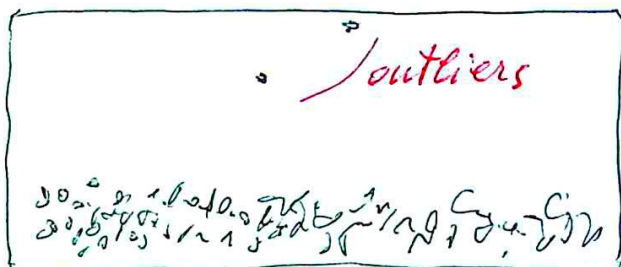
un expected data range



non std distrib.



mixed dataset



outliers

Quantitative quality measurement

- le metriche + utili tendono a essere specifiche dell'applicazione e talora specifiche del set di dati - Vediamone un elenco limitato:
 - sample size - molti campioni → risultati statistici + affidabili
molti campioni → diminuire rumore, parametri della popolazione + accurati. Il sample è migliore qualità non è comunque automatico. Quando un campione è "grande" → Cap. 17
 - Error rates - qui ci stiamo riferendo ai valori piccoli
picoli → misurare la quantità dei rumori → sample size
Troppi errori → dubbio sulla qualità dell'esperimento
 - Data range - ci riferiamo agli estremi. Nella pratica difficilmente i dati di un set di dati fuor'area sono.
 - Variance - se troppo piccolo o troppo grande → problema
troppa variance → potrebbe oscurare i dati reali
curare la normalizzazione x permettere confronti

7.4 - Improving data quality through transformations

- distribuzioni fortemente sfegliate da un lato \rightarrow diventare simili a $Geom$ con una trasformazione log
- la conoscenza del dominio aiuta a scegliere la trasformazione
- Testare prima e dopo a vedere se ci sono miglioramenti

7.5 - What are outliers? = valore relativamente insolito

- decidere se "outlier" ha una componente di soggettività.

Il valore anomalo è un H_0 non valido

- ci può essere rumore, errori, malfunzionamento dei sensori
(es. 0,0254 \rightarrow può diventare 254 se la virgola non è festiva bene)
- una data sbagliata (come gli bettini)

Il valore anomalo è valido es. dati sul patrimonio netto di uomini maschi negli USA, con nome di bettino Jeff - Qui incontriamo

Jeff Bezos con un patrimonio $> 100e9$

es. assicurazioni contro le inondazioni - in Texas \rightarrow Utah

es. carte di credito - banca tutto l'anno, ma un milione a Natale

es. transazioni con carte di credito relativamente piccole in una data area \rightarrow un addebito particolarmente elevato potrebbe essere fraudolento

- non-rappresentativo, anomalo, estremo, deviante

- sono difficili da gestire

Outliers richiedono decisioni

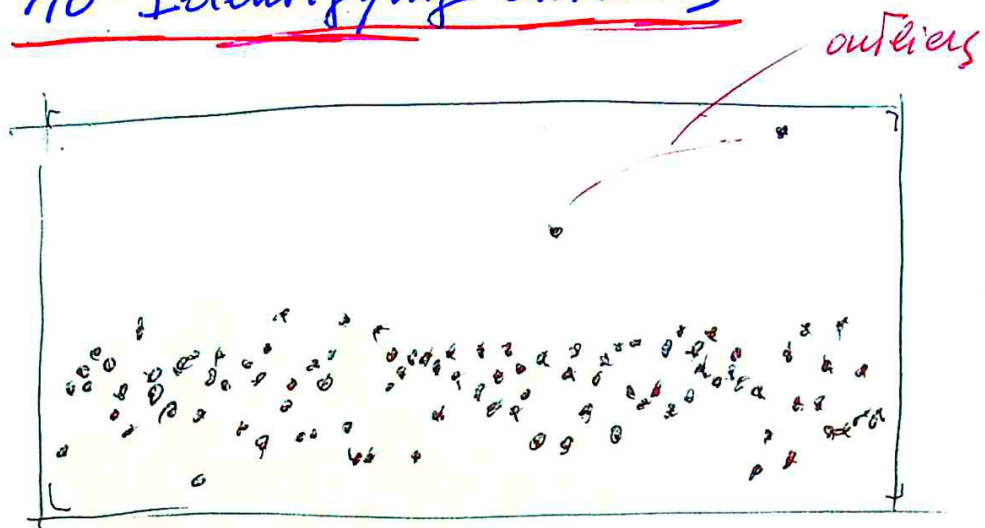
- Rimozione outliers - se non riflettono il sistema in oggetto

- Tenere outliers

- 1) applicare trasformazioni a ridurre loro impatto (es. log)
- 2) usare analisi robuste rispetto a essi (statistiche ^{non} parametriche, regressioni robuste (ponderate))

3) analizzare i sottogruppi → separazione dei dati in gruppi con caratteristiche simili, poi analisi statistica sui singoli gruppi

7.6 - Identifying outliers



La sola ispezione visiva non è un metodo ottimale.

- 1) non è scalabile, funziona bene su dati piccoli o unidimensionali;
- 2) rischio di soggettività bias soggettivi; tendenza escludere
- 3) se le prove diverse → probab, scelte diverse → risultati diversi

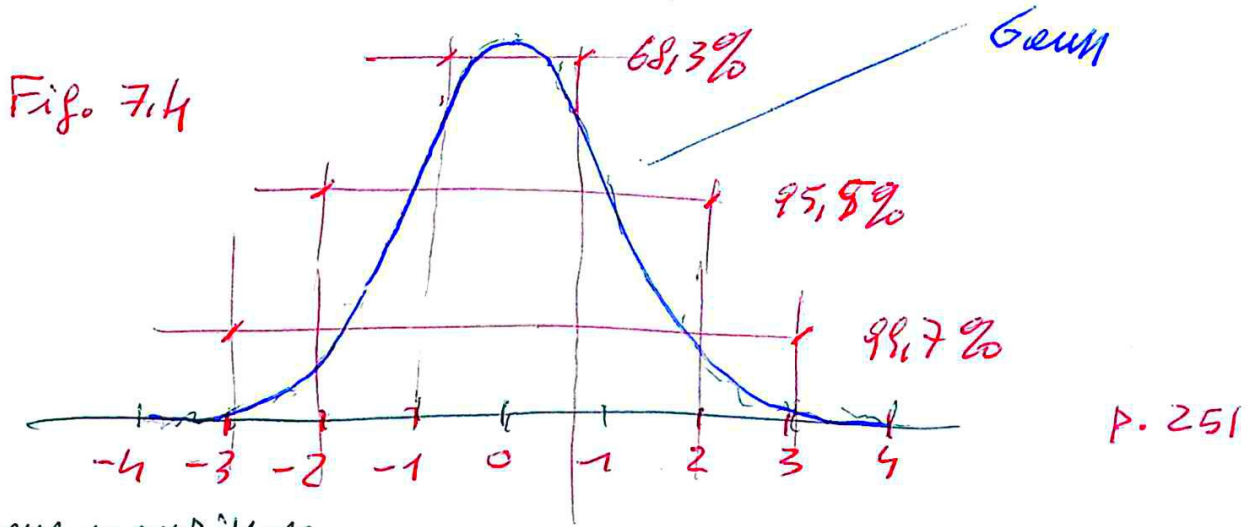
Absolute Threshold Detection - ciò che è fuori da un certo intervallo è escluso. Qui si vuole molta conoscenza del fenomeno. È una scelta personalizzata x quel set di dati.

The z-score method

Per rimediare al con precedente → relatività, ecco x che' z-score

z-scores and data proportions

• poniamo valutare la quantità 'di dati' entro diversi limiti di z-score
vediamo meglio con una figura:



• σ bene memorizzare

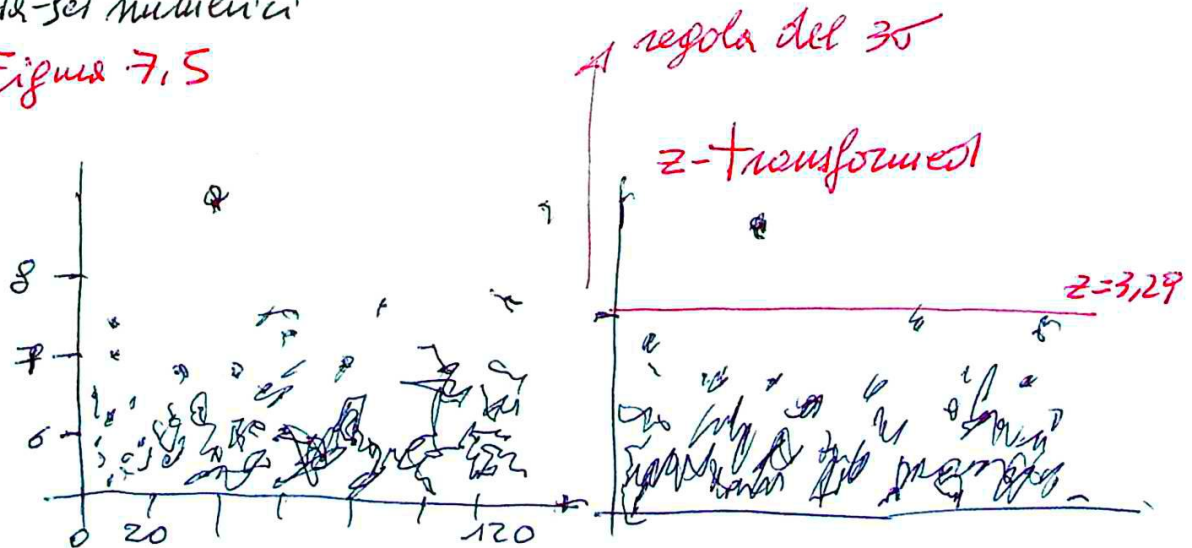
σ	\rightarrow	68,3%	
2σ		95,5%	X valore di dettaglio \rightarrow Cap. 10
3σ		99,7%	

- possiamo usare z-scores x identificare dati relativamente grandi come outliers

The z-score method

- E' il metodo + comune x identificare outliers. Si convertono i dati con z-score, si considerano quelle che superano dette soglie.
- Il punto di forza e' che il metodo funziona x un numero di data-sets - indipendentemente dalle unita' di misura
- Vale x i data-set numerici

Esempio - Figura 7.5



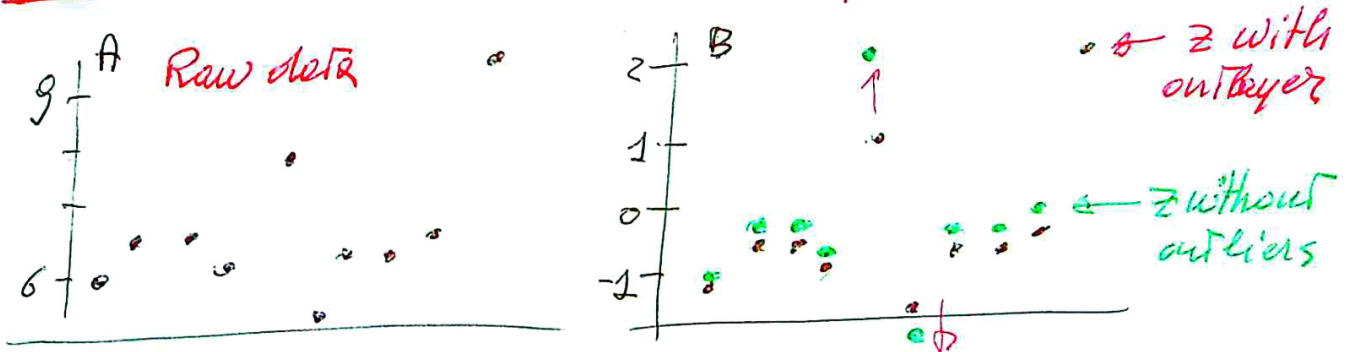
in generale: $|z| > 3,29$

- z-score è riproducibile
- il metodo è "relativo"
- ha un parametro **la soglia** (una certa soggettività)
 - ↳ Valori + comuni % 2.3, 3, 3.29
- Se il fenomeno in oggetto è disperso molto disperso →
 - ↳ anche 5, 8
- i dati "non-normali" potrebbero richiedere soglie molto alte →
 - ↳ z-modificato

Modified-z

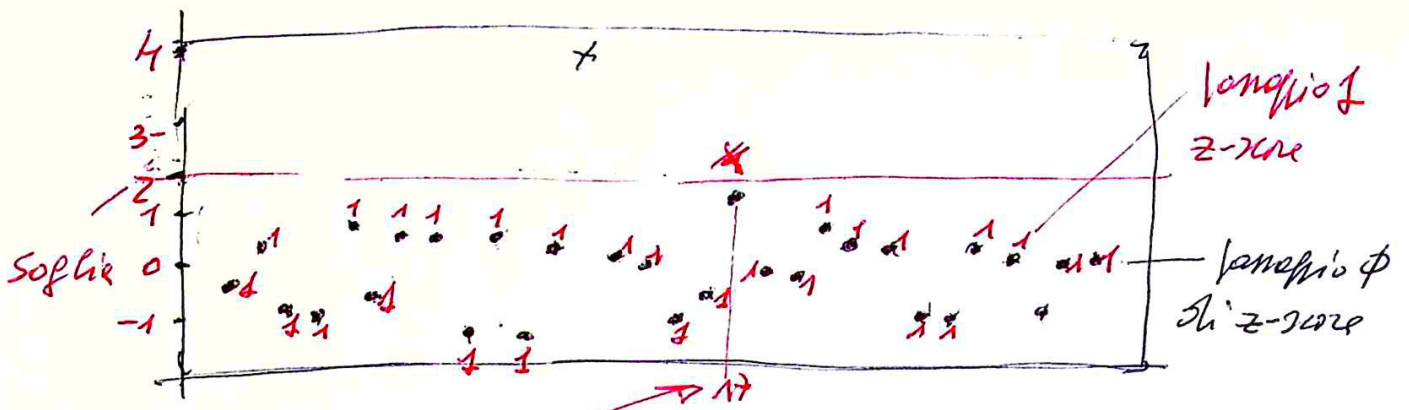
- ma la mediana e la MAD, sono le + appropriate x fenomeni centrale e dispersione - Poi approccio numerico (soglia)

Iterative z-score method - vediamo problema con figura 7.6



z-transformation modifica i dati: senza outlier li espande rispetto allo zero (e' cambiata la media)

- 1) Implementiamo z-score come prima descritto
 - diamo una etichetta ai dati + estremi che superano la soglia classificandoli come outliers
- 2) Ricalcoliamo z-scores senza gli outliers - Seguiamo con una label i nuovi dati che superano la soglia
- 3) Ripetere fino a che non ci sono + outliers Vedi fig. 7.7



Nota - fatto uno z-score, eliminiamo outliers, \rightarrow risultato

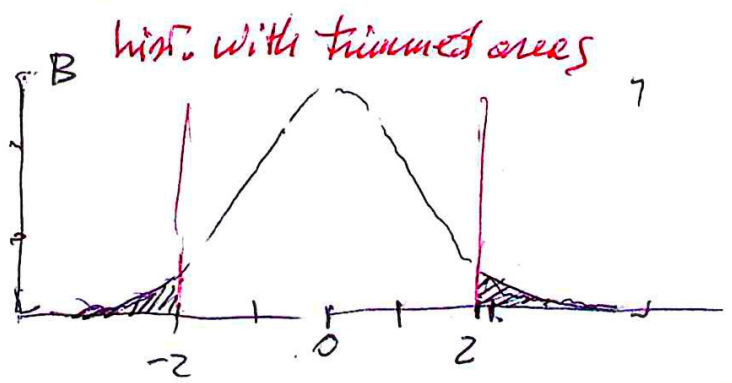
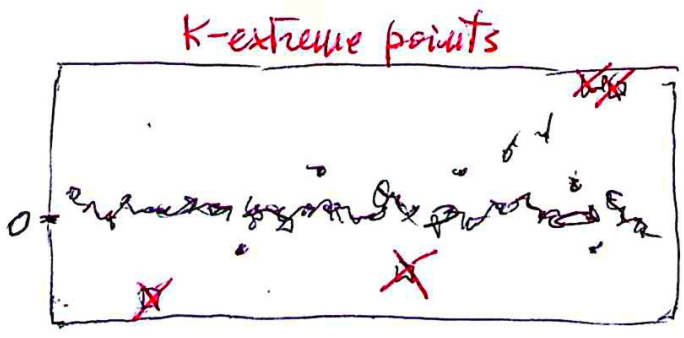
- che dice del ϕ 17? non sembra insolitamente grande (soglia $z=2$) \Rightarrow
- metodi iterativi sono ^{no} rimuovere troppo
- fidarsi ciecamente e' non consigliabile - infezionare i dati
- ma sono outliers grandi e piccoli \rightarrow quelli grandi possono oscurare quelli piccoli
- non e' di utilizzo comune questo metodo

È utile sapere che \exists queste possibilità

Removing data by trimming

Simile al metodo della "soglia", qui però si rimuovono i punti estremi, oppure si rimuovono le "code" della distribuzione.

- Es. possiamo impostare $z=2$ x ogni coda Fig. 7.18
- L'autore non e' un fan di questo metodo
- i vantaggi sono che il metodo e' semplice, ripetitivo



Manual, automatic, and semi-automatic cleaning

- se automatico → implica fiducia nel metodo
- se manuale → ispezione visiva. Esperienza del dominio
- semi-automatico → mix, magari si merita in automatico i pesanti outlier, si procede a mano nella eliminazione
 - alla fine del processo, si può dare (bisogna essere esperti) ok all'uso in automatico in quelle circostanze
 - migliora il problema della scalabilità - rende efficiente le evoluzioni in termini di tempo - si può concentrare solo sui dati segnalati
 - molto utilizzato nella elaborazione di dati finanziari e finanziari (es. può segnalare un grosso versamento su un deposito, lo specialista viene informato x eventuali azioni)

What happens to rejected outliers?

- potrebbe essere fattibile rimuoverli (come se \bar{X})
- si possono presentare in una "feature" e non in altre, potrebbe essere demerito cancellare una intera riga, solo perché lo contiene (es. outlier in una sola colonna)
- possono essere "accoppiati", allora devono essere abbinate alle corrispondenti righe in una o più colonne diverse, o in data sets separati
- possono essere rimpiazzati con "NaN" = not a number - Nella computazione sono ignorati.

Hybrid methods

- si possono mescolare e abbinare diverse strategie di rilevamento.
 - es. soglia evoluta (si eliminano quelli eclatanti) /
/ metodo iterativo (x identificare outlier non-erozi) //
 - z-modificato (x i dati non sono Gauss)
- molti metodi ottimizzati sono specifici

7.7 - Analysis-based solutions to outliers

- talora gli outliers sono significativi, li vuoi mantenere, ma non vuoi che alterino i risultati.
- ci sono metodi di analisi progettati x minimizzare il loro impatto.

Nonparametric analyses

- molti metodi es. t-test, correlazione, ANOVA hanno alternative non-parametriche che implicano la trasformazione in base al rango o al test delle mediane. Sono robuste rispetto a outliers.
- tuttavia queste analisi non-parametriche possono avere sensibilità ridotta nel rilevare effetti sottili.

Permutation-based tests

- creano distribuzioni statistiche "null hypothesis"
 - ↳ si ipotizza che Z relazione tra due serie di dati o variabili qualificate.
 - ↳ se vero \rightarrow l'effetto osservato è dovuto al caso

- questa "null hypothesis" distribution si ottiene eseguendo 369
do numeri casuali - I valori euomeli in una condizione
vengono eseguiti casualmente ed altre condizioni -
- cio' significa che i valori euomeli vengono incorporati nel
processo di inferenza statistica. → Cap. 16

ESEMPIO - H_0 = ipotesi nulla

- se $(x_1 - x_n)$ e' di un certo tipo → rifiuto H_0
- l'ipotesi alternativa e' H_1
- H_0 la si ritiene vera fino a prova contraria
- il test di ipotesi ci dara' una regola di decisione (o rifiuto)

• Se $H_0 = \text{True}$ → molto improbabile che il campione sia fuori
al rifiuto

REGIONE di RIFIUTO

$R = x_1 - x_n$ porta a rifiutare se

statistica
osservata

NON
COMPLETO

$$R = \{ (x_1 - x_n) : t(x_1 - x_n) \in I \}$$

intervallo
regione di rifiuto

• dare un test = dare sua regione di rifiuto

↳ ogni test e' enunciato ellen sua R

$I \subset \mathbb{R}^m$

meglio esplicitare
sua definizione Cap. 16

Weighted analyses

- se il problema è che gli enomali hanno un impatto sproporzionato sui risultati → perché non de-ponderarli?
- de-ponderare = valore $\times r_i$ $0 < r_i < 1$ in base al loro impatto
- Es. $[1, 2, 3, 40]$ → vettore di ponderazione $[1, 1, 1, \dots]$ →
→ $[1, 2, 3, 4]$
- è un metodo complicato → Cap. 15

Subgroups analysis

- se ci sono molti enomali → potrebbe essere che rappresentino un gruppo qualitativamente distinto

Es. consumo di carburante di 100 auto

$$80 = 18 \div 25 \text{ miles/gallon}$$

$$20 = 40 \div 60 \text{ —————}$$

(questi escono come outliers, in realtà sono auto ibide electric-fos)

↳ l'analisi statistica si farà x gruppi separati

7.8 - Missing data - Fig. 7.9

- i dati mancanti possono avere un grande impatto. Possono mancare x vari motivi:
 - Drop out - es. se traccia gli stessi individui nel tempo (studio longitudinale) - Alcuni partecipanti possono interrompere il test → dati perpici
- i dati perpici sono utili in alcune analisi, ma non in altre

- Malfunctionamento attraverso

- perdita irrimediabile, corruzione, inutilizzabili

- Errori umani

Che fare con i dati mancanti?

Vediamo varie modalità:

- Rimozione x riga (appropriata x dati accoppiati)

↳ es. misure prestazioni sportive prima e dopo una settimana di allenamento

- se ci sono dati mancanti → nessun dato di quell'individuo è utilizzabile

- Analysis-specific row removal

- però i dati di quella riga potrebbero essere validi, utilizzabili

↳ mantenere la riga nel data set, escluderla dalle analisi che necessitano proprio di quei punti

- utile nel caso di dati limitati

- si può sostituire il dato con NaN

- Imputation (replacement)

- sostituire i dati con valori stimati, basati sui dati disponibili

↳ spesso si usa la media

Time 1	Time 2
5	7
6	6
4	?
8	7

Fig. 7.18

6,67 = media di Time 2

Predictive modeling

- concettualmente simile a sostituzioni
- "fitting" the model parameters richiede molti dati

7.9 - Exercises

1) implementare "iterative z-score" outlier removal →

→ ottenere Fig. 7.7

2) data trimming - creare un dataset casuale $y = \exp^{\sin(x)}$
dove $x \sim \mathcal{N}(\mu, \sigma)$ con $N = 10000$ - Fare un copia, per
essere originale con trimmed.

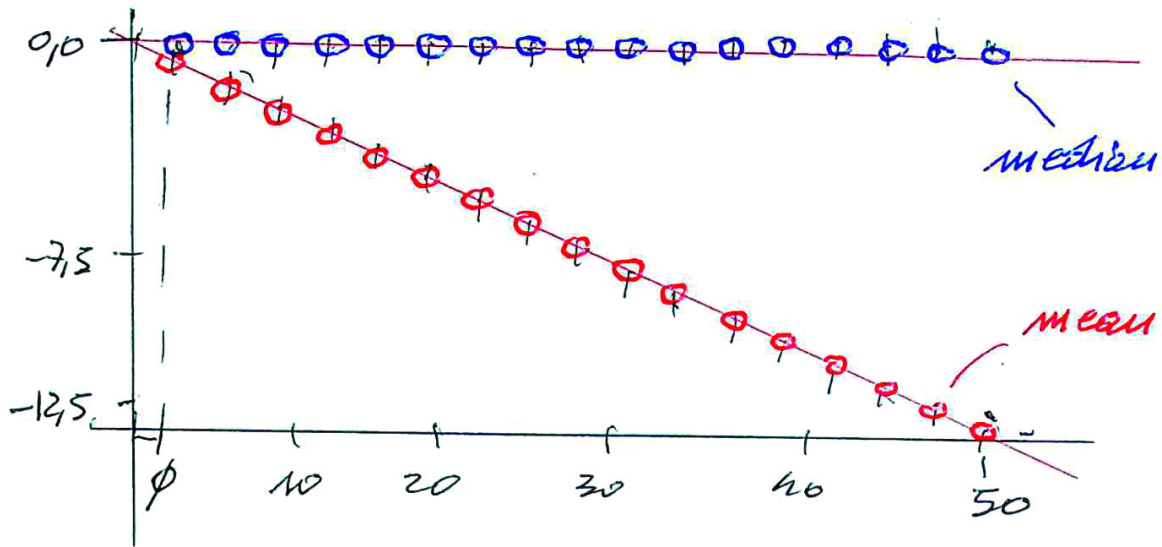
$\mu_{trim} = 4\%$ dei dati (2% di ogni coda)

- sostituire i dati rifiutati con N&N → 9600 dati validi
- calcolare media, mediana, prima e dopo troncamento →

$\bar{x} = 1,222$		$\tilde{x} = 0,989$		<u>mp.mqm</u>
$\bar{x}_{trimmed} = 1,054$		$\tilde{x}_{trim} = 0,989$		<u>mp.arg50%</u>

3) mettersi codice 2) in un loop con $K = 1\% \rightarrow 50\%$

- calcolare media, mediana, e ogni cosa
- moltiplicare lo stesso set di dati casuali (→ descrittivi confrontabili)
- testare K es. se $K = 12 \rightarrow 8800$ punti validi
- visualizzare la caratteristica descrittiva con un grafico Fig. 7.10 dove si mostra la variazione % di media, mediana vs. $K\%$ to trim -



Commento

- la mediana è inalterata (n° di rifiutati è uguale nella coda di destra e in quella di sinistra)
- la mediana cala
- x questo se il trimming a due code non è appropriato
- provare con trim in una sola coda - In questo caso come influisce su st01?
- i risultati cambiano qualitativamente con una diversa distribuzione di dati?

p. 263

4) Vediamo l'impatto della rimozione degli outliers sulla forma della distribuzione - valutare se tale procedura si fonda sull'algoritmo di calcolo del n° dei bai -

$N = 1000$ casuali, estratti da una distribuzione F , nell'intervallo 5 e 100 quasi di libertà \rightarrow Cap. 111

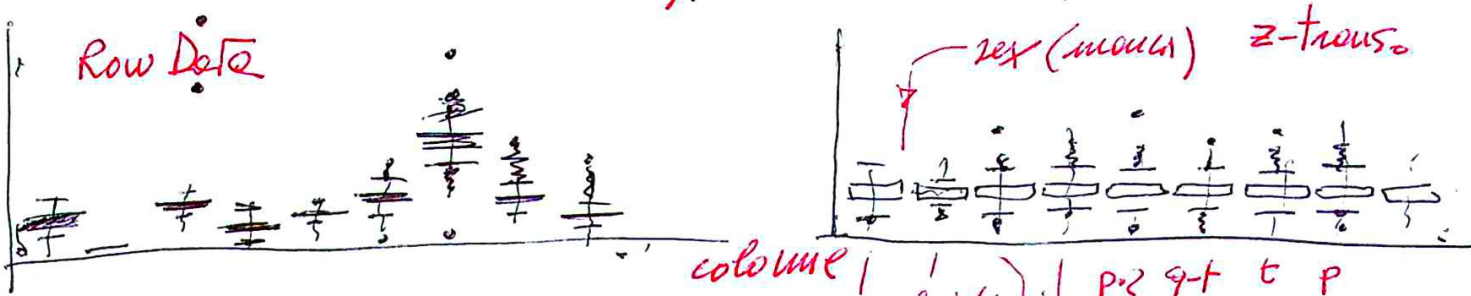
mp. random. $f(5, 100, n_{i,j} = N)$

- mi è venuto z -score e poi mi è venuta la soglia $z = 3$

L calcolare % di eliminati $N = 1000$
 $Eliminati = 983$
 $\% rimasti = 1,70\%$

%

- facciamo una copia del dataset, applichiamo z-score a tutte le colonne (eccetto "row"), otteniamo un "boxplot"



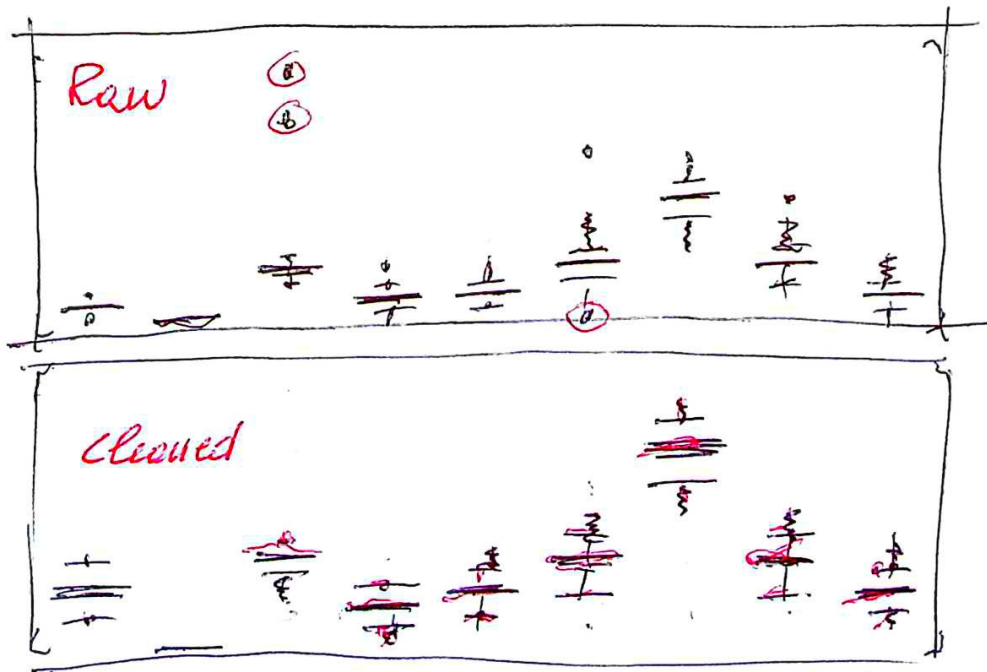
vediamo che ci sono molti outliers
 L x ogni variabile criteri diversi
 L bisogna essere cauti x poter scegliere

NOTA

qui e' sufficiente un solo criterio applicabile a tutte le variabili.

6) Riconoscimento outliers e loro rimozione

- soglia $z = 3,29$
- nei row data, gli outliers sono diventati NaN



age
 - outliers positivi e negativi sono stati eliminati.

• stampiamo le medie di raw e cleaned

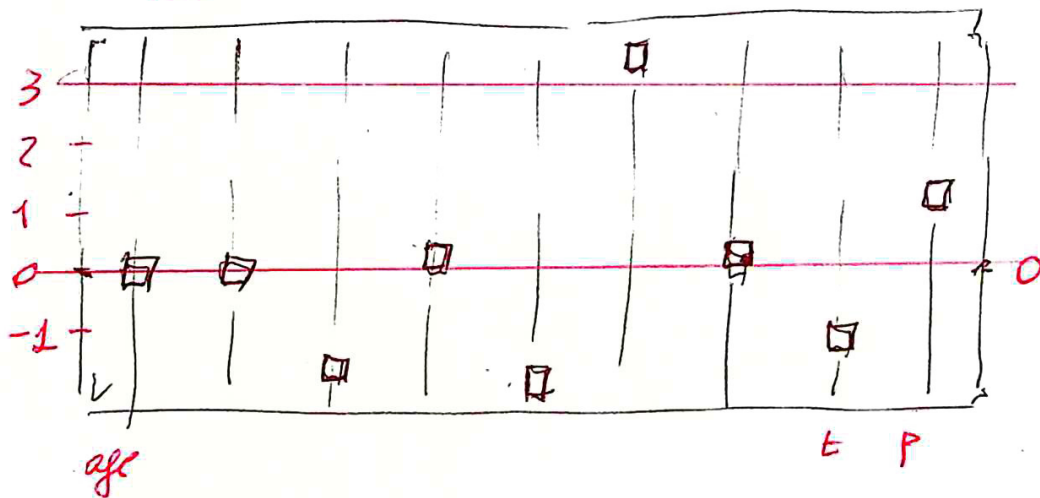
age	46,47	46,47
sex	0,55	0,55
height	166,19	163,84
weight	68,17	68,33
qvs	88,92	87,48
p-2	155,15	160,38
q-t	367,21	368,05
t	169,95	168,28
p	90,00	91,14

prende un po' il femminile

ora ne facciamo un grafico in %

$$\frac{\text{cleaned} - \text{raw}}{\text{raw}} \cdot 100$$

raw cleaned



• qui abbiamo visto che i dati anomali sono presenti nei dati reali -

p. 268

— FINE CAPITOLO 7 —

8 - Probability Theory

8.1 - From descriptive to inferential statistics

- Questo è un filo di molto
 - le statistiche inferenziali usano i dati x fare previsioni "fermate" (educated guesses) sul mondo al di là del mio dataset -
- es. Occhiali inferiore
 - descrittive: dimensione, forma, colore, caratteristiche di)
 - inferenziali: con questi occhiali guardare il mondo
- L'acqua meglio partiamo dalla probabilità
- L'oro un aiuto a prendere decisioni in un mondo incerto
- L'la probabilità quantifica l'incertezza

"Pure" vs. "computer" probabilities

- La teoria della probabilità è molto orientata matematicamente
- E però la statistica applicata → quel prerequisito non è più misurabile - Questo sarà il taglio del filo -
- loro riferimento ad approcci tipo "computer", "empirici" -
- La statistica computazionale comporta implicazioni di calcolo, errori di arrotondamento, algoritmi astratti - Tutte queste problematiche possono essere affrontate con Python e R -

8.2 - What is probability?

Nella realtà concreta nulla è certo. La probabilità viene
tro $\phi \frac{1}{2}$, può essere espressa come percentuale

8.2.1 Il problema con la probabilità

• Allora la probabilità è anche contraddittoria - Vediamo
di riconoscere fatti corrette e non corrette:

- 1) piovera' x il 40% della giornata
- 2) — nel 40% della città, il minuenente sarà asciutto
- 3) c'è una probabilità su 5 che finirà oggi a un certo punto, — **OK**
da qualche parte a Osaka
- 4) i meteorologi sono sicuri al 40% che finirà oggi

L'no, la confidenza è differente dalla probabilità ^{40%}

L'non essere confidente al 99% di una probabilità di fatto del

Esempi di probabilità

- lancio della moneta - ogni faccia ha probabilità = 0,5
- lancio dei dadi - ha 6 facce \rightarrow probabilità = $1/6 \approx 16,66\%$
- gioco delle carte - 52 carte \rightarrow — = $1/52 = 1,92\%$
- prob. di una carta rossa = 0,5
- — — — refina (di ogni seme) = $1/52 = 1,92\%$

- la probabilità di piova nel deserto non è 50, anche
se 2 sono le probabilità

- la probabilità somma = 1 - potremmo visualizzarlo come "torta"
L sul PC potrebbe non tornare esattamente
L vedremo nel seguito esempi di ciò

When do we need probabilities?

41

- Alcuni eventi sono con' probabili che non serve il concetto di probabilità -
- Alcune situazioni (es. un referto medico) rich'edono invece la probabilità -

8.3 - Probability vs. proportion

- Probabilità e proporzione sono concetti correlati; ma distinti -

Es. Lanciamo 10 monete

L'esperimento di 6 monete con teste \rightarrow

$$\text{probabilità} = 0,5 = 50\%$$

$$\text{proporzione} = 0,6 = 60\% \rightarrow$$

PROBABILITÀ = la possibilità che un evento accada

PROPORZIONE = una frazione del tutto

Es.

- impiego 5,1 min/day a lavarsi i denti su 17 ore di veglia -

- proporzione della parte di veglia usata a lavare i denti:

$$5,1/1020 = 0,5\%$$

- probabilità che "un minuto scelto a caso" implichi il lavaggio dei denti = 0,5%

In questo caso numericamente uguali

(1020 min)
"

- accade che il valore della probabilità "dipenda da come è posta la domanda -
- la probabilità riguarda un evento futuro
- la proporzione è una semplice conta -

8.1 - Computing probabilities

Dove nasce la probabilità?

Intuizione dalla conoscenza semantica

Conoscenza semantica = l'informazione sul mondo che conosciamo.

- La maggior parte delle probabilità che conosciamo sono inutili nell'analisi statistica formale. Potremmo inoltre avere divisioni che in realtà sono fallaci.

Una formula matematica

- Impareremo la meccanica della determinazione della significatività statistica =

- prob. che una statistica descrittiva di un campione rifletta le caratteristiche di una popolazione con parametri noti.

L deriva da una formula

Misure empiriche

p. 277

La probabilità si misura da misure empiriche.

Nei casi + semplici, matematicamente:

$$P(x_i) = \frac{N(X=x_i)}{N(X)} \quad \text{(8.1)}$$

↳ set degli eventi possibili

Es. Barattolo con 90 biglie / 40 blu
30 gialle
20 arancione

quale la probabilità di finire su un colore?

42

- questa risposta può essere soddisfatta esclusivamente dai dati empirici

$$blu = \frac{40}{90} = 44,4\%$$

$$giallo = \frac{30}{90} = 33,3\%$$

$$\rightarrow 44,4 + 33,3 + 22,2 = 99,9\%$$

$$\underline{\hspace{10em} 22,2\%}$$

More complicated analytical probabilities

- Se prendo un generico set di dati empirici \rightarrow ho bisogno di calcolo numerico

Calcolare probabilità empiriche

Vale sempre la 8.1 - però conosciamo meno, es. sappiamo che nel barattolo ci sono 40 biglietti

Requisiti sui dati di cui calcolare le probabilità

- due sono i requisiti

- i dati possono essere numerici - discreti; ordinati; categoriali **Es.**

- supponiamo di voler conoscere la probabilità empirica che un pinguino abbia un certo peso

$N=100$ deve accettare un intervallo di peso

Es. tra 3 e 3,1 Kg

- i dati di intervallo devono essere convertiti in intervalli discreti

- bisogna dunque contare gli eventi favorevoli (il numeratore della 8.1)

%

— i dati devono avere etichette o contenitori mutuamente esclusivi — Es.

se un pingüino pesa tra 3 e 3,1, non può pesare tra 2,8 e 2,9

Contro-esempio

— dove le persone prendono le notizie? — cap. 3

— possiamo calcolare la probabilità di quella tabella? NO

↳ x che $\sum \text{prob}_i > 100\%$

Interpretare le probabilità empiriche

Si basano sui campioni / campioni diversi possono avere probab. empiriche diverse —

Per superare queste difficoltà — ^{misurezione di campioni + grandi} calcolo di intervalli di confidenza

— le probabilità sono solo uguali alla ^{proporzione} $\frac{\text{numero}}{\text{totale}}$, che è un conteggio di cose già accadute

8.5 - Probability functions, mass, and density

Probabilità = la possibilità che un evento accada

Funzione di probabilità = una espressione matematica che collega ogni elemento di un set a un valore di probabilità numerico

Mass di probabilità = una funzione che descrive le probabilità

x un set di eventi discreti esclusivi

Densità di probabilità = una funzione che descrive le probabilità per eventi continui esclusivi

in formule:

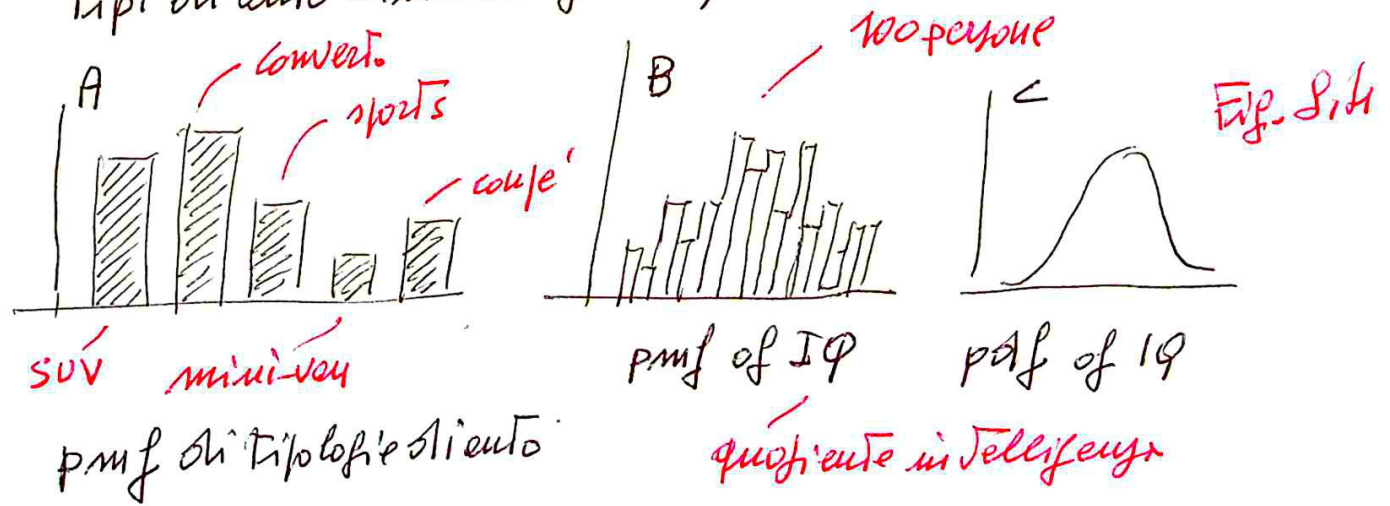
$f(x) = P(X=x)$ (8,5)

set di valori / un particolare valore
permette di eseguire a quell'evento una probabilità

$f(x_i) = P(X=x_i)$ 8,6/8,7/8,8/8,9
 $P(X=x_i) \geq \phi$ / $P(X \neq x_i) = \phi$ / $\sum_{i=1}^m P(X=x_i) = 1$

- Massa di probabilità x eventi discreti (invisibili con barre)
- Densità di probabilità x continui (con linee)

Es. (non reale) barre in Finjeco
Tipi di auto - sono categoriali; non hanno ordinamento



SUV mini-van pmf di tipologie di auto

pmf of IQ pdf of IQ
quoziente in Jellifera

- Consideriamo B
- se proporzione -> stai considerando quel particolare campione di dati
↳ non puoi fare inferenza sul resto della popolazione
- se massa di probabilità -> si possono fare inferenze su persone ed
di fuori del tuo campione

deve essere casuale, e rappresentativo p. 287



col computer le grandezze sono solitamente discretizzate (binning)
 → funzioni di densità empiriche → essendo discretizzate, sono tecnicamente rappresentate come mappe (non come densità) -

• è utile pensare alle mappe di probabilità come a stime della densità di probabilità - **Es.**

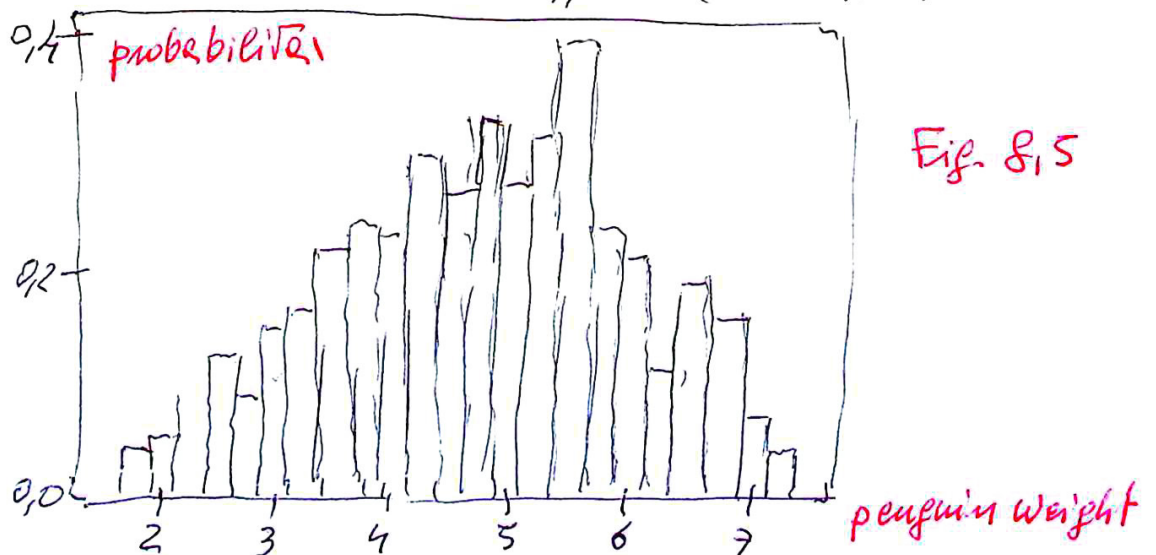
- ho viaggiato in Antartide x pesare pinguini
- ero dotato di una bilancia ad altissima risoluzione (misura da 10 a 16 Kg)
- voglio calcolare la funzione di probabilità empirica che descrive i pesi dei pinguini
- discretizzo in contenitori di larghezza 0,25 Kg ⇒

⇒ nel mondo reale sono presenti 2 valori numerici discreti:

• l'analogo matematico è che non ha senso calcolare l'integrale in un po'. Necessità di un intervallo - **combinate or bias dell'integrazione** - ⇒

⇒ le funzioni di probabilità sono continue

pinguin = mp.erctanh (mp.random.uniform(size=473) * 1,8 - 0,9) * 2 + 4,5
 bin_edge = mp.linspace (mp.min (pinguin), mp.max (pinguin), step = 0,25)



8.6 - Cumulative distribution function (cdf) 114

Sono la somma cumulativa delle masse di probabilità o della funzione di densità. (qui invece di $\Sigma \rightarrow \int$)

- Ogni ptb della cdf dice la probabilità di avere un valore minore o uguale al suo valore. *In formule*

$$F(x) = P(X \leq x) \quad \leftarrow \text{prob. di essere } \leq x \quad (8.10)$$

$$f(x) = P(X = x) \quad \leftarrow \text{probab. di un valore specifico} \quad (8.5)$$

• Nota la sottile distinzione tra le due

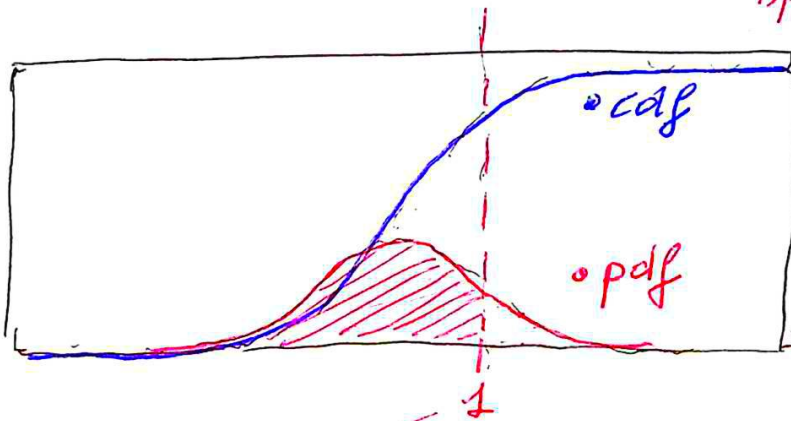


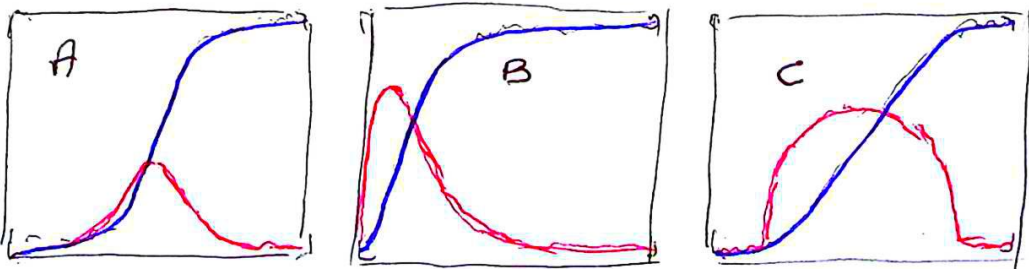
Fig. 8.6

il valore della cdf $\times x = 1 =$ area e minima della pdf

Nella Fig. 8.7 alcuni esempi di pdf e corrispondenti cdf

L cdf = somma cumulativa \rightarrow sempre crescente

Fig. 8.7



• cdf è utile a calcolare un valore p

L interpretabile solo a dati che possono essere ordinati numericamente mette
in modo significativo

8.7 - Expected Value

- Valore atteso, come si ottiene, come lo si interpreta, come lo si relaziona con la media e la varianza?

$$\bar{X} = \sum_{i=1}^n \frac{x_i}{n} \quad (8.11)$$

$$E[X] = \sum_{i=1}^n x_i p_i \quad (8.12)$$

Valore atteso

↳ è una media ponderata

Se $p = \frac{1}{n}$ → le due equazioni sono identiche

- il valore atteso è analitico - calcolo teorico da una distribuzione
- la media è empirica - statistica descrittiva empirica di un campione di dati

Es. dado "perato"

faccia	prob
1	1/4
2	1/4
3	1/8
4	1/8
5	1/8
6	1/8

$$E(x) = \frac{1}{4} + \frac{2}{4} + \frac{3}{8} + \frac{4}{8} + \frac{5}{8} + \frac{6}{8} = 3$$

x un dato non perato

$$E(x) = 3,5$$

Fig. 8.18

- calcoliamo il valore medio facendo lanci, supponendo di aver ottenuto:

$$x = [1, 3, 4, 4, 4, 3, 2, 5] \rightarrow \bar{x} = 3,25$$

notiamo che non coincide con $E(x)$ — questo è da mettere in conto

Computing expected value — nel caso di distrib. uniforme

$$E(x) = \int_a^b \frac{x}{b-a} dx \rightarrow \frac{(b+a)}{2} \quad (8.18)$$

8.13

Stima dei valori attesi

45

Molte distribuzioni di dati e pdf sono sconosciute \rightarrow

• per avere $E[\cdot]$ possiamo esaminare l'istogramma, fare una ipotesi con la somiglianza a una distribuzione nota \rightarrow calcolare il valore atteso

• es. molte distribuzioni empiriche sono \approx Gauss

• per avere stime empiriche ricche, e varie, $N \rightarrow \infty$, in pratica non possibile

\rightarrow campione con N finito \rightarrow una semplice stima

Valori attesi e momenti statistici

$$E[X] = 1^{\text{st}} \text{ moment}$$

$$E[X^2] = 2^{\text{nd}} \text{ "}$$

$$E[X^k] = k^{\text{mo}} \text{ momento}$$

potremmo calcolare la varianza attesa, se il tuo 2nd momento è centrato sulla media \rightarrow

\rightarrow varianza = 2nd mean-centered moment:

$$\sigma_X^2 = E[(X - E[X])^2] = \text{---} \quad (8,19)$$

$$= E[X^2] - (E[X])^2 \quad \text{---} \quad (8,20)$$

• Esaminiamo il caso della distrib. uniforme

con la 8,18 $E[X] = \int_a^b \frac{x}{b-a} dx \quad \Rightarrow$

per calcolare x^2

$$E[X^2] = \int_a^b \frac{x^2}{b-a} dx = \frac{1}{b-a} \frac{b^3 - a^3}{3}$$

- combiniamo i due momenti e otteniamo la varianza (8,18)

$$\sigma_X^2 = \frac{b^3 - a^3}{3(b-a)} - \left(\frac{b+a}{2}\right)^2 \rightarrow \text{calcolando} \rightarrow \frac{b^2 - a^2}{12} \quad (8,28)$$

ripetiamo $\sigma_x = \frac{b^2 - a^2}{12}$ distrib. uniforme (8,28)

- Softmax = trasformazione matematica,
 - converte un insieme di numeri in una mappa di probabilità
 - L'output in apprendimento automatico
 - non è direttamente una probabilità
 - 2 fornire una distribuzione di probabilità che codifica previsioni utili x prendere una decisione

$$\sigma(\alpha_i) = \frac{e^{\alpha_i}}{\sum_{j=1}^n e^{\alpha_j}} \quad (8,29)$$

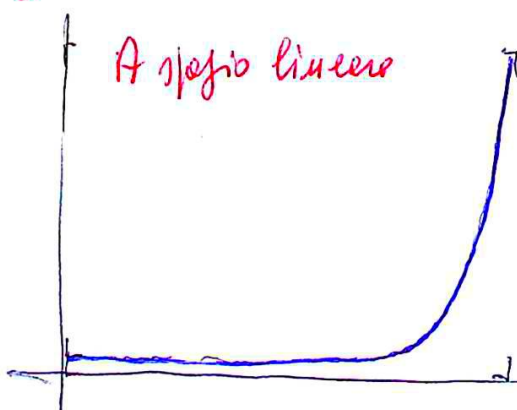
Softmax non confondere con deviazione std = σ_x pedice

Es. numerico

Row	Softmax
4	0.042
5	0.114
7	0.844

prevedendo un intervallo + ampio, graficando: Fig. 8,9

Softmax



Softmax

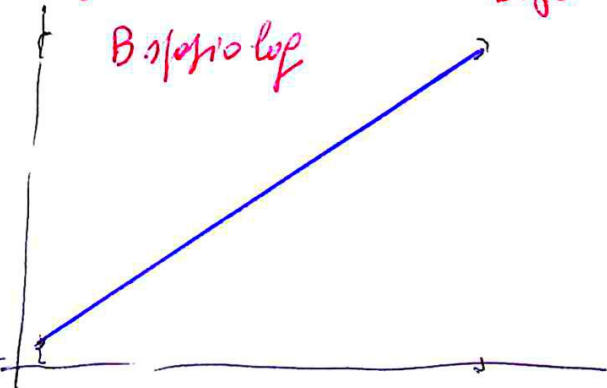


Fig. 8,10

$\sigma(x)$ produce veramente una funzione di mona? 46

Si perché soddisfa i requisiti di una distribuzione di probabilità:

- eventi reciprocamente esclusivi
- $0 < \sigma(x_i) < 1$
- $\sum \sigma(x_i) = 1$ (8,30)

A cosa serve softmax?

- classificazione in ambito apprendimento automatico
- Es. foto cane-gatto
- trasforma il set di pixel dell'immagine in una label
- [cane, gatto]
- questi valori sono trasformati in una lista di probabilità mediante l'uso di softmax
- in base ai valori di probabilità \rightarrow si selezionano etichette

Es. predire se: no tumore, tumore benigno, tumore maligno

- per un particolare paziente \dots una tabella simile a **Figo 8.19** - viene processata con softmax

4	0,042	— 4% no tumore
5	0,114	— 11% benigno
7	0,844	— 84% maligno

8.9 - Exercises

1) esploriamo le unite' dei pdf dati da "saipy"
 — distrib. Gauss

— creiamo pdf e partire da Gauss

$x_4 = \text{mp.linspace}(-4, 4, 400)$

$\text{pdf}_4 = \text{stats.norm.pdf}(x_4)$

normalizzato

$\text{pdf}_4N = \text{pdf}_4 * (x_4[-1] - x_4[0])$

$\text{mp.sum}(\text{pdf}_4N) \text{ ————— } = 1$

2) come prima, ma differente intervallo

————— $\text{mp.linspace}(-2, 2, 300)$

[n'eto, verifico che = 1

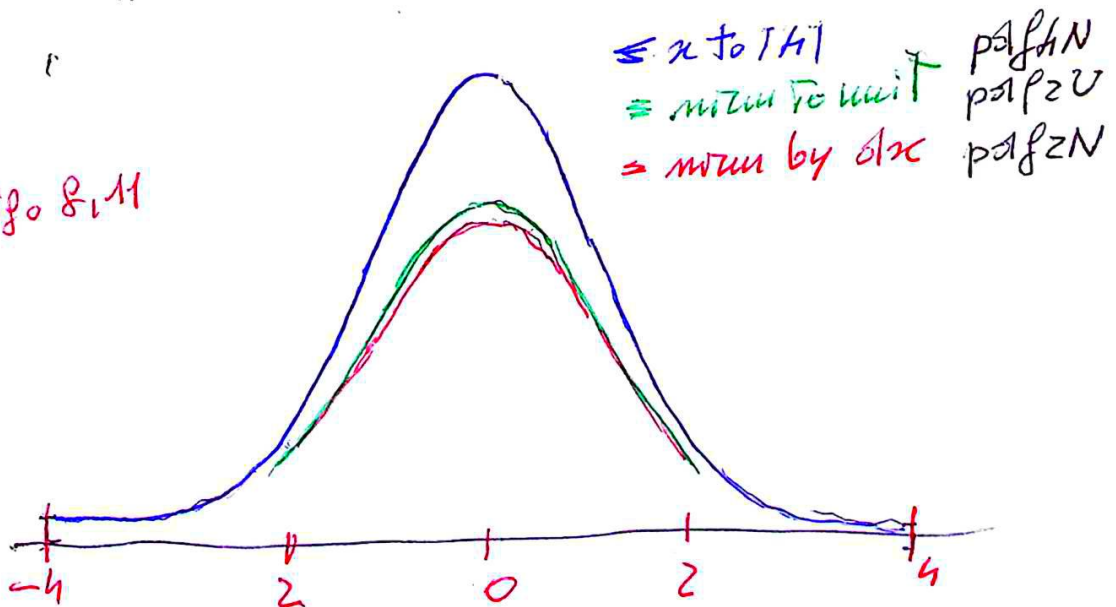
— in realtà ottenuto .955

Provo una diversa normalizzazione

$\text{pdf}_{2U} = \text{pdf}_2 / \text{mp.sum}(\text{pdf}_2)$

$\text{mp.sum}(\text{pdf}_{2U}) \text{ ————— } = 1$

Fig 8.11



le tre distribuzioni sono tutte correlate!
 L semplicemente abbiamo diviso x numeri diversi

3) due pdf, entrambi $-4 \div 4$ $\left\{ \begin{array}{l} 100 \text{ punti} \\ 1000 \text{ punti} \end{array} \right.$

uno un ciclo for (2 steps) x grafico

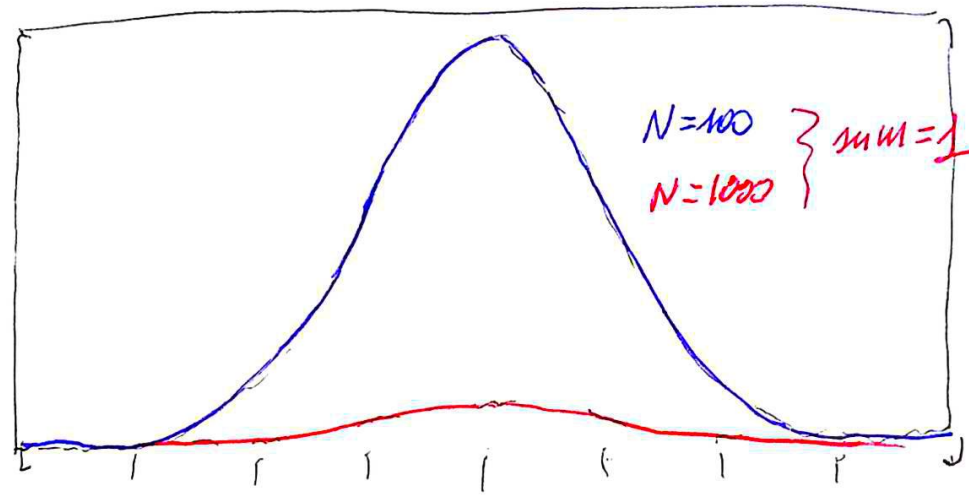


Fig. 8.12

4) poniamo e "cdf"

te -4 e 4 , 300 steps - creiamo pdf

- facciamo la somma cumulativa (1)
- sulla normalizzata (2)

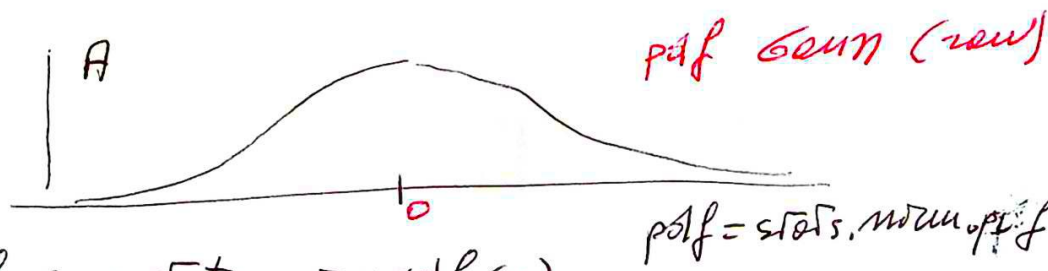


Fig. 8.13

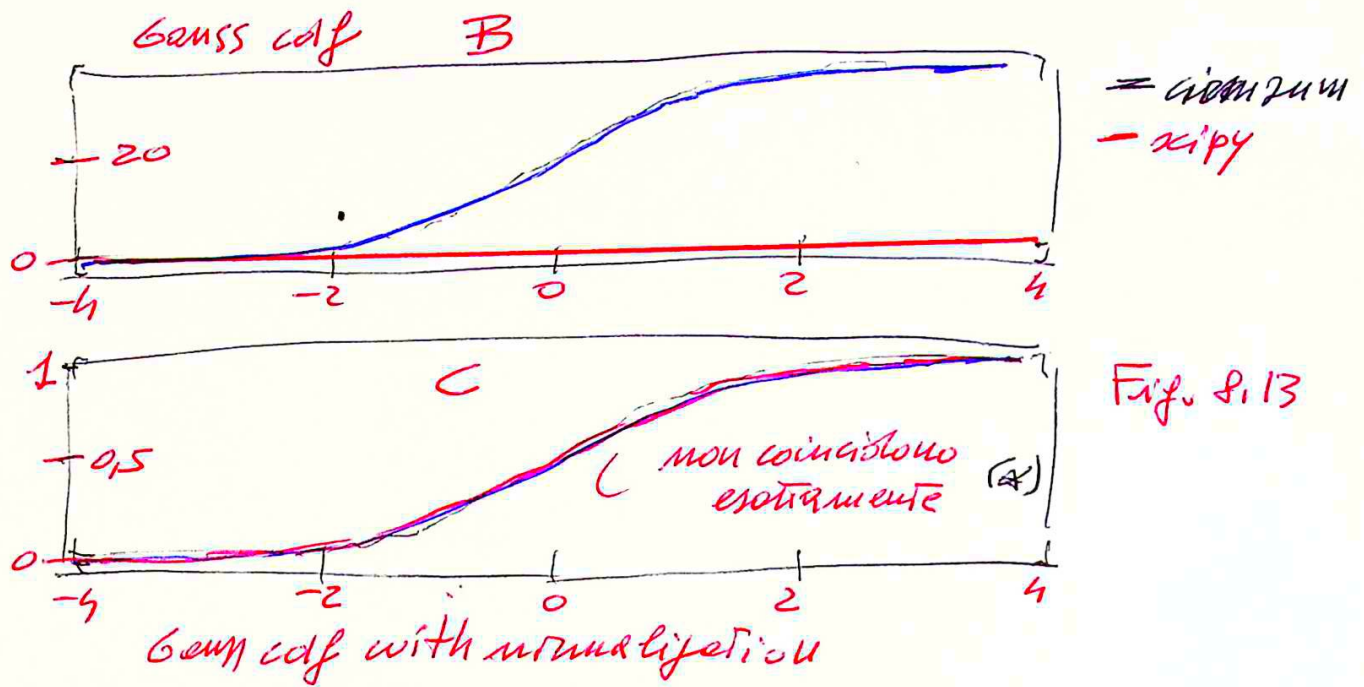
pdf-sp = stats.norm.pdf(x)

calcolo manuale

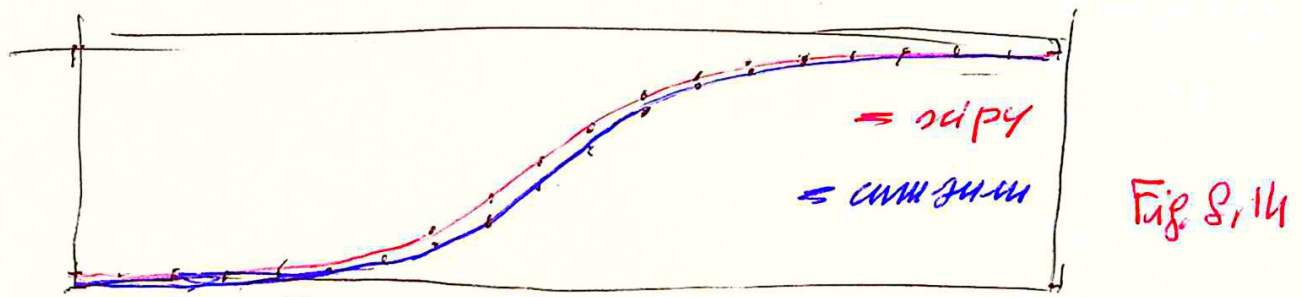
cdf-my = mp.cumsum(pdf)

cdf-myN = mp.cumsum(pdf) * (x[i] - x[0])

grafico:



5) per curve meglio (*) mi sono misurato con risoluzione minore (20 punti) - ricostruisco il grafico C -



la verità è che "scipy" è + accurata - Le differenze sono tanto maggiori quanto migliore è Δx

6) Calcoliamo cdf empirica non usando Python

$$x = \text{np.linspace}(-6, 6, 1001)$$

$$pdf = \text{stats.norm.pdf}(x - 2.7) + \text{stats.norm.pdf}(x + 2.7)$$

normale

$$cdf = \text{np.cumsum}(pdf) * \text{np.mean}(\text{np.diff}(x))$$

migliore normale

$$pdfN = pdf / \text{np.sum}(pdf)$$

$$cdfN = \text{np.cumsum}(pdfN)$$

grafico

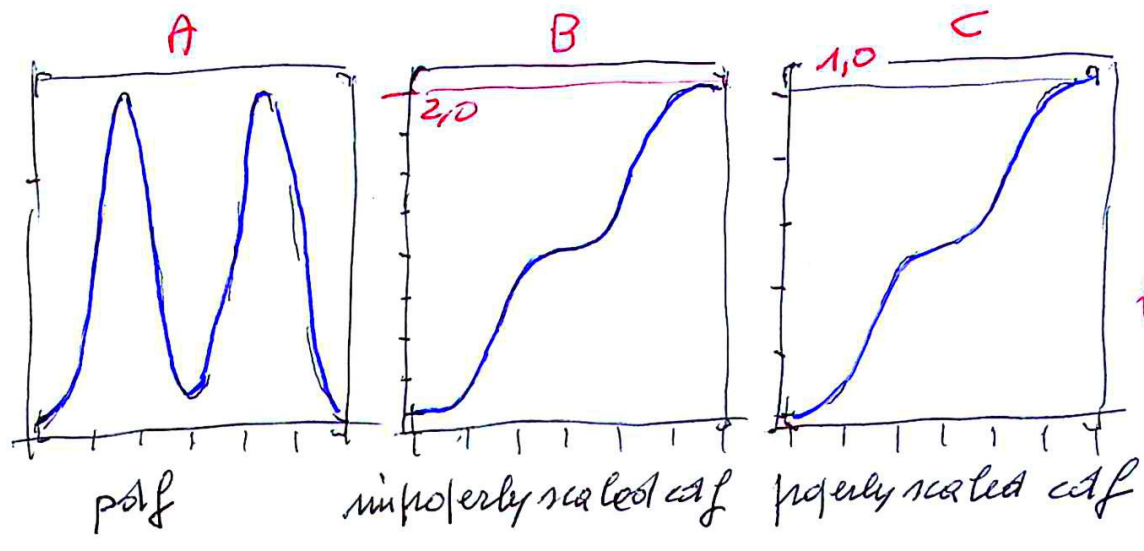


Fig. 8.15

Aveendo sommato 2 pdf, il risultato non e' un vero pdf, e' utile se scalato con dx -

- meglio quindi preventivamente calcolare pdfN, e solo dopo calcolare cdfN

7) ancora cols - TEST: pdf fuo' essere calcolato come derivata di cols -

$$n = \text{mp} \cdot \text{linstace} (0, 10, 200)$$

$$\text{cdf} = \text{stats} \cdot \text{lognorm} \cdot \text{cdf} (n, 1, 1/2)$$

$$\text{pdfD} = \text{mp} \cdot \text{diff} (\text{cdf}) \quad \text{---} \text{ da derivata}$$

$$\text{pdfA} = \text{stats} \cdot \text{lognorm} \cdot \text{pdf} (n, 1, 1/2) \quad \text{] analitica}$$

$$\text{pdfA} * = n [1] - x [\varphi]$$

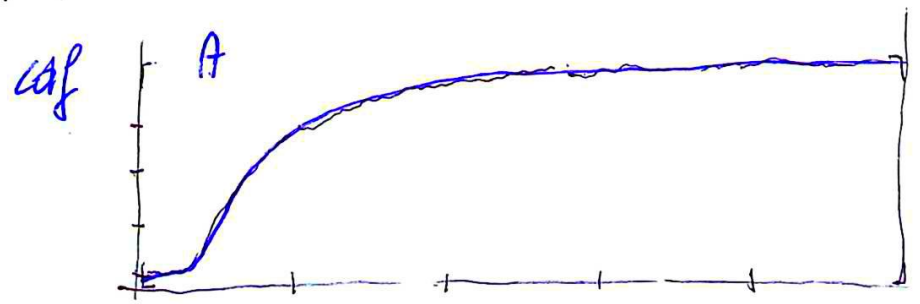
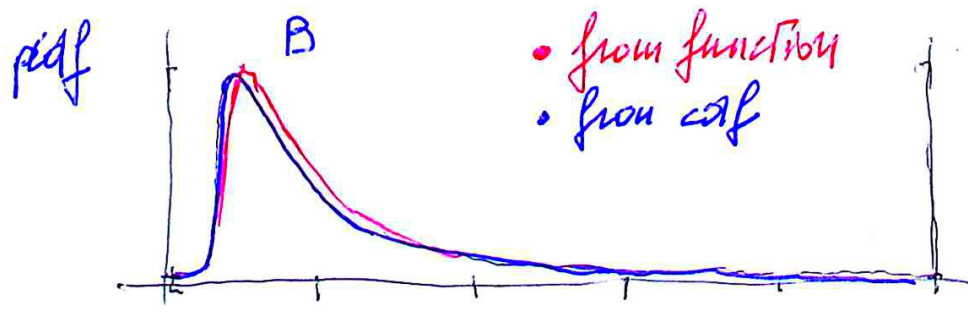


Fig. 8.16



8) differenza tra probabilità empiriche e analitiche -

40 biglie blu / 30 gialle / 20 arancioni / $\Sigma = 90$

- avremo un vettore di 90 elementi

40 elem. = 1 / 30 elem = 2 / 20 elem = 3

- scegliamo con estrazione casuale una biglia nell'urna

- ripetiamo l'estrazione 500 (con sostituzione)

L registriamo i totali x ogni colore

- convertiamo in proporzioni

- tracciamo la proporzioni empirica e quella analitica

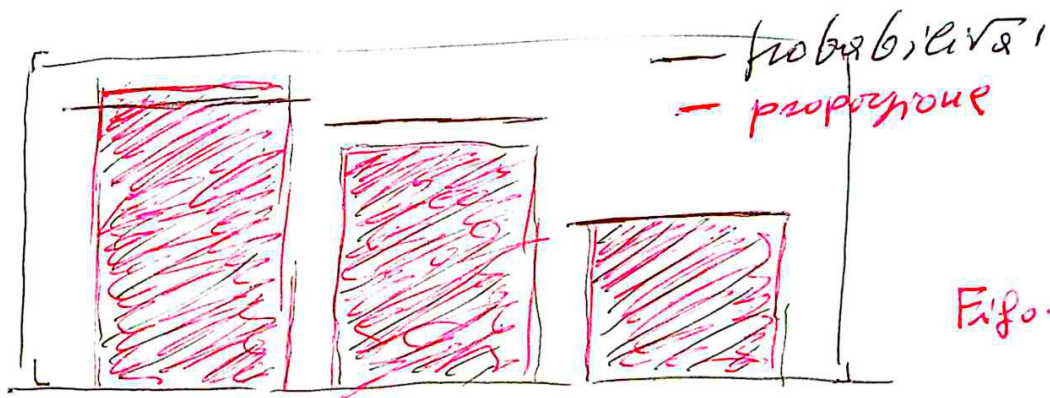


Fig. 8.17

9) nel precedente probabilità e proporzioni non sono uguali ma simili - Aumentando i campioni convergerebbero (legge dei grandi numeri) -

Vediamo di misurare con successo variando n° campioni

L x quantificare questo misuro la formula della

root mean square - una misura tipica di concordanza

in statistica. p = probab. analitica / \hat{p} = probab. empirica

$$RMS = \sqrt{\frac{1}{c} \sum_{i=1}^c (p_i - \hat{p}_i)^2} \quad (8,31)$$

*RMS = ϕ
concordanza
perfetta*

indice del colore $\rightarrow c=3$

- prendiamo il codice dell'esercizio precedente e lo mettiamo in un loop "for", variando la dimensione del campione da 20 a 2000 in passi di 10
- memorizzare ogni RMS x ogni dimensione del campione
- graficare

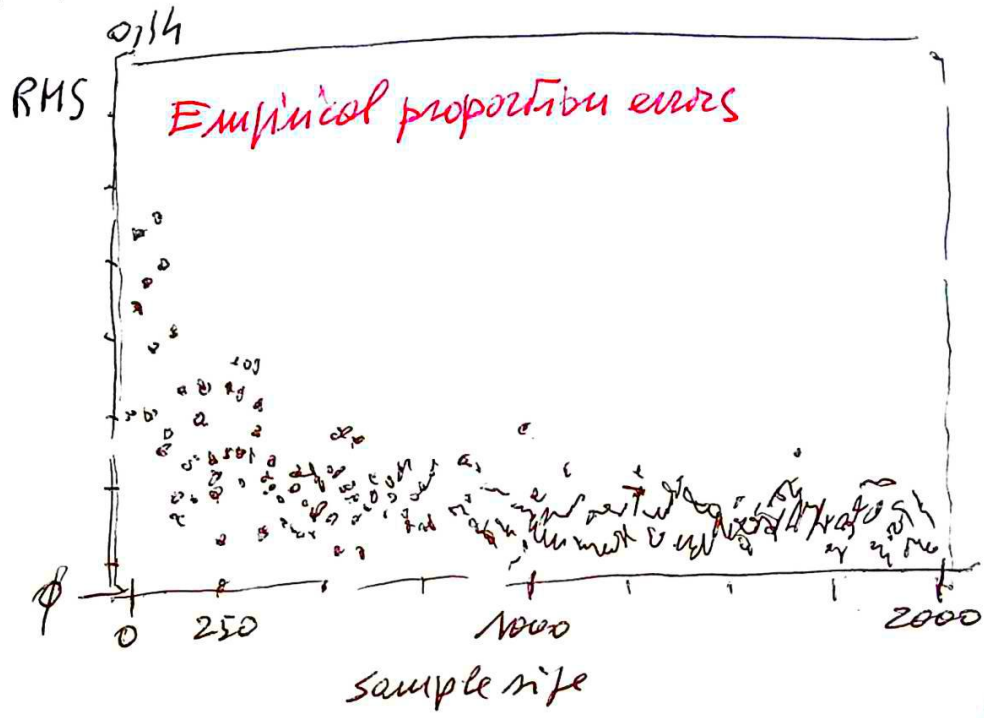


Fig. 8.18

- prova con 20 e 20'000 e passi di 100
- aumentare il n° di biflie di un fattore 10

legge grande numeri

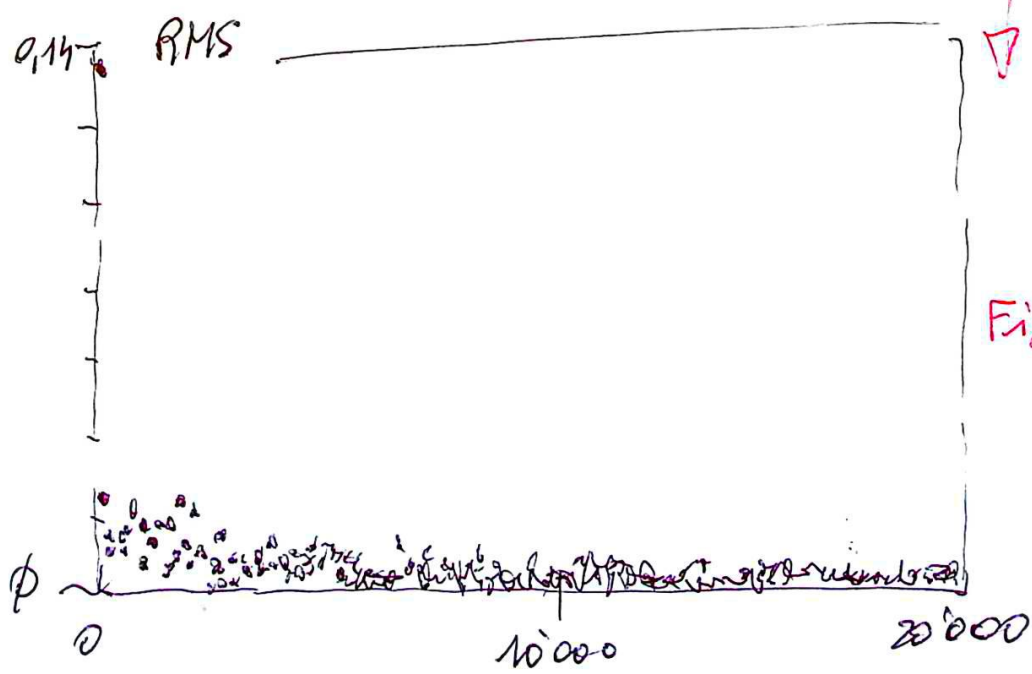


Fig. 8.18-bis

10) calcolare funzioni di proporzione cumulative empiriche, utilizzando dati simulati -

- copieremo i valori "p" e una e due code - **Capo 10**

$N = 1000$ casuale gaussiani

creare vettore di 41 numeri spazati linearmente, tra $-3, 3 \rightarrow$ **servono come "cut-off" boundaries**

ci riferiamo ad esso con la lettera greca zeta ζ

- essendo **60%** \rightarrow la maggior parte dei valori sono sopra $\zeta = -3$ // la metà dei dati sopra $\zeta = 0$ // la maggior parte dei dati $\zeta < 3$

- facciamo un ciclo sopra ogni valore di ζ (1) calcolando la proporzione di dati sotto ζ

- calcoliamo la proporzione di dati sopra ζ (2)

- Visualizzare **Fig 8.19**

$N = 1000$

$K = 41$ - **bins**

$data = mp.random.randn(N)$

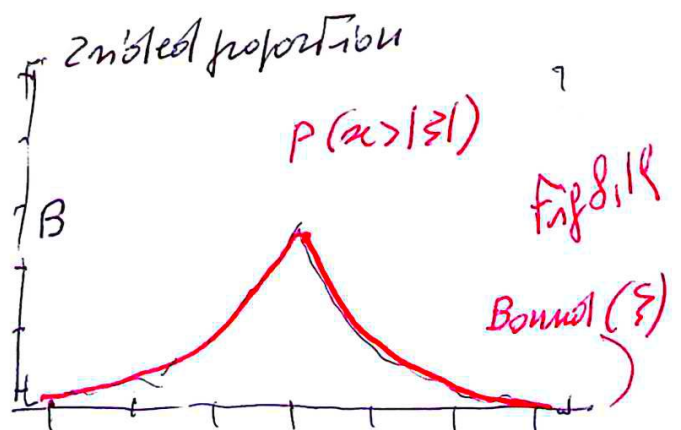
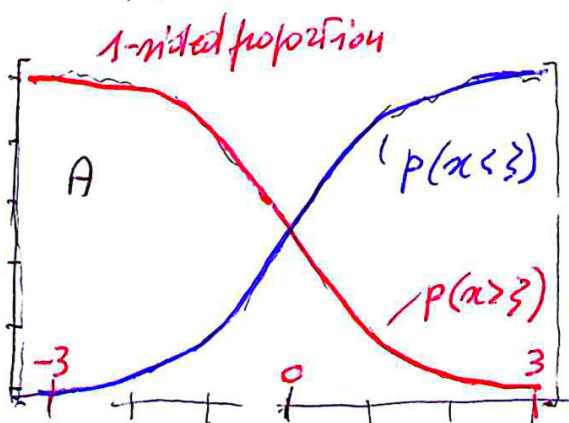
$bounds = mp.linspace(-3, 3, K)$

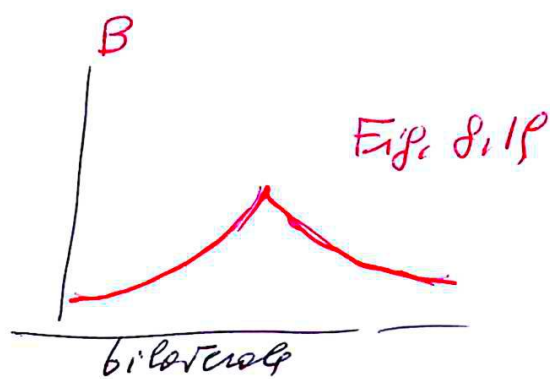
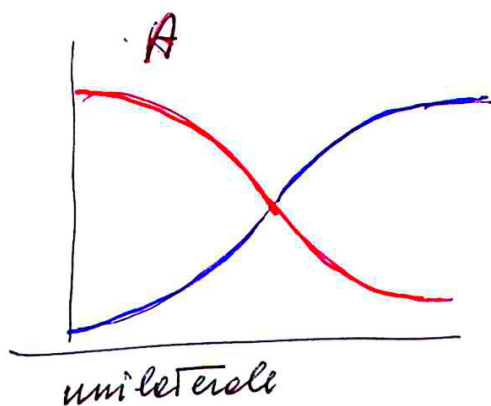
\overline{for} z in $bounds$ enumerate($bounds$):

$emp_prop_GT = mp.sum(data > z) / N$

$emp_prop_LT = mp.sum(data < z) / N$

$emp_prop_z_tail [z] = mp.sum(data > mp.abs(z)) / N$

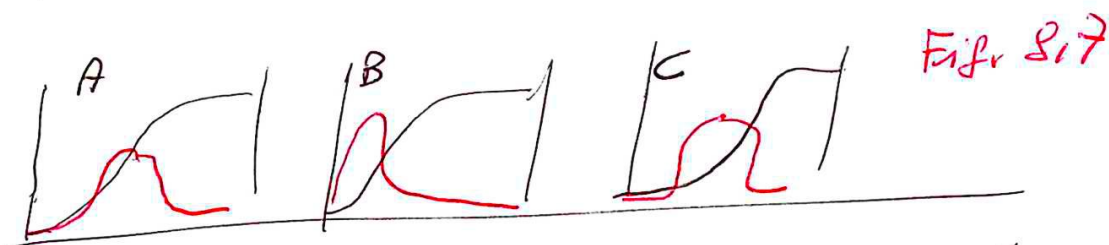




ogni funzione di proporzione unilaterale è essenzialmente una stima empirica di una *cdf* (o del suo *qf*) (A)

ogni funzione di proporzione bilaterale non sembra una *cdf*, ne sembra il *pdf* di una *Cauchy*

11) vedere la 8.7



usando distribuzioni prese dalla libreria di "scipy" e che offrono i riferimenti. Ne troviamo un elenco qui <https://docs.scipy.org/doc/scipy/reference/stats.html> #

continuous distributions

— FINE CAPITOLO 8 —

9 - Sampling and distributions

9.1 - Sampling variability and its annoyance

- Sei uno studente dei dati, ma spesso si fida e vuole sapere se uomini o donne sono più propensi a fare esercizi in determinati momenti della giornata -
- le scelte qui sarebbero troppo personali, quindi non generalizzabili
- allora decidi di campionare 1000 uomini, 1000 donne --- ma vedi i dati \rightarrow ricampioni altri 1000+1000 ma su persone diverse
- ci possiamo aspettare che le due distribuzioni siano simili, ma non identiche
- quindi abbiamo una variabilità nei campioni (non la variabilità dell'intero di 1 campione)
- questa distinzione ha impatto su statistiche inferenziali, su intervalli di confidenza p. 310
- la differenza tra le due variabilità è la base di ANOVA
- **La variabilità del campionamento è fastidiosa** che influenza la variabilità a un livello diverso da quello che abbiamo considerato finora.
- la motivazione a misurare un campione - invece che un individuo - è che quest'ultimo potrebbe non essere rappresentativo della popolazione. **Ma il campione ha una sua variabilità** \rightarrow non è garantito che sia un parametro affidabile

- pertanto ci fidiamo di fin' della media di un campione, che non di un singolo μ_0
- questo e' vero ma... - vedi veridicitá del campione -
- se i campioni sono grandi es. $N=10'000$, il problema si risolve ^{meno}

Un esempio con dati casuali

$N=50$, ciascuno dei quali comprende 500 valori casuali

$N=500$

$nSamples=50$

$kHistBins=20$

$edges = np.linspace(-3, 3, kHistBins)$ bin x info

dichiaro le matrici

$allHistY (50, 20)$

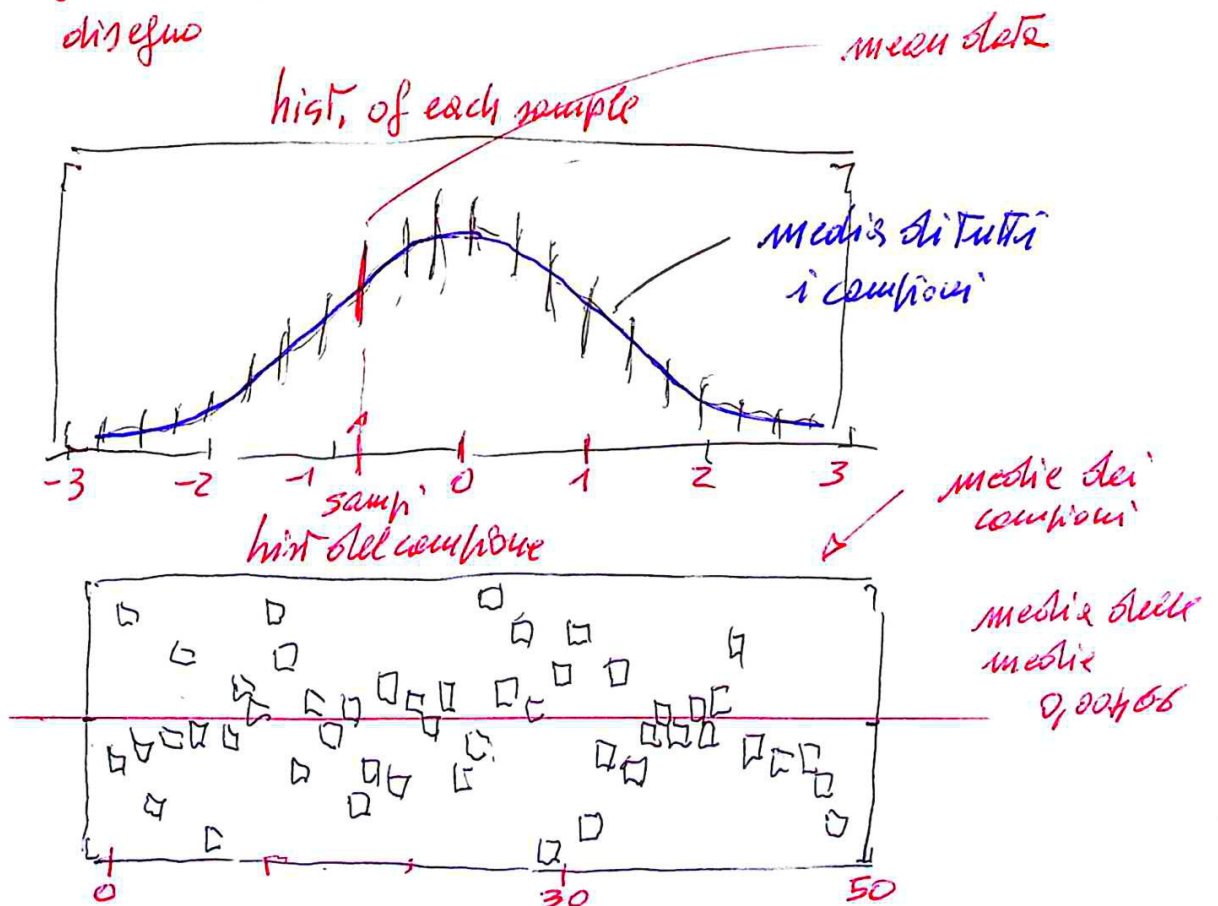
$allMeans (50)$

for sampi in range (nSamples):

$data = np.random.normal(loc=0, scale=1, size=N)$

facio histo, calcolo la media μ_i

disegno



A cosa è dovuta la variabilità del campione?

È inutile cercare x tentare di ridurre.

Natural variation

È comune nei dati biologici, nella neurofisiologia, nella psicologia, nei trattamenti medici, nelle pratiche culturali.

Es. persone diverse reagiscono in modo diverso a pari stimolo →

- (1) è probabile che campioni diversi abbiano caratteristiche diverse
- (2) è ancora + probabile che campioni di età, sesso, stile economiche, background culturali diversi, siano diversi

Nota Variabilità anche in sistemi non biologici:

magnitudo di terremoti, n° di stelle in una galassia, altezza degli edifici in una città.

Rumore da misurazione i sensori sono imperfetti

Dinamiche e cambiamenti la natura è dinamica - Stesse misure in tempi diversi possono venire **Es. cambiamenti geopolitici**

Sistemi complessi la maggior parte degli effetti di un sistema scientifico sono complessi **Es. con retroazioni**

Stocasticità (casualità) c'è casualità nell'universo che non possiamo misurare, che non comprendiamo **Es. moto browniano, eventi**

In concreto, limitare la variabilità, senza compromettere la generalità è un processo complicato.

9.2 - Creating sample estimate distributions

Es. K campioni indipendenti, presi da una popolazione

$S =$ i singoli campioni nella loro interezza

$S_i =$ i mo campione

- ciascun campione ha il suo valore medio $= \bar{S}_i \rightarrow$

K campioni producono K medie - Fig. 9,2 A

A) Data & means

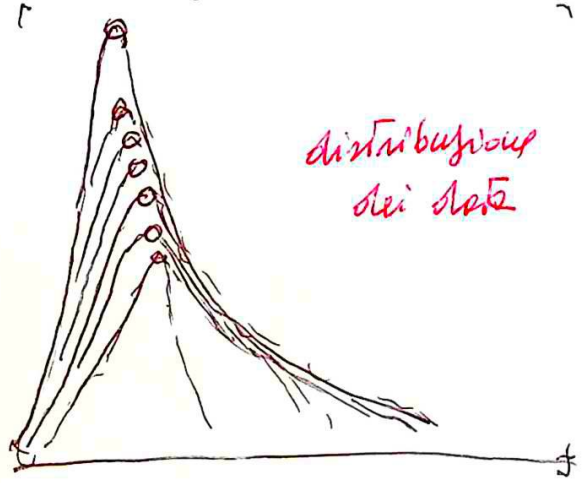
S_1	S_2	...	S_K
0	0		0
0	0		0
0	1		1
1	0		1
1			1
0			0
⊕	⊕	...	⊕

data points

medie

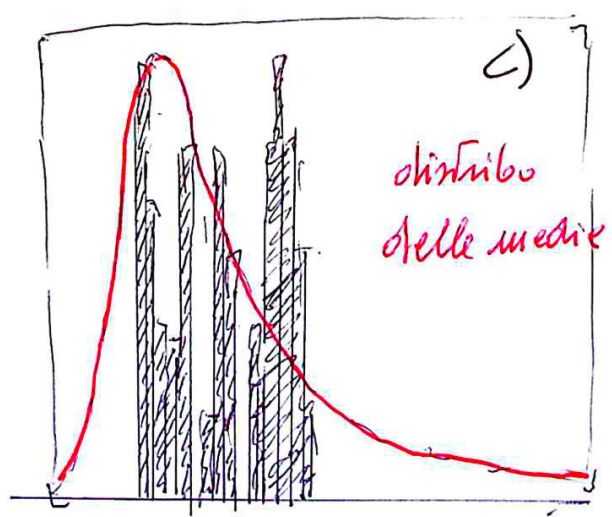
ciascun S_k ha la sua distrib.

B)



distribuzione dei dati

Le medie dei campioni sono una collezione di K valori, che producono un loro istogramma. Questo non è calcolato dei "raw data", ma viene dalla statistica descrittiva che abbiamo calcolato.



distribuzione delle medie

Importante la distribuzione delle medie da K campioni

53

non è la stessa della media degli n inf. campione
(vedi Teorema del limite centrale)

Perché le distribuzioni di stima del campione sono importanti?

- esse rivelano la variabilità intrinseca della popolazione
- esse consentono di calcolare gli intervalli di confidenza

Si possono calcolare le distribuzioni di stima del campione e
statistica descrittiva (es. std, mediana, coeff. di correlazione...)

9.3 - Standard error of the mean

- La "sample mean" è una stima della media di una popolazione
L se faccio più "sample mean" ottengo valori diversi →

Quanto è precisa la mia stima?

Io dico l' "errore std della media" = SEM = ^{statistica} descrittiva

SEM = quantità di variabilità che ci si può aspettare nelle
("medie di campioni casuali" ripetuti con la stessa
dimensione, ricovero della popolazione)

- viene anche detto deviazione std della popolazione

in formula

$$SEM = \frac{\sigma}{\sqrt{N}} \quad (9.1)$$

- campione

- in concreto però non conosciamo σ , al suo posto usiamo
la std del campione

$$SEM \approx \frac{s}{\sqrt{N}} \quad (9.2)$$

$$SEM \approx \frac{s}{\sqrt{N}} \quad (8,2)$$

s = std del campione $e' \approx \sigma$ *Se è*

- il campione è random
- _____ è rappresentativo
- N è sufficientemente grande

Std error of the mean vs. std deviation del campione

Entrambi quantificano la variabilità associata a un set di
 Anche matematicamente sono simili, ma riflettono aspetti
 diversi dell'incertezza \rightarrow *interpretazioni distinte*

Conceptual meaning

std deviation = dispersion dei valori attorno alla media del campione

SEM = precisione della media del campione, *intesa come stima della*

della popolazione
 = quanto è probabile che s sia vicina a σ

Calculation

std = si ottiene moltiplicando ogni s in un campione

Applications

- come misura della variabilità in un "set sample"
- in trasformazioni di tipo z-transf
- a generare intervalli di confidenza
- a valutare la significatività statistica delle statistiche di test
 (es. t-values, coeff. di regressione)

Impact of sample size

54

La deviazione std non varia con la dimensione del campione
SEM diminuisce all'aumentare della _____,
ma std no.

9.4 - Random and representative sampling

Lo scopo delle statistiche inferenziali è valutare ^{se} le caratteristiche di un campione si generalizzano a una popolazione.

Campioni rappresentativi:

ES. Abbiamo dati su quanti soldi delle famiglie americane finiscono in cibo spazzatura.

Potremmo da qui prevedere la morte delle stelle? No, il campione non è rappresentativo.

ES. Dati sulle opinioni degli studenti universitari sul reddito di base universale. Potremmo da qui inferire come la popolazione concepisce UBI? **No** x che gli studenti universitari non rappresentano la popolazione.

Forse potrebbero esserci correlazioni con studenti univ. di altre nazioni.

- Consideriamo ancora questi studenti.

Ma ci riferiamo a un compito di percezione visiva in cui bisogna avere due forme 3D motate. Lo possiamo generalizzare? Potrebbe essere ragionevolmente sì.

Campioni casuali

È una delle modalità x avere campioni rappresentativi -

Ritorniamo a UBI // supponiamo che gli intervistati siano tutti isaiti & "equite economica" - Le loro opinioni saranno diverse da quelle di altri studenti -

Gli studenti isaiti & "tone equite" anche em' enorme istee diverse -

→ raccogliere dati casualmente su tutte le popolazioni universitarie

→ dati raccolti casualmente da tutta la popolazione sotto studio

Attenzione offri denaro a chi ti lascia intervistare? questo crea una selezione di caso -

offri crediti formativi? // stai chiedendo fieno in

supermercato? probabilmente selezioni femminili o senza lavoro // compri numeri telefonici? comunque un sotto insieme

Raccogliere campioni casuali è cosa complicata -

Se possibile confrontare le statistiche descrittive del campione con statistiche note della popolazione - Es. età, equilibrio di genere, anni di istruzione del proprio campione con quelli ISTAT -

Se qui c'è corrispondenza → + fiducia che i dati possano essere generalizzati -

Buona notizia se i campioni individuali hanno un bias, ma differenti campioni hanno un bias diverso - mediando le medie - in qualche misura si compensano ("tutte grandi numeri") -

Independent and identically distributed def

IID (*i.i.d.*, *i.i.d.*) è un termine tecnico usato nella matematica stats.

IID = fornire una definizione più approfondita di "random part" di dati casuali e rappresentativi. Secondo IID:

- Indipendenti, cioè "non correlati" tra loro *Eg.*

- campionare in base all'ultima cifra del n° di fruizione sociale *garantisce indipendenza*
- campionare nello stesso nucleo familiare = *dispendente*

- Identica distribuzione nei dati empirici questo è comunque approssimativo

9.5- Law of Large Numbers

Il campionamento è una necessità pratica. È ragionevole che campioni più grandi riflettano + accuratamente la popolazione.

$$LLN = \lim_{n \rightarrow \infty} P(|\bar{x}_n - \mu| > \epsilon) = 0 \quad (9.3)$$

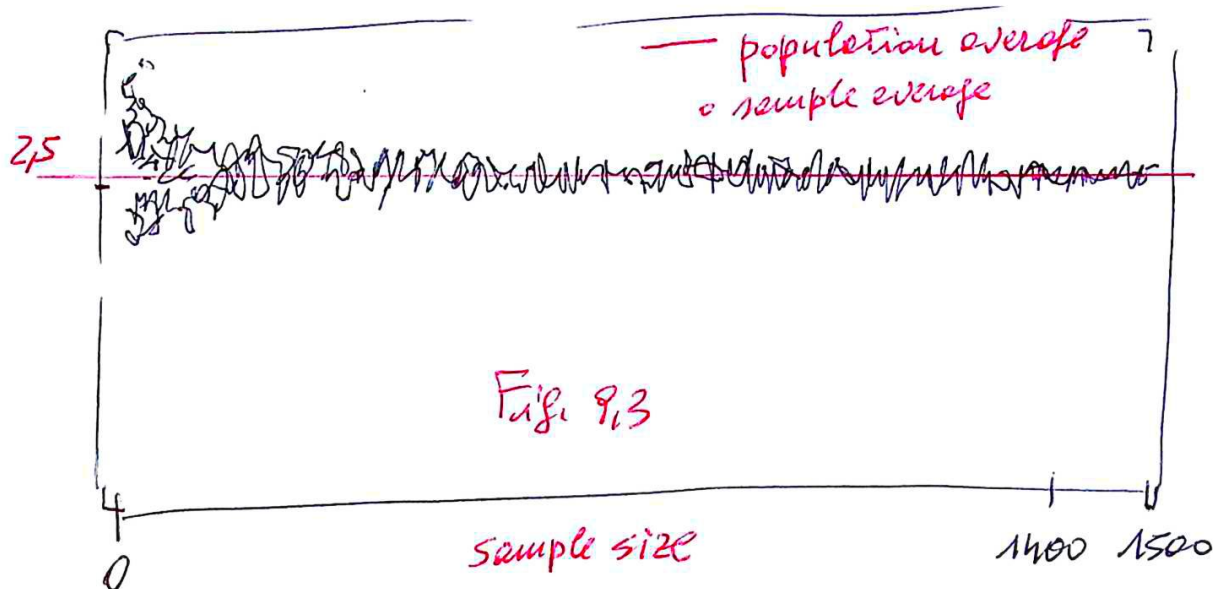
sample size / *sample mean* *true population mean*

LLN and sample size (LLN demo #1)

Facciamo una dimostrazione usando i dati simulati, con che possiamo riprodurre tutta la popolazione, quindi conoscere la vera media

- ripetiamo [1, 2, 3, 4] molte volte → 4'194'304 dati simulati
- media esatta = 2,5

- prendo campioni casuali da 1 a 1500 Fig. 9.3



- (1) le medie dei campioni non sono esattamente $= \mu$ - sembrano vicine
- (2) sembra che non ci siano bias
- (3) la variabilità diminuisce, aumentando le dimensioni
- (4) anche a 1500 non sono esattamente uguali

LLN and repeated samples (LLN demo #2)

L'idea è che ogni campione (ogni esperimento) sia sensibile alla variabilità ^{alea} casuale (rumore, altre variazioni non sistematiche) →

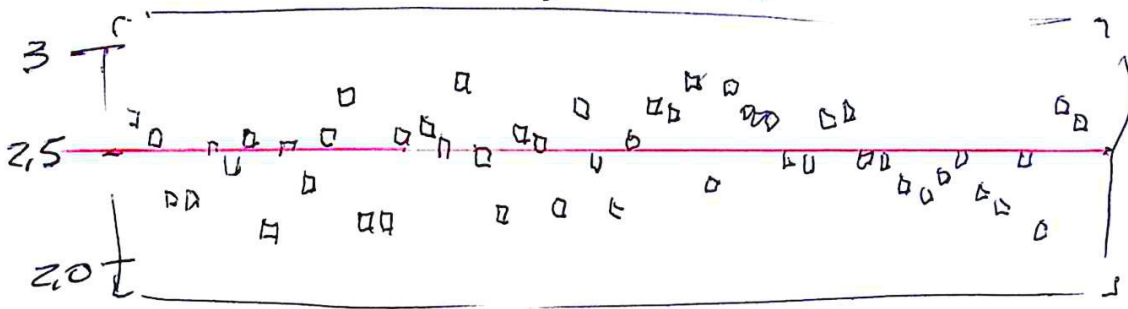
È improbabile che un campione possa fornire una buona stima di μ .
→ consideriamo la media delle medie dei campioni.

Per mostrare la forza di questa procedura, facciamo questa simulazione:
Della popolazione precedente prendo 50 campionamenti, ciascuno di 30 elementi. Calcolo la "ammoletive average" Fig. 9.4B.

I campioni vanno da 1 → 1500

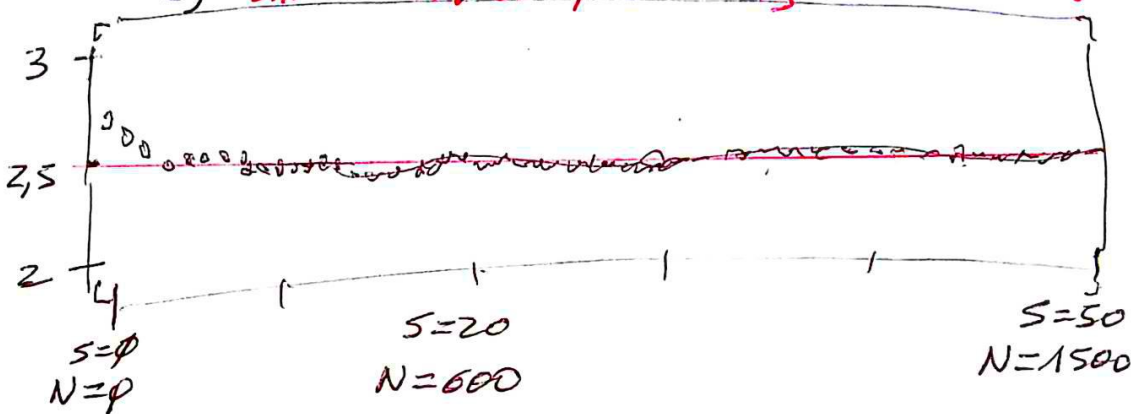
In un ciclo for memoria tutte le medie - Grafico medie e medie cumulative -

A) each sample mean



B) Cumulative sample means

Fig. 9.14



- (1) Le medie dei campioni individuali non hanno bias (*)
- (2) La media cumulativa si è approssimata a un solo dato un piccolo numero di campioni → i campioni individuali hanno una variabilità maggiore di quella delle loro medie

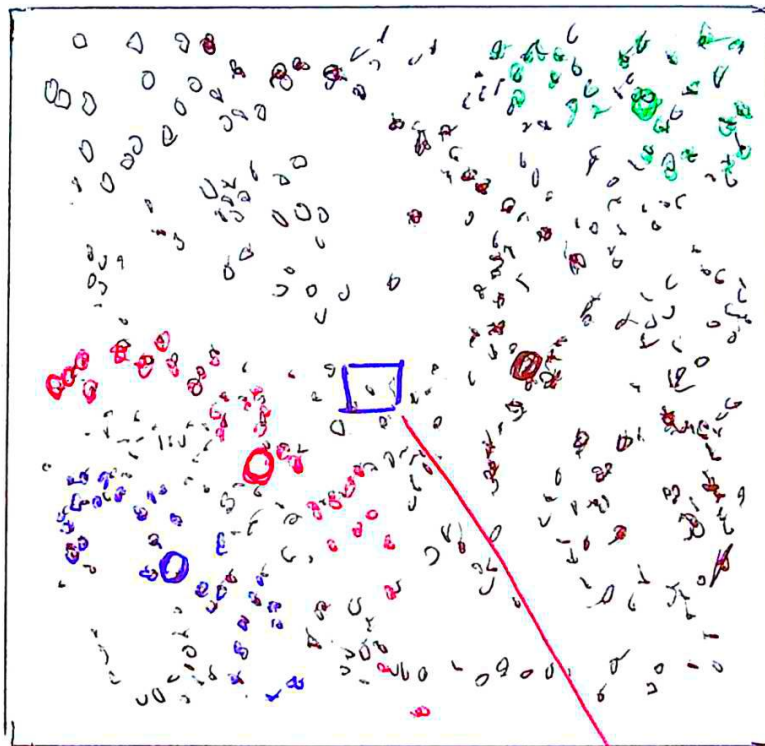
Completiamo ora la Fig. 9.5 che rappresenta una popolazione di dati su due variabili. Da questi dati estraiamo due campioni, la rosa e la blu. Entrambe hanno la media fissata al centro della figura. Questa è una rappresentazione vicina del concetto (*) -

Nella realtà i campioni sono sempre fogli (Tempi e costi) → combinare i dati di più studi indipendenti porta a una valutazione + accurata della realtà (meta-evalui). È particolarmente importante nella ricerca medica.

- i campioni dei diversi frutti possono contenere distorsioni (bias) statistiche e campioni non rappresentativi, come in profetie false, stime e regioni maldefinite ---

- Se le distorsioni non sono le stesse \rightarrow LLN comunque funziona + accurate - **Es.**

Immaginiamo di raccogliere 6 campioni, ma ogni campione ha un bias -
La media delle medie è + accurata -



con bias

Fig. 9.6

media delle medie

9.6- The Central Limit Theorem

CLT = una distribuzione di medie di campioni \rightarrow **Good**
anche se le distribuzioni originali non lo sono

CLT = utile a garantire la validità delle ipotesi che usiamo nelle statistiche inferenziali

CLT part 1: sampling distributions

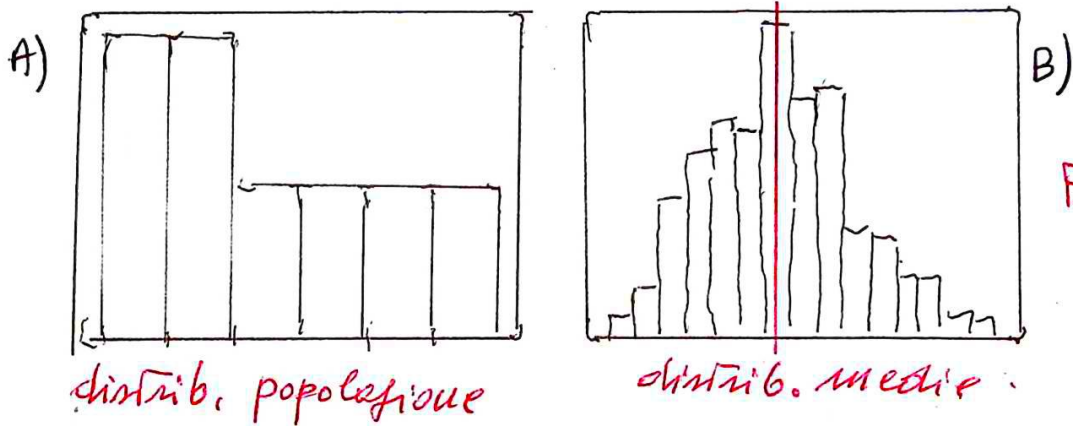
57

Popolazione di interi (1-6) moltiplicando il peso ponderato
del **Cap. 8** - 1,2 $\rightarrow p=1/4$ // gialli $p=1/8$

Fig. 9.7A mostra istogramma della popolazione

- conduco 500 campioni casuali, ciascuno con $N=30$

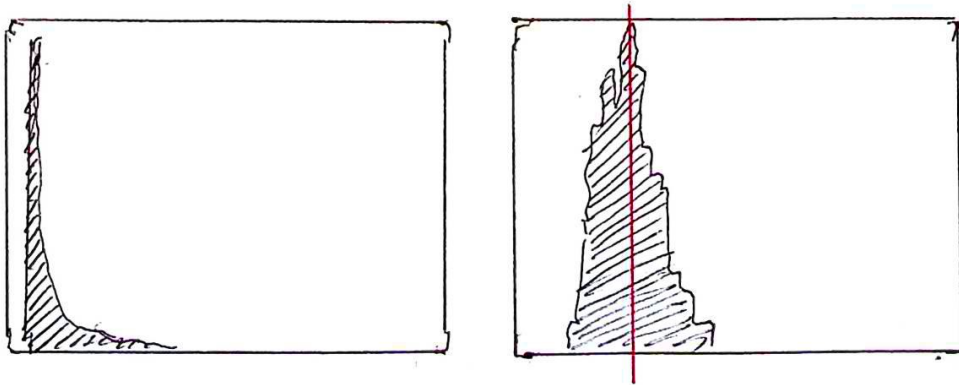
Fig. 9.7B mostra istogramma delle medie (**Green**)



Es.

Dati con legge di potenza - Fig 9.8A - Estraggo 500 campioni con $N=30$
Faccio istogrammi delle medie Fig 9.8B

Fig. 9.8



LLN e CLT formano un duo potente

↳ miglioriamo la predizione, diventiamo **Green**

CLT part 2: mixing variables

CLT = miscela casuale di variabili \rightarrow tende a CLT

L' scegli cose casuali nell'universo e mettile in pila \rightarrow verso Gauss

Fig. 9.8 illustra questo concetto

- due variabili - un'onda sinusoidale Dati 1

- un rumore uniforme Dati 2

P. 329

le due set di dati sommati insieme hanno una distrib. \approx Gauss

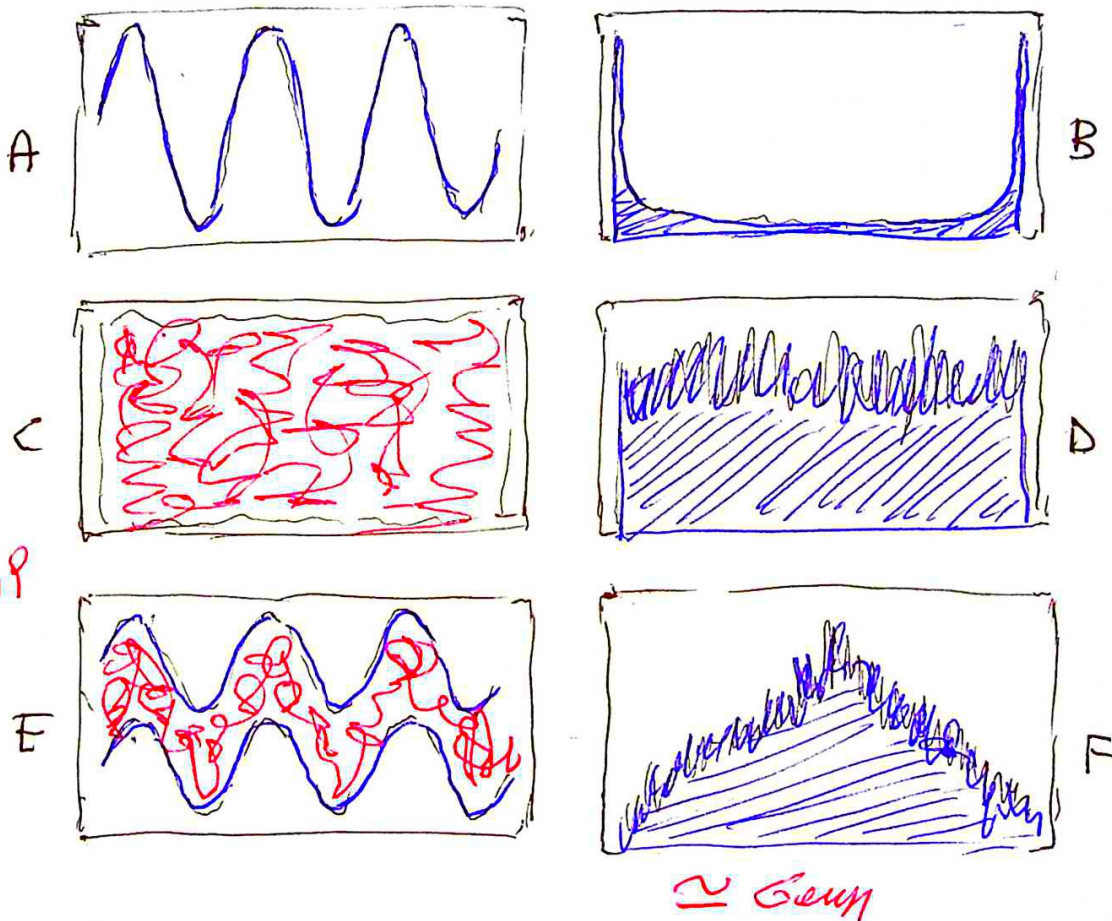
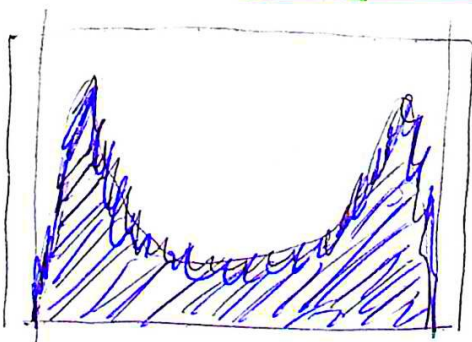


Fig. 9.8

Il risultato non è banalmente generato - Es. la similarità di scala delle variabili - Se riduco il numero di un fattore 10 \rightarrow



F' Fig. 9,10

The distributions of sample means

58

CLT = le medie dei campioni sono GOMM, media μ , std = SEM
in formule

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{N}}\right) \quad (9,4)$$

\bar{x} = media del campione, μ = media della popolazione, σ = std della popolazione, N = numerosità del campione

se $N \rightarrow \infty \rightarrow SEM \rightarrow 0$

NOTA μ è indipendente dal valore di N

Implications of the CLT

- Estimating population means from samples

CLT in combinazione con LLN dice che rappresentando i dati su molti campioni, otteniamo una buona stima della popolazione

- Assumption of normal distribution

Molte inferenze statistiche si basano su GOMM - Ma, nel mondo reale non sempre si ha GOMM. \rightarrow medie dei campioni come base della inferenza statistica \rightarrow intervalli di confidenza.

- Demixing multivariate signals

con immagini e segnali, spesso fin'ora questi sono mischiati.

\exists un metodo di equalità x de-mixere, chiamato metodo delle componenti principali indipendenti, basato sul presupposto che i segnali siano non-Gaussiani.

9,7 - Exercises

1) Esplorazione della LLN

- manipoliamo il mondo

- 200 punti estratti casualmente da $\text{Geom} \text{ o } \mathcal{N}(\phi, \sigma^2)$

varia da 0.1 a 10
in 10 steps

- iteriamo sui valori di σ^2

L calcoliamo, memorizziamo medi, std del campione

L x generare una distribuzione di stima, iteriamo 20 volte x ogni σ^2

tau2 levels = np.linspace(0.1, 10, 10)

sample size = 200

numsamples = 20

results = np.zeros((numsamples, len(tau2 levels), 2))

loop over: multisamples, tau

for m, tau2 in enumerate(...):

for sampi in range(numsamples):

data = np.random.normal(phi, np.sqrt(tau2), size=sample size)

results[sampi, m, 0] = np.mean(data)

results[sampi, m, 1] = np.var(data, ddof=1)

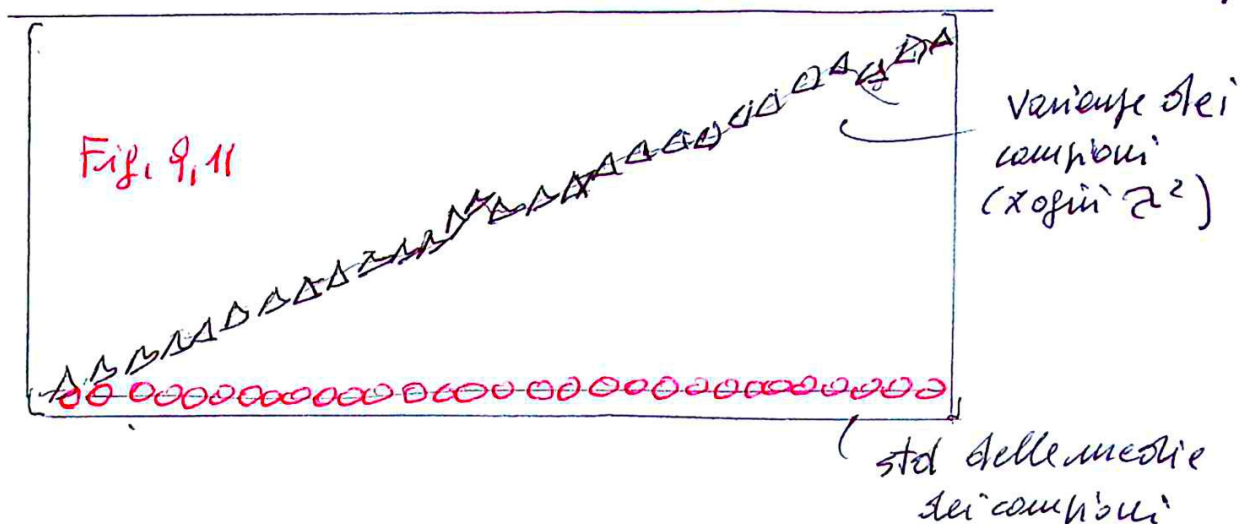
for grafico

stampo results / mean of results [i, j, phi]] 1° grafico

mean [i, j, phi]]
var of results [i, j, phi]] 2° grafico

Fig. 9.11





- 2) Modifichiamo il codice della Fig. 9.3 (campioni e medie al venire di N)
- calcoliamo std (invece che le medie)
 - veniamo random beam con $\sigma = 2,4$
 - LLN e' formalmente definita x le medie, ma n'ottien endue ad altre statistiche descrittive

```
sample sizes = np. range(10, 1001)
```

```
pop_std = 2,4
```

```
population N = 1000000
```

```
population = np.random.randn(population N) ↗ std=1
```

```
_____ = population / np.std(population, ddof=1)
```

```
_____ = _____ * pop_std - freq std
```

```
sample stds = np.zeros(len(sample sizes))
```

```
for sampi in range(len(sample sizes))
```

```
    pick a random sample
```

```
    sample = np.random.choice(population, size=sample sizes[sampi])
```

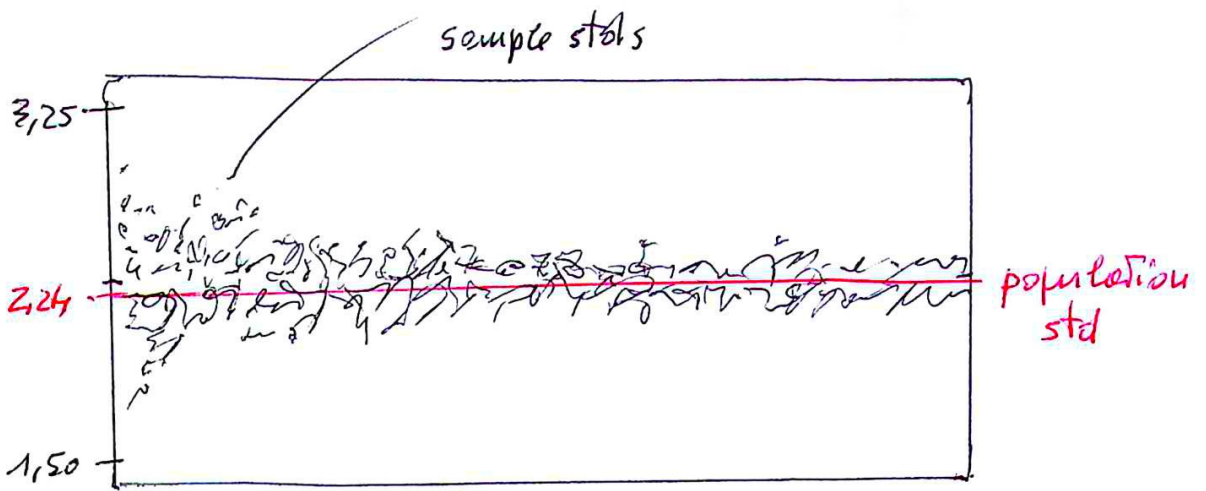
```
    sample stds[sampi] = np.std(sample, ddof=1)
```

plot

```
    sample stds
```

la nra diff. e' costante

Fig. 9,12



3) Calcoliam qui la differenza media tra due campioni, tratti da distribuzioni separate. È un esempio che percorre le statistiche inferenziali (calcoleremo il numeratore di t-test)

$$\text{pop Mean 1} = 3$$

$$\text{2} = 3,2$$

$n = 10^6$

$$\text{population 1} = \text{np.random.randn}(\text{population 1})$$

$$\text{---} = \text{population 1} - \text{np.mean}(\text{population 1}) + \text{popMean 1}$$

analogamente x population 2
un campione

$$s1 = \text{np.mean}(\text{np.random.choice}(\text{population 1}, \text{size}=30))$$

$$s2 = \text{---} \quad \text{popMean 1} - \text{popMean 2} \quad \text{---} \quad -0,200$$

$$s1 - s2 \quad \text{---} \quad -0,585$$

Ripetiamo la procedura

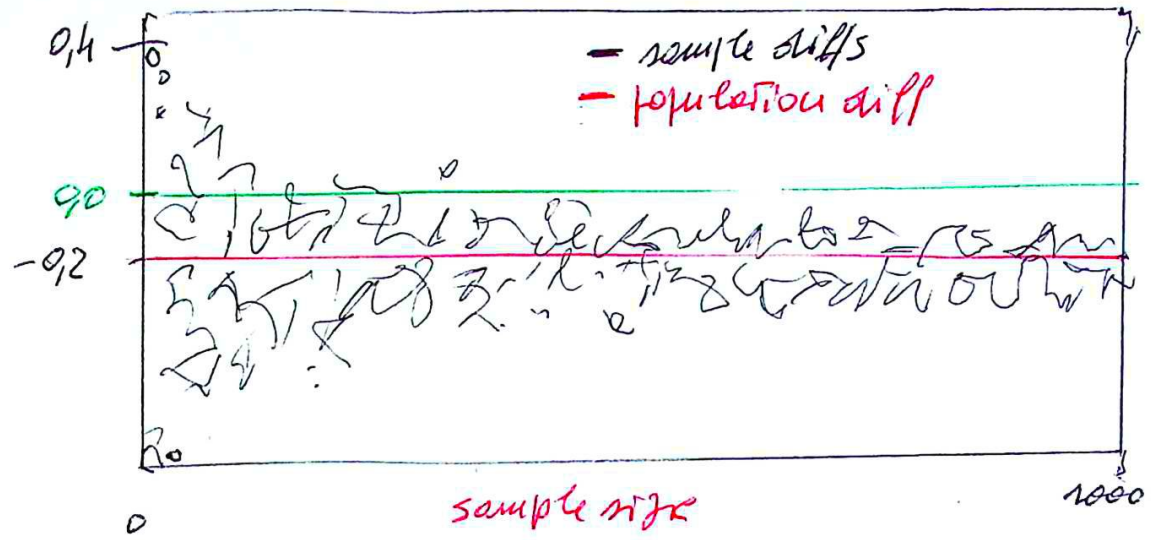
$$\text{sample diffs} = \text{np.zeros}(\text{len}(\text{sample size}))$$

for sampi in range(len(sample size)):

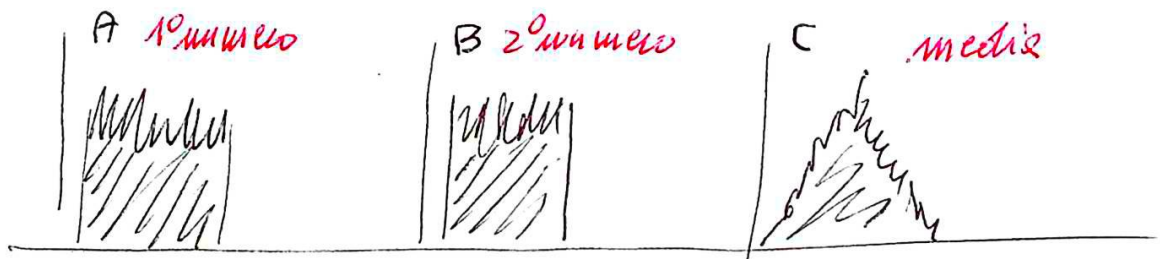
pick a random sample

$$s1 = \text{np.random.choice}(\text{population 1}, \text{size}=\text{sample size}[\text{sampi}])$$

$$s2 = \text{---} \quad \text{sample diffs}[\text{sampi}] = \text{np.mean}(s1) - \text{np.mean}(s2)$$



4) Scegliamo due interi casuali tra ϕ e 100, calcoliamo la media. Ripetiamo 1200 volte - Ad ogni passo memorizziamo questi tre numeri - Facciamo inferenza -



5) CLT funziona se i campioni sono abbastanza grandi -
 Quando e' "abbastanza grande" ?
 Ripetiamo il codice di Fig. 9.8 (dati con legge di potenza)
 N=5 (con n° campioni = 500) poi
 N=100 con osserviamo ?

$$6) N_{pop} = 10^6$$

$$x^2 / \mu \sim \mathcal{N}(91)$$

population = np.random.randn(Npop) * x2

sample sizes = np.arange(5, 500, 8)

number of sampls = 1000

sample means = np.zeros(number of sampls)

fw_hms = _____ (len(sample sizes))

peak_vals = _____

line colors

c = np.linspace(0.9, 0.9, 9), (0, 0, 0), len(sample sizes)

for Ns in range(len(sample sizes)):

calcolo di media di molti campioni

for expi in range(number of sampls):

sample means [expi] = np.mean(np.random.choice(
population, size = sample sizes [Ns])

calcolo dell'istogramma e della media

yy, xx = np.histogram(sample means, np.linspace(0.4, 1.6, 41))

yy = yy / np.sum(yy)

calcolo FWHM

ym = yy / np.max(yy) *normalizzo*

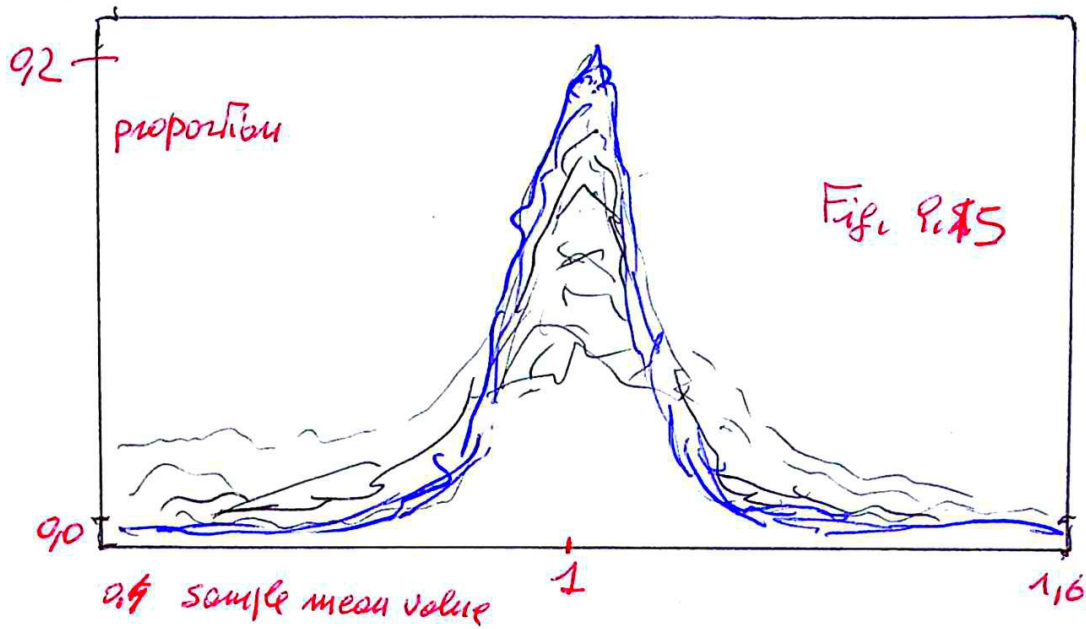
idx = np.argmax(ym) *trovo peak index*

peak_vals [Ns] = xx [idx]

fw_hms [Ns] = xx [idx - 1 + np.argmax(np.abs(ym [idx - 5 :] - 0.5))] -
xx [np.argmax(np.abs(ym [0 : idx] - 0.5))]

grafico xx, yy

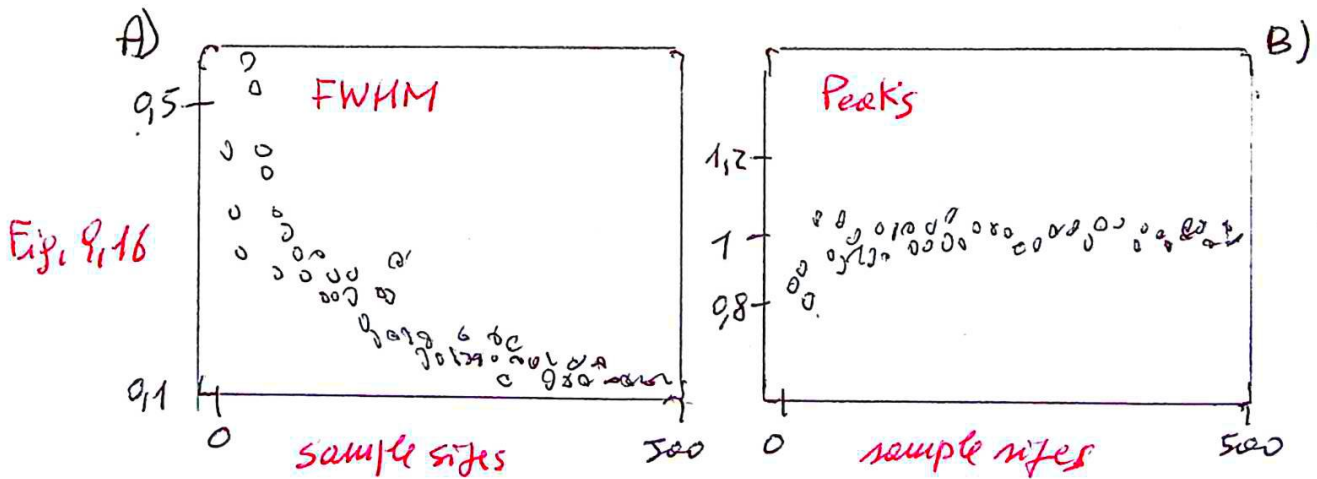




- più spesse le linee → maggiore il n° di campioni (da ϕ a 500)
- FWHMs empirici, e la distribuzione dei picchi possono essere visualizzati con grafici raster.

```

axs[0].plot(samplersizes, fwhms, 'ks', _____)
axs[1].plot(_____, peakvals, _____)
    
```



7) confronto di SEM analitico con due stime empiriche

$$SEM = \frac{\sigma}{\sqrt{N}} \approx \frac{s}{\sqrt{N}}$$

calcoliamo il vero SEM x 25 dimensioni del campione, spaziate
log tra $N=10 < N=10\%$ della popolazione.

Successivamente stimiamo SEM in due modi

(1) std di campione casuale / $\sqrt{\text{una dimensione}}$

(2) std di più medie campionarie empiriche

↓ 50 campioni casuali x ogni nze di campione, calcolare SEM
del campione — calcolare medie del campione

Graficare

$$N_{pop} = 10^8$$

$$\text{population} = \text{np.random.randn}(N_{pop}) * x_2$$

$$\text{samplesizes} = \text{np.linspace}(\text{np.log}_{10}(10), \text{np.log}_{10}(N_{pop}/10), 25, \text{dtype=int})$$

$$\text{numExps} = 50$$

$$\text{theory} = \text{np.std}(\text{population}) / \text{np.sqrt}(\text{samplesizes})$$

$$\text{standard error} = \text{np.zeros}((\text{numExps}, \text{len}(\text{samplesizes})))$$

$$\text{samplemeans} = \text{np.zeros}((\text{numExps}, \text{len}(\text{samplesizes})))$$

for expi in range(numExps):

for idx, ssize in enumerate(samplesizes):

generate a random sample

$$\text{rsample} = \text{np.random.choice}(\text{population}, \text{size}=\text{ssize})$$

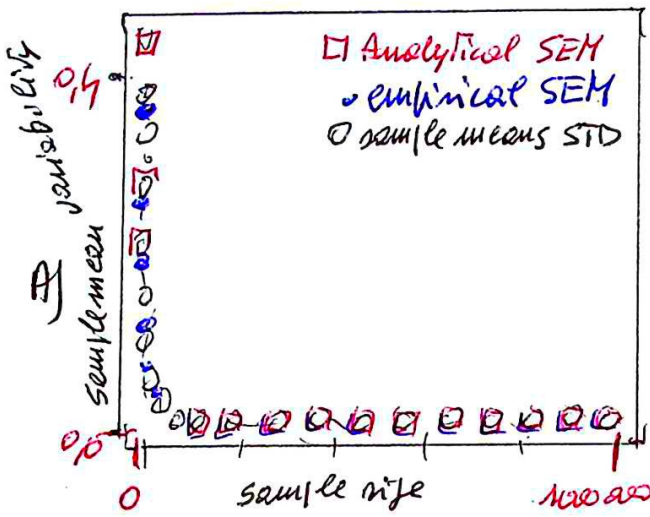
comput std error and sample mean

$$\text{standard error}[\text{expi}, \text{idx}] = \text{np.std}(\text{rsample}, \text{ddof}=1) / \text{np.sqrt}(\text{ssize})$$

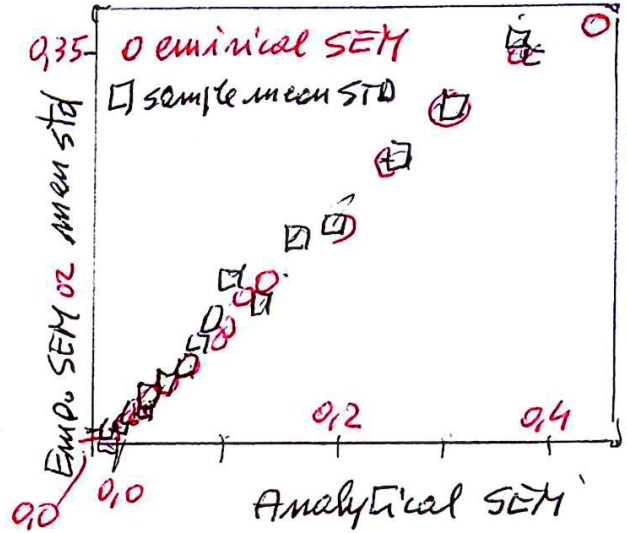
$$\text{samplemeans}[\text{expi}, \text{idx}] = \text{np.mean}(\text{rsample})$$

Plot $\text{axs}[\phi].\text{plot}(\text{samplesizes}, \text{theory}, \text{label}=\text{Analytical SEM})$
 $\text{axs}[\phi].\text{plot}(\text{samplesizes}, \text{np.std}(\text{samplemeans}, \text{axis}=\phi, \text{ddof}=1), \text{label}=\text{Sample means STD})$

$\text{axis}[1].\text{plot}(\text{theory}, \text{mp.mean}(\text{standard error}, \text{axis}=\phi) - \text{empirical SEM}) - 62$
 $(\text{mp.std}(\text{sample means}), \text{mp.mean}(\text{standard error}))$
 sample mean std



Estimates by sample size



Consistency of estimates

Fig. 9.17

B)

- la variabilità delle medie campionarie diminuisce drasticamente all'aumentare del campione Fig. 9.17A
- le stime teoriche e empiriche di SEM corrispondono bene Fig. 9.17B
 I valori di SEM grandi \rightarrow campioni piccoli

P. 339

— FINE CAPITOLO 9 —

10 - Hypothesis Testing

10.1 - Hypotheses

How to specify a hypothesis

- Spesso si parte da una domanda \rightarrow si pone una ipotesi
- Il caffè aiuta a studiare? \rightarrow 1-2 tappe/die in un momento
 - alcune ipotesi possono essere simulate o pensate matematicamente
 - A noi interessano ipotesi che possono essere valutate con i dati

(Why) do we need hypotheses?

- 1) le ipotesi migliorano la progettazione degli esperimenti
 - 2) da idee vaghe pensiamo a idee concrete
 - 3) aiutano a sviluppare nuove teorie
 - 4) consentono una valutazione quantitativa
 - 5) portano alla comprensione dei meccanismi sottostanti
- ! forse la voce + importante

Strong and weak hypotheses

L'ipotesi deve essere chiara, specifica, falsificabile

- il farmaco X aiuta i pazienti? \leftarrow è solo una domanda
- _____ ha un effetto \rightarrow ipotesi debole
- _____ riduce i sintomi della malattia Y in base alla fisiologia
- L'ipotesi forte
- L'ipotesi debole non fornisce dettagli

- non è la debole se non è troppo poco \rightarrow una nuova legge? solo dopo qualche caso di applicazione possiamo formulare ipotesi forte

Esempi veri, da valutare

- 1) \exists altri universi con leggi fisiche - non falsificabile
- 2) indomare bianchenia intima viola migliora l'umore - F
- 3) le piante crescono in modo diverso nell'acqua zuccherata - W
- 4) i libri di Max Cohen sono fantastici - opinione
- 5) lavarsi le mani x 20 secondi riduce la diffusione di malattie infettive - F
- 6) una mela al giorno toglie il medico di turno - proverbio, domanda
- 7) le persone sono + creative dopo aver guardato una "stand-up comedy"?

10.2 - IVs, DVs, models, and other stats things

Dependent variables (DVs)

- DV include l'evento di una malattia in uno studio medico
- il peso perso dopo un regime dietetico
- il n° di click in uno studio di marketing
- si possono avere anche + DV

Independent variables (IVs)

- e' una variabile che si usa x spiegare le variazioni nelle DV
- alcune possono essere manipolate sperimentalmente
- altre possono essere misurate
- variabili regressive, esplicative, predittive

Ricerca	DV	IVs
canco	mentalita'	eta', esercizio, tipo canco
peso perso	Delta peso	tipo dieta, esercizio, rudi sonno
E-commerce	profit	ovetime, budget <u>Twitter</u> <u>google</u> <u>Facebook</u>
psicologia	umore rispondito	ore sui social, genere

- alcune variabili sono numeriche, altre categoriali

Residuals

- IV raramente spiega completamente DV

↳ la differenza tra ciò che si è previsto e il risultato effettivo = residuo

↳ chiamato anche errore, deviazione, imprecisione

Model

è un framework x interpretare i dati - I modelli vengono in profondità, specificità, capacità esplicativa, dettaglio matematico -

- in questo libro Model = formula matematica

Es. Altezza^y persone dipende da altezza madre e padre, da alimentazione infantile (intercetto di ϕ e ω) →

$$y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

con analisi di regressione calcolo β_i , la deviazione residua la chiamano con ϵ

- vedremo come trasformare un'ipotesi → modello

Test statistic

È una statistica descrittiva di un set di dati, oppure è un'ipotesi -

statistica inferenziale → determinazione probabile che la statistica di test sia stata generata a caso (casuale...), invece che essere un effetto della popolazione -

statistica di test → coeff. di correlazione, t-test, F-value, chi-square value

Null and alternative hypotheses

Le statistiche inferenziali includono l'assunzione ("fit") di un modello ai dati -

In molti casi un risultato importante (es. p-value) riflette la compatibilità tra due modelli derivanti da due ipotesi - Le due ipotesi sono definite **null** e **alternative** -

Una "ipotesi nulla" è spesso abbreviata in H_0 oppure $H\phi$ 65

- la **null** dice che nulla di interessante accade (DV correlato da IV).
 - ↳ l'equazione matematica di H_0 ha alcuni dei β_i nulli.
 - la **alternative hypothesis** è quella con cui specifichiamo la ricerca.
 - È una alternativa alla **null**.
 - ↳ la si dice anche **effect hypothesis**, nel senso che specifica che c'è un effetto
 - si usano anche H_1, H_2 quando ci sono + **alternative**
- H_A = le persone acquistano + widgets dopo aver visto la pubblicità X rispetto a quelle Y
- H_0 = il tipo di pubblicità non ha effetto sui widgets

10.3 - Can you prove a hypothesis?

- Confutare una ipotesi può essere semplice. Basta un caso contrario.
- Molte ipotesi in biologia non possono essere confutate in pratica - Questo xche' la comprensione della realtà è limitata -
 - accettiamo queste imperfezioni, dati non coerenti; possono ∇ essere utili a il progresso della scienza.
 - Raccolgo i dati e noto che H_A mi sembra meglio ai dati di H_0 , questo non implica che H_A sia dimostrata - Potrebbero \exists modelli migliori \rightarrow con una certa probab. H_A è migliore di H_0 .
- Es. supponiamo di avere $x_1=1, x_2=2$ & $y=3$
quale modello ti fa? certamente $y=x_1+x_2$, ma \exists un altro modello? Sì, es. $-5x_1+3x_2$
- Quale modello scegliere? x esempio quello + semplice
(Rasoio di Occam)
- Non riusciamo a dire se abbiamo scelto il "modello vero", possiamo solo valutare H_A vs. H_0 .

10.4 - Sample distributions under H_0 and H_A

La migliore dispersione SEM può fornire supporto statistico a favore di una ipotesi.

Es. La caffeina migliora l'umore

H_A : 1 Tazza di caffè aumenta lo stato emotivo auto-riferito dopo un'ora dal consumo

H_0 : la Tazza di caffè non ha impatto

- per testare ipotesi conduciamo un esperimento:

200 partecipanti (100 partecipanti in doppio cieco), entrambi dopo un'ora riportano lo stato emotivo. migliore

- lo stato emotivo è un costrutto complesso (lo valutiamo da 1-10)

modello
$$\delta = \beta c + \epsilon \quad (10,1)$$

$c = \phi$ se decaffeinato
 $= 1$ " con caffeina

$\beta =$ impatto della caffeina

$\epsilon =$ cambiamenti di umore dovuti ad altre cause

$\bar{\delta} =$ variazione auto-dichiarata / pedice C, D / caffeina / deca

H_A : $\bar{\delta}_C > \bar{\delta}_D$

questi i termini dell'esperimento

H_0 : $\bar{\delta}_C = \bar{\delta}_D$

• abbiamo bisogno di un gruppo di controllo

└ può essere che umore ↑ solo aver bevuto bevanda calda
 senza questo gruppo $\bar{\delta}_C = \phi$ diventa un'ipotesi debole/minima

• "cambin l'umore" potrebbe anche essere scritto come $\bar{\delta}_C \neq \bar{\delta}_D$
 (differenza tra ipotesi unilaterale e bilaterale)

• la variazione media è riferita a 100 campioni - Non è necessario che tutti abbiano umore ↑ - Basta che la media ↑

• scriviamo le due ipotesi:

$$\bar{\delta}_C - \bar{\delta}_D > \phi \quad \text{---} \quad \bar{\delta}_C - \bar{\delta}_D = \phi$$

- in relazione al pto precedente, in alcune analisi come la regressione, le ipotesi sono + direttamente collegate al modello -
 Ricordiamo la 10,1: $S = \beta C + E$

$H_0: \beta = \phi$ $H_A: \beta > \phi$ diciamo $\Delta = \bar{S}_0 - \bar{S}_0 \rightarrow$

$H_0: \Delta = \phi$ $H_A: \Delta > \phi$

- Come testiamo l'ipotesi?

di Δ ?

Supponiamo che la verità sia $H_0 \rightarrow$ quale il valore atteso

Δ in pratica non sarà esattamente ϕ

- ripetiamo + volte l'esperimento con 200 partecipanti:

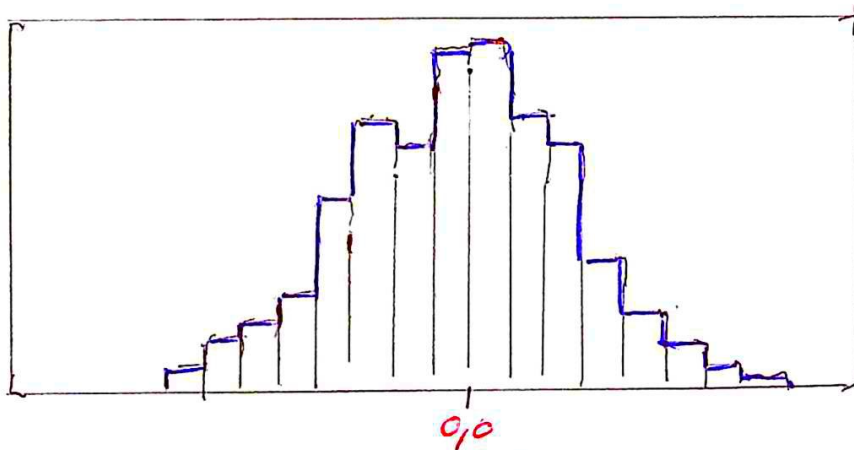
L vogliamo la distribuzione di Δ

LLN ci dice che questo valore convergerà a ϕ , se l'ipotesi è ^{corretta}

Visualizziamo una ipotetica distribuzione di Δ .

Esperimento ripetuto 1000 volte.

Fig. 10,2



Questo grafico dice che è valida H_0

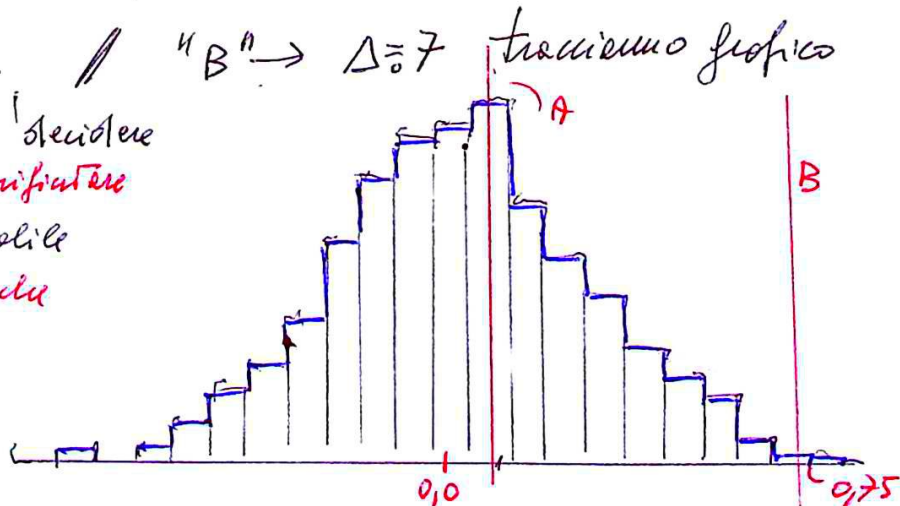
- Supponiamo un altro risultato

"A" $\Rightarrow \Delta = 0,1$ // "B" $\rightarrow \Delta = 0,7$ tracciamo grafico

A da' poca info x decidere
 L con A non sono sufficiente

B è con' poco probabile
 L può decidere da

H_0 è vera



• vediamo di tradurre in forma quantitativa

67

10.5 - Where do H_0 distributions come from?

Come si crea un'ipotesi nulla? Necessitiamo di una distrib. di H_0 per valutare la "p" che non sia dovuta al caso.

- H_0 distribuzione qualitativa \rightarrow distribuzione ottenibile sempre dati
- H_0 empirica \rightarrow randomizziamo i dati x o creiamo un set di dati artificiali (mescolati, permutati, sottocampi) che potrebbero essere sotto H_0 . \rightarrow creiamo dati falsi con le stesse caratteristiche di quelli reali, ma in cui non è vero $H_0 = \phi$ e vero.
- H_0 empirica richiede di avere già dati. \rightarrow Cap. 16

10.6 - P-values: Definition and misinterpretations

p = probabilità di ottenere per vero H_0

\hookrightarrow p piccolo = poco probabile

p = probabilità di ottenere una statistica del test estrema, o estrema della statistica di test osservata

\hookrightarrow supponendo che sia vero H_0

NB p piccolo non implica che H_A sia vero

possiamo calcolare la probabilità della statistica osservata $\&$ H_A

\hookrightarrow se p piccolo \rightarrow dati più coerenti con H_A

\hookrightarrow non significa che H_A sia vero

• i valori di "p" sono sfumati \rightarrow vedi più avanti!

P-values and statistical significance

significance = rifiutiamo H_0 sulla base di "p".

i valori di "p" sono continui da $\phi \rightarrow 1$

• bisogna prendere una decisione \rightarrow dobbiamo binarizzare

es. scegliendo come soglia 0.5

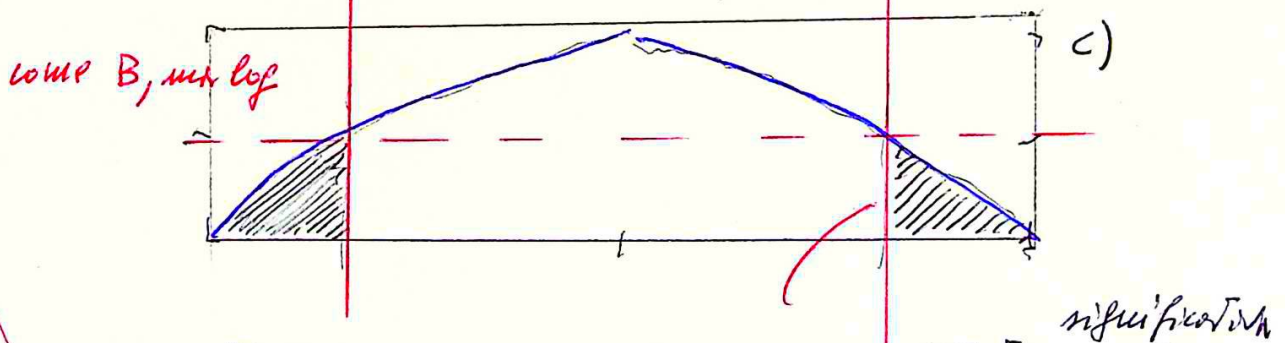
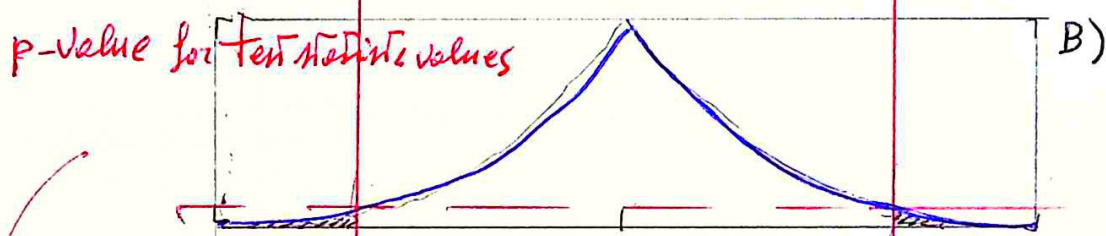
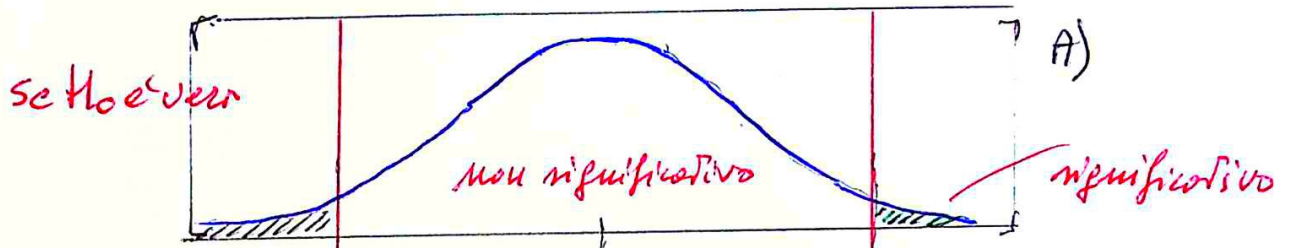
L'insieme anche altre soglie 0.01 / 0.001

- ho statistica di test se $p < \text{soglia} \rightarrow$ statistica significativa

- di solito la soglia si indica con α

es. $\alpha = 0.05$ o $\alpha = 5\%$

• dividiamo quindi in "regioni significative"



una statistica con p associato < 0.05 è statisticamente

Nei grafici "test a due code"

qui si vede meglio cosa

il grafico B) può interpretare la curva, ma viene dal fatto che
col $z < \phi$ e invertito per $z > \phi$ — curva

se $z < \phi$ allora p corrisponde a cdf

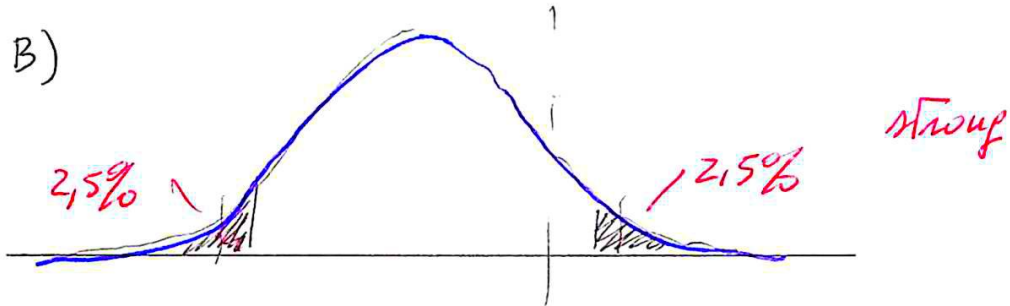
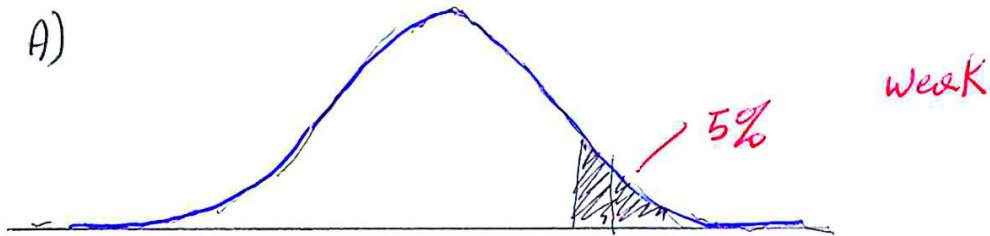
se $z > \phi$ " " " " $1 - \text{cdf}$

P-values and distribution tails

x specificare una ipotesi: effetto positivo, negativo, $\neq \phi$

strong una coda $\left\{ \begin{array}{l} \Delta < \phi \\ \Delta > \phi \end{array} \right.$ due code $\Delta \neq \phi$ weak

le due code insieme x una ipotesi debole



Area x rigettare e' la stessa - In pratica si usa sempre la bilaterale.

Talora anche una coda - Es. **FTEST** nelle ANOVA e nelle regressioni.

Oppure con distribuzioni di potenza

p si origina dagli stessi dati che definiscono H_0

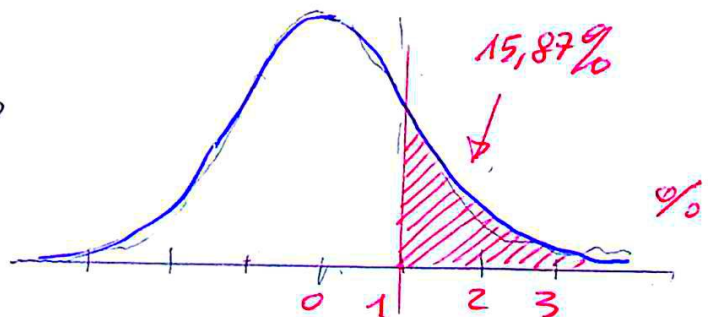
p-values espletici mediante formula

Esempio

punteggi z, come moltiplicare p? z sono Geom , media ϕ
 quale la probabilita' di $z > 1$? \rightarrow dobbiamo calcolare l'area
 a destra di $z=1$ nella pdf corrispondente a Geom .
 Come si vede da figura:

$z > 1 \rightarrow p = 0,158$

in unita' di σ



empirical p-values

Se usiamo statistiche basate su permutazioni x avere H_0 empirico

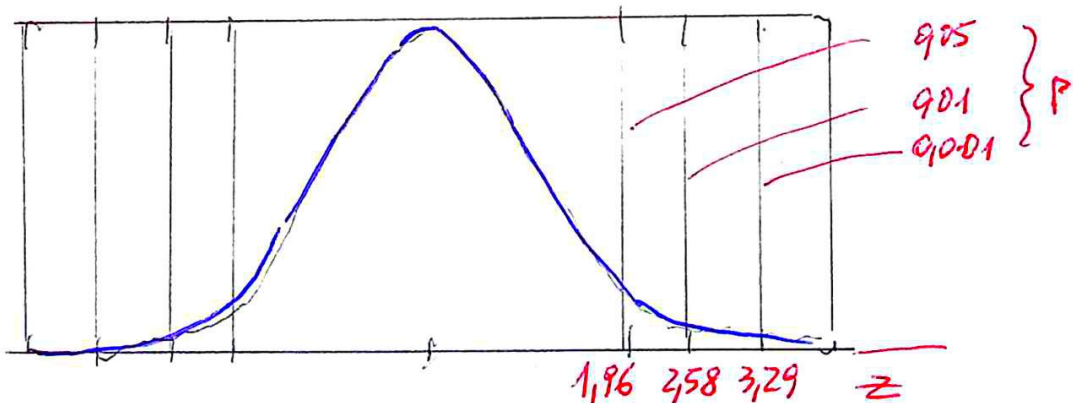
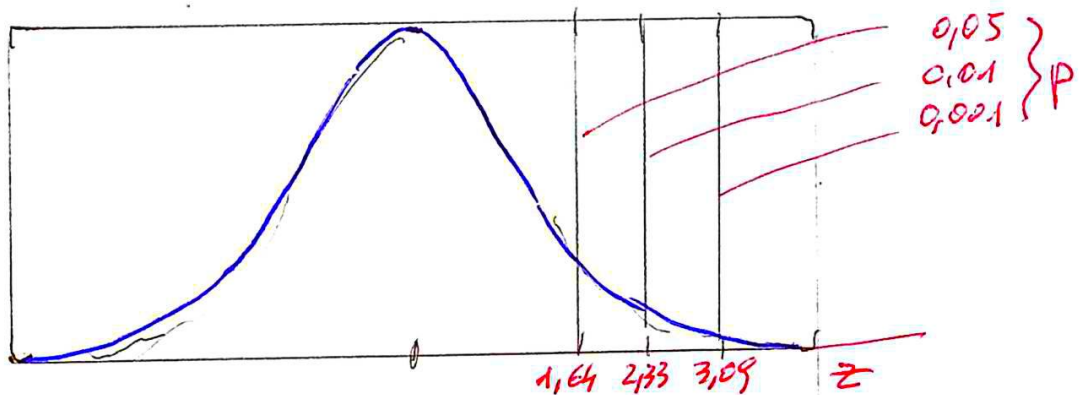
→ p-value = distanza normale rispetto al centro della distrib. di H_0 → Cap. 16

Precision of p-values - si può calcolare con precisione arbitraria -

In pratica i numeri sono troncati - Per valori di p piccoli si usa la notazione scientifica -

P-z combinations to memorize

alcune di queste combinazioni meritano la memorizzazione



Misinterpretations riferite a $p=0.02$ → se $p=0.02$:

- 1) l'effetto è presente x il 98% della popolazione topologica
- 2) 98% di probabilità che la mia statistica di test = parametro della H_0
- 3) $p < soglia$ → l'effetto è reale

4) $p < \alpha$ → una variabile causa l'altro

5) se $p = 0.08$ → $p > \alpha$ → H_0 è vera

Risposte e correzioni

- 1) c'è una $p = 2\%$ che non ci sia alcun effetto/ e la mia statistica era dovuta a variabilità, rumore, piccola dimensione del campione, distorsione sistematica.
- 2) c'è una $p = 2\%$ che la mia statistica sia dovuta a rumore
- 3) $p < \alpha$ → è improbabile che l'effetto sarebbe stato osservato
- 4) $p < \alpha$ → variabili correlate
- 5) $p > \alpha$ → l'effetto avrebbe potuto essere osservato se H_0 vera
 è improbabile che mio H_A specificato sia vero - vedi ipotesi alternative che siano migliori di H_0 .

Discutibile sui 5 punti!

- 1) $p \approx$ effetto medio all'interno del gruppo in relazione alla variabilità!
 non dice nulla sulla presenza o meno dell'effetto x ogni individuo -
 Si può determinare se un effetto è significativo per ciascun individuo,
 ma non è questo che indica p .
- 2) p non dice nulla sulla relazione tra campione e le caratteristiche della popolazione - È semplicemente la probabilità che la statistica di test avrebbe potuto essere osservata se $H_0 = vera$ -
 Si possono usare gli intervalli di confidenza x quantificare la relazione tra le caratteristiche del campione e i parametri della popolazione.
- 3) $p < \alpha$ non dimostra che l'effetto è reale, dice che l'effetto osservato è improbabile data $H_0 = vera$.

L in senso colloquiale lo diciamo, ma non è appropriato)

h) p che solo non stabilisce la causalità - Le interazioni causali possono essere determinate tramite manipolazioni sperimentali o altri tipi speciali di causalità -

p significativo \rightarrow relazione tra le variabili, non causalità -

$p > soglia$ \rightarrow non dimostra H_0 \rightarrow dice che l'effetto è probabile se $H_0 = vera$ -

Problems with p-values

p non dimostra una ipotesi, non dimostra causalità, non dice la probabilità di popolazione che prova l'effetto -

$p < 0,001$ oppure $p > 0,3$ è relativamente facile da interpretare, perché siamo su valori estremi \rightarrow elevata confidenza su presenza/assenza dell'effetto -

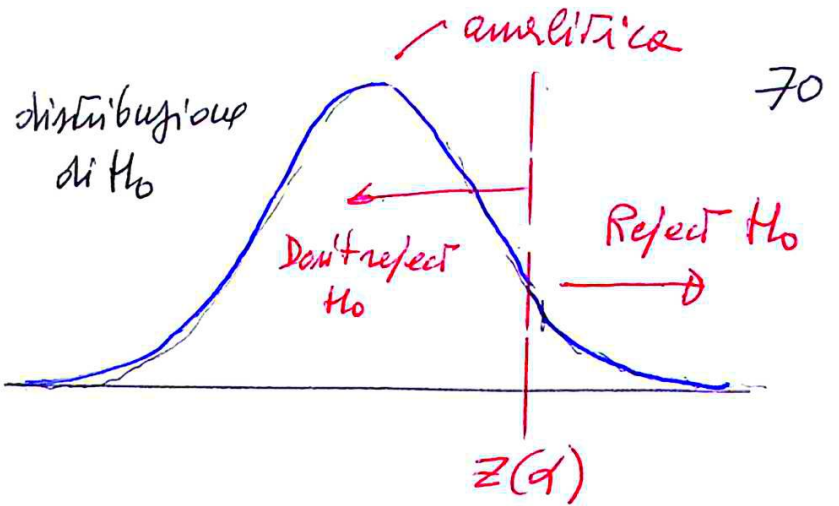
p tra $0,05$ e $0,3$ è + difficile da interpretare, anche se siamo vicini alle soglie di $0,05$, e un po' di rumore potrebbe cambiare tutto -

10.7 - P-values and significance categorization

Per capire ciò che ci circonda raccogliamo dati, e valutiamo se una ipotesi è in linea con essi - Il risultato è statistico \rightarrow
 \rightarrow ci affidiamo a p a decidere in modo informato -
Come si forma lo spazio bidimensionale della decisione?

- lo vediamo con una figura:

Vediamo come Tabella
 le decisioni statistiche
 corrette



Realta' su H_0

Vera	Falsa
Vero $< \phi$ ($1-\alpha$)	(β)
(α)	Vero $> \phi$ ($1-\beta$)

falso negativo

Non rigetto H_0

Rigetto H_0

statisticamente
 significativo

- puoi rifiutare H_0 quando e' davvero falsa
 - non rifiuti H_0 " " " vera
 - la decisione implica incertezza
- falso positivo

10,8 - Type I and Type II errors

- un giudice:

- dice e' un innocente sei colpevole
- dice e' un criminale sei innocente

Where do errors come from?

- non possiamo provare che H_A vera
- rifiutiamo H_0 se il nostro teste e' $< 5\%$

la decisione è statistica \rightarrow possono esserci false scelte

Type I errors rifiuto H_0 quando è vero

la probabilità di questo è la "p-value significance threshold"

L'errore di α \rightarrow falso positivo

Type II errors non rifiuto H_0 quando è falso

la probabilità di questo è correlata alla potenza statistica del test

$(1-\beta)$ \rightarrow Cap. 17

L'errore di β \rightarrow miglior con campioni grandi

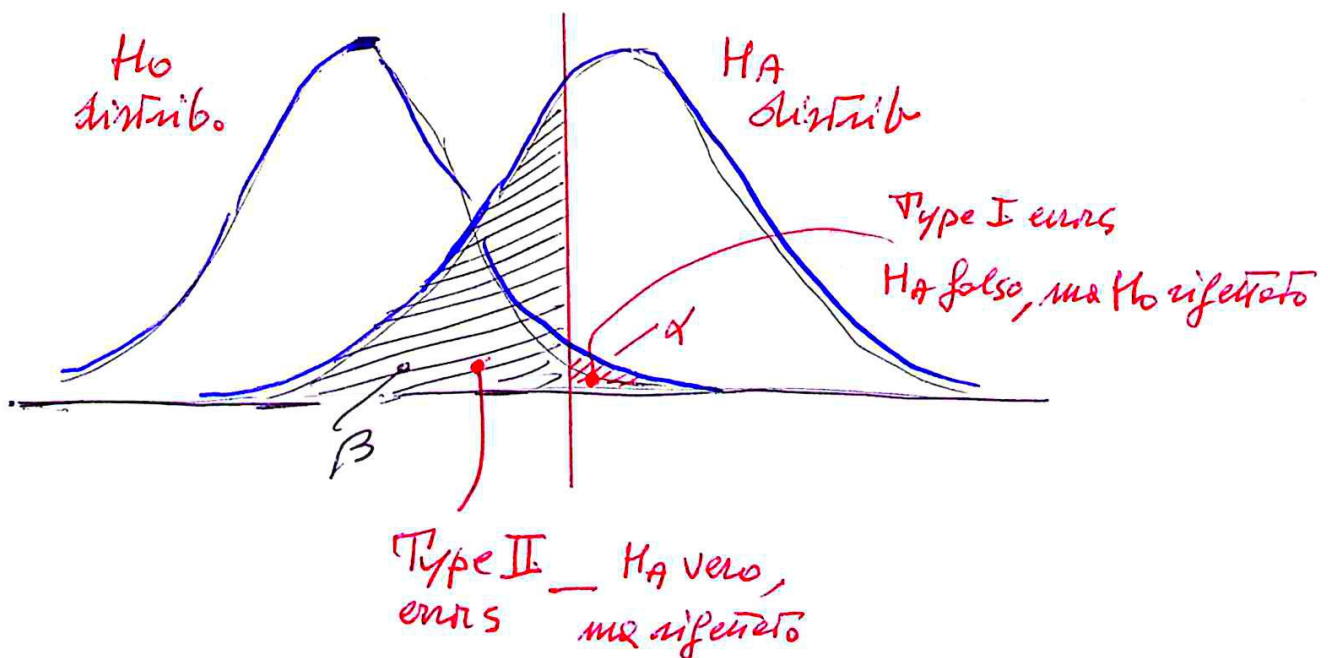
The balance of Type I and Type II errors

I due tipi di errori sono correlati. Se diminuisco l'uno, l'altro aumenta

Esperimento mentale

H_A vero al 100% — poniamo ripetere molti esperimenti \rightarrow
avremo un'intera distribuzione di statistiche di test in cui $H_A = \text{vero}$

Potremo visualizzare le distribuzioni sotto H_0 e H_A



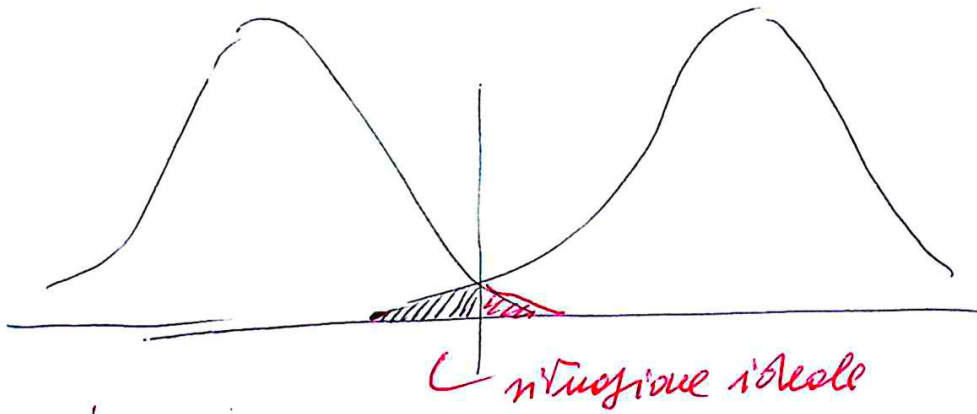


Fig 10.12 C

L'ora si diminuisce la variabilità
 o si aumenta la dimensione dell'effetto } → Cap. 11

- però non conosciamo la realtà → non conosciamo la distrib. di H_A - Abbiamo un solo valore di H_A , e ci domanderemo se quel valore della statistica di test empirica ha un valore insolitamente alto - tratto da H_0 - o se ha un valore tratto da H_A .

10.9 - Various interpretations of "significant"

Se p abbastanza piccolo → statisticamente significativo

L non è l'unica interpretazione

teoricamente

Significante: statisticamente, teoricamente, clinicamente, socialmente

Statistical significance = prob. che una statistica di test empirica venga osservata con $H_0 = \text{vero}$

Theoretical significance - ci stiamo riferendo a una teoria scientifica

L non c'entra nulla con quella statistica

Clinical significance - risultato rilevante x la pratica clinica

L non correlato con statistica

Practical, societal, educational, etc.

La statistica vive al confine tra scienze dure e soft →
come utilizza "p" calcolato e compito che richiede competenze
specifiche

10.10 - Multiple comparisons

= Testare più ipotesi x uno stesso set di dati

- con Test appropriati Type I ↑

- non si deve pensare ai singoli Test, ma al loro insieme

↳ + membri → + probabile che qualcuno sia etichettato erroneamente

Esempio - soglia $p < 0.05$ - facciamo 3 Test → (FWE)

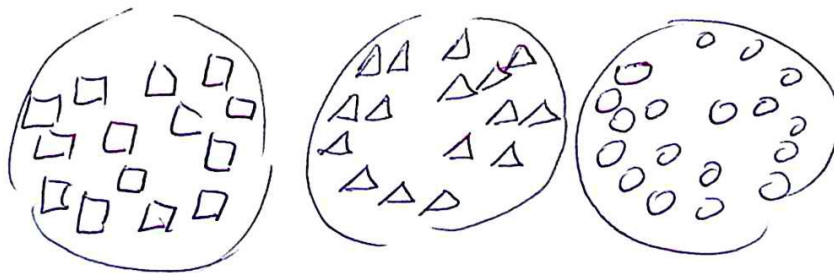
→ probabilità combinata di Type I = $0.15 = 15\%$ →

ogni singolo $0.05 \rightarrow \sum_1^3 = 0.15 =$ "family-wise error" rate
in formula $FWE = m\alpha$ anche non esattamente accurato

Dipendenza tra i dati può diminuire FWE

altra formula $1 - (1 - \alpha)^m$

Una dipendenza può essere vista confrontando
↳ qui i teste sono correlati x che entrambi contengono A



H_1 □ vs. △

H_2 ○ vs. △

H_3 □ vs. ○

Solutions to the multiple comparison problem

72

- bisogna regolare la soglia di significatività / corrisponde a FWE
L e non alla soglia individuale
3 metodi x fare questo → Ch. 14 in ANOVA

Bonferroni = impostare FWE = soglia desiderata = es. 0,05

e la soglia dei singoli test = $d/m = 0,01667$

- 1) questa correzione può essere troppo riposta (presuppone la indipendenza dei test)
 - 2) si basa solo su m e non sulle caratteristiche dei dati
 - 3) aumenta la probabilità di Type II → impedire la scoperta
fittoria di ~~dati~~ risultati veri
- Conclusione: buona se m non troppo grande, se i test sono indipendenti

False discovery rate (FDR)

utile se molti valori di p da test correlati (serie temporali), p di dati geografici, genetica - FDR fornisce soglia critica q che viene calcolata dalla distribuzione di p .

$$q = \frac{\text{n° previsto di false ndente}}{\text{vere}} \rightarrow \text{Esercizio 11}$$

$q = f(p)$ e' un limite → un risultato potrebbe essere etichettato come significativo o non significativo solo in funzione di p

Cluster correction = buona x dati che hanno struttura di correlazione significativa

Test significativo se è all'interno di un insieme di test contigui

L Fig. 10,14 = soglia - colore corrisponde a valore di p

- la soglia minima del cluster può essere determinata in base a considerazioni a priori (es. un cluster deve contenere almeno 100 ms di punti contigui, oppure almeno 6 pixel contigui significativi)
- la soglia minima in base a "permutation testing" x decidere una dimensione empirica del cluster sotto H_0

Questi metodi controllano Type I, ma Type II ↑

L

10.11 - Degrees of freedom

$$\begin{aligned} &\rightarrow \text{elem}3 = 10 \\ \text{elem}1 &= 2 \\ - 2 &= 3 \end{aligned}$$

Se dico: 3 elementi - media = 5

↳ quanti gradi di libertà?

$df = n^{\circ}$ gradi di libertà opure ν (μ)

una statistica descrittiva viene calcolata dai dati - potrebbe essere vincoli!

df è usato come parametro nelle distribuzioni t →

df utilizzati x inferenza sulla popolazione sottostante

df = utile x controllare errori, correggere interpretazioni dei risultati nelle regressioni e nelle ANOVA

È unica formula x calcolare df -

Possiamo pensare al df di una equazione come il n° delle osservazioni meno il n° di parametri in formula

$df = N - K$ ————— e' una semplificazione

N : n° di elementi o il n° di gruppi o le condizioni in esperimento

Esempio df associato a t-test → $\begin{cases} N-1 \\ N-2 \end{cases}$ dipende se c'è

un solo gruppo oppure due — $K = n^{\circ}$ di IVs — vedi dopo

I gradi di libertà possono essere anche non interi -

W, Z - Exercises

1) nella Fig. 10.5 il codice assume $p < 0,05$ - Modificare in modo da utilizzare un arbitrario p .

In seguito modificare codice

2) valutare il tono di falsi allarmi - \rightarrow uso t-test x determinare se la media di un campione devia significativamente da un valore prefissato - la media con attesa = ϕ , ma un campione ^{casuale} la può avere un po' diversa in modalità "statisticamente significativa" - Questo sarebbe un Type I \rightarrow la vera media = ϕ , la media del campione $\neq \phi$ con $p < 0,05$ -

$N=20$

$X = np.random.randn(N)$

$p = stats.ttest_1samp(X, \phi)[1]$

p-value from t-test = 0,0185

media = np.mean(X)

\rightarrow media = 0,14

simulo i dati e t-test

alpha_range = np.linspace(0,001, 9, 15)

$M=100$

type1_errors = np.zeros(M, len(alpha_range))

means_and_pvals = np.zeros(M * len(alpha_range), 2)

loop sopra alpha

for q_i, α in enumerate(alpha_range):

for exp_i in range(M):

$X = np.random.randn(20)$

$p = stats.ttest_1samp(X, \phi)[1]$

loop su esperimento

memorizzo

type1_errors[exp_i, q_i] = $p < \alpha$

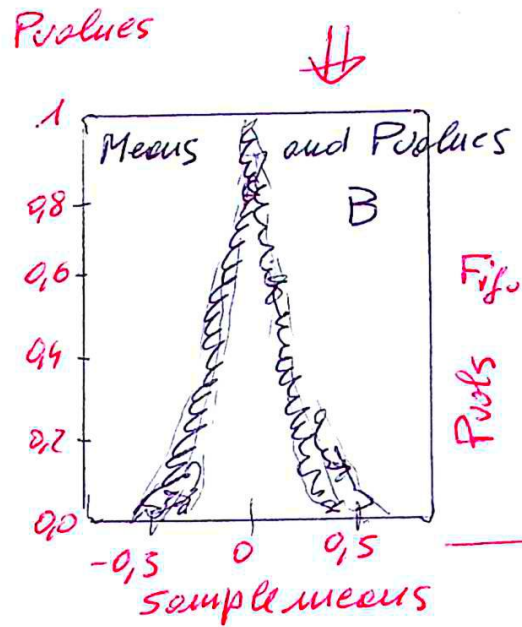
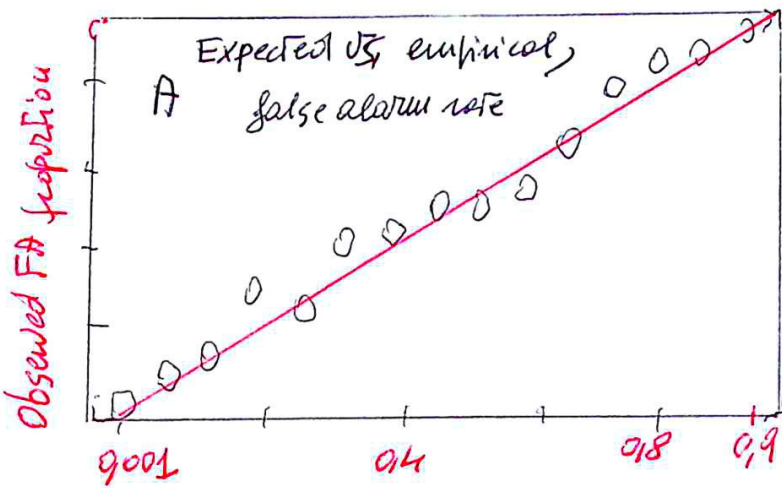
ritraslo media e p

means_and_pvals[$exp_i * \text{len}(\alpha_range) + q_i, :$] = [np.mean(X), p]

grafico

\Rightarrow

$\alpha \times 5 [1] \cdot \text{plot}(\text{meansAndPvals}[\phi], \text{meansAndPvals}[\phi, 1], 'k')$ 74



↑↑ Predicted FA proportion

$\alpha \times 5 [\phi] \cdot \text{plot}(\text{alpha Range}, \text{np.mean}(\text{type I errors}, \text{axis}=\phi), 'ks')$
 $\alpha \times 5 [\phi] \cdot \text{plot}(\text{---}, \text{alpha Range}, \text{---})$

Ho ripetuto t-test 100 volte, registriamo il n° di volte che p è (150%)
 sotto soglia (cioè quando il test è significativo)

Con una soglia $\alpha = 0.05 \rightarrow$ si tratta di circa 5 campioni

↳ risultato con $p < 0.05 = \text{Type I}$

Metto un secondo ciclo con $\alpha = 0.001 \div 0.9 \rightarrow$ rappresento
 graficamente il tasso di falsi allarmi empirici in funzione di α
 FIG. 10.15A.

Nella FIG. 10.15B vediamo che - più la media è vicina a $\phi \rightarrow$

↳ migliore è il valore di p - Non è una relazione perfetta e che
 t-test dipende anche da σ - In ogni caso è una relazione ben
 consolidata - Qui la media rappresenta il numeratore della statistica t
 e il denominatore è \approx cost in tutte le simulazioni. \rightarrow Cap. 11

3) Le distribuzioni H_0 di questo capitolo non sono pdf Gauss -
 Gauss sono usate spesso x visualizzare le distrib. H_0 x che 'conosci'
 in stats, e non richiedono parametri.

Ma molte altre distrib. x H_0 sono 'importanti' -

Esploriamo la distribuzione $\therefore H_0$ di un t-values -

La distribuzione t richiede un parametro: gradi di liberta'

La forma della t-distrib. dipende da df

- Prendiamo una "famiglia" di distribuzioni supponendo $H_0 = \text{vero}$.

`tvals = np.linspace(-4, 4, 1001)`

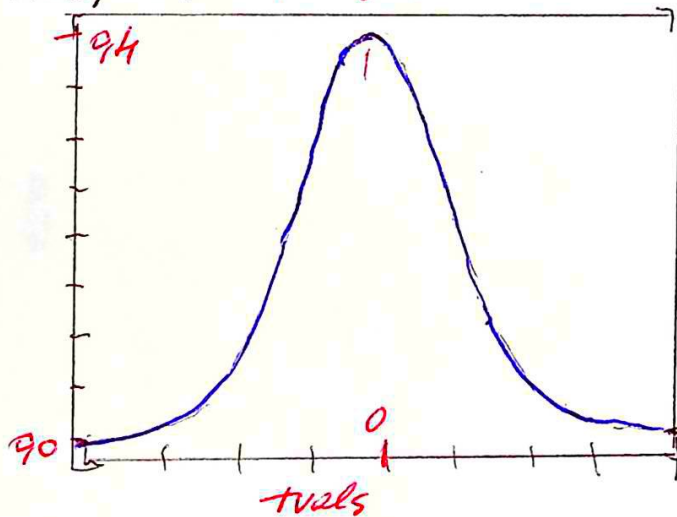
calcolo pdf

`tpdf = stats.t.pdf(tvals, 20)`

`plt.plot(tvals, tpdf)`

= df

lor' grafico



One vogliamo usare df
 come parametro

`dfs = np.arange(1, 41)`

for df in dfs:

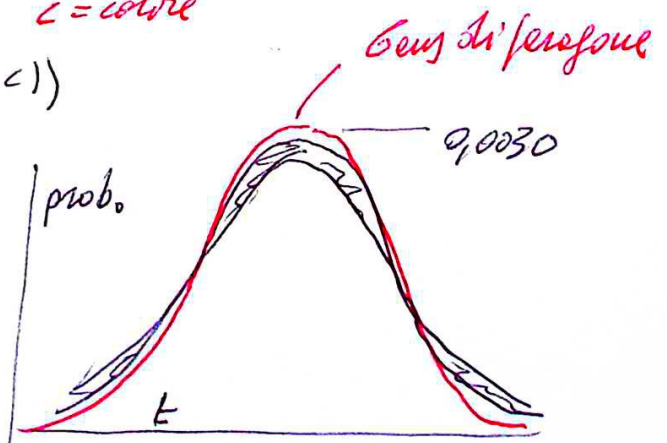
calcolo pdf

`tpdf = stats.t.pdf(tvals, df) * np.diff(tvals[1:2])`

nono colore x ogni curva, stampa c = colore

`plt.plot(tvals, tpdf, color = (c, c, c))`

df ha un aspetto blando - E ha sempre
 code + spene di Gauss



4) Due modi di usare la correzione FDR

75

1) accettare i valori "new" considerando significativo $\forall p / n_{\text{new}} < \alpha$

2) derivare una soglia di significatività e considerare significativo $\forall p$ empirico $< \alpha$ empirico

Es. FDR derivato da α empirico, corrispondente a

$$p < .05 \text{ potrebbe essere } q < .008$$

$\forall p\text{-value} < .008$ sarà considerato significativo

$\forall \text{ --- } > .008$ " " non-significativo

importo *for corrections*

from statsmodels.stats.multitest import fdr corrections

valore di soglia corretto

$$p\text{Thresh} = .05$$

pvalues

$$k = 40$$

$$pvals = np.random.uniform(low=.001, high=.3, size=k) ** 2$$

ordine p-values

$$pvals\text{Sort} = np.sort(pvals)$$

linear interpolated distribution

$$pvals\text{Interp} = np.arange(1, k+1) / k$$

adjusted p-values

$$pvals_adjusted = pvals\text{Sort} / pvals\text{Interp}$$

questa funzione dà un tuple: ϕ rigettate, 1 adjusted p-values

$$qq = fdr\text{corrections}(pvals\text{Sort}, p\text{Thresh}) [1]$$

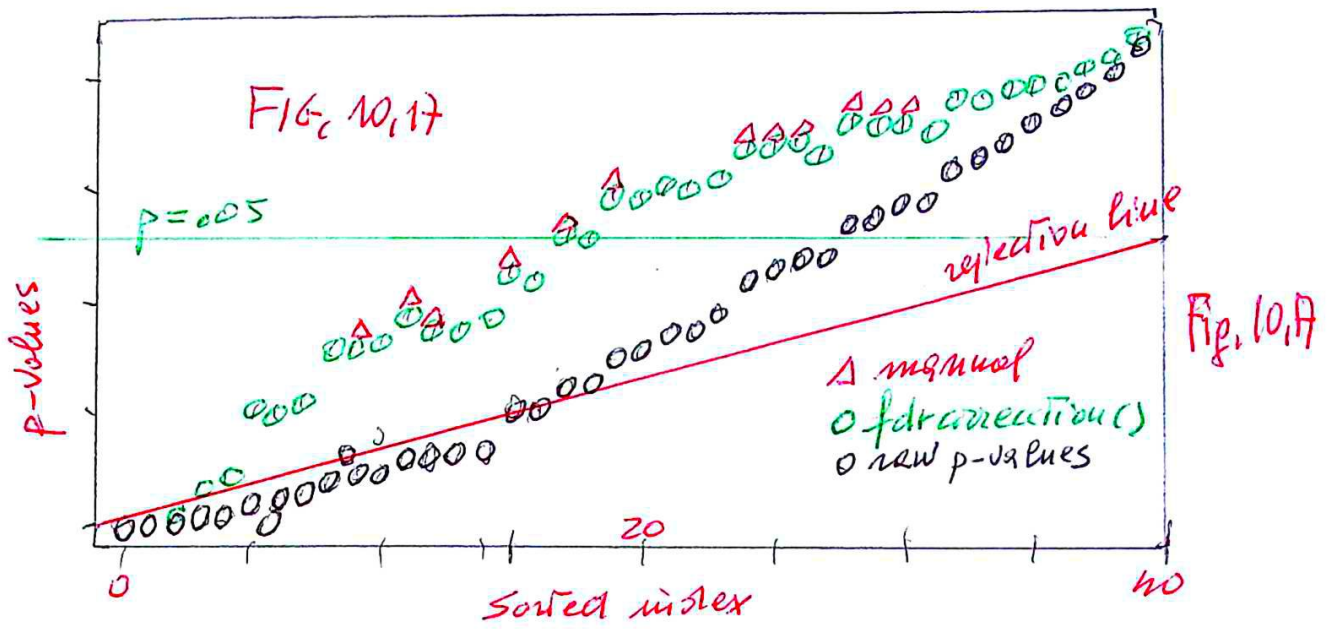
grafico

$$\text{plt.plot}(qq, 'ko', \text{---}) \text{ fdr corrections ()}$$

$$\text{---} (pvals_adjusted, 'k', \text{---}) \text{ p-values}$$

$$\text{---} (p\text{Thresh} * pvals\text{Interp}, \text{color}='gray', \text{---}) \text{ rejection line}$$

vedi Plot a pagina successiva



abbiamo creato 40 random p-values come u^2 con $u \in U(,0.01, 0.3)$

$K=40$ Ordinati, interpolazione lineare con v/k

dove v = vetore da $1 \rightarrow K$

calcolati adjusted p-values dividendo p-values / pvals Interp
abbiamo usato la funzione `fdrCorrection(s)`

- FDR correction considera che \forall raw p-value sotto la "rejection line" sono significativi -
 - FDR correction " " \forall adjusted p-values sotto la linea originale $p=q=0.05$ sono significativi
 - con FDR i valori di p sono diversi da quelli raw
- **Uso complementare:** determinare di approssimare senza modificare p_2 \rightarrow esercizio successivo

Esempio - t Student test

tra medie

Vogliamo fare il test x il confronto di due dati affiatati -

6 coppie di gemelli, un gemello fumatore, l'altro no
andiamo ad esaminare una variabile distribuita $N(\mu, \sigma^2)$, relativa
alle particelle in bronchi -

	F	\bar{F}
1	60,6	47,5
2	12,0	13,3
3	56,0	33,0
4	75,2	55,2
5	12,5	21,9
6	29,7	27,9

non fumatori $\int = \mu_x - \sigma^2 + \mu_y + \sigma^2$ YouTube

Ipotesi $H_0 \rightarrow \mu_F = \mu_{\bar{F}} \rightarrow \mu_F - \mu_{\bar{F}} = \phi$

Ipotesi altern. $H_1 \rightarrow \mu_F - \mu_{\bar{F}} \neq \phi$

supponiamo $\alpha \rightarrow 0,10$ $n = 6$

I dati non sono indipendenti x che i soggetti sono gemelli -
 \rightarrow non possiamo usare il t-test x gruppi indipendenti -

ma: **differenza tra medie = media della differenza** \rightarrow

scritto in modo diverso l'ipotesi nulla:

$H_0 \rightarrow \mu_d = \phi$ / $H_1 \rightarrow \mu_d \neq \phi$

calcolo la differenza

	d
1	13,1
2	-1,3
3	23
4	20
5	-9,4
6	1,8

Quindi mi sono ridotto al test su
un solo gruppo \rightarrow test su una media

\rightarrow uso t-test

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

che sotto H_0 e' distribuito come

gradi di liberta' $(n-1)$

\rightarrow nel nostro caso

$$t = \frac{\bar{d} - \phi}{\frac{s_d}{\sqrt{n}}}$$

radice della varianza campionaria

d	
1	13,1
2	-1,3
3	23
4	20
5	-8,4
6	1,8

$$\bar{d} = 7,87$$

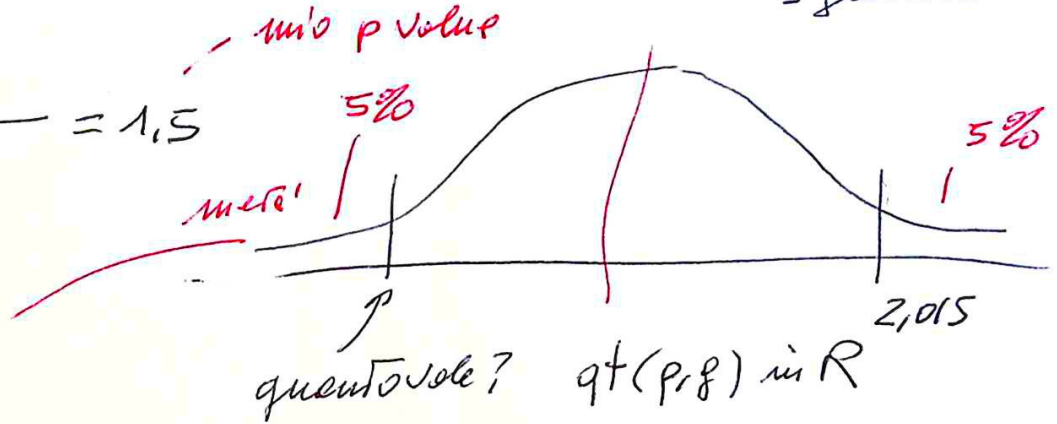
$$gdf = n - 1 = 5$$

$$s_d^2 = 164,5 \rightarrow s_d = 12,8$$

il nostro valore di t vale

t -student con 5 gradi di libe.

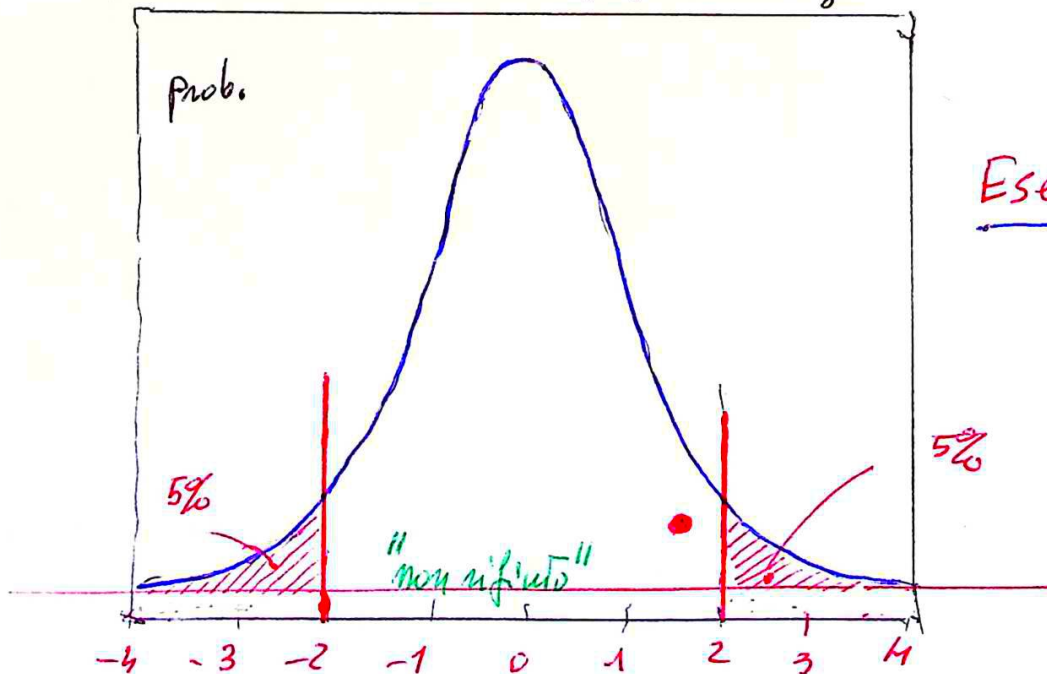
$$t = \frac{7,87}{\frac{12,8}{\sqrt{6}}} = 1,5$$



$$2\alpha = 0,1$$

Eserciziario di STATISTICA Vol. 1
 100 esercizi multi di statistica descrittiva - *soliti argomenti*

t -Student 5 gradi di libe



Esercizio 39

5) determinare una soglia critica da applicare ai p-value zero. 77

Vedi Fig. 10.17 di esercizio 4.

E' il p-value grezzo + grande che si trova sotto la soglia di rifiuto
Nel nostro grafico e' il n° 18 osservato, e il suo valore e' $p = 0,0191$.
Quelli sotto soglia sono i p-value di "non rifiuto".

Scrivere codice x identificare questo p.
(usare `fdrcorrection`)

6) compariamo tra loro Bonferroni e FDR (usare `fdrcorrection`) -

Dobbiamo creare un set di valori p, contare quanti di questi sono
significativi, usando vari metodi di correzione.

N=100 casuali come $p = u^2$ per $u \in [0, .25]$ con $\alpha = .05$
Calcoliamo la percentuale dei p che sono sotto soglia.

`pvals = mp.random.uniform(.001, .25, N) ** 2`

`bon_thresh = .05/N`

`q = fdrcorrection(pvals, .05)`

FDR significativi, Bonferroni significativi

`print(100 * mp.mean(q[phi]))` → 86%

`(pvals < bon_thresh)` → 11%

Facciamo un loop di 100 sets di p-values, calcoliamo media, std dei
p-values sotto soglia

`sigTest = mp.zeros((100, 2))`

for `expi` in `range(100)`:

`pvals =` _____ (*)

`bon_thresh = .05/N`

`q = sigTest[expi, phi] = 100 * mp.mean(q[phi])`

`1 =` _____ (`pvals < bon_thresh`) →

FDR → 80,45% — significativo

Bonferroni → 8,89% — " "

FDR + indulgente di
Bonferroni
esperienza x
scoglie

⇒ Bonferroni si basa sul n° di test N
 FDR si basa sulla distribuzione di P

Mettiamo il ciclo di esercizio ϕ in un altro loop, che varia i valori di P da 2-0500 a passi di 25 ripetuti loop -

$N_s = \text{np.linspace}(\text{np.log}_{10}(2), \text{np.log}_{10}(500), 25, \text{dtype}=\text{int})$
 distribuzione log dei P ↗

$m \text{Repetitions} = 100$

matrice del risultato finale (non x ogni ripetizione)

$\text{sigTest} = \text{np.zeros}((\text{len}(N_s), 3))$

for m_i, m in enumerate(N_s):

for $_$ in range($m \text{Repetitions}$):

p-values & corrections

$\text{pvals} = \text{np.random.uniform}(0.001, 0.25, m) * m$

$\text{bon_thresh} = 0.05/m$

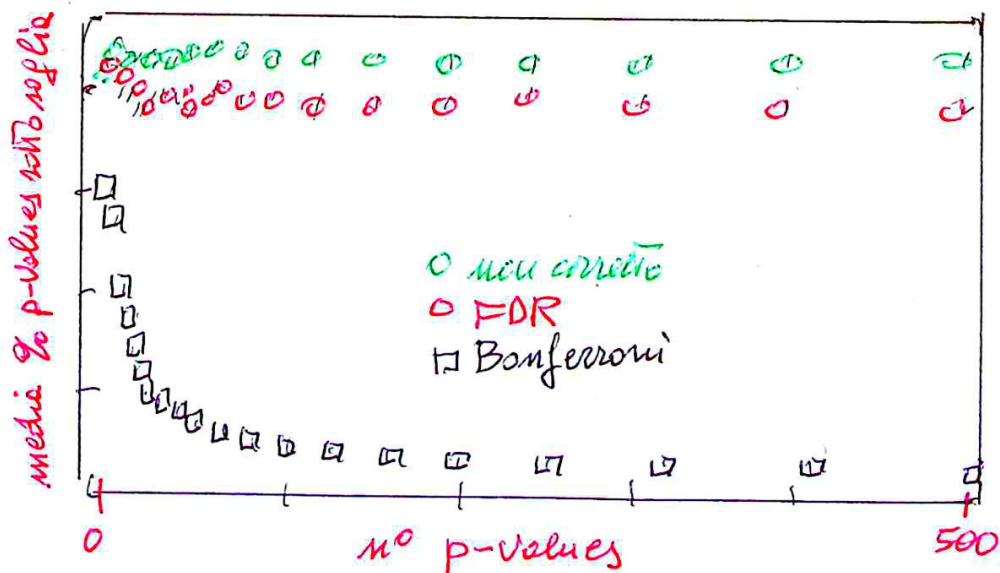
$q = \text{fdr_correction}(\text{pvals}, 0.05)$

$\text{sigTest}[m_i, \phi] += 100 * \text{mean}(\text{pvals} < 0.05)$ *non corretto*

_____ 1 += _____ ($q[\phi]$) *FDR corretto*

_____ 2 += _____ ($\text{pvals} < \text{bon_thresh}$) *Bonferroni corretto*

La distribuzione dei P non viene al variare di $N \rightarrow$
 \rightarrow FDR \approx costante



8) Dimostreremo "empiricamente" il rischio di "errore flessibile", cosa che può succedere eseguendo ripetutamente un test t -test, inferenziale, aumentando ogni volta il campione. \rightarrow Cap. 18

E' una pratica non etica, che aumenta il rischio di ottenere $p < .05$ anche quando non c'è alcun effetto.

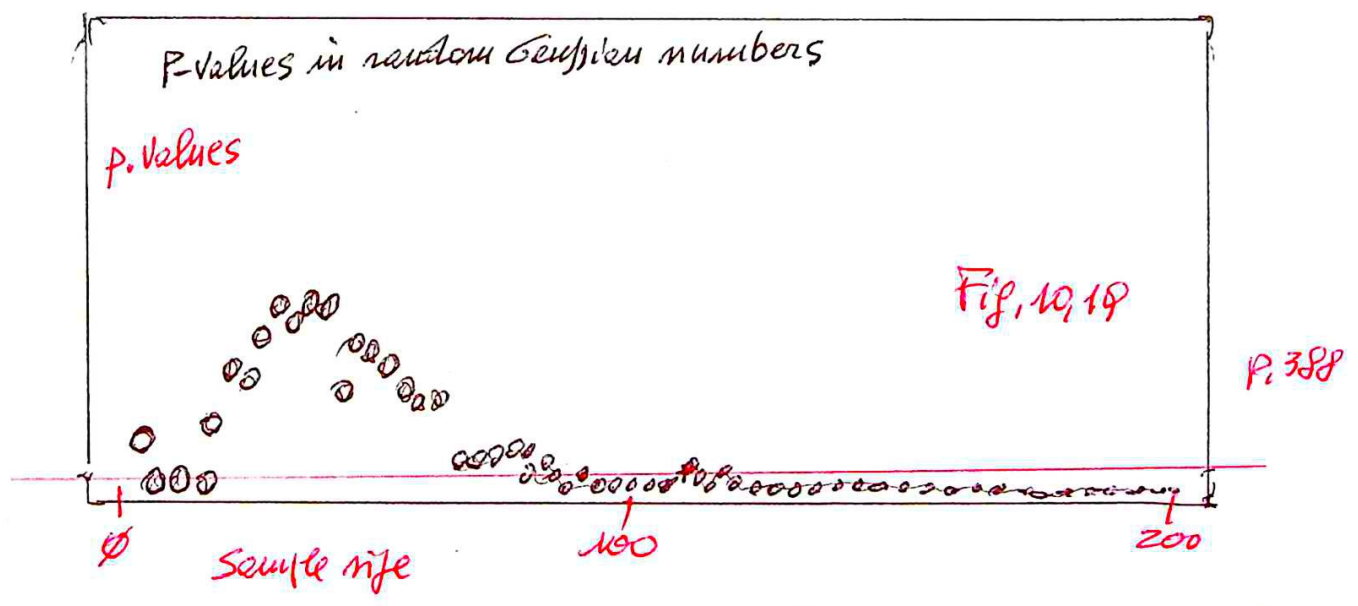
5 numeri da $N(0,1)$ - eseguiamo t-test contro $H_0: \mu = \phi$
naturalmente il test può dare risultato non significativo.

Ad ogni iterazione aggiungo 3 numeri casuali al campione.

L'insieme è un t-test

Fare grafico di p vs. dimensione campione

NOTA l'errore flessibile viene in funzione della casualità dei p



Stiamo generando numeri casuali da una distribuzione casuale con $\mu = \phi$
Tuttavia il test diventa significativo solo $\approx 2/100$

Se aumento i campioni, alla fine confermo l'ipotesi, ma ciò non corrisponde alle realtà.

9) incorettamento

— FINE CAPITOLO 10 —