

MODERN STATISTICS

INTUITION, MATH, PYTHON, R

MIKE X COHEN

 SINCXPRESS

Modern Statistics - Kindle ebook

1- Introduction to this book

La statistica riguarda l'uso dei dati per aiutare a prendere decisioni di fronte all'incertezza.

Molte procedure statistiche alla fine danno una probabilità "A".

1.2 - Statistica, scienza dei dati, apprendimento - -

- statistica = collezione, organizzazione, analisi, interpretazione, presentazione dei dati
- apprendimento = uso di "pattern" contenuti nei dati per predire o fornire in modo una classificazione
- scienza dei dati = riguarda prevalentemente le applicazioni più recenti - fin qui ignorate - come immagini, serie temporali, testi

1.3 - Target audience

Libro non orientato esclusivamente alla matematica della statistica (statistica teorica), ma focus su concetti e applicazioni pratiche (Python).

1.4 - Prerequisites

- high school math // - Calculus & Linear Algebra // programming
- ChatGPT trasforma con buona efficienza linguaggi in SAS, MATLAB, Julia

1.5 - Exercises

Sono forniti i sorgenti Python e R.
Al capitolo 20 brevi descrizioni degli esercizi.

1.6 - Learning from simulated data

- a partire da equazioni nella letteratura, →

Spesso ciò proviene da letteratura "classica", non "computer-oriented", questo approccio può essere utile a statistici teorici.

- usare ben dati la letteratura, ma spesso servono esempi poco correlati con ciò che ci interessa.
- usare Tutorial SW x implementare, questo va bene se già si conosce la statistica
- SIMULAZIONE dei DATI - è l'approccio di questo libro - con esempi e buon senso questo conduce a risultati che non sono troppo remoti dalla realtà.

1.7 - Using the code in this book

- è bene personalizzare i codici forniti. Se si producono cose nuove, utili, farele di blog dell'autore.

1.8 - Online resources

- vedi info x slowbooks
- ChatGPT-4 - il suo utilizzo è contenuto nel libro.

— FINE CAPITOLO 1 —

2 - What are (is?) "data"?

- plurale o singolare? meglio plurale.

2.2 - Where data come from, what do they mean?

- i dati in forma numerica che usiamo "non sono la realtà", ma una sua rappresentazione. Sono misure che abbiamo fatto sulla realtà.
- la qualità e la precisione dei dati sono fondamentali.

2.3 - What do data look like?

- i dati sono memorizzati in un computer. Spesso sono organizzati come tabelle (righe, colonne) -

riga = operazione colonna = misura

esempio:

lunedì	Iced-coffino	10/10
martedì	Espresso	8/10
mercoledì	Iced-mocha latte	6/10
-	-	-
-	-	-
Giorno	Coffee	Grading

	misura A	B	C		
operazione 1	15,94	23,60	15,76	-	-
-	2	18,52	32,15	4,15	-
-	3	37,96	-11,67	43,46	-
-	-	-	-	-	-
-	-	-	-	-	-

2.4 - Limitations of data

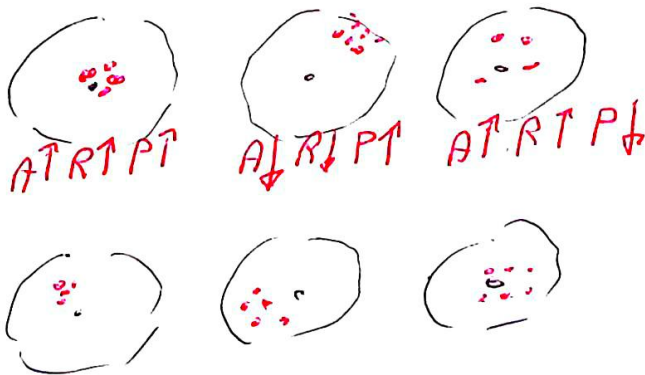
- i dati provengono da misure imperfette (necessariamente), es. la temperatura di un camion di nostro di un motore non può essere misurata direttamente.

- nelle malattie virali conosciamo solo i casi dichiarati in modo ufficiale, il resto ci sfugge.
- i dati sono soggetti a "rumore"
- ci possono essere dati "fuori dalla ragionevolezza", "skew statistics" (quelle fuori del 3 σ)
- i dati sono espressi secondo unità di misura, \exists metodi per trasformarli in numeri adimensionali.

2.5- Accuracy, precision, resolution, range

- accuracy = molto legata alla qualità dello strumento di misura
A
- precisione = capacità di avere misure simili a fronte di ripetizione della misura
P
- risoluzione = la distanza numerica tra due successive misurazioni. Es. un termometro misura 200°, 225°, 250° ma non gli intermedi intermedi.
R

esempio grafico:



2.6- Data Types

- non confondere con tipo dato del SW
- possono essere "numerical" o "categorical" (etichette)

esempio:

Family	Type	Description	Example
Numerical	Discrete	No arbitrary precision	Population
	Interval	Meaningful intervals, arbitrary precision	Temp. in °C
	Ratio	Interval but with meaningful zero	Height in cm
Categorical (labeled)	Nominal	Non-sortable, discrete	Movie genre
	Ordinal	Sortable, discrete	Education level

- circostanza di questo cosella, dentro il PC e' trattata in numeri - Es. BMW=1, Mercedes=2, ...

2.7 - From anecdotes to populations

- speculation, theory - $N = \phi$
- anecdote $N = 1$
- case report $N = 1$ - in medicine
- sample - se voglio per medici di un animale, non lo so per tutti quegli animali esistenti, ma un sottoinsieme
- pilot study N piccolo
- small scale / large scale
- observational study - cose dati senza manipolazione o forzatura esterne
- convenience sample - dati raccolti con poca fatica (es. nucleo familiare, amici)
- population - cio' che si studia non si fa solo il singolo

2.8 - Data management

Sistemi di documentazione, archiviazione, memorizzazione - Backups -

2.9 - Ethics of making up data

- e' facile generare "fake data" - Quello che conta e' saperli riconoscere -

— FINE CAPITOLO 2 —

3 - Visualizing data

3.1 - Why visualize data?

- per meglio capire, per sintetizzare, per trovare errori -

3.2 - How visualize data

- per dati poco numerosi, li si possono visualizzare tutti

Dom	Sab	Ven	Gio	Mer	Mar	Dom
18,28	21,10	21,25	29,04	18,47	19,12	19,30

3.3 - Bar plots

es. dove la gente prende le sue informazioni?

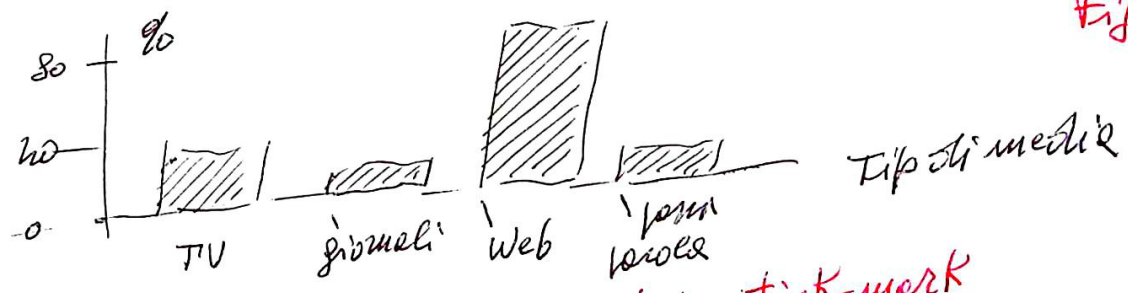


Fig. 3.3

non equidistanti
ordine arbitrario

tick-mark
non necessita il colore

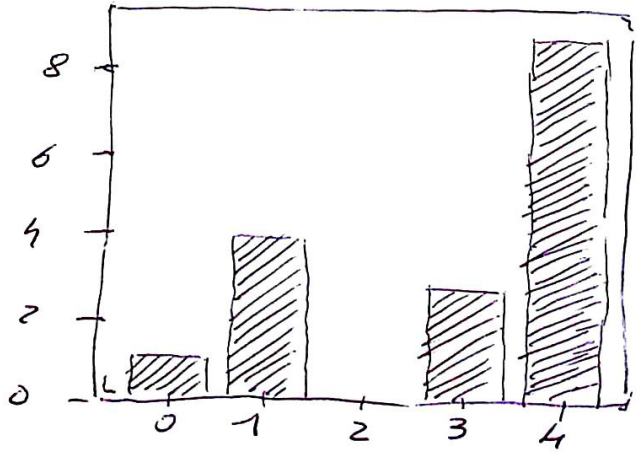
```

import matplotlib.pyplot as plt
Y = [1, 4, 3, 8] # altezza delle barre
X = [0, 1, 3, 4] # bar locations
plt.bar(X, Y)

```

Fig. 3.4

stats_ch03_visualization



%

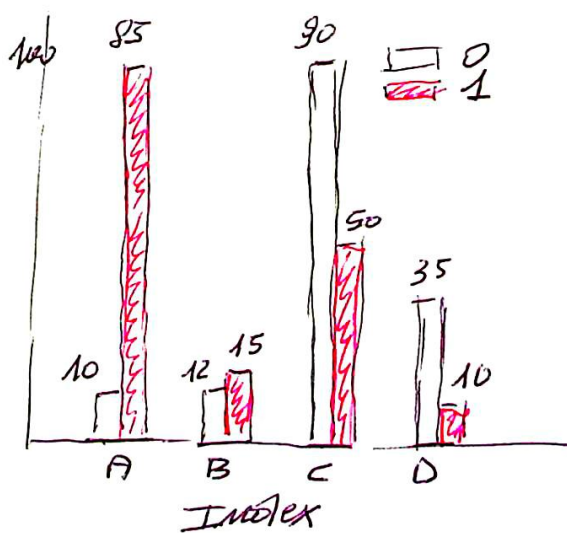
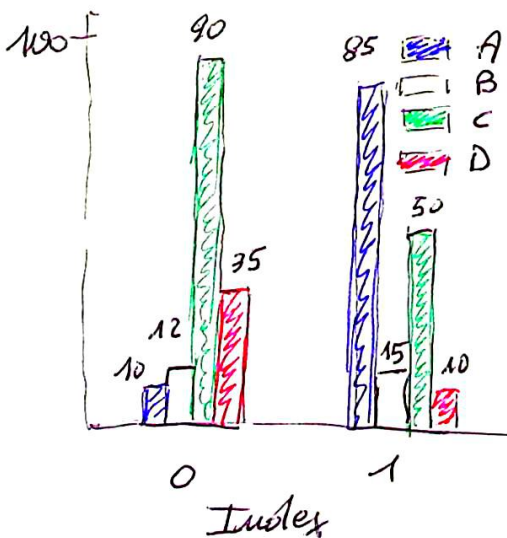
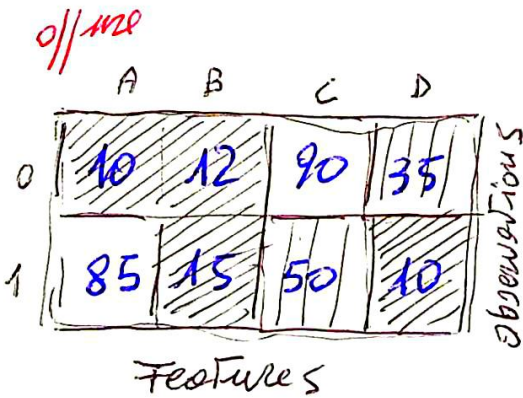
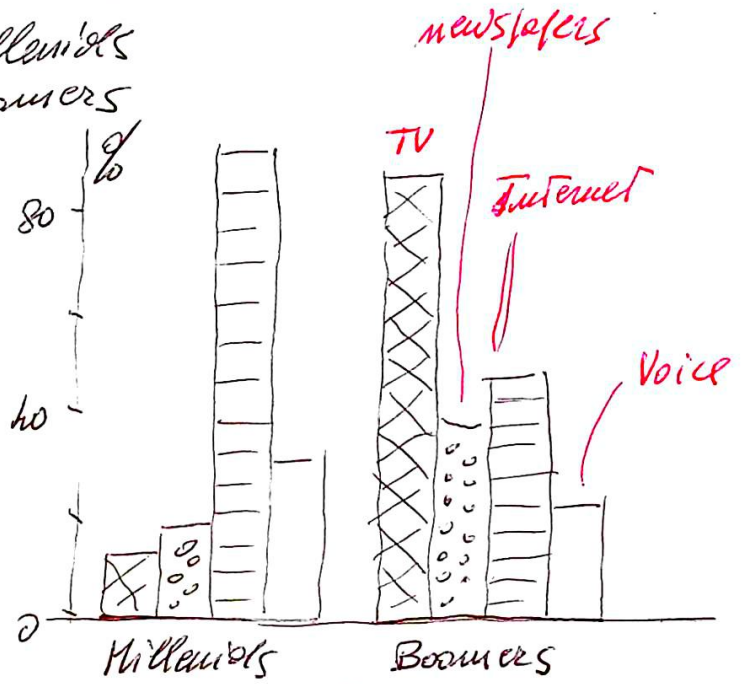
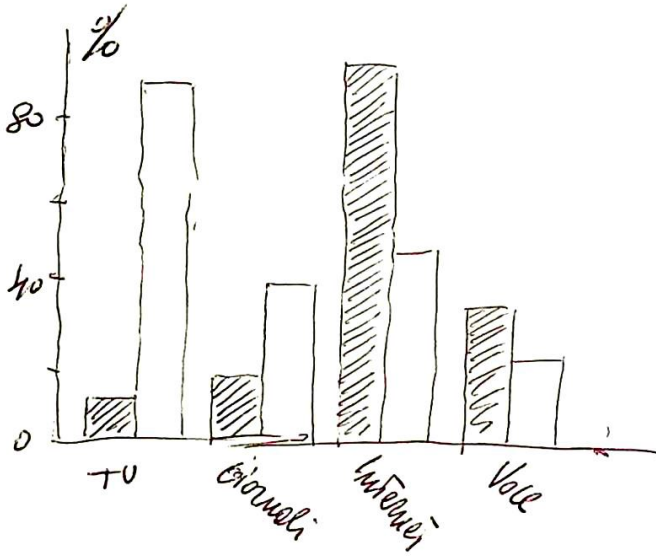
- i dati possono essere rappresentati in modi differenti

news_sources = mp.winey ([[12, 17, 95, 35], [90, 40, 50, 25]])

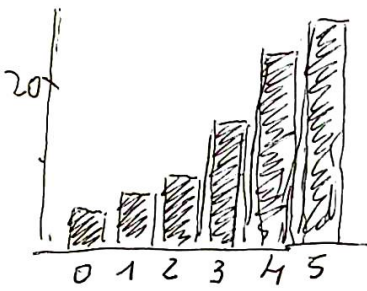
agegroups = ['Millennials', 'Boomers']

Rappresentazione x news source

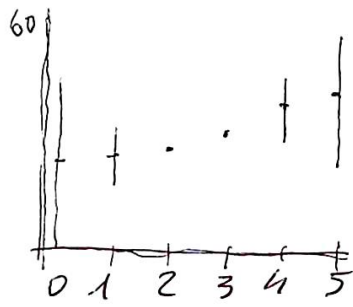
Millennials
Boomers



A) Bar plot



B) Error plot



C) Error bar plot

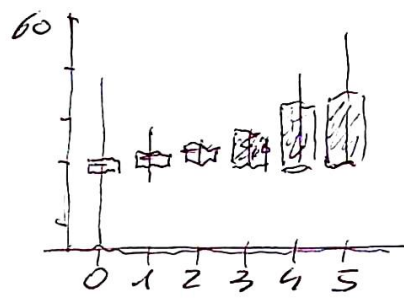
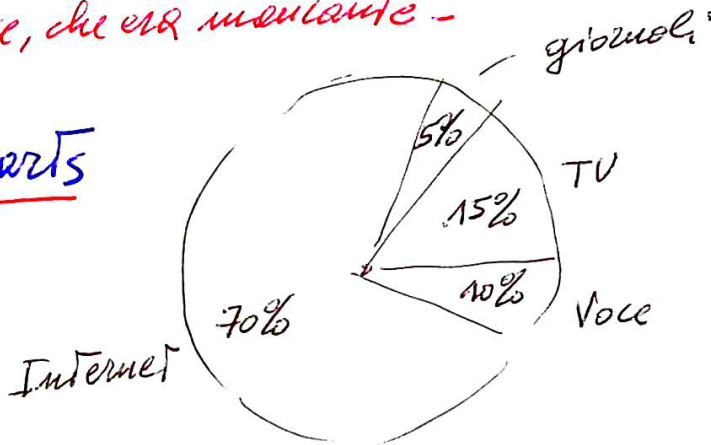


Fig. 3.7

Stessa non ben qui' una rilevazione finita, e' il risultato di una formula sviluppata codice, che era mancante -

3.4 - Pie charts

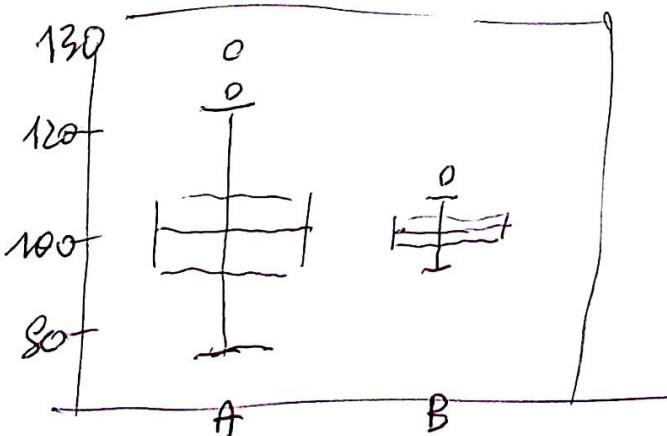
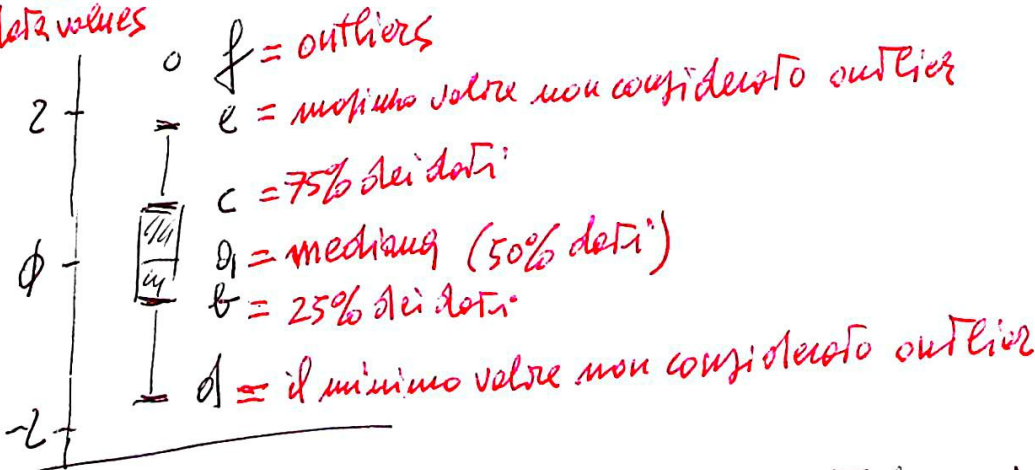
NB $\sum \% = 100\%$



3.5 - Box plots

pet. tight layout (c)

data values



due distribuzioni di dati con uguale medie e mediana, ma differente variabilità -

Possano essere confrontati

3.6 - Histograms

Altre sono simili a grafici a barre, ma sono molto diversi:

- bar plots = "categorical data"
- histograms = "numerical"

es. $X = [1, 2, 2, 3, 3, 5, 5, 5, 5, 6, 7, 7, 7, 8, 8, 9]$

```
plt.hist(X, bins=len(set(X)), color='gray', edgecolor='k')
plt.xticks(np.arange(np.min(X), np.max(X)+1))
plt.tight_layout()
```

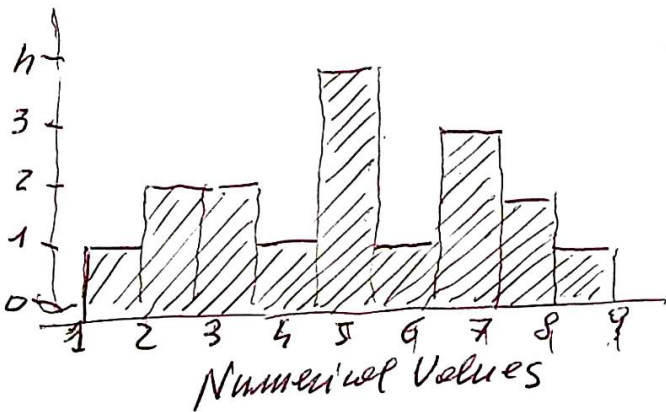


Fig. 3.11

i ticks mark non sono al centro della barra →
facciamo una modifica

```
plt.hist(X, bins=np.arange(0.5, 9.5, step=1), color='-', edgecolor='-')
```

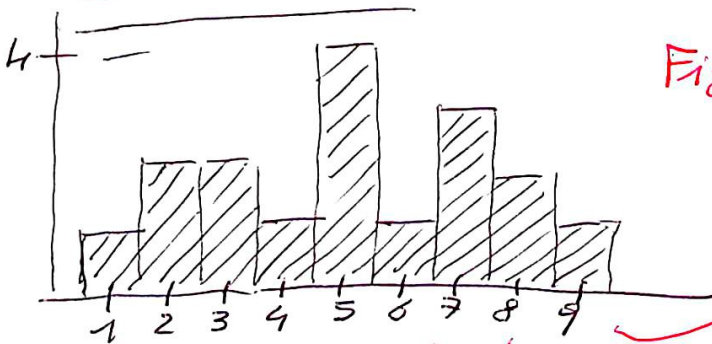


Fig. 3.12

molto importante
x istogrammi "multiplici"

es: istogrammi delle lunghezze delle manufatti

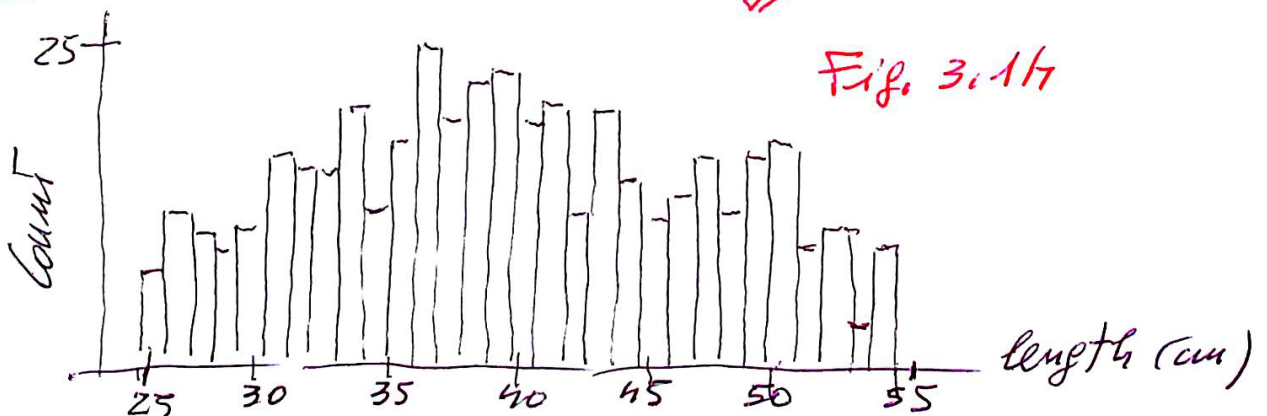


Fig. 3.17

Istogrammi della monguete

monguete = np.arange(np.random.uniform(-.75, .75, size=500) * 15 + 40)
 plot.hist(monguete, bins=30, color=)

Istogrammi con Bins differenti

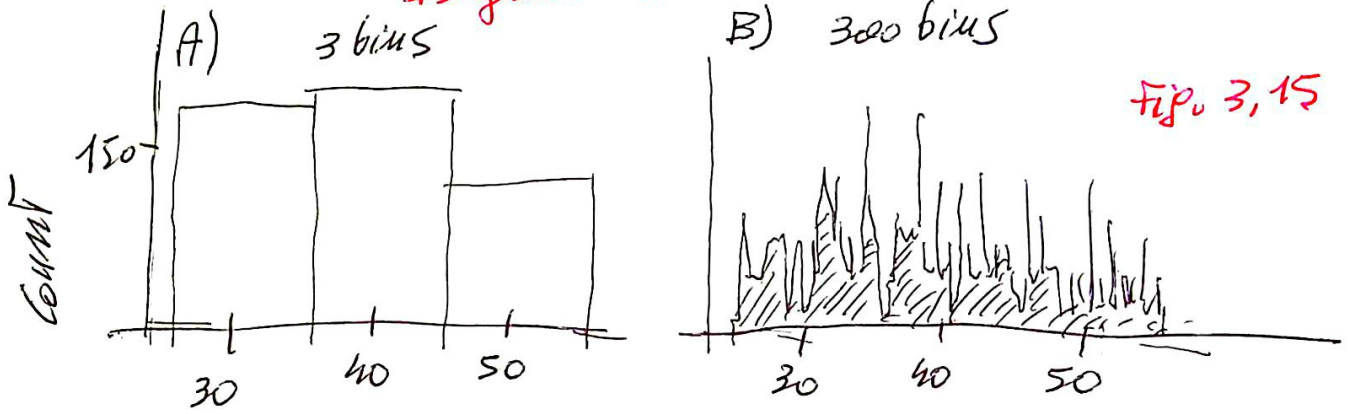


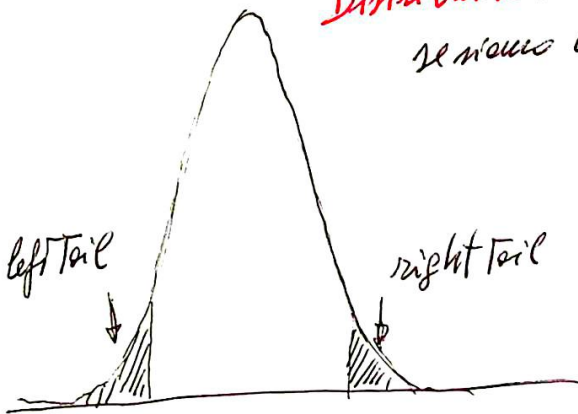
Fig. 3, 15

quale il n° giusto di bins? - in concreto $3\phi - 4\phi$

Distribution tails

se siamo attorno allo zero →

negative tail
 positive



one tailed
 two tailed

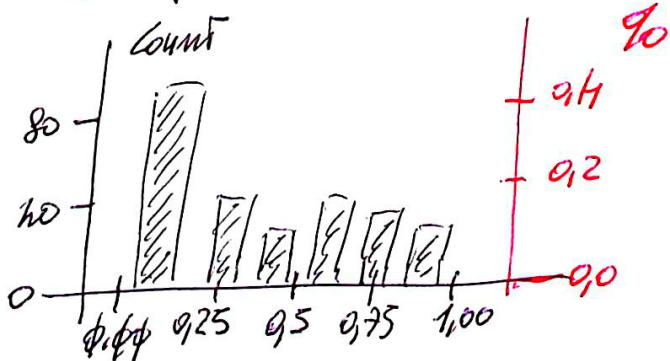
La "shape" di un istogramma - ordinando i dati lungo x - riflette proprietà intrinseche della distribuzione.

- l'ordinamento non ha senso in un "boxplot" (labels)

- se ho due distribuzioni, e metto su y i dati grezzi →

→ in generale i due grafici saranno difficilmente confrontabili

→ meglio convertire in % (con questo cambio y)



%

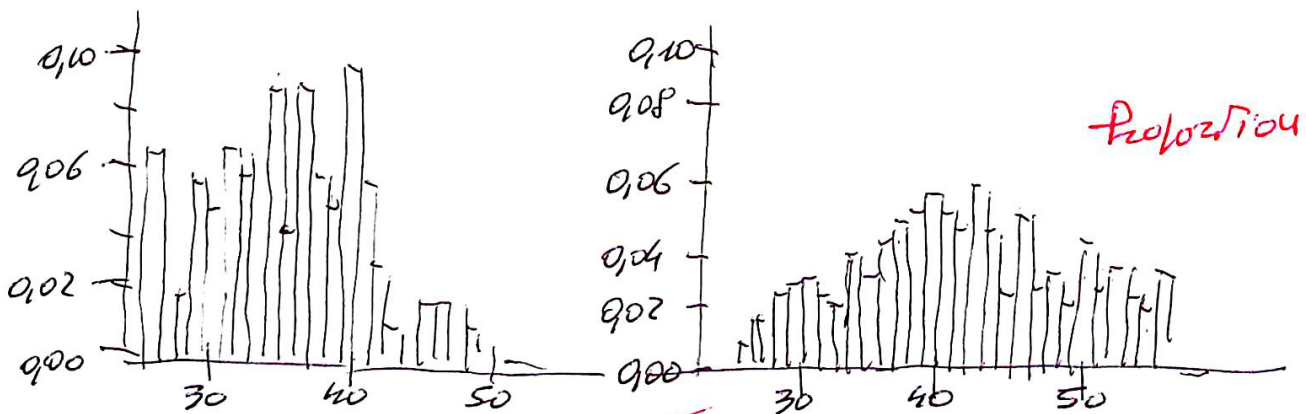
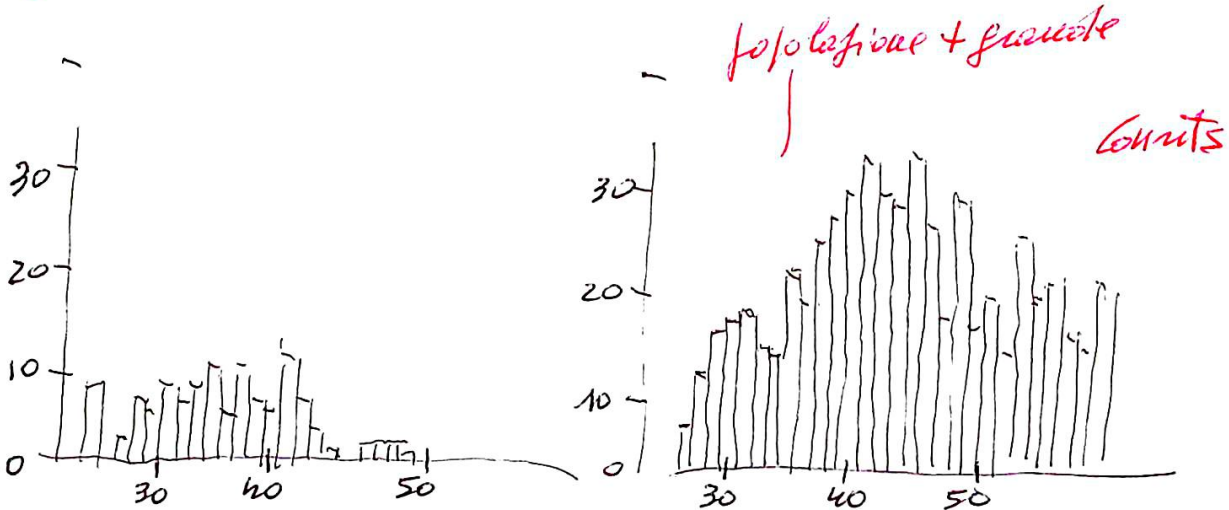
Raw counts	Proportion
Facile interpretazione	Sfuso x ricadurre ai dati
Difficile comparare con altri	Facile da comparare
Σ arbitrario	$\Sigma = 1$ (o 100%)
buono x ispezione qualitativa	buono x ispezione quantitativa

Fig. 3.18

Come convertire "counts" a "proportion"?

$$\tilde{b}_j = \frac{b_j}{\sum_{i=1}^K b_i} \quad (3.2)$$

$\Sigma = 1$ — comportamento dei due tipi di grafici

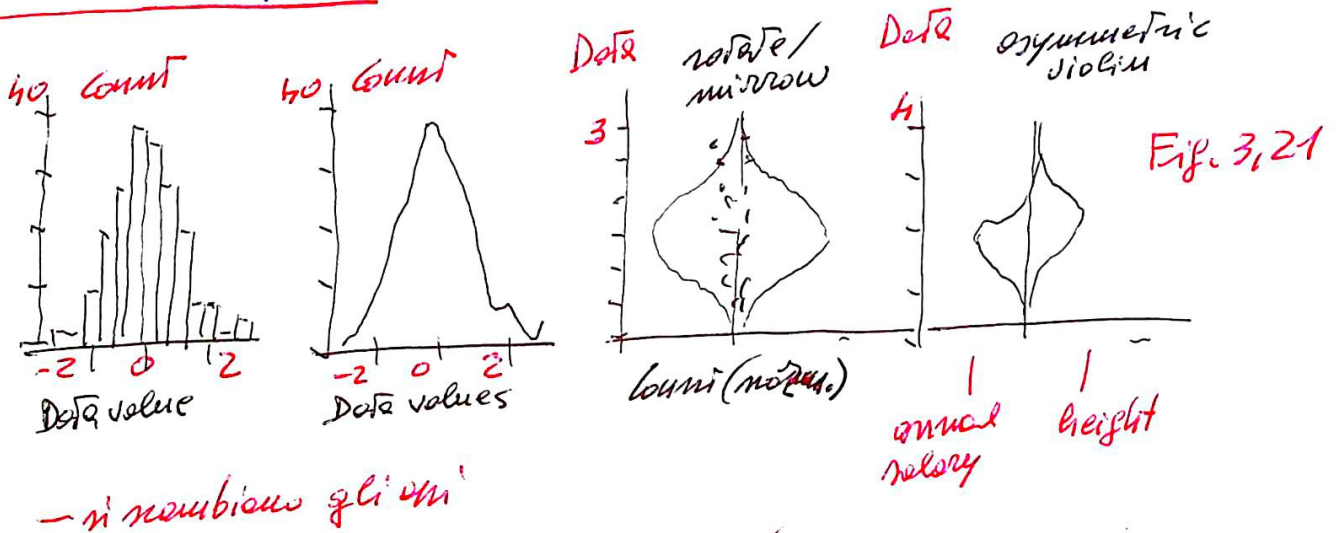


comparabili $\Sigma = 1 \rightarrow$ forma della distribuzione
indica le differenze di struttura

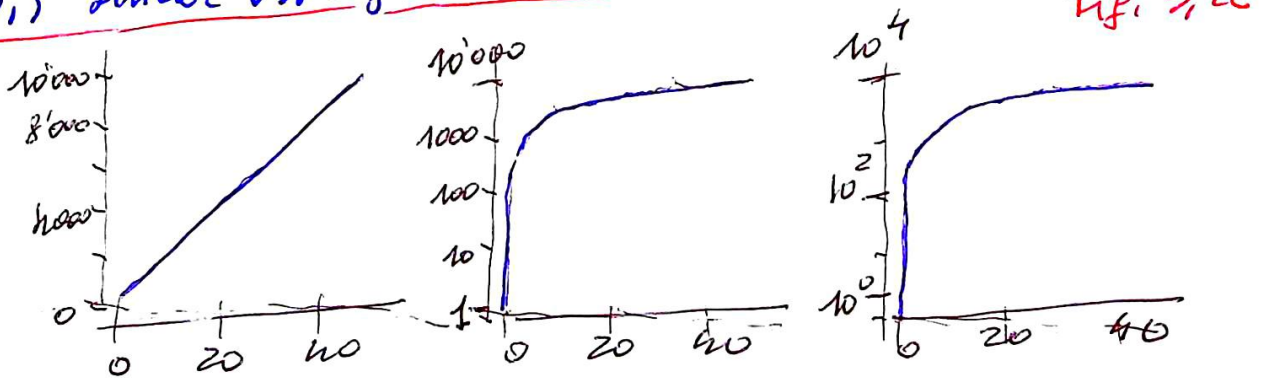
3,7 - lines vs, bars in a histogram

- se le barre sono sufficienti → ha senso tracciare una linea tra bins

3,8 - Violin plots - ci sono defezioni



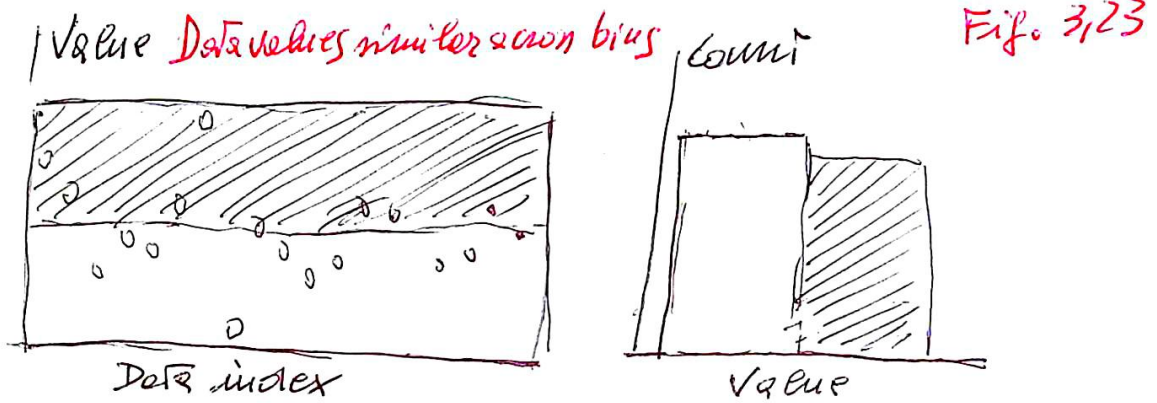
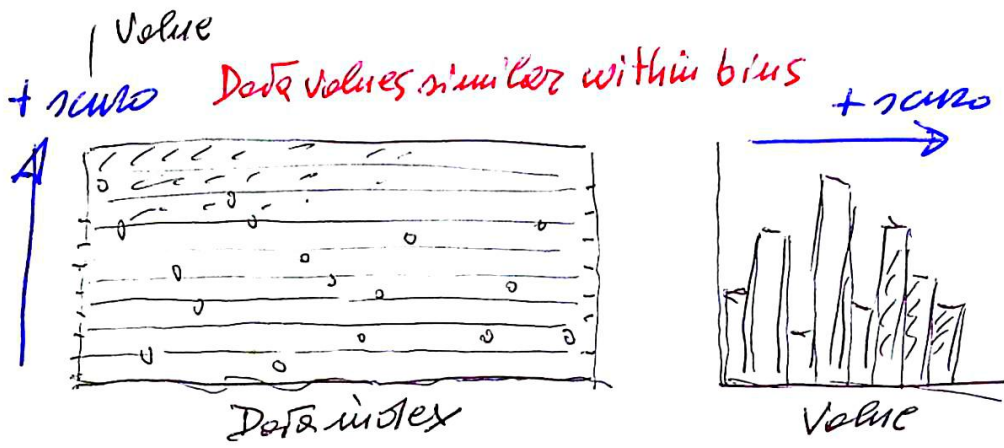
3,9 - linear vs. logarithmic axis scaling



3,10 - Discretizing continuous data

E' importante nel decidere se perseguire ANOVA o regression, e nella visualizzazione in una "regression analysis" -

- per capire meglio visualizziamo la figura 3,23



- campionare e' quindi non problematica complicata

3.11 - Radial plots

- utili quando sono "wrapping around"
- la distanza dal centro indica la magnitudine (equivalente a altitudine)

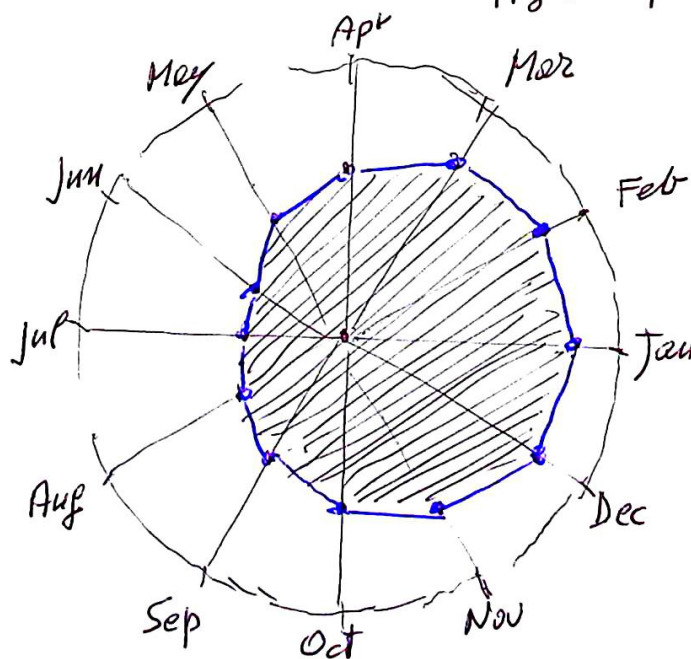


Fig. 3,24

```

ax = plt.subplot(111, polar=True)
ax.plot(theta, tempC, 'ko')
ax.set_xticks(theta [-4])
ax.set_xticklabels(months)
ax.set_yticks([10, 20, 30])
ax.set_ylim([0, 30])
ax.set_title('_____')

plt.tight_layout()
plt.show()

```

3.12 Color

in matplotlib \rightarrow color = (R, G, B) tra 0 e 1

- colore può dare diversità, ma anche confondere
- usare con cautela
- colori / se venuti in "famiglie" di figure "non riconoscibili"

3.13 - Exercises

	0	1	2	3
A	0	3	6	9
B	1	4	7	10
C	2	5	8	11

disegnare x righe e per colonne, stare in
 usare random seed x error bar plot
 30 osservazioni, 6 misure, date da:

$$Y_i \sim N(\mu_i, \sigma_i^2)$$

mu, sigma:

$$\mu_i = (i+1)^2$$

distrib. norm. media μ

$$\sigma_i = 30(2i/5 - 1)^2$$

" " dev. std = σ

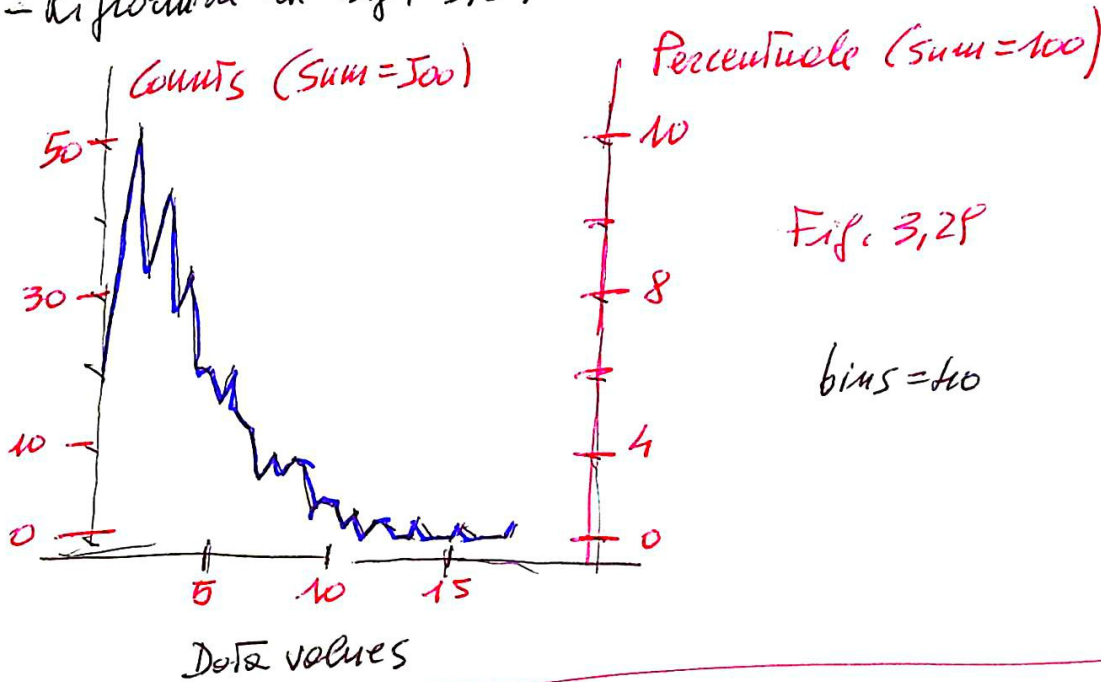
np.random.normal(mu, sigma)

ottenso la Fig. 3.7

Graficare - 60 persone che esprimono un gusto
 - e' appropriata una torta?

Vaniglia
 Strawberries
 Pistachio
 cioccolato

- Riprodurre la Fig. 3,28



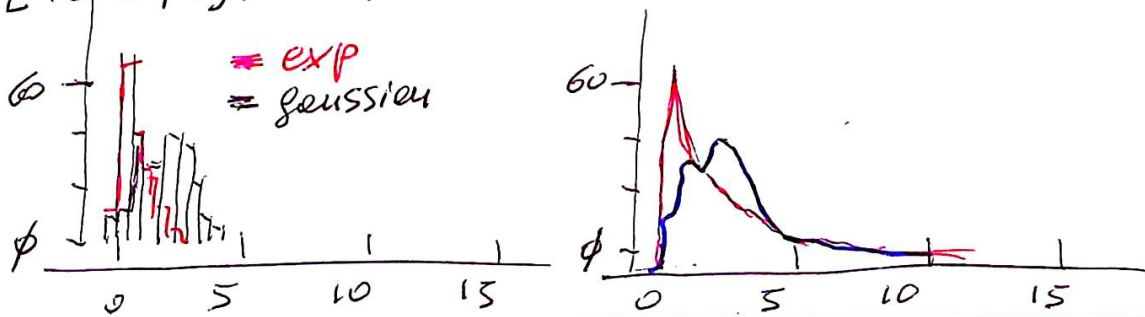
- Bone o linee? $N=200$ $\left\{ \begin{array}{l} normale \\ esponenziale \end{array} \right.$

$G \sim N(2,1)$

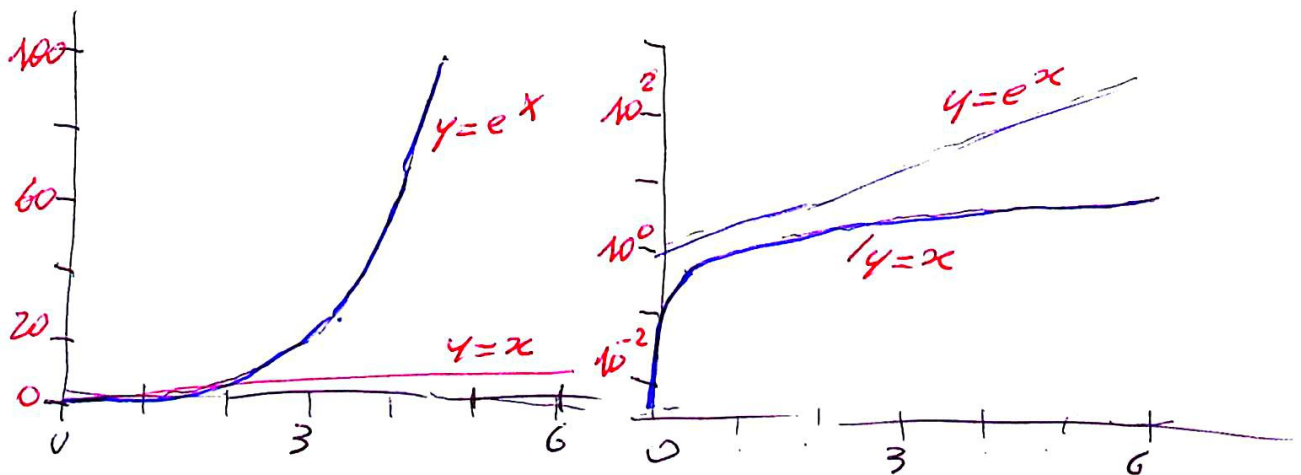
$E \sim \exp(N(\phi,1))$

bins = 30

usare mp. histogram ()



- linear vs. logarithmic



Riproduzione la ff. 3,32 - 123 numeri random
(normale, uniforme)

memorizzati su 123 righe, 2 colonne

vic' due
date sets / 123

mp.random.uniform(ϕ , 1, size=123)

df[[ϕ], [ϕ]]

sns.violinplot (date= ϕ , palette='gray', ax=axis[ϕ])

sns.stripplot (date= ϕ , ax=axis[ϕ], palette='dark:w')

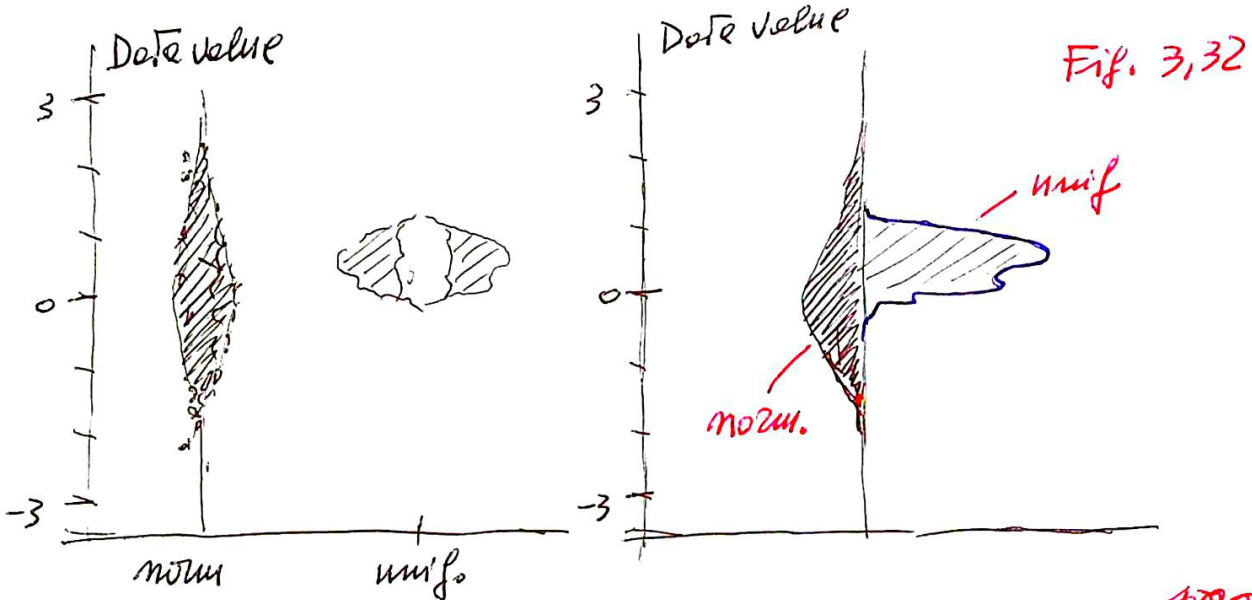


Fig. 3,32

tipi di distribuzione

pc 100

df = pd.DataFrame (mp.hstack ((mp.random.randn (123, 1),
mp.random.randn (123, 1))), columns = ['norm', 'unif'])

intra gli error in sequenza orizzontale (colonna)

0	-0,145	0,261	
1	-0,845	0,888	(123 righe, 2 colonne)
2	=	=	(123, 2)
	norm	unif	
	i		
	(123,)		

non solo label, definiscono il tipo di distribuzione

dobbiamo combinare le colonne in una

```
df_all = pd.DataFrame(pd.concat((df['norm'], df['unif'], series=...),
```

```
columns = ['y']) -> df_all (246, 3)
```

nome della colonna di 246 righe

0	0,34	norm
1	0,61	norm
2	-	norm
3	-	-
	y	dist2



```
df_all['dist2'] = 'unif'
df_all['dist2'][:len(df)] = 'norm' ignora uniform
df_all[''] = '-'
print(df_all)
```

ora creiamo "split violin plot"

```
sns.violinplot(data=df_all, x='', y='y', palette='gray',
ax= axs[1], split=True, hue='dist2')
```

```
plt.tight_layout()
plt.show()
```

— Fine capitolo 3 —

4 - Descriptive Statistics

4,1 - Descriptive vs. inferential statistics

- Statistica descrittiva = numeri che caratterizzano un set di dati (media, mediana, varianza, inclinazione, spettro, covarianza)
 - Statistica inferenziale = algoritmi efficaci a uno o più set di dati per verificare se e' probabile che le statistiche descrittive di quel set di dati si generalizzano in altri set di dati (F-value, t-value, ANOVA, regression)
- La stat. inf. utilizza i dati per fare affermazioni che non abbiamo.

CONVENZIONE: $\left\{ \begin{array}{l} \text{Stat. desc.} = \text{"descriptive statistics"} \\ \text{Inf. stats} = \text{"statistics"} \end{array} \right.$, YouTube

4,2 - Data distributions

"Gaussian style"?

300 persone cui ti chiedo: quante volte hai visionato

- il video puo' essere visionato in parte
- anche lo 0 e' ammesso
- non hanno senso numeri < 0

Data type Ratio

come visualizzare?

$\text{timesWatched} = \text{mp. round}(\text{mp. abs}(\text{mp. random}(\text{randm}(500) * 20)))/2$

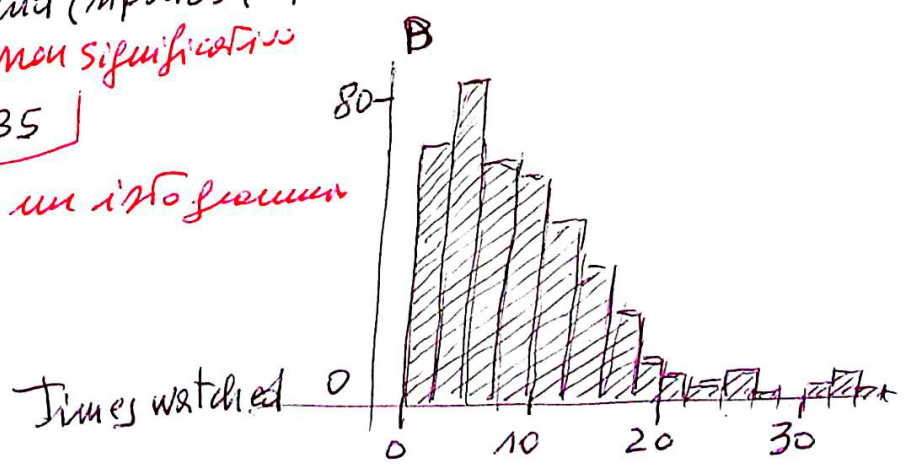
(500,) p = 103

fu' un "outlier" non significativo

$\text{timesWatched}[300] = 35$

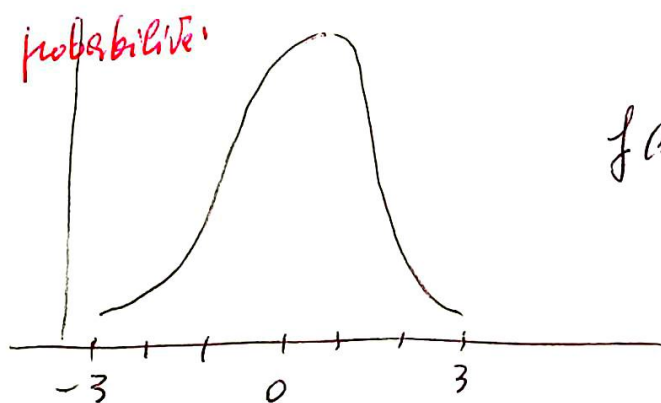
un grafico semplice e un istogramma

$\text{bins} = 'fd'$



† dati possono provenire da misure, o da formule -

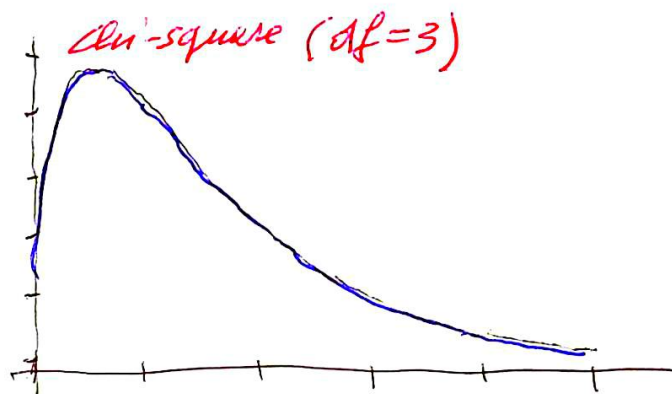
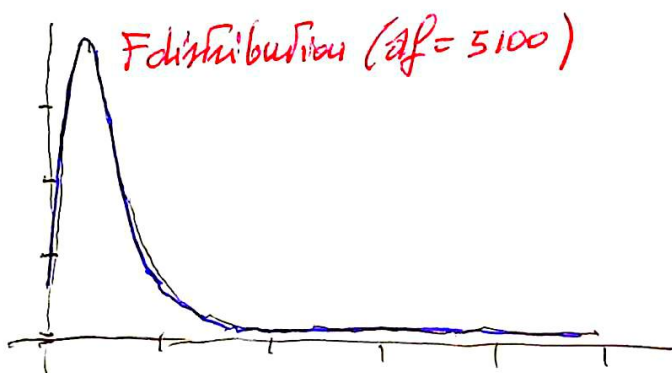
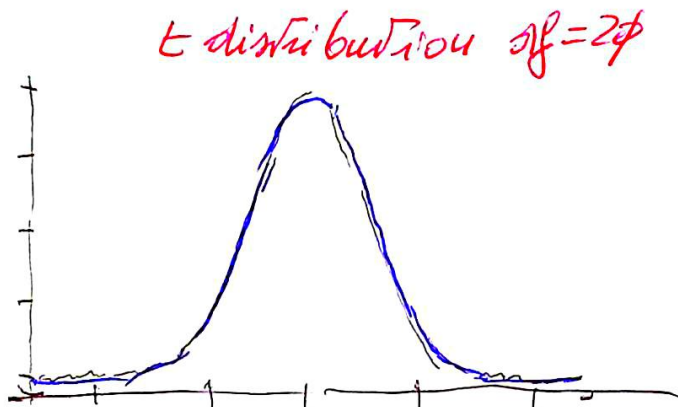
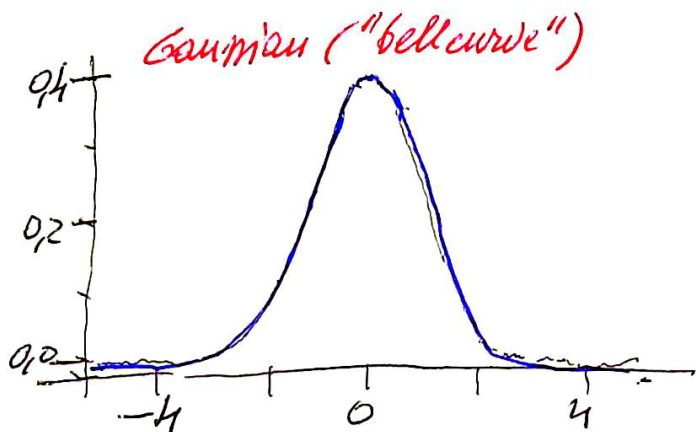
- ogni dato sets ha la sua distribuzione
- se genero due sets di numeri random, ne faccio due inferenze, i vede che non sono identici, ma molto simili -
- una distribuzione qualitativa non viene creata raffigurando i dati misurati, ma vedendo una formula matematica -
- Se un modello \rightarrow probabilità \rightarrow es. Bernoulli



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-x^2/(2\sigma^2)} \quad (4.11)$$

- ci sono miliardi di distribuzioni statistiche p. 107
- danno una visione qualitativa delle caratteristiche dei dati
- la mappa delle statistiche si basa su ipotesi relative alle distribuzioni sottostanti \rightarrow e' bene conoscere le distribuzioni
- inferenza statistica = ~~defezionamento della probabilità dei dati~~ ~~defezionamento statistico~~ ~~inferenza statistica~~ ~~descrittiva di un campione~~ = interpretazione dei dati; che permette di trarre indicazioni su fenomeni non osservati direttamente — stella osservazione di una parte dei dati "campione" — selezionato usualmente mediante un esperimento casuale (aleatorio), si indicano le caratteristiche della intera popolazione -
- Esempi pratici in Biologia, Fisica, informatica - Es. Al suo nessi volti da campionamenti casuali di determinate distribuzioni
- lo scarto tra le previsioni e i valori osservati: "residuals" permette di valutare la qualità del modello
- legge dei grandi numeri, Teorema del limite centrale - Per comprendere i concetti statistici fondamentali -

Esempi di distribuzioni



Esempi di distribuzioni storiche

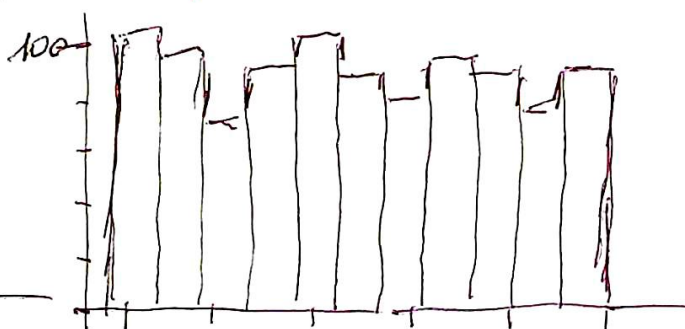
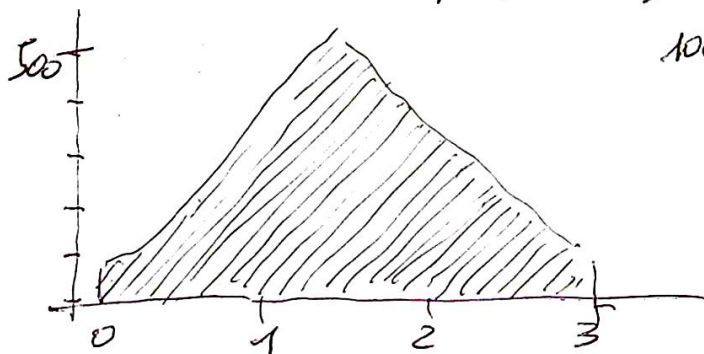
$x = \text{np.linspace}(-5, 5, 10001)$

$y = \text{stats.norm.pdf}(x)$ *Gaussian*

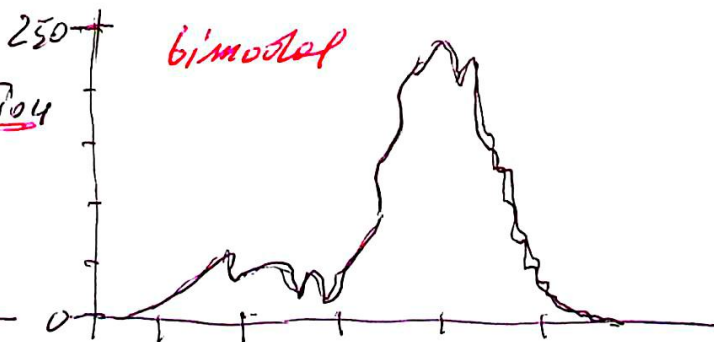
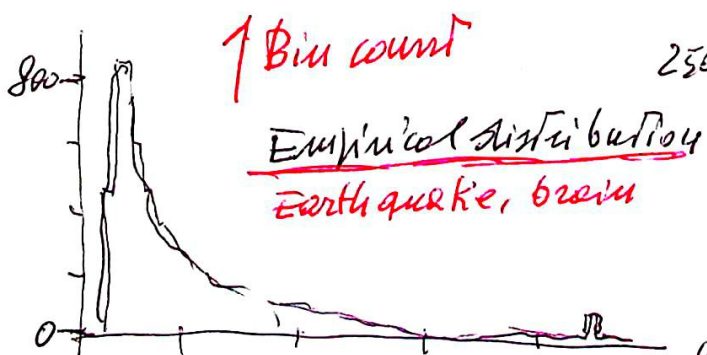
$= \text{stats.t.pdf}(x, 20)$ *t-distribution*

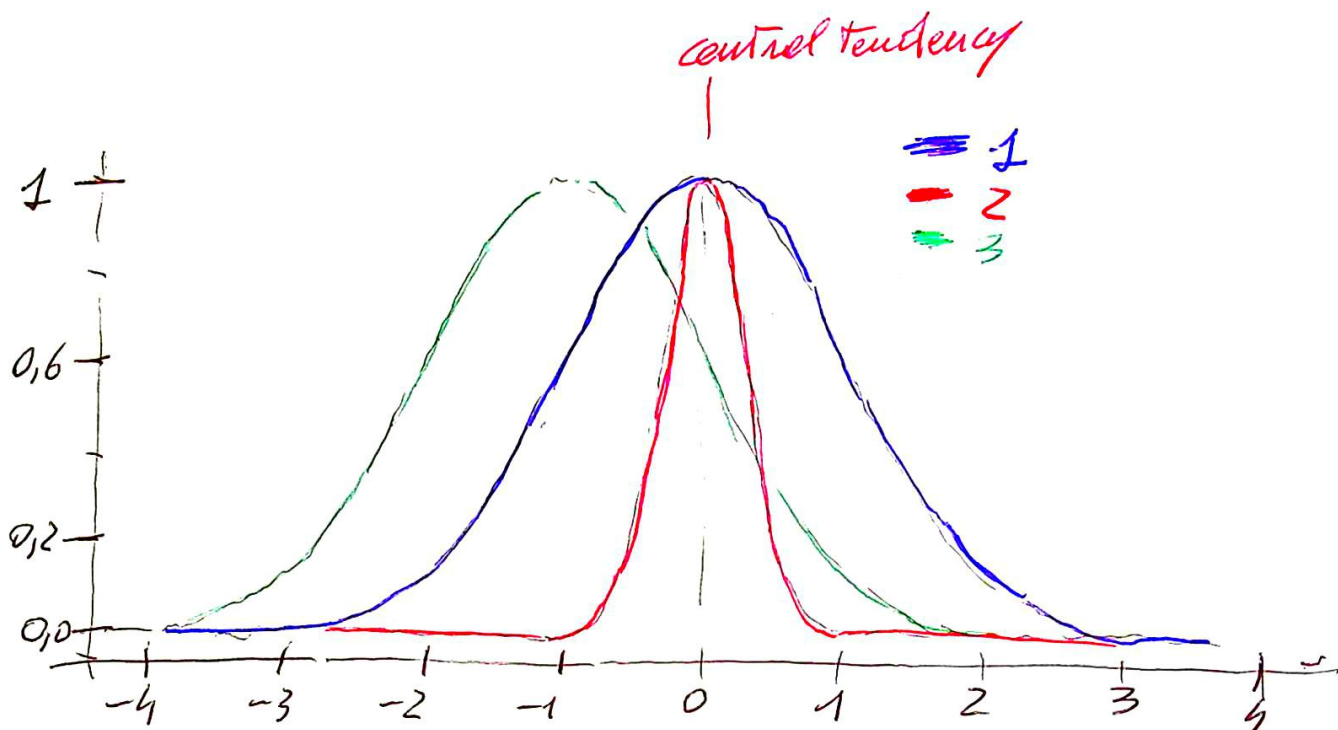
$= \text{stats.f.pdf}(x, 5, 100)$ *F-distribution*

$= \text{stats.chi2.pdf}(x, 3)$ *Chi-square distribution*



→ data value





Questi grafici stanno un'idea qualitativa su "central tendency" e su "dispersion".

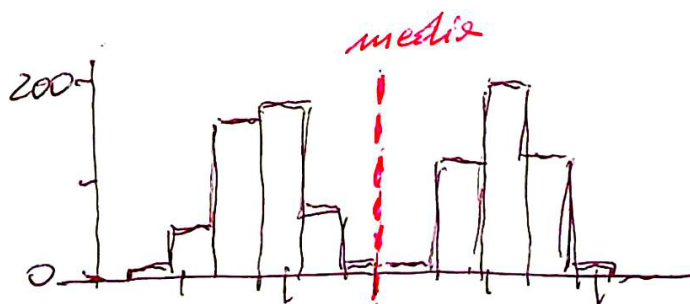
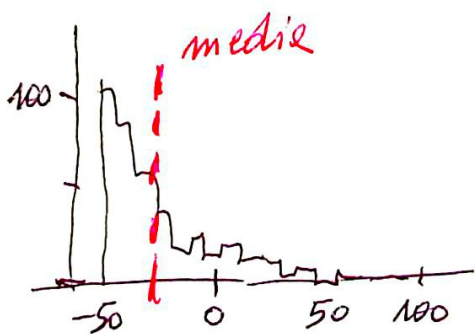
h.3 - Central tendency

Ci sono molte modalità x quantificare "central tendency"

mean, median, mode

— media ————— $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

La possiamo visualizzare come un segmento verticale sovrapposto alla distribuzione. Si presta bene con distribuzioni unimodali e simmetriche.



qui non utile

Media \neq expected value — sensible a outliers

Mediana divide in due metà ugualmente numerose i componenti

$$\text{med}(x) = \tilde{x}_i, \quad i = \frac{n+1}{2} \quad (h14)$$

sorted data
i = indice che corrisponde alla mediana

può essere un numero decimale

È facilmente interpretabile?

mediana ≠ valore estremo

Moda una misura della "central tendency"

è il valore che compare più volte

la moda può essere calcolata nei dati sono discreti

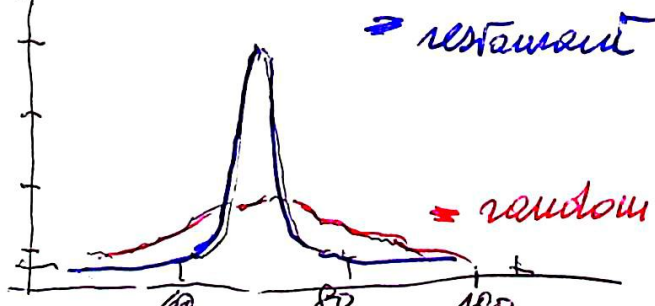
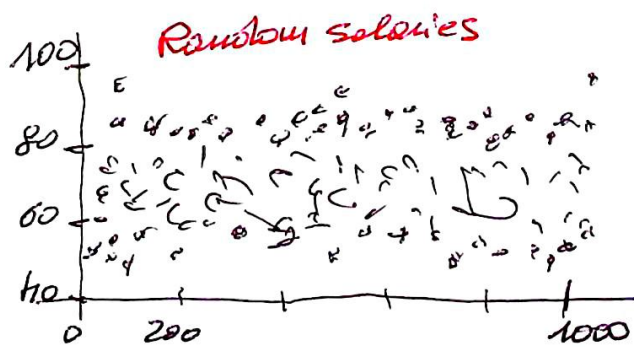
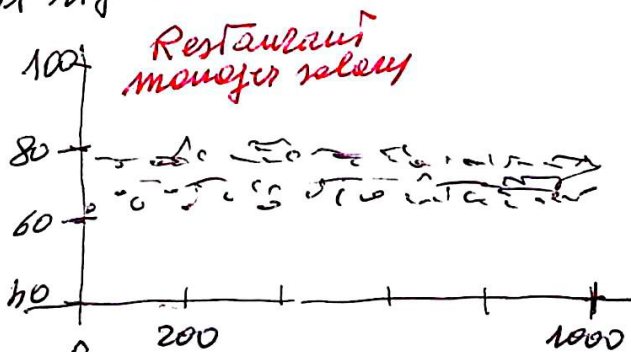
ci può essere + di una moda

$[\underline{\phi}, \underline{\phi}, \underline{1}, \underline{1}, 2, 7]$ $\left\{ \begin{array}{l} \text{moda 1} = \phi \\ \text{moda 2} = 1 \end{array} \right.$

in un istogramma la moda sta nella barra alta

h1h - Measures of dispersion

Si riferisce all'ampiezza della distribuzione



Varianza σ^2 - sample variance s^2

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (4,5)$$

MAD

$$MAD = \frac{1}{n-1} \sum_{i=1}^n |x_i - \bar{x}| \quad (4,8)$$

absolute difference

- 1) s^2 è sensibile (-permette di individuare-) outliers
L ~~è~~ usata x financial risk
- 2) s^2 è continua e ha derivata morbida
L ottimo x ottimizzazioni
- 3) s^2 è strettamente legata alla distanza Euclidea
- 4) s^2 è momento di 2° ordine
- 5) s^2 è legata al "least-squares-algorithm" for fitting regression model to data

a) MAD è una buona misura della dispersione

b) MAD meno influenzata da outliers

c) MAD usata in machine-learning

d) MAD usata in "optimization methods"

$n-1$ = degrees of freedom

in la stessa unità
dei campioni

Standard deviation

$$STD = s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

(4,8)

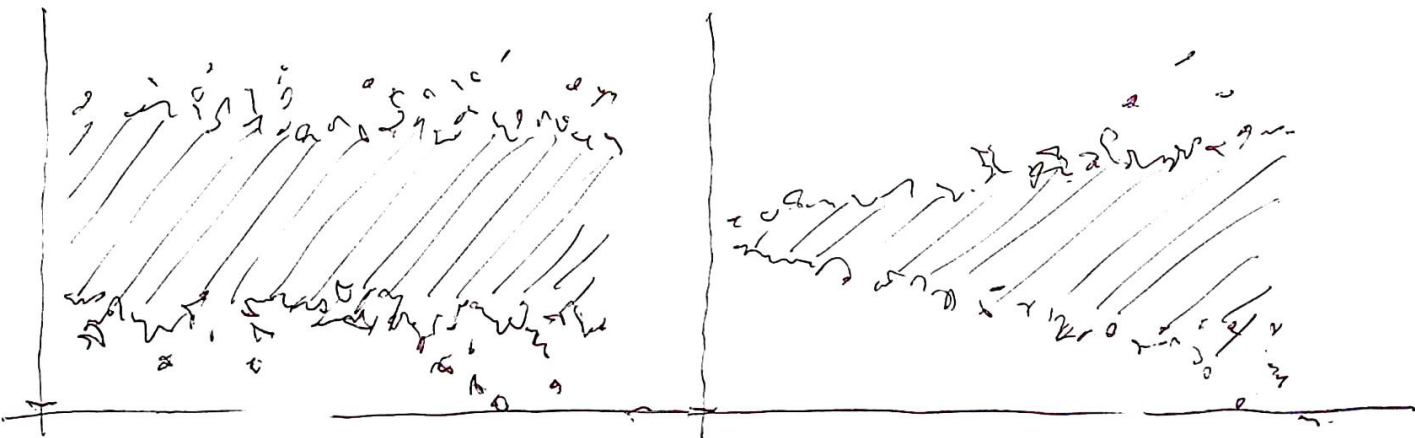
Una comune moltiplicazione:

Z-scoring → dati in unità di deviazione st. ^{interpretabile}
es. la regressione in queste unità è più facilmente

Heteroscedasticity and Homoscedasticity

Homoscedasticity — è una variabile che ha uguole varianze per tutti i valori

— in alcuni casi la varianza cresce secondo la x che esmente



Homoscedasticity

Heteroscedasticity

importante x correlazione, regressione

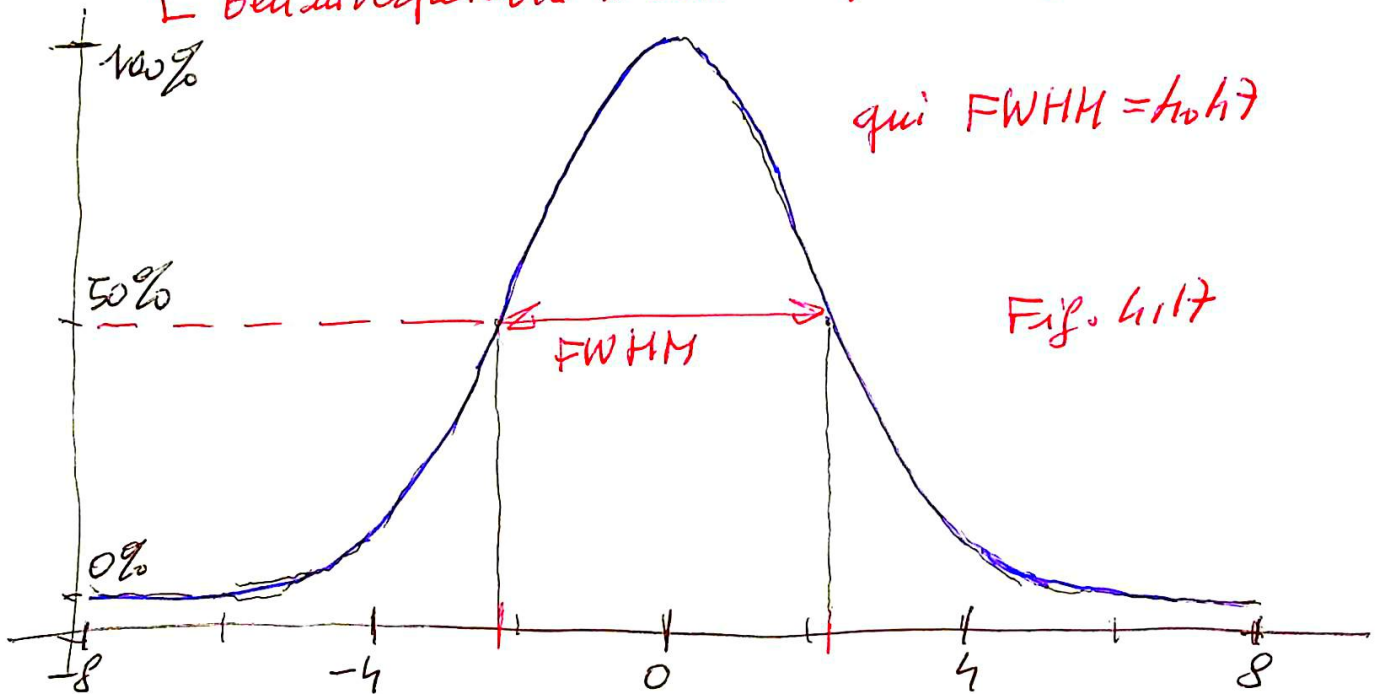
Esempio di Hetero — : ricchezza e spese

se esmente la ricchezza, la varianza degli acquisti esmente

FWHM - Full width at half maximum

FM FWHM è una misura dell'ampiezza della Gaussiana

- ↳ può essere calcolata analiticamente o empiricamente
- ↳ ben interpretabile x distribuzioni gaussiane



nel caso della Gaussiana:

$$g(x) = \exp\left(\frac{-x^2}{2\sigma^2}\right) \quad (4.10)$$

$$\text{FWHM}(g(x)) = 2\sigma \sqrt{2 \ln(2)} \quad (4.11)$$

per una distribuzione empirica non è direttamente calcolabile - Si usano algoritmi per fare questo calcolo. P. 128

↳ vedi Esercizio 10

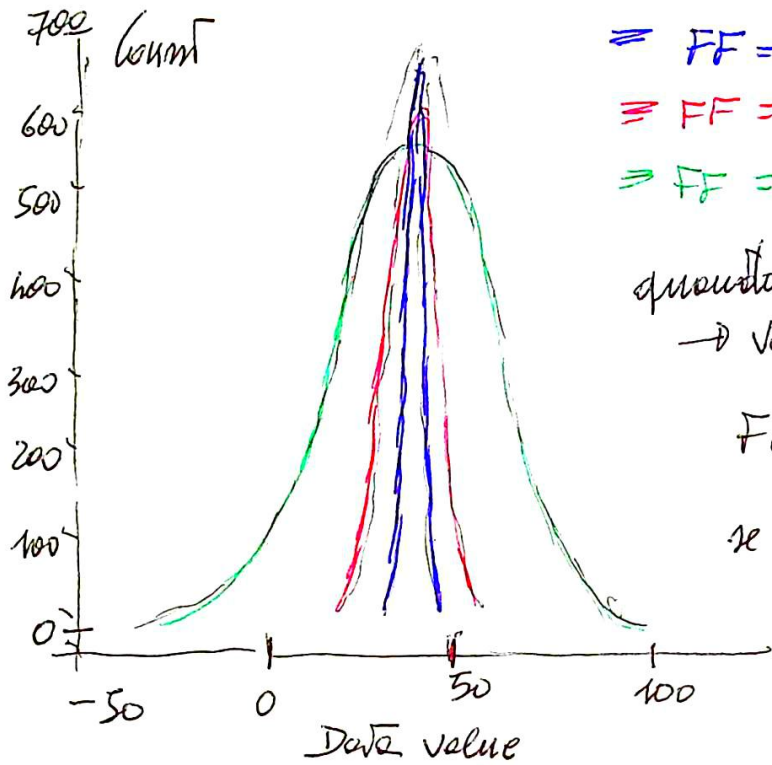
Fano factor and CV

solo x deviazioni > ϕ

$$\text{FF} = \frac{s^2}{\mu} \quad (4.12)$$

$$\text{CV} = \frac{s}{\mu} \quad (4.13)$$

vedi grafico seguente



$\Rightarrow FF = 0,10$

$\Rightarrow FF = 1,02$

$\Rightarrow FF = 9,97$

quando questo parametro e' piccolo \rightarrow
 \rightarrow varienza piccola rispetto alla media

$FF = \frac{s^2}{\bar{x}}$; $CV = \frac{s}{\bar{x}}$;

se aumentato \rightarrow la varienza diventa
 maggiore della media

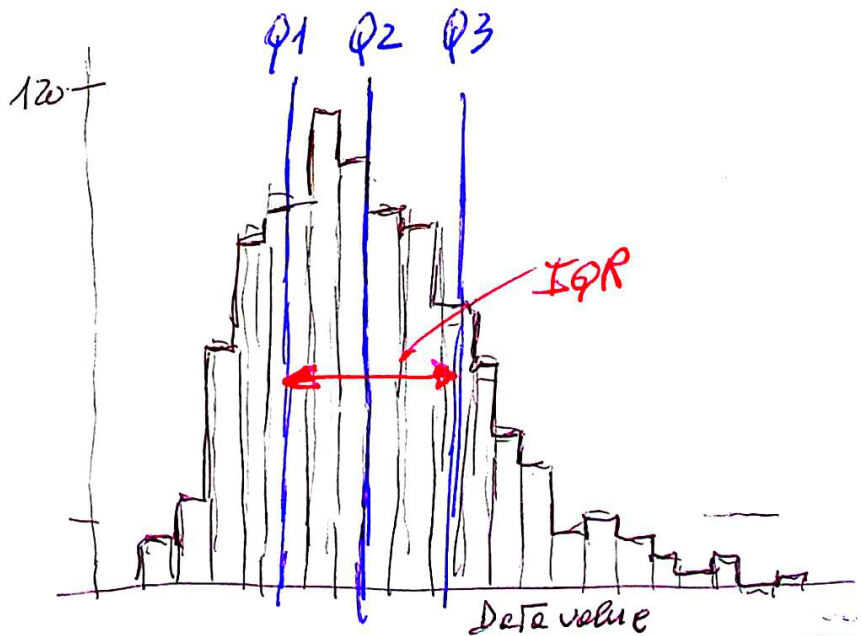
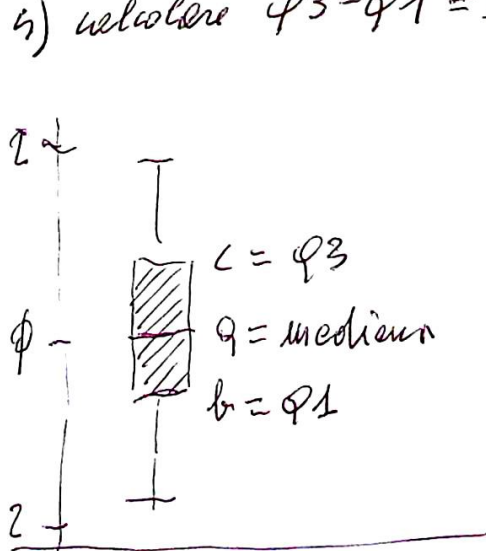
Fano e' nelle unita' dei campioni — CV e' adimensionale

MIS - Interquartile range (IQR)

IQR = distanza numerica tra 25% e 75% dei dati

\rightarrow può essere calcolato con i seguenti passi:

- 1) calcolare la mediana "quartile 2" o "Q2"
- 2) _____ del subset < Q2 (a sinistra) "Q1"
- 3) _____ > Q2 (a destra) "Q3"
- 4) calcolare $Q3 - Q1 = IQR$



questi confini sono anche detti "p25", "p75" dove p = percentile

Se IQR piccolo \rightarrow distribuzione stretta

IQR è una misura non parametrica della variabilità (basata su mediana) \rightarrow insieme agli outliers

QQ plots $q =$ quantile

Les: IQR diviso in h quantili:

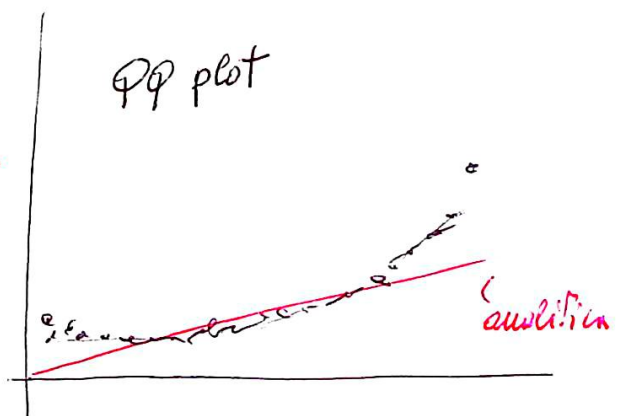
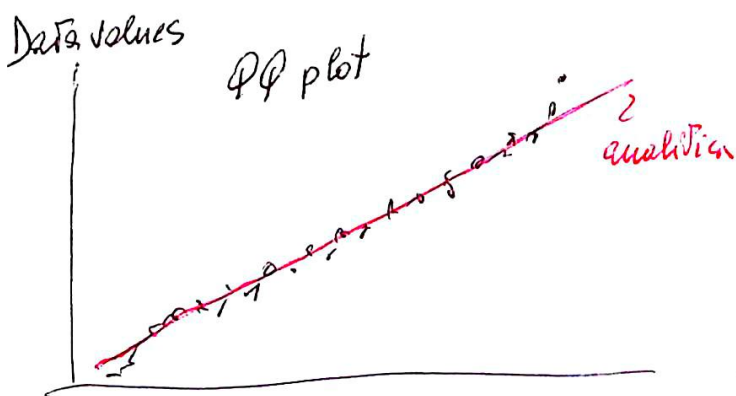
- QQ mostra la relazione tra i dati empirici e con normale
- QQ mostra se possibile ANOVA, se possibile regressione
- QQ aiuta a capire se ci sono problemi con i dataset

- abbiamo due dataset e provengono dal campionamento di una Gamma?

Le loro forme integrano un sottopopolo Gamma

Una QQ ce lo dice meglio

- su x metto Gamma empirica
- su y — dati empirici



4.7 - Statistical "moments"

momenti = numeri che definiscono lo "shape"

$\in 1^o, 2^o, \dots$

Momento non standardizzato $m_k = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^k$ (4.14)

il valore del momento non standardizzato dipende dalle unità di misura \rightarrow standardizzato :

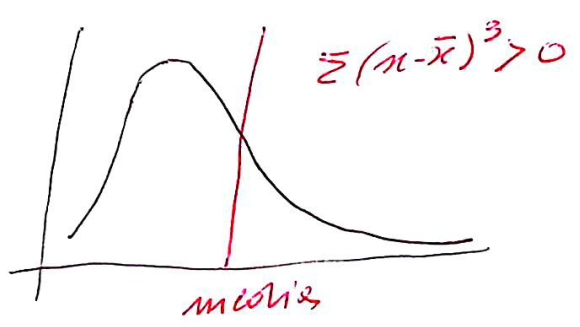
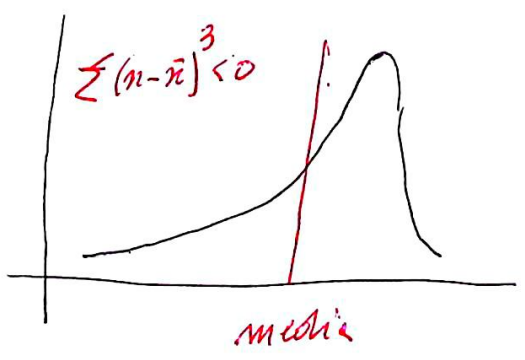
$$m_k = \frac{1}{N \sigma^k} \sum_{i=1}^N (X_i - \bar{X})^k \quad (4.15)$$

N° momento	Nome	Descrizione	Formula
Primo	Mean	Average	\circ
Secondo	Variance	Dispersion	
Terzo	Skew	Asymmetry	
Quarto	Kurtosis	Tail fatness	

First moment = media = central tendency of the distribution
= centro di massa dei dati

Second moment = variance = dispersione attorno alla media

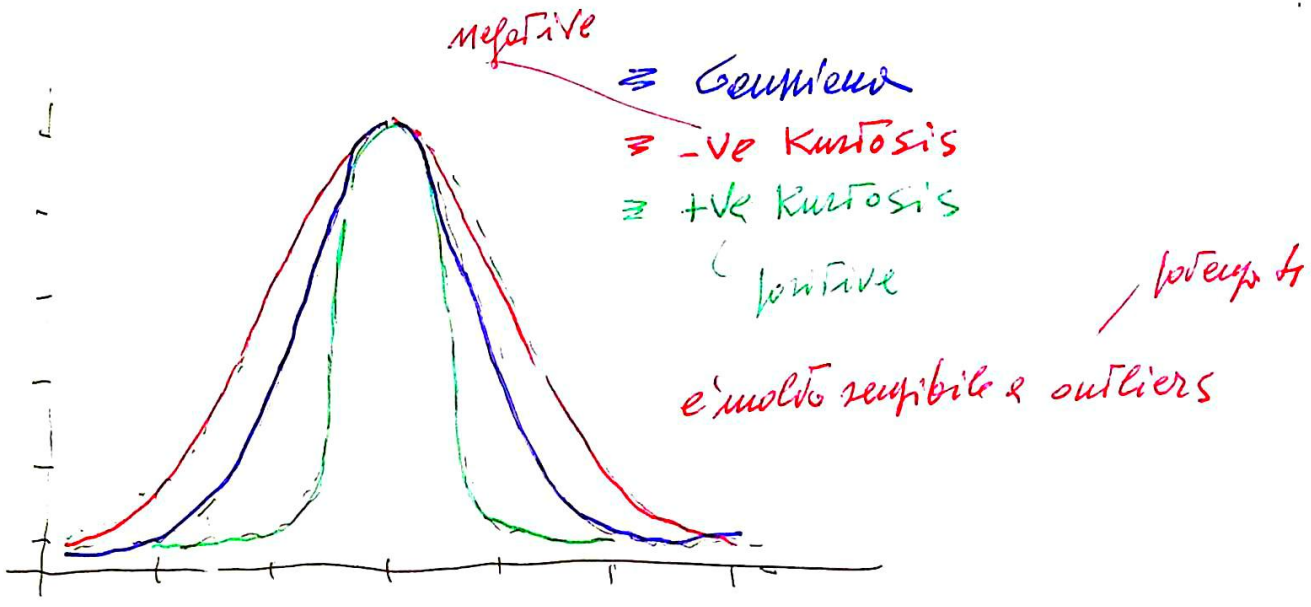
Third moment = skew - momento di 3° ordine in modo standardizzato = asimmetria della varianza rispetto alla media



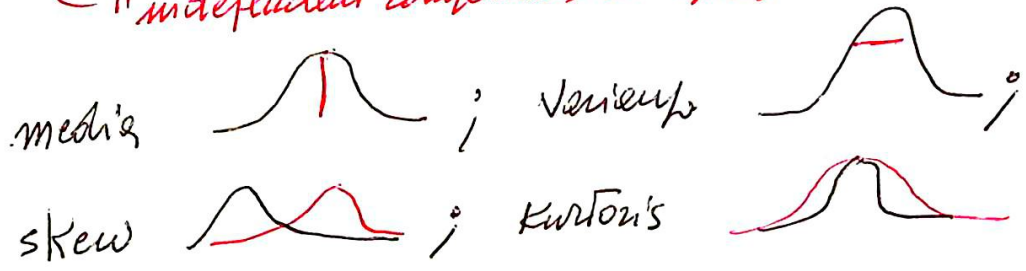
Fourth moment = Kurtosis stats. Kurtosis (x)

Kurtosis = $\frac{\mu_4}{\sigma^4}$ - momento di 4° ordine Gaussian Kurtosis = 3

L = allungamento di una distribuzione attorno al suo valore medio



usate x stimare il rischio nel mondo finanziario
 usare il numero, mediante una tecnica denominata
 "independent components analysis"



H.8 - Histograms part 2: Number of bins p. 139

Relazione tra "bin count" (variabile K) e larghezza del bin (w) -

$$K = \left\lceil \frac{\max(x) - \min(x)}{w} \right\rceil \quad (H.17)$$

rounded up

ci sono vari metodi x determinare n° bins

Method	Formula	Advantage
Arbitrario	$K = h \cdot \phi$	simple
Sturges	$K = \lceil \log_2(N) \rceil + 1$	depends on count
Friedman-Discovis	$w = 2 \cdot (IQR / \sqrt{N})$	" on count & spread

= FD = F-D

Inter-quantile range
 plt. hist (data, bins = 'fd')

n'pu' avere anche un w. verificabile, ma e' di difficile interpretazione.

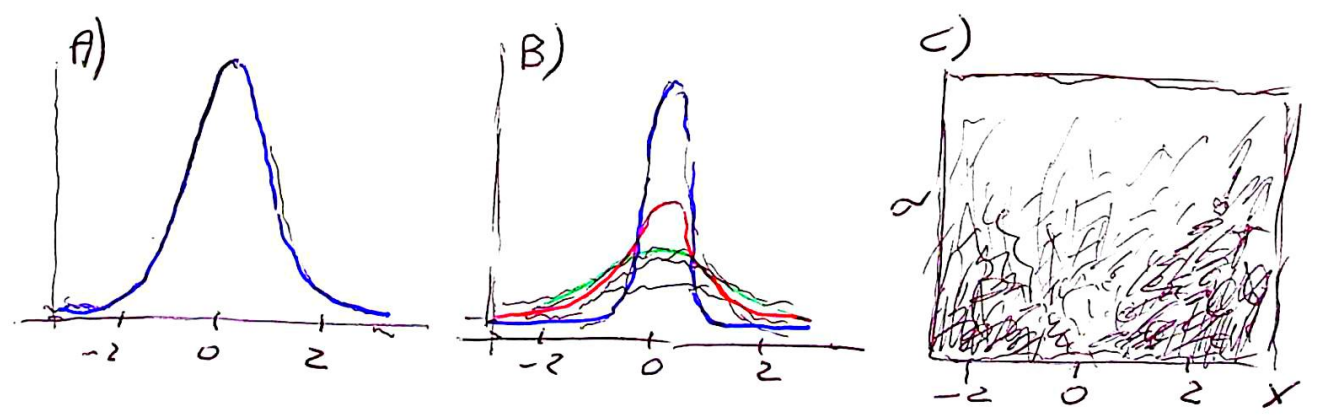
Other descriptive stats

∃ molte altre distribuzioni (es. Hurst exponent)
 Le serie temporali hanno "spettro" e "autocorrelazione"
 I set di dati multivariati → covarianza
 Le matrici → rango, "condition number", "singular value spectrum"
 - in ogni caso - anche a quante - n' applicano i concetti esposti in questo

4.9 - Esercizi

La Gaussiana (o.k.d. distribuzione normale, curva a campana) e' la piu' importante.
 - crea una famiglia di gaussiane in una matrice (dimensione x sigma),
 es. 50 gaussiane con $0.1 \leq \sigma \leq 3$
 - crea una immagine come in fig. 4.29C

Fig. 4.29



```

plt.imshow(G, extent = [x[0], x[-1], sigma[0], sigma[-1]],
           cmap = 'gray', aspect = 'auto', origin = 'lower', vmin = 0, vmax = 255)
N = 50
x = mp.linspace(-3, 3, 111)
sigma = mp.linspace(0.1, 3, N)
G = mp.zeros(N, len(x))
for i in range(N):
    s = 1/(sigma[i]*mp.sqrt(2*mp.pi))
    eTerm = x**2/(2*sigma[i]**2)
    G[i,:] = mp.exp(-eTerm)
  
```


L'integrale delle densità (qui sono normali) = 1, noi facciamo l'integrale solo in un sottointervento \rightarrow integrale < 1

- per un risultato migliore \rightarrow estendere x

- se siamo vicini al fronte moltiplicativo \rightarrow il picco = 1, non finì l'integrale

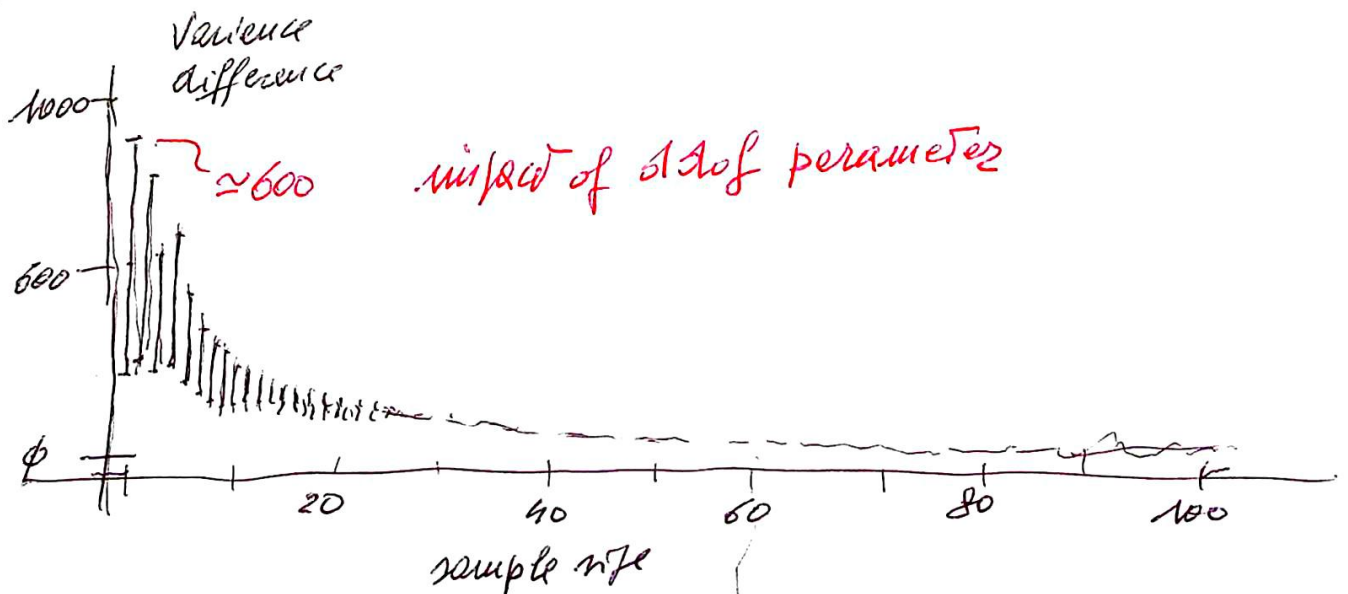
Scriviamo la media senza usare "numpy" - calcolare 3 statistiche

"media", "mediana", "varianza"

Perché dividere per "N-1"? due differenze?

- numeri casuali tra -100 e 100, calcoliamo varianza due volte (una con dof=1, l'altra con dof=1 - *ne calcoliamo la differenza in valore assoluto* - ripetiamo il calcolo del n° di campioni -

- mando numeri casuali \rightarrow poniamo ripetere esperimento e fare la media - *vediamo una bella errore* - mando 25 runs



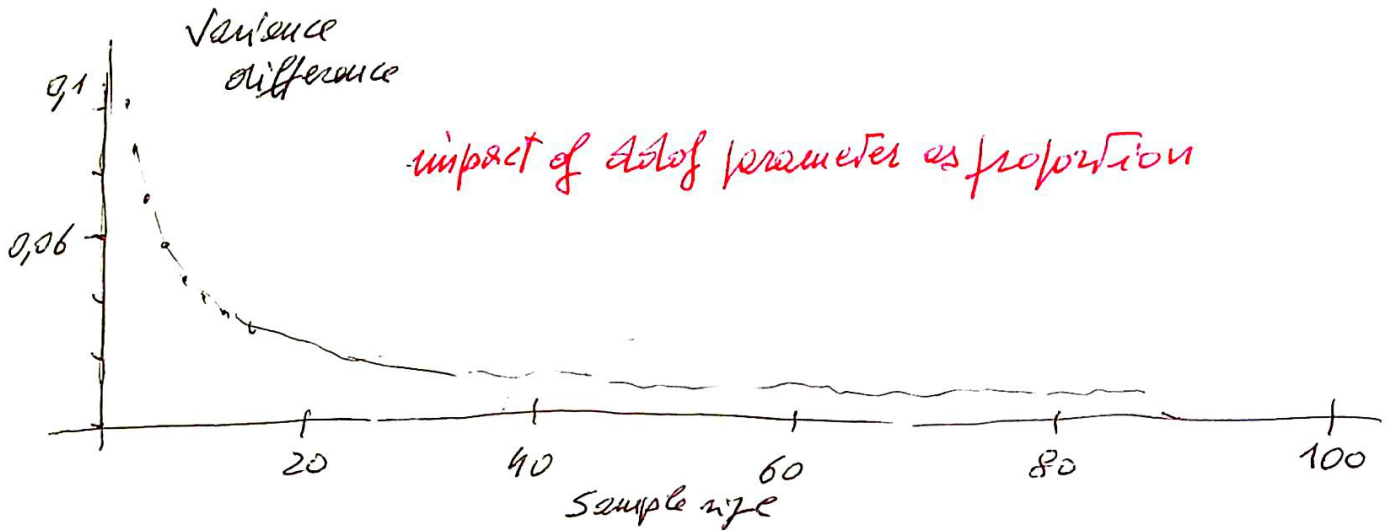
600 è un errore molto grande? \rightarrow provo con valori $-10 \div 10$ e dopo due questo valore \downarrow \rightarrow difficile valutare

Questo è l'effetto del valore dei dati \rightarrow normalizziamo e risolviamo.

Calcoliamo nuova valore di media

$$\frac{V_1 - V_0}{V_1 + V_0}$$

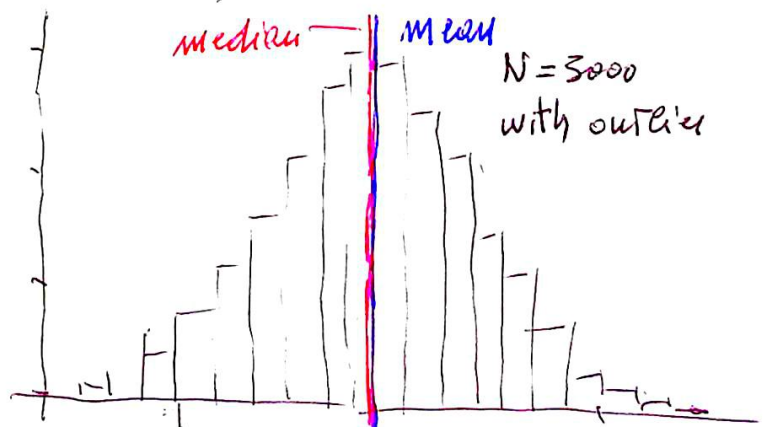
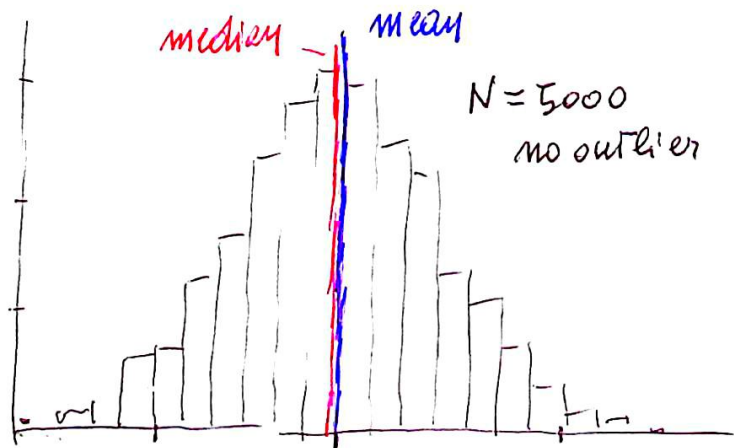
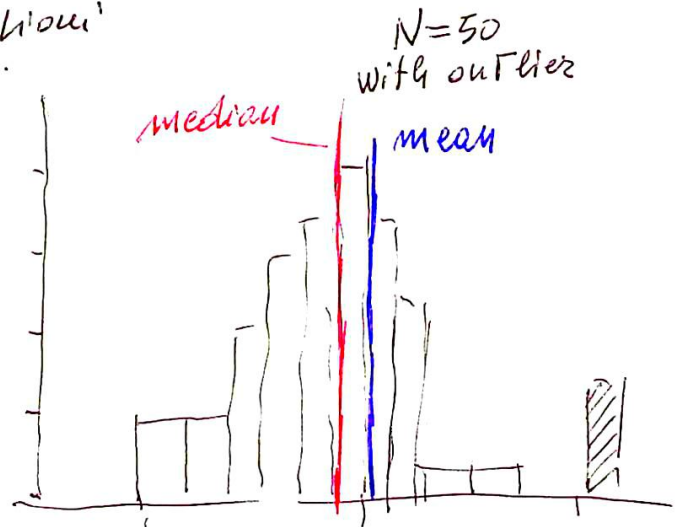
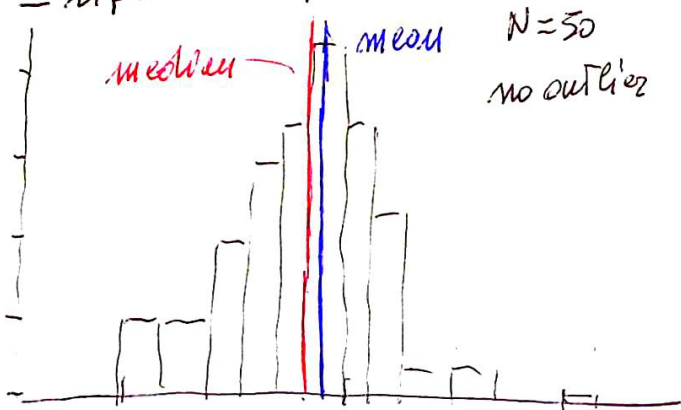
varianza con dof=1



il miglioramento è dato da un fattore $1/(2N-1)$ → non dipende dal valore dei dati

Esploriamo l'impatto di un singolo outlier sulla media, rispetto alle mediane (x campioni 'veri' e pochi) mediane

- creiamo 50 numeri random da 0 a 100, vogliamo media
- visualizziamo su istogramma
- creiamo un set anomalo sostituendo il valore maggiore x 4
- calcoliamo media, mediana
- ripetiamo procedura x 5000 campioni



$N=50$ mean increased of 0.53 | $N=5000$ ————— 0.02
 median ————— 0.00 | ————— 0.00

'most'

Calcolare i momenti statistici

mean, variance, skew, kurtosis = stats.moments (loc=1, scale=2, moments=)

Average = 1
 Variance = 4
 Skewness = 0
 Kurtosis = 0

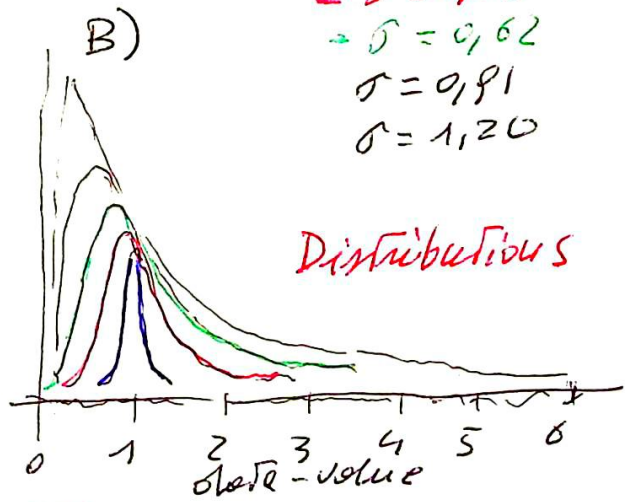
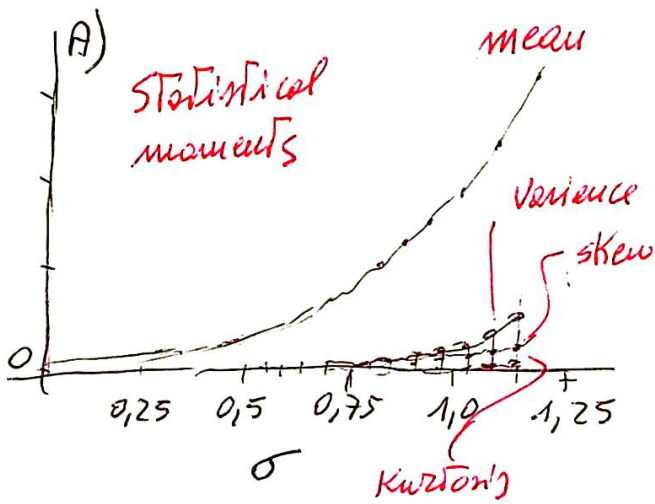
uniform
 log normal
 exp normal → vedi Python x i
 diversi parametri

Come variano i momenti cambiando la distribuzione?

scipy.stats.skew()

$N = 13'524$ casuali, gamma → trasformare in 20 log-normal,
 il gamma definito come $\exp(X\sigma)$ con $1 \leq \sigma \leq 20$ in 20 steps

$\sigma = 0.10$
 $\sigma = 0.33$
 $\sigma = 0.62$
 $\sigma = 0.91$
 $\sigma = 1.20$

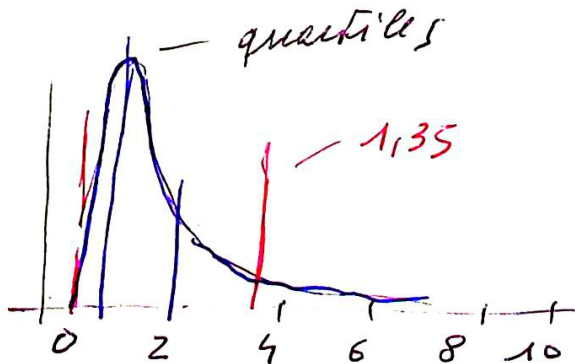
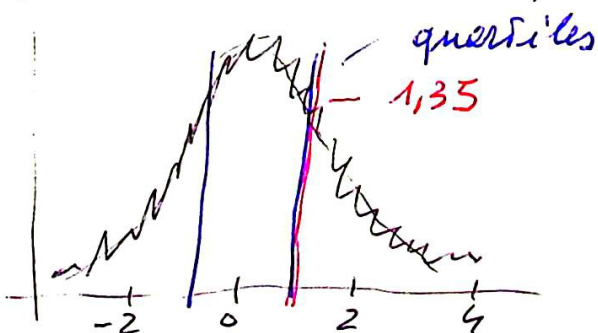


IQR simile a deviazione std
 errore & media

$x_{gamma} : IQR = 1.350$

$N=10'000$ — verificare // con chi è e^x , quanto vale IQR?

L'hoza del PC: IQR = 1.360; IQR = $1.481 \times e^x$



calcolare FWHM empirico x Geant

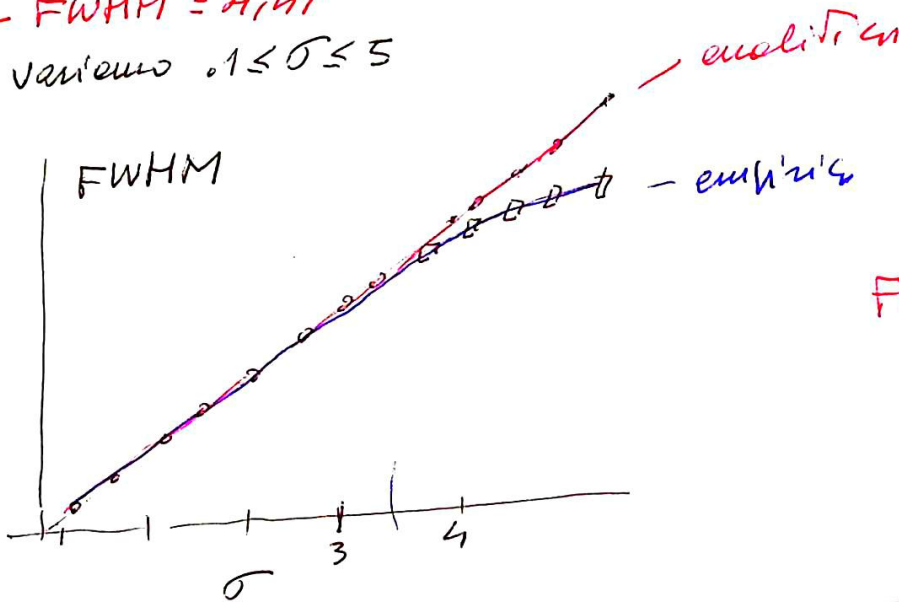
- definire la funzione emp FWHM(x,y)
- normalizzo i dati $\tilde{y} = (y - y_{min}(y)) / (max(x(y)) - min(x(y)))$
- trovo il picco come indice

L'evento normalizzato devo cercare attorno a 0.5 - pre-peak - post-peak = 4,46

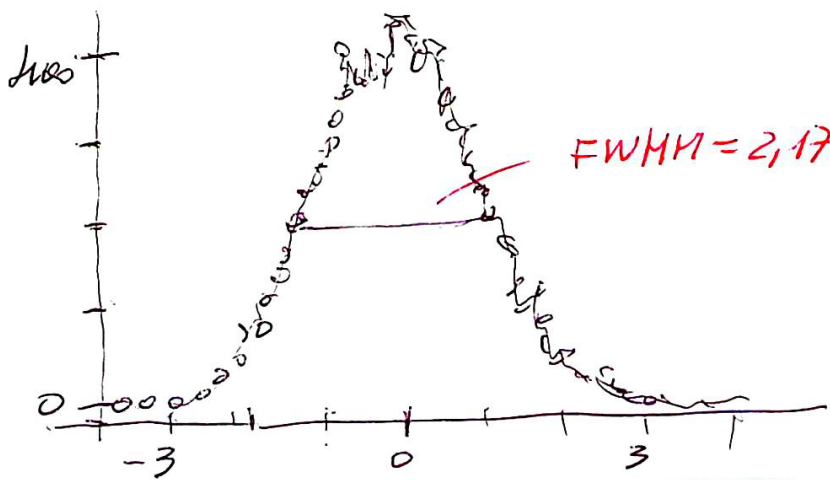
- calcolo FWHM - distanza tra pre e post
- calcolare FWHM x Geant formula $\sigma = 1,9, -8 \leq x \leq +8, N = 1001$

L FWHM = 4,47

- variando $0.1 \leq \sigma \leq 5$

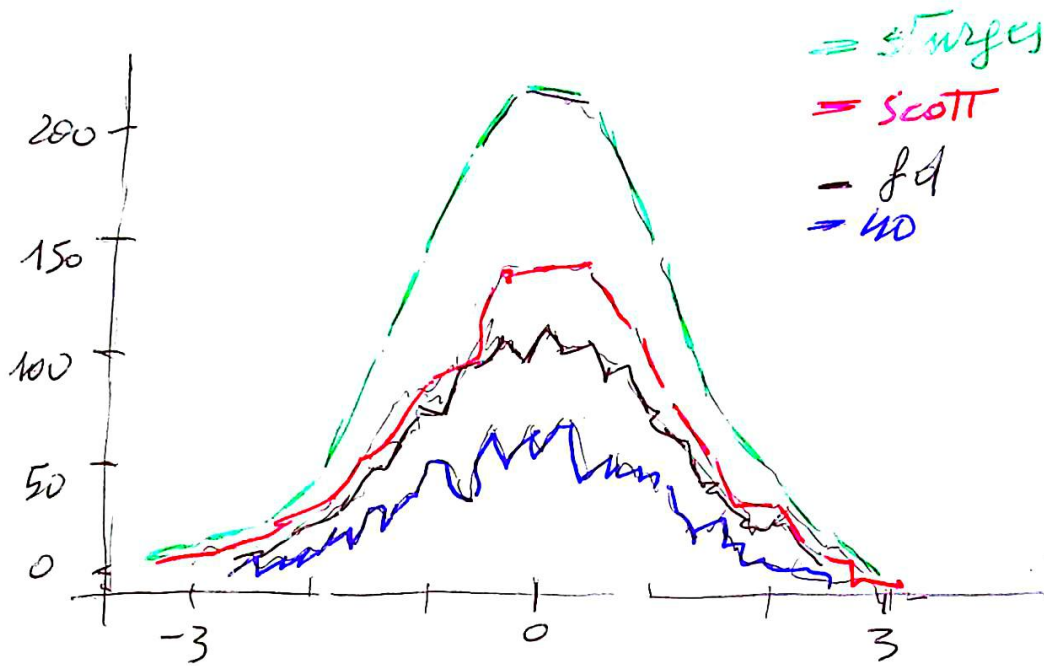


N = 12345 Geant, bins = 100, disegnano FWHM

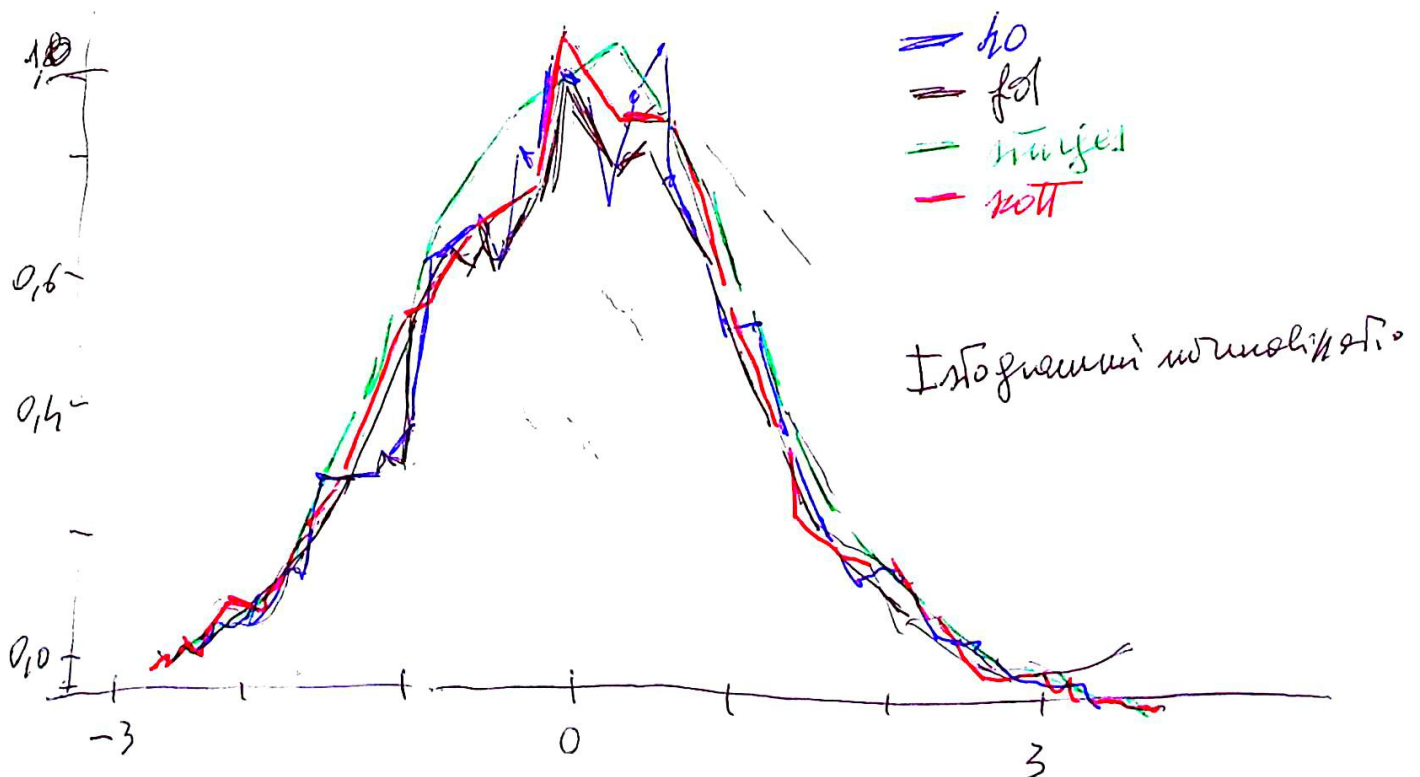


Influenza delle varie regole di calcolo di bin su istogramma

N = 1000 Geant



i valori sono elemento diversi, ma gli esperimenti sono simili
 → non impetisco qualitativamente sui risultati - p. 156
 variando il n° dei bins, varia la quantità di dati inclusi
 L → etepe differenti
 L → n' fu' normalizzare plotando $y / \max(y)$



— Fine ESERCIZIO 4 —

5 - Simulating data

19

5.1 - Why simulate data?

- Validare i metodi di analisi
 - ↳ È una gamma di metodologie, ciascuna con i suoi parametri che possono influenzare il risultato. Simulando si possono confrontare metodologie diverse. Questo è particolarmente utile quando i dati simulati hanno caratteristiche simili a quelli reali.
- Capire vantaggi e limiti dei vari metodi - consente di gestire l'influenza del rumore - con non fattibilità con i dati reali.
- Consente di capire in modo più approfondito come funzionano i metodi dei segnali.
- Capire meglio i dati di cui si dispone.
- Pensare in modo + critico e obiettivo ai dati - come selezionare il metodo di analisi.
- Usare un approccio "proattivo".
- Statistica computazionale (o empirica). Richiede significatività statistica, intervalli di confidenza - Questo diventa necessario se i dati violano le assunzioni richieste dalla statistica parametrica.
- Molte diverse simulazioni richiedono pensiero critico.
- ---

migliorano le capacità di scrittura e codifica.
- Consente molte prove senza necessità di esperimenti.

5.2- Random data from distributions

- meglio partire da istogrammi

5.2.1- Normally distributed random data

$$X \sim \mathcal{N}(\mu, \sigma^2) \quad (5.1)$$

nota: *tilde*

$\mathcal{N}(\phi, 1)$ $\mu = \phi; \sigma^2 = 1$ = normal std distribution

• σ^2 allarga o restringe la distribuzione

• μ = media

5.2.2- Uniformly distributed data

$$X \sim \mathcal{U}(a, b) \quad (5.2)$$

$\mathcal{U}(\phi, 1)$ = std uniform distribution

se scrivo $Y = 2\pi X - \pi$ (5.4)

↳ quale sarà il suo intervallo? $[-\pi, \pi]$

∀ a, b con $a < b \Rightarrow Y = a + (b-a)U$ (5.5)

dove $U \sim \mathcal{U}(\phi, 1)$ (5.6)

in questa distribuzione

$$\mu = \frac{a+b}{2} \quad \sigma^2 = \frac{(b-a)^2}{12} \quad (5.7)$$

$$(5.8)$$

vale

$$Y = \mu + \sqrt{3}\sigma (2U - 1) \quad (5.9)$$

mp. random. uniform $(a, b, \text{size} = N)$

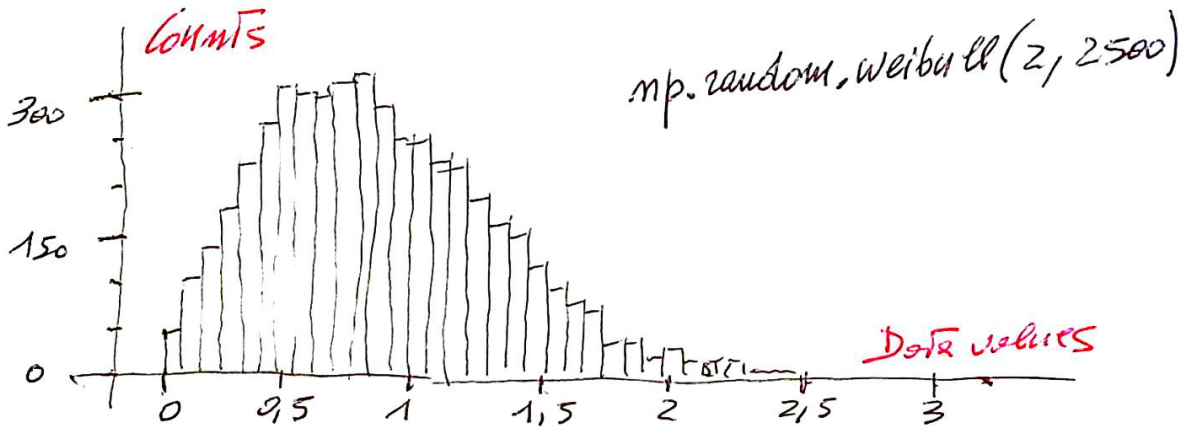
5.2.3 - Random data from other distributions

20

- dati casuali da una distribuzione di Weibull

`mp.random.weibull` in formula

$$Y = \exp(X^\sigma + \mu) \quad (5.11)$$



5.2.4 - Random integers

`mp.random.randint(1, 5, size=10000)` / esempio
`mp.random.lognormal(1, 0.5, 5000)`

5.3 - Random elements of a set

es. (1, 2, 3, 6, 7, 8)

$S = [1, 2, \text{mp.pi}, 10]$

`mp.choice(S, 1) = 2`

$t = ["a", "b", "hello"]$

`mp.random.choice(t, 1)`

es. `mp.random.choice(S, 4)`

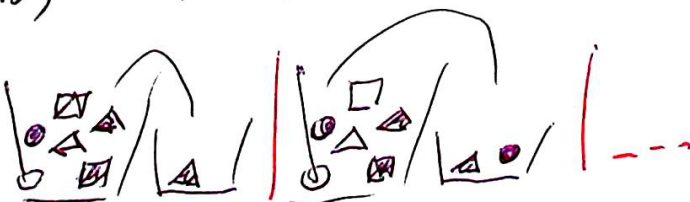
↳ array([10, 1, 1, 3.14])

A) Sampling without replacement



`mp.random.choice(S, 4)`

B) Sampling with replacement



`(S, 4, replace=True)`

NOTA con replacement = creiamo un nuovo dataset più grande,
 perché gli elementi di partenza possono essere selezionati più volte
 ↳ la statistica del nuovo insieme sarà diversa
senza replacement = produce un dataset identico, tranne l'ordine

5.4 - Random permutations

↳ è un modo di distribuire casualmente gli elementi.

`l = np.arange(5)`

`print(l)`

`print(np.random.permutation(l))`

⇒

`[0 1 2 3 4]`

`[2 0 1 3 4]`

↳ non ordinare (sort) gli elementi

`theData = np.arange(-3, 4) ** 3`

`newIdx = np.random.permutation(len(theData))`

⇒

`shufData = theData[newIdx]`

`[-27 -8 -1 0 1 8 27]` theData

`[3 4 1 6 2 0 5]` newIdx

`[0 1 -8 27 -1 -27 8]` shufData

- randomizzazione di dati "scorrelati" es. altezza e peso
 - rimescolo altezza, ma non peso → non più correlate
- calcolare la correlazione tra due set di dati ci può dire quanto sono indipendenti

5.5 - Reproducing randomness

- genero due volte dati casuali: mi serve due nuovi uguali?
 n' quando sto facendo, in modo da avere risultati confrontabili

`np.random.randn(3,3)` - qui le matrici sempre differenti

ma se mi do

rs = np.random.RandomState(17)

rs = randint(3, 3) — qui ottengo lo stesso risultato

perche' Python ha salvato il valore di "seed"

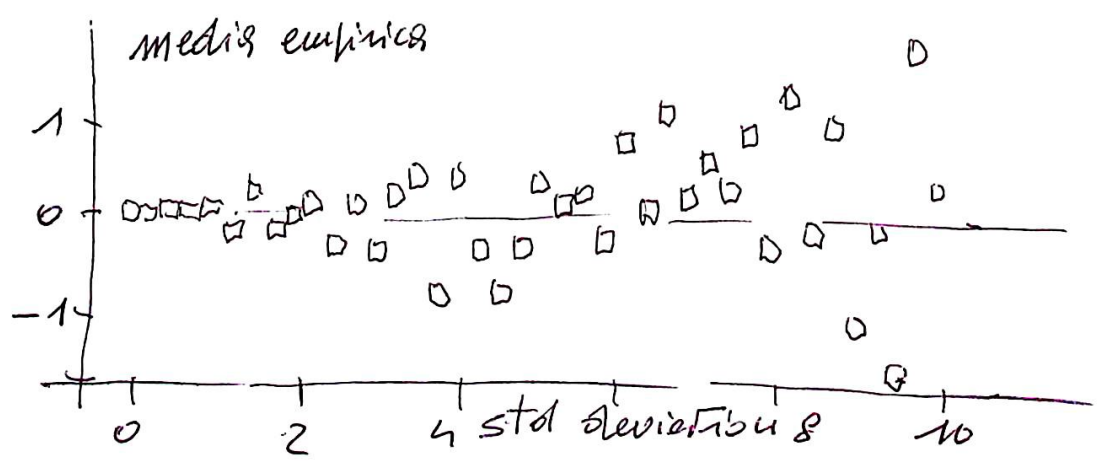
usare con cautela

numeri casuali e' buona informazione sulla "robustezza" di alcune ipotesi o algoritmi.

5.6 - Running experiments with random numbers

- genero i valori di una popolazione casuale, Gauss
- faccio la media
- o mi posso la domanda: in che modo σ^2 influenza sulla correlazione con un'altra variabile?
- decidiamo come varia σ^2 e il grado di correlazione
- Les. ANOVA non robusti rispetto al valore di σ^2
- es. impatto di variazione σ^2 sulla media x std. Gauss
- in teoria non ha nessun impatto

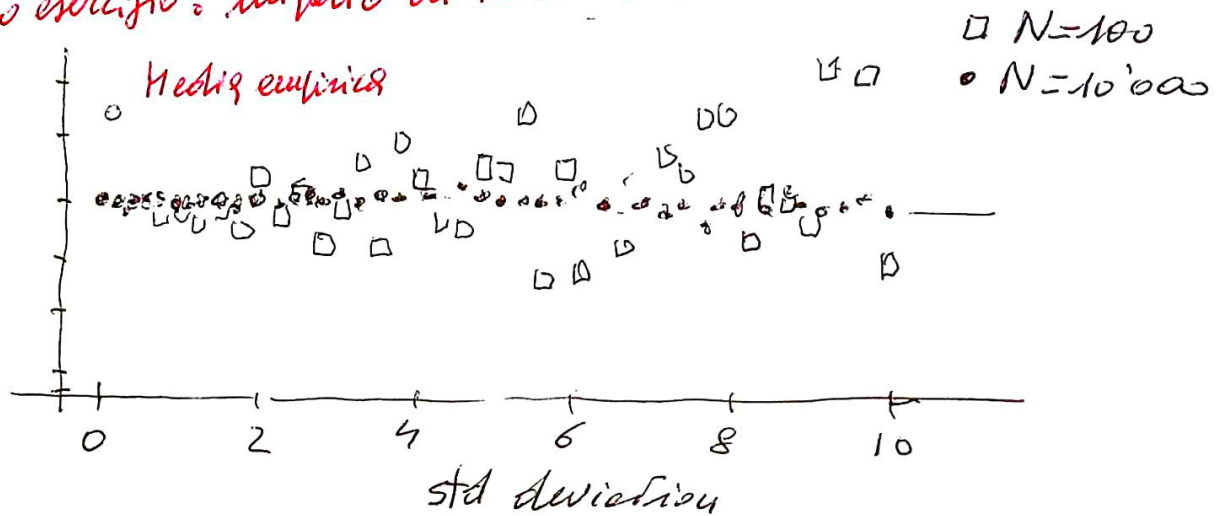
In concreto: N=100
 in un ciclo for faccio venire σ^2 , calcolo i campioni Gauss,
 ne faccio la media **results**
grafico σ^2 vs results; **grafico media**; perche' / tutti i campioni



queste serie di dati empirici \rightarrow statistiche inferenziali

- poiché le variazioni della media sono simmetriche (circa) rispetto allo zero \rightarrow la deviazione non è sistematica -
- Questo effetto appare in una distribuzione log-normal

Altro esercizio: impatto di N sulla media



5.7 - The emerging world of data - simulations

truffico

- serie temporali, dati finanziari, meteo, processi bio, fisica,
- "Generative deep learning" = intelligenza inventata, video falsi
- \rightarrow esamineremo poco di questo, ma questi sono fondamenti sugli argomenti del libro

5.8 - Finding publicly available real datasets

- sono aperti come Kiplogh, sono specifici
- alcuni archivi sono free, alcuni free + restrizione
- \rightarrow su questo \neq standardizzazione
- due siti popolari specifici
 - UCI machine learning repository ⁷ // Kaggle ⁸
 - <https://archive.ics.uci.edu/ml/datasets.php>
 - <https://www.kaggle.com/datasets>

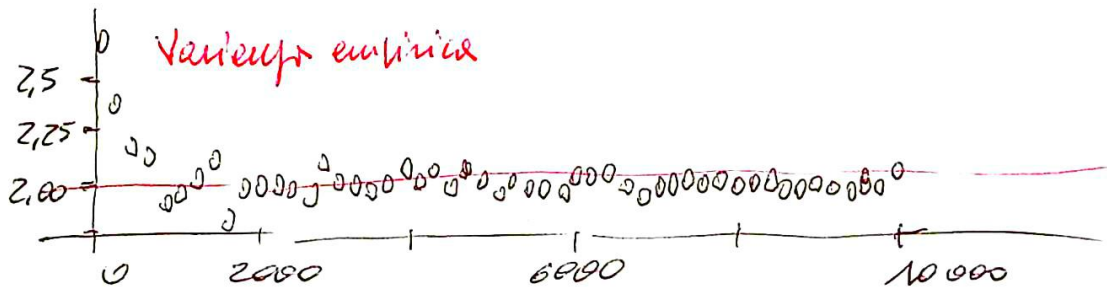
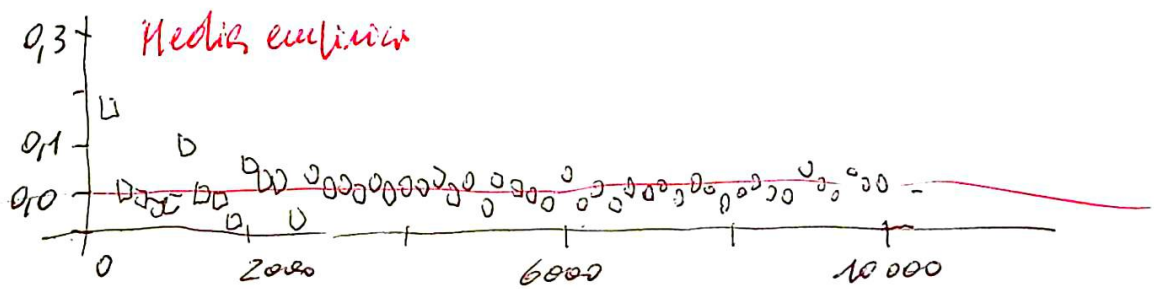
Esercizio 1 - relazione tra N e accuratezza della media empirica \bar{Z}
 e sulla varianza - $N=10$ $N(0,2)$

media empirica = $-0,184$ (intervallo ± 2)
 varianza $\hat{\sigma}^2 = 1,688$

in un loop for $N = 10 \div 10'000$ step 200

L in ogni loop neo dataset (stima media, varianza)

faccio grafico



- le dispartenze sono via + che -

- n' intravede la legge dei grandi numeri

- Teorema sul limite centrale

p. 184

grandi numeri: se $N \rightarrow \infty \Rightarrow$ media empirica $\rightarrow \mu$

n non distribuzione arbitraria

limite centrale: $\sum_{i=1}^n x_i \rightarrow$ e' una distribuzione Gaussiana

L es.

risultato n più volte su condizioni diverse presi da una popolazione di distribuzione arbitraria - ripetendo molte volte $\sum x_i \Rightarrow$ Gauss

in formule:

dato X_j (famiglia di n) aleatorie indipendenti, identicamente distribuite \Rightarrow

$n; X_j; \text{iamo } E[X_j] = \mu, \text{ Var}[X_j] = \sigma^2$
 con $j = 1, \dots, n; \phi < \sigma^2 < +\infty$ *poniamo*

LIMITE
CENTRALE

$$Y_n = \frac{\frac{1}{n} \sum_{j=1}^n X_j - \mu}{\sigma/\sqrt{n}}; \Rightarrow Y_n \xrightarrow{D} Y \sim \mathcal{N}(\phi, 1)$$

con altre descrizioni

definim X_i con media $= \mu_x$, varianza $= \sigma_x^2 \Rightarrow$ *consideriamo*
 la variabile $(\bar{X} - \mu_x) / \sigma_{\bar{X}}$ dove $\sigma_{\bar{X}}^2 = \sigma_x^2 / n$
 dove $\bar{X} = \frac{\sum X_i}{n}$
 \rightarrow se $n \rightarrow \infty = \mathcal{O} \sim \mathcal{N}(\phi, 1)$

Esercizio 2 *confermeremo empiricamente*

$$\mu = \frac{a+b}{2}$$

$$\sigma^2 = \frac{(b-a)^2}{12}$$

$Y = \text{mp. random. uniform } (a, b, \text{ size} = 132K)$

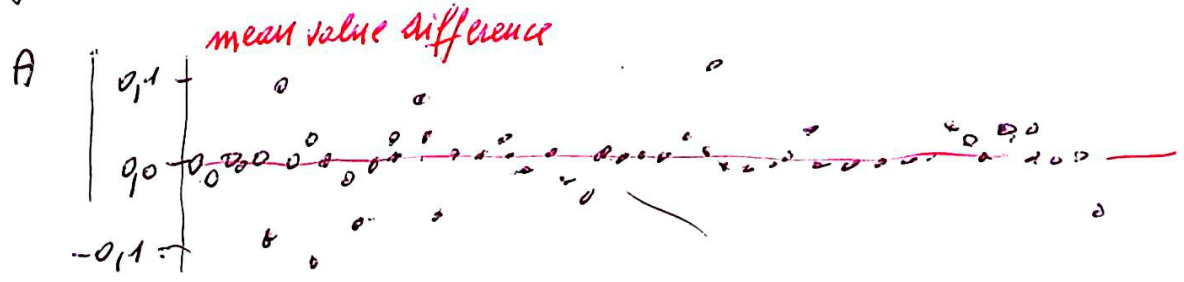
mean Diff = mp.mean(Y) - $\frac{(a+b)}{2}$ $\rightarrow -0,040$

var Diff = (mp.var(Y, dof=1) - $(b-a)^2/12$) $\rightarrow 0,000$

ora invece $N_5 = \text{mp. sample } (10, 10200, \text{ step} = 200)$

in un ciclo "for" randomico i, N in N_5

genero ed ogni passo una serie di dati, ne calcolo media e varianza
 disegno la discrepanza



il pannello B mostra valori solo $> \phi \rightarrow$ perché $\sigma^2 \Rightarrow$

non si hanno bias sistematici

• le regressioni e le ANOVA si basano su disubsempo quadratico

Esercizio 3

abbiamo visto come da dati non gaussiani si possa trovare una

gaussiana - partiamo da $Y = \exp(X\sigma + \mu)$ con $\mu = 2, \sigma = 1,5$

- calcoliamo la media

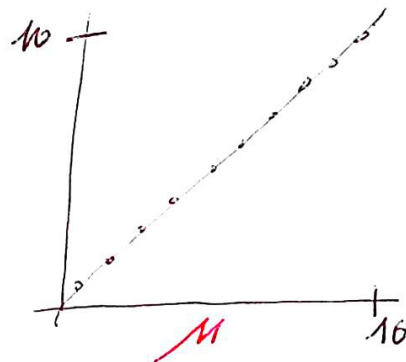
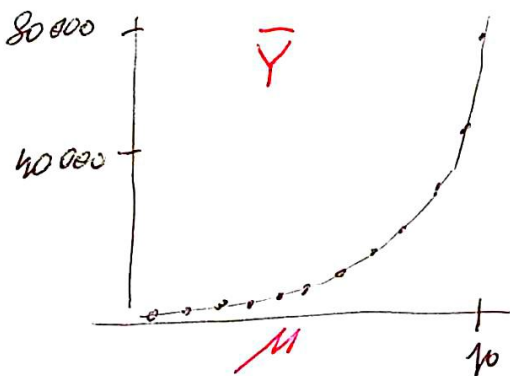
media attesa $\bar{Y} = e^{\frac{\mu + \sigma^2}{2}}$ (5.13)

- confermeremo numericamente che $\mu \approx 2$

- ciclo "for" - ricalcoliamo μ ad ogni passo (13 valori di μ tra 1 e 10)

- compareremo μ e \bar{Y}

$\ln(\bar{Y}) - \frac{1}{2}\sigma^2$



Esercizio 4 - relazione tra media e deviazione std di una

distribuzione uniforme nell'intervallo (a, b)

$Y = \text{mp. random. uniform}(a, b, n \times N)$, $N = 1001$

calcolo le medie con due formule \rightarrow
 / empiriche \rightarrow 5,0161
 / 4,9839

calcolo le medie formali \rightarrow 4,9895
 / 5,0000 media delle medie

calcolo della deviazione std con due formule

$\sigma_{-a} = 0,5713$
 $\sigma_{-b} = 0,5713$
 $\text{std}(Y) = 0,5867$

$Y = \mu + \sqrt{3}\sigma(ZU - 1)$
 $U \sim U(0,1)$

creiamo un dataset con data media e deviazione:

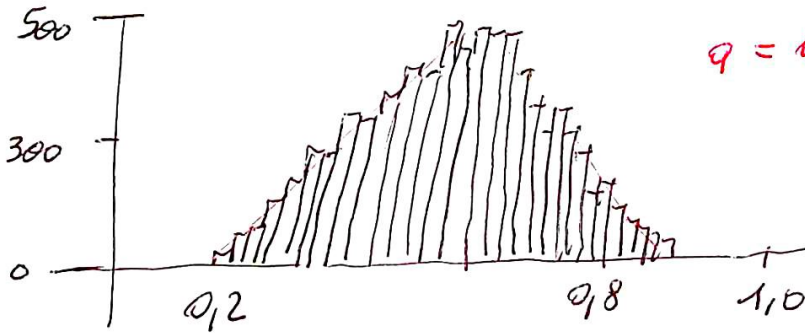
$\mu = 3,5, \sigma = 1,8, N = 100000$,
 calcoliamo media \rightarrow mean = 3,5
 deviazione teorica \rightarrow std = 1,798

Esercizio 5 - quale mediana e moda x distrib. uniforme (a, b) ?

Esercizio 6 - distribuzione triangolare [a, b] ∈ ℝ, c ∈ [a, b]

$$f(x) = \begin{cases} \frac{2}{b-a} \cdot \frac{x-a}{c-a} & \text{se } a \leq x < c \\ \frac{2}{b-a} & \text{se } x = c \\ \frac{2}{b-a} \cdot \frac{b-x}{b-c} & \text{se } c < x \leq b \end{cases}$$

N = 10 000



a = 0.2, b = 1.0, c = 0.6

$$F_c = (c-a)/(b-a)$$

$$U = \text{mp.random}(N)$$

$$Y[U < F_c] = \text{---}$$

$$Y[U > F_c] = \text{---}$$

plt.hist(y, 'fd',
 xlim=[a-0.2, b+0.2])

facciamo lo stem con cui valori diversi -

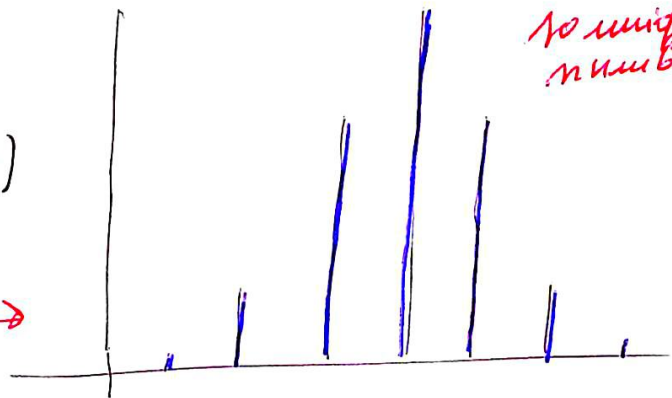
Esercizio 7 x = mp.random.normal(loc=0, scale=1, size=100000)

$$x = \text{mp.round}(x)$$

$$\text{print}(\text{mp.unique}(x))$$

$$\text{plt.hist}(x, \text{bins} = 'fd')$$

e' corretto, ma lascia a bocca asciutta



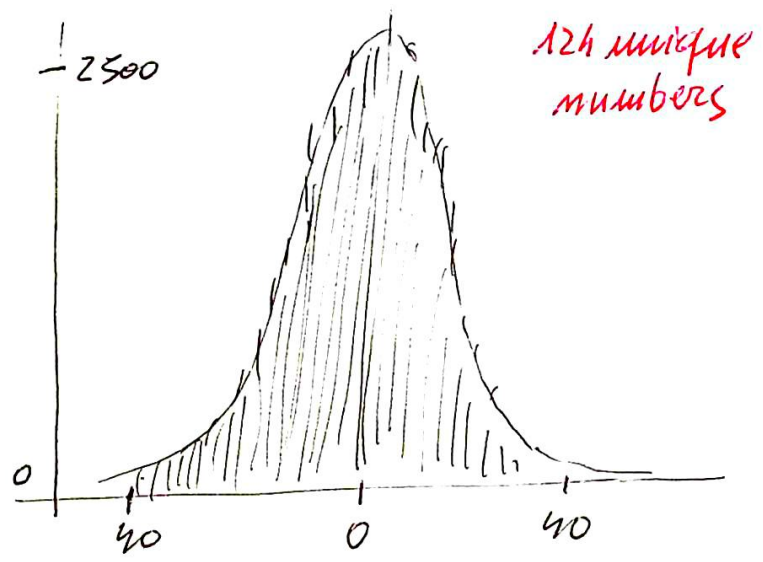
↳ vedo i valori!

$$x = \text{mp.random.normal}(loc=0, scale=15, size=100000)$$

$$x = \text{mp.round}(x)$$

$$\text{plt.hist}(x, \text{bins} = 'fd')$$





Esercizio 8 - combineremo diverse tecniche, e insieme un nuovo metodo (correlazione, permutazione)

- matrice 100×2 con numeri presi da Campiura STD
- la seconda colonna = 2ª colonna + la prima

$N = 100$

$M = \text{mp. random. random}(N, 2)$

$M[:, 1] += M[:, 0]$

calcoliamo il coefficiente di correlazione

$r_{\text{real}} = \text{mp.corrcoef}(M, T)[1, 0]$ — matrice con diagonale nulla
 / — transposta

tra $-1 \leq r \leq 1$ - se $r = 1$ correlazione perfetta
 se $r = 0$ " " " nulla $\rightarrow r = 0.6$ nel nostro caso

ora facciamo 2 "shuffle", metto in memoria di lavoro
 un paio numeri casuali interi (portione) x riorganizzare le righe
 della prima colonna senza modificare l'ordine della seconda —
 I dati non sono cambiati, è venuto l'ordine \rightarrow dev. std, medie
 non sono cambiate —

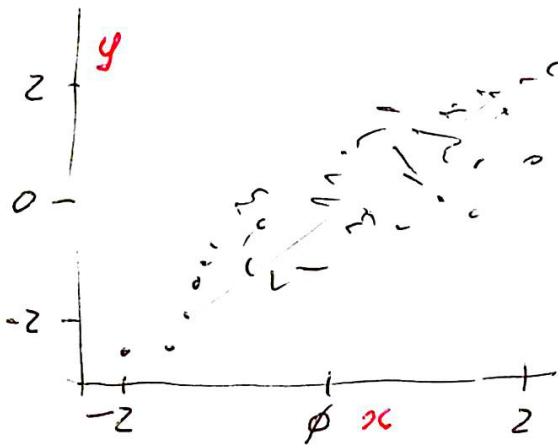
Ricalcoliamo il coefficiente di correlazione - Grafichiamo -

$\rightarrow r_{\text{real}} = 0,666$

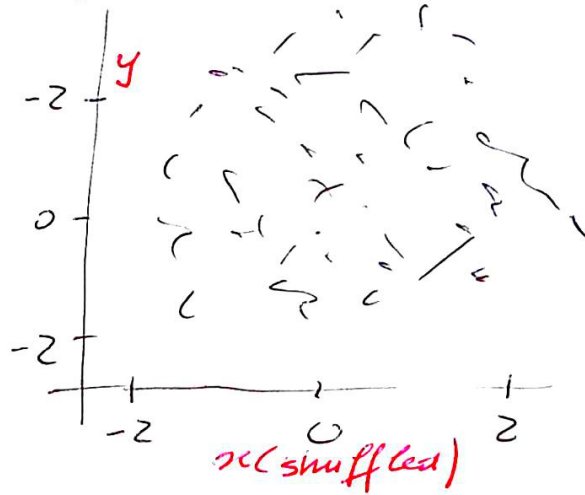
shuffled correlation = $-0,076$

Vediamo meglio con un grafico

A real correlation $r = 0,73$



B shuffled correlation $r = 0,05$



Esercizio 9

$N = 3000$

multimedia con una di queste

data = stats.expon.rvs(3, size=N)

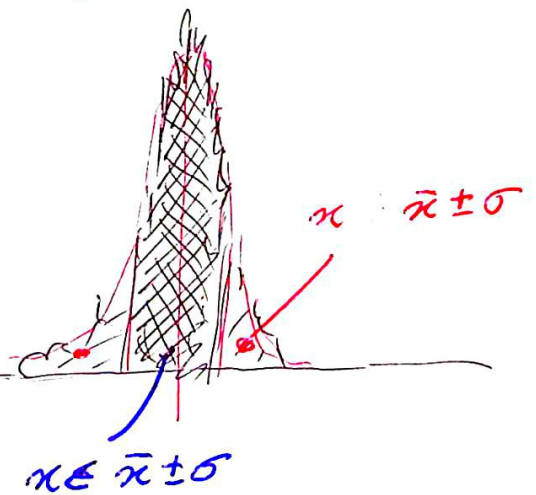
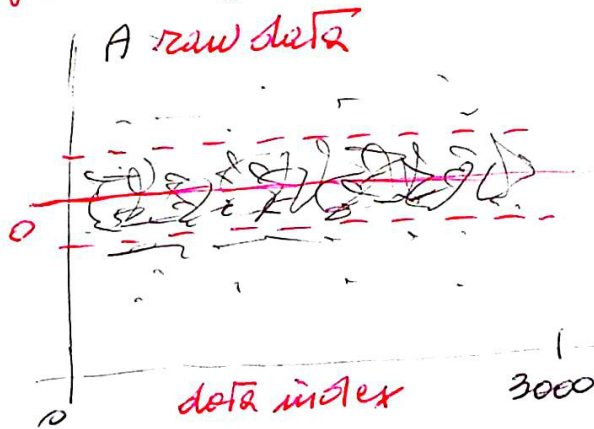
— = — . norm.rvs(loc=3, scale=1, size=N)

data = stats.laplace.rvs(size=N)

— = — . gambel.rvs(size=N)

calcolo media, teorica, dev. std teorica
facio into grammi

produce x le
vere distribuzioni



Esercizio 10 $\mu = \phi$, $N = 10000$, σ in 10 km, $\sigma = 10$

log-normal — la media e la dev. std sono una funzione non lineare di μ e σ — $X = e^N$ con $N = \log X \sim N(\mu, \sigma^2)$

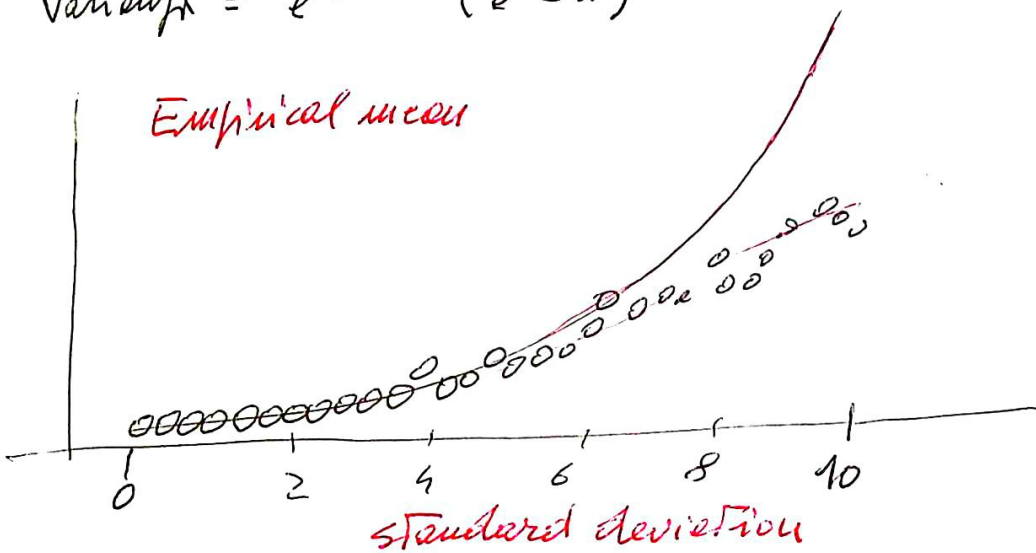
$\Rightarrow X = e^N$ segue una log normale $\log f(\mu, \sigma^2)$

$$f(x) = \frac{e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}}{x\sqrt{2\pi}\sigma} \quad \text{con } x > \phi$$

in modo formale $\mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_\phi^+$

valore atteso = $e^{\mu + \frac{\sigma^2}{2}}$

varianza = $e^{2\mu + \sigma^2} (e^{\sigma^2} - 1)$

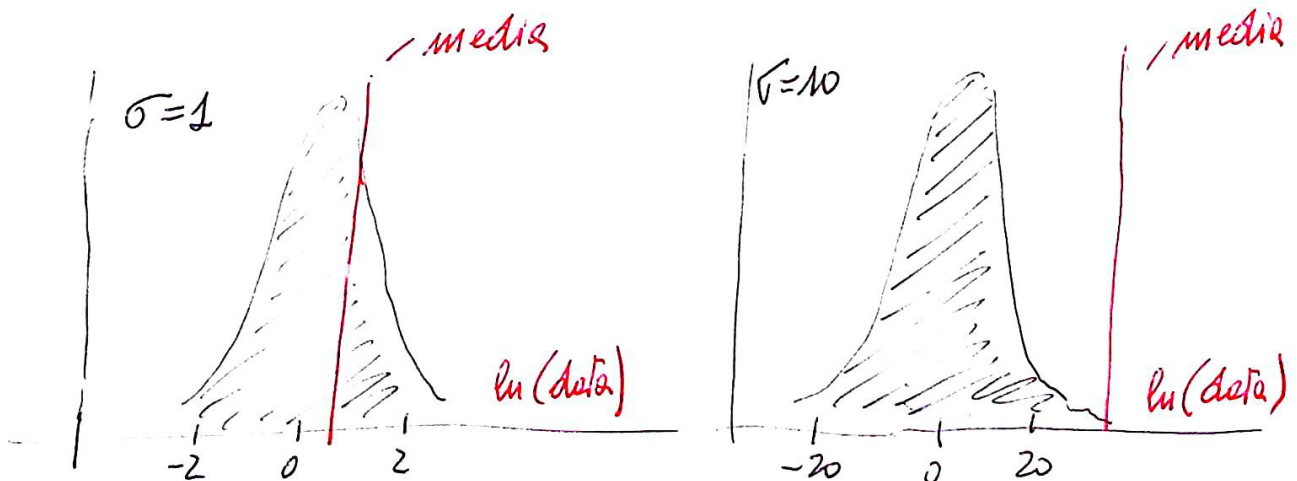


la media attesa varia in funzione della deviazione std.

• si incrementa una sottostima e partire da $\sigma_0 = 5,11$
 ↳ perché? e^x sale molto rapidamente - Inoltre σ controlla direttamente la $\ln(\text{data})$ - Questa è una situazione critica da comprendere.

Inoltre è inclusa anche la media in questi valori così grandi - Questi valori in concreto mancano nelle distribuzioni empiriche → **divergenza** -

• per illustrare il concetto vedere fig. 5.17 che mostra iogrammi di log-transformed data x due valori di σ . È imponente apparire normali, ma anche la media è log-transformata.



In realtà (concreti), empirici, la probabilità si stima con
quanti è il minimo.

• questa difficoltà si presenta anche nelle statistiche computazionali (es. quando si calcolano i p-values su shuffled data).

- per queste problematiche, riflessione, ragionevolezza.

— FINE CAPITOLO 5 —

6 - TRANSFORMATIONS

6.1 - What, why, and how of data transformation

• Cosa sono le trasformazioni di dati?

- sequenze di operazioni matematiche sui dati

• Perche' trasformare i dati?

- moltiplicare x rendere confrontabili le scale

- passare da una distribuzione \rightarrow Gauss

- normalizzare il range di variazione

- eliminare influenza dei valori troppo estremi (fuori dei 3σ)

- z-score

- x memorizzare + facilmente

• Come trasformare i dati?

- scipy.stats

- usando funzioni scritte da altri

- funzioni di nostra implementazione

• Quali tipi di trasformazioni?

- lineari, non lineari

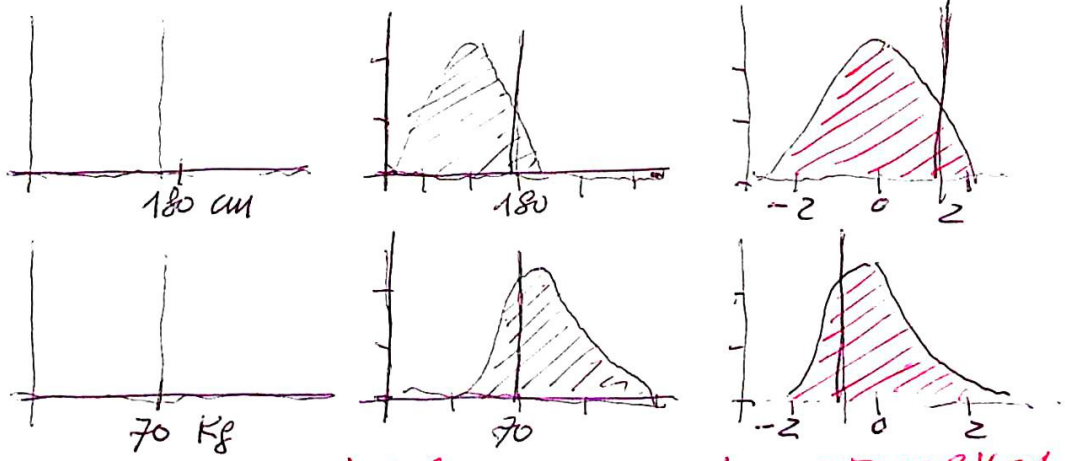
- lossless, con perdita di dati

- iterative, non iterative

altre, peso di
 A_{1-2} una persona
 B_{1-2}
 C_{1-2} distribuzioni normalizzate (misura in std deviation)

6.2 - z-score standardization

es. cambiare peso e altezza delle persone



raw values normalized

• Z-score math

- Hard & Soft assumptions -

1) shift the data $\rightarrow \bar{x} = \phi$

2) scale the data to the std deviation

$$z_i = \frac{x_i - \bar{x}}{s} \quad \text{-----} \quad (6,1)$$

Vediamo i risultati applicati a un caso numerico:

ERA	^{mean} 0,60	3,36 - std
DIVENTA	900	1,00

• Interpretazione

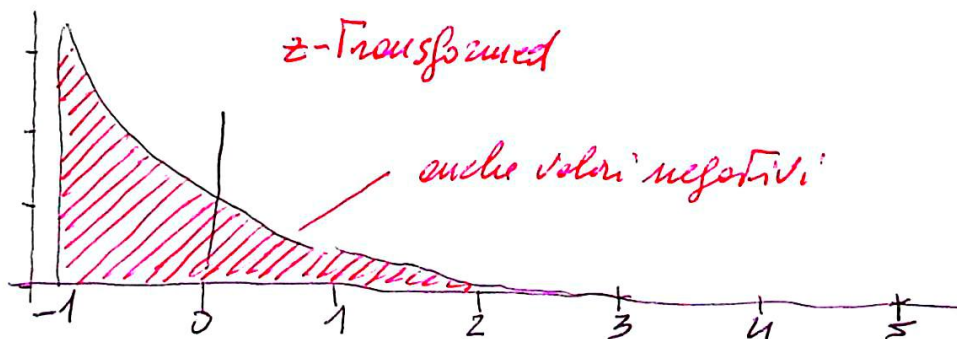
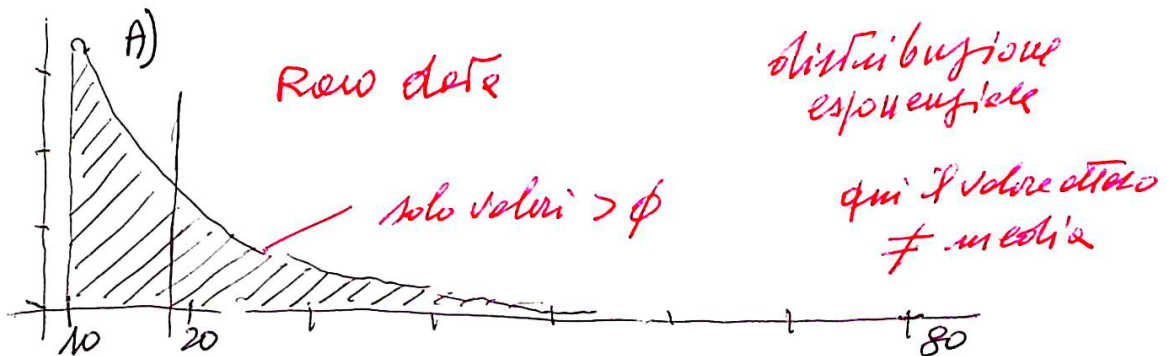
- la z-transformation sposta e modifica i dati, ma \rightarrow
 \rightarrow non viene la distribuzione

• Hard & Soft assumptions

- la std deviation $\neq \phi$ hard

- media e std hanno significato reale in questa distribuzione (fu' anche non essere vero \rightarrow soft)

L' beni applica & beni

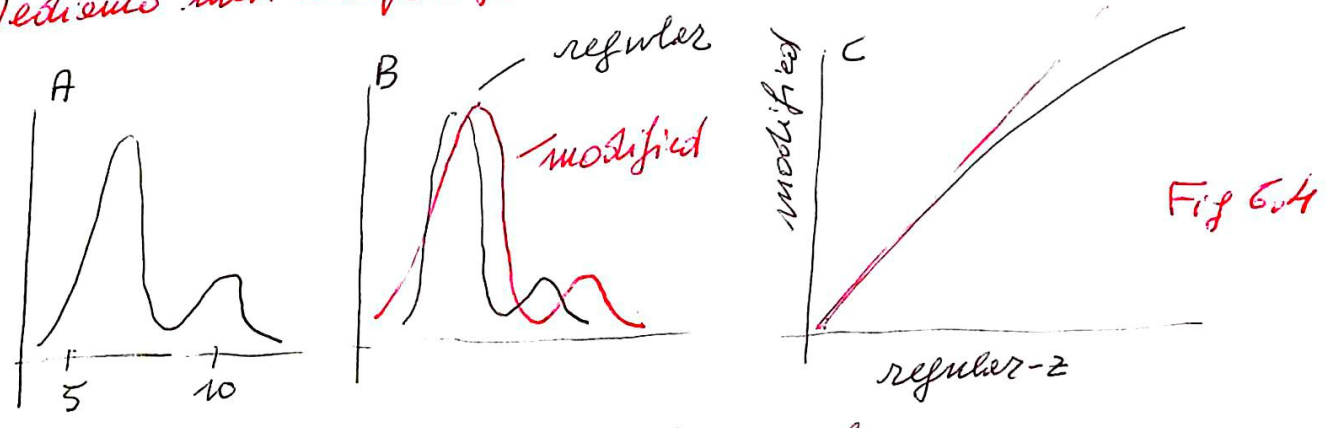


The modified z-score method

- si usa quando la distribuzione e' "fortemente non gaussiana" →
- regular z-score
- modified z-score → sottrae la mediana
divide per la differenza assoluta della mediana
- ↳ in concreto sottrae una misura di tendenza centrale

$$M_i = \frac{x_i - \tilde{x}}{1.4826 \cdot MAD} \quad \text{dove } \tilde{x} = \text{mediana}, \quad MAD = \text{mediana}(|x_i - \tilde{x}|)$$

Trovo quantile della distrib. Gauss
vediamo una comparazione:



- sono tra loro abbastanza simili in generale
- ma allora, quando diventa utile?
- intendo nella modified gli outliers contano poco →
- modified e' utile a testare se ci sono valori anomali durante la pulizia dei dati

p. 213

6.3 - Min-max normalization

tipicamente tra $[0, 1]$, $[-1, 1]$, $[\phi, z\pi]$

range minimo

$$\tilde{x}_i = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (6.9)$$

Tra i valori relativi e quelli originali, coefficiente = 1
per un range arbitrario

$$x^* = a + (b - a) \tilde{x}_i \quad (6.10)$$

6.4 - Z-normal vs min-max scaling

Z-normal produce una media e una deviazione std note
la scala min-max non preserva le caratteristiche distributive,
tranne i limiti sup e inf - Min-max si presta meglio se
i dati non hanno distribuzioni con lunghe code -

Min-max si presta bene a machine-learning - In questo
ambito si fa distinzione tra "normalizzazione" e "standardizzazione"
min-max \downarrow L Z-normal

Trasformazione in %

$$\text{percent} = 100 \cdot \frac{\text{ref} - \text{new}}{\text{ref}} \quad (6.11)$$

La % permette di paragonare dati molto diversi tra loro -
La misura con valori > 0 , con i negativi l'interpretazione è
problematica

6.6 - Nonlinear data transformations

Qui l'intento è cambiare la forma della distribuzione -
Necessità molta cautela nell'utilizzo - In quelle che presenteremo
la relazione monotona tra i punti è conservata -

cioè se $x_1 < x_2 \rightarrow \tilde{x}_1 < \tilde{x}_2$

questo è molto utile a interpretare i risultati di analisi
statistiche come i "test t", correlazioni, regressioni -

esempio di non monotonia: $\tilde{x} = \sin(x)$

06 90h

Rank-Transform

modificare i dati numerici con una scala rifinita (pollici, euro) in posizioni ordinali - Es.

$X = [1, 2, 3, 9348753945, 2.01]$ ----- (6,12)

$\tilde{X} = [1, 2, 4, 5, 3]$ ----- (6,13)

prima ordinio, *nessuno un profomido, con \tilde{X}*

qui abbiamo una perdita di informazione

esempio $X = [10, 2, 4, 5, 5]$ ----- (6,14)

$\tilde{X} = [4, 1, 2, 3, 3]$

fractional ranking $\tilde{X} = [5, 1, 2, 3.5, 3.5]$ "tied rank"

questa procedura è usata in molte statistiche inferenziali non parametriche (es. test di Wilcoxon, che è un'alternativa non parametrica al "t-test") (es. la correlazione di Spearman, che è una correlazione non parametrica).

logarithm and square root transformation

log-transform si prende il log₁₀ dei dati

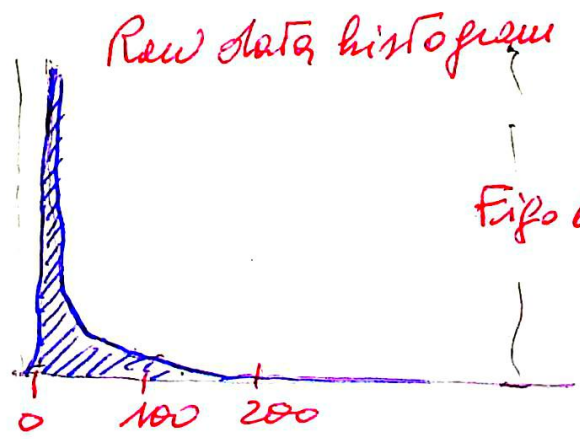
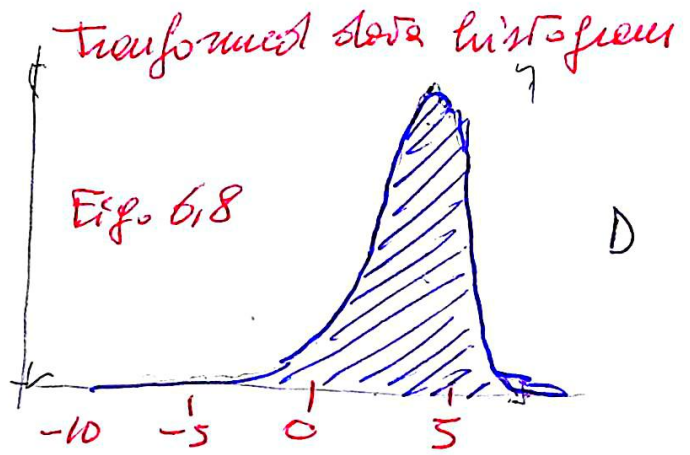
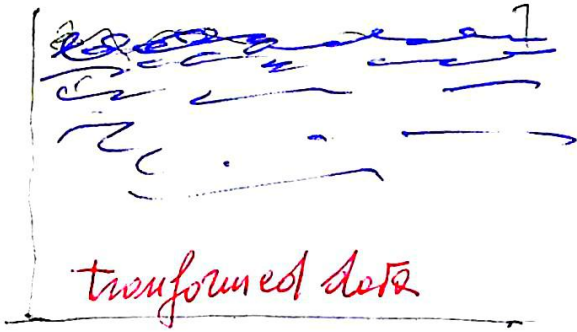


Fig 6.8



notare la distorsione relativa nel grafico D, però la distribuzione è unimodale ebbene non è una funzione.

Si presta bene se i dati seguono una legge di potenza \rightarrow la variabilità è \propto ai valori (eteroschedasticità, vedi ch. 4)

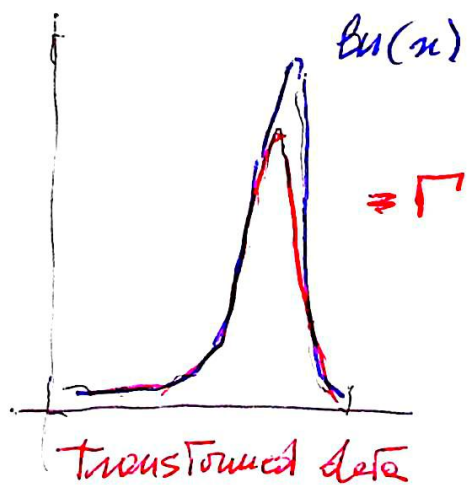
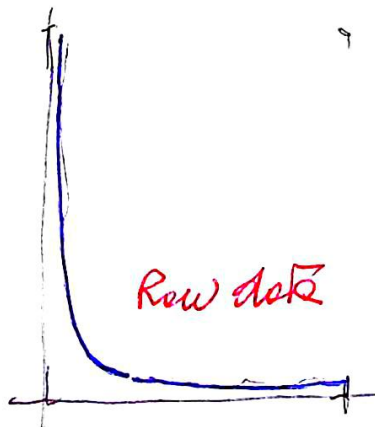
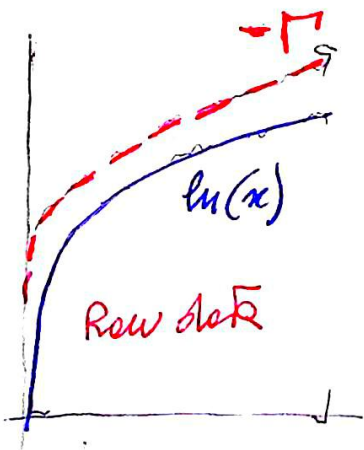
Si può applicare solo per valori $> \phi$

Se si hanno valori $< \phi$, si può procedere preventivamente a una traslazione.

Radice quadrata

eteroschedasticità

valida x valori $> \phi$ - Questa trasformazione può ridurre

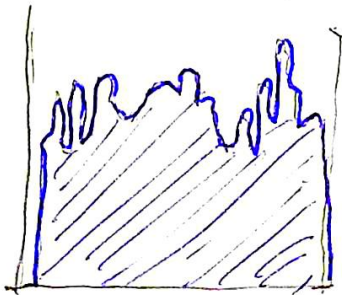


Fisher - Z

x deformare una distribuzione uniforme $(-1, 1)$ in una approssimativamente Gauss - Verfosco "allungati" i valori prossimi agli estremi dell'intervallo -

$$x_z = \frac{1}{2} \ln \left(\frac{1+x}{1-x} \right) \text{ con } -1 < x < 1 \quad (6,17)$$

Raw data histogram



Fisher data h_0

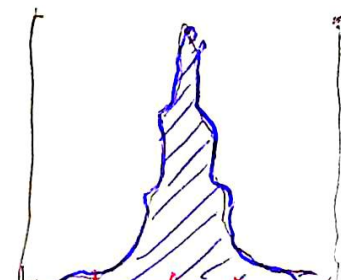
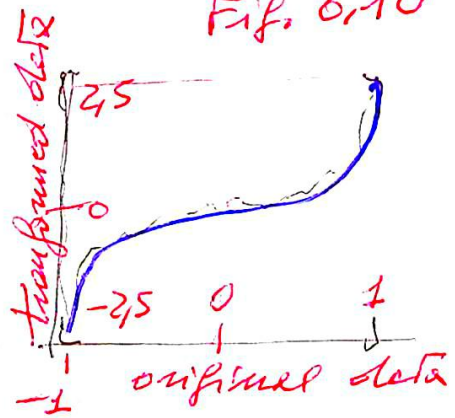


Fig. 6.10



se $x \rightarrow \phi$ la funzione $\rightarrow 1$
 $x \rightarrow 1$ \rightarrow molto pesante > 0
 $x \rightarrow -1$ $\rightarrow \phi \rightarrow \ln(\cdot)$ diventa $< \phi$

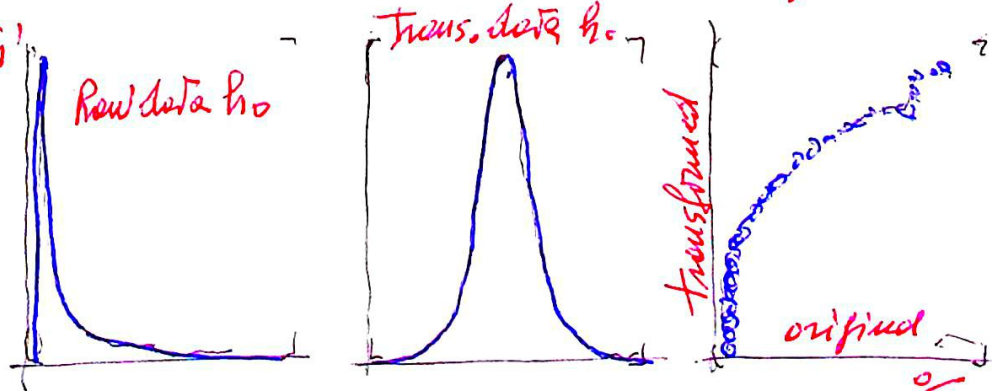
ma la 6,17 equivale alla tangente iperbolica \rightarrow tenifuo' ottenere con $\text{mp. or } \tanh(\cdot)$

Trasformare \forall distribuzione in Gaussiana

- 1) trasformare il tempo
- 2) min-max in $[-PPP, PPP]$
- 3) Fisher-Z

Gaussiane (emotona)

e' con perdita di dati



Interpreting transformed data

Z-transform e' facile da interpretare - Ogni aumento della deviazione std in A \rightarrow incremento di st nella deviazione std di B

es. trasformo log, quindi eseguo analisi di regressione -

Potrei avere un effetto lineare di una variabile su un'altra, ma tutto questo si basa su trasformazioni non lineari \rightarrow

in realtà l'effetto e' non lineare - **TENERE PRESENTE**

- min-max offre una interpretazione + chiara, anche se non necessariamente la + precisa -

When to transform your data

Lo z-normalizing puo' essere utile a confrontare variabili su scale diverse -
Usare le trasformazioni solo se necessario -

6.8 - Exercises

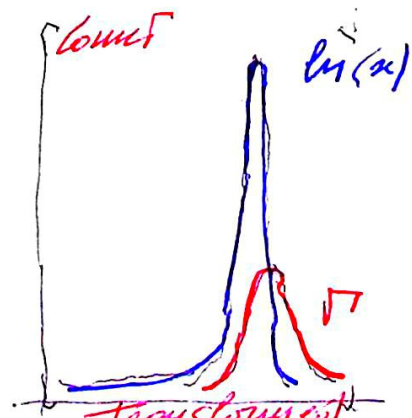
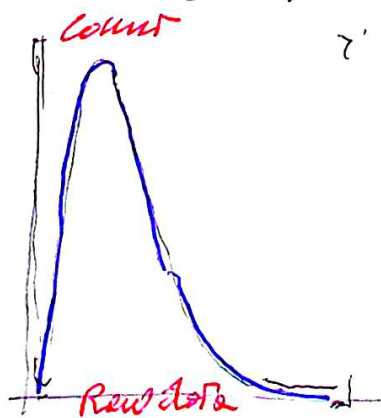
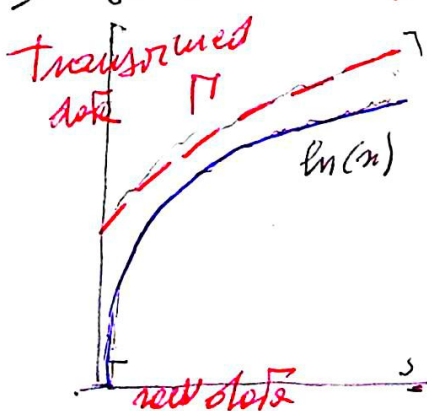
1) combinare insieme 6.8 e 6.10 - tradurre in Python -

testare su 10 numeri Gauss (14,3, 34) **valore 2** \rightarrow

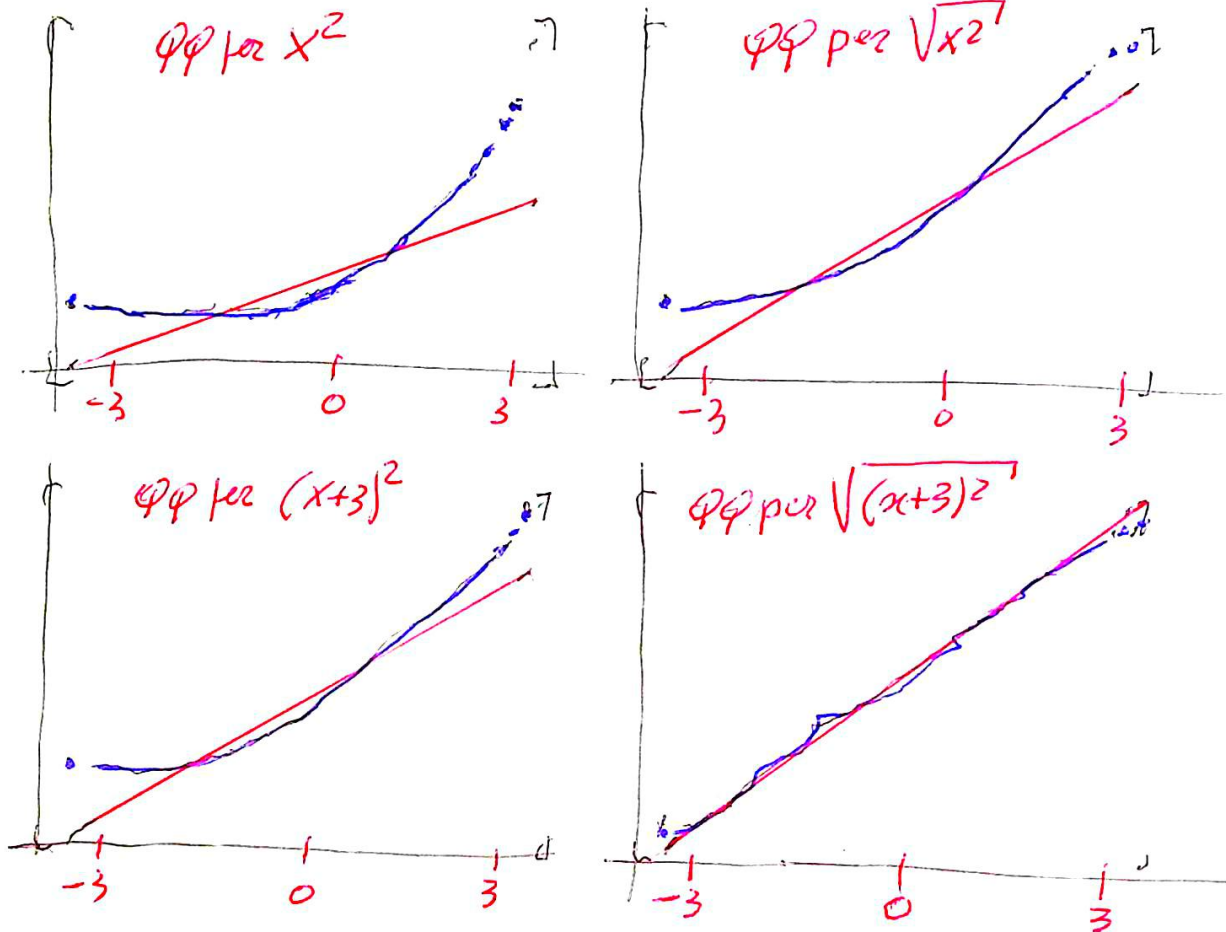
min value = 14,3 ; max value = 34

su Web 7 funzioni che lo fa, trovarle.

2) Fig. 6.8 - Modifica mando $(X+3)^2$ invece che X^2



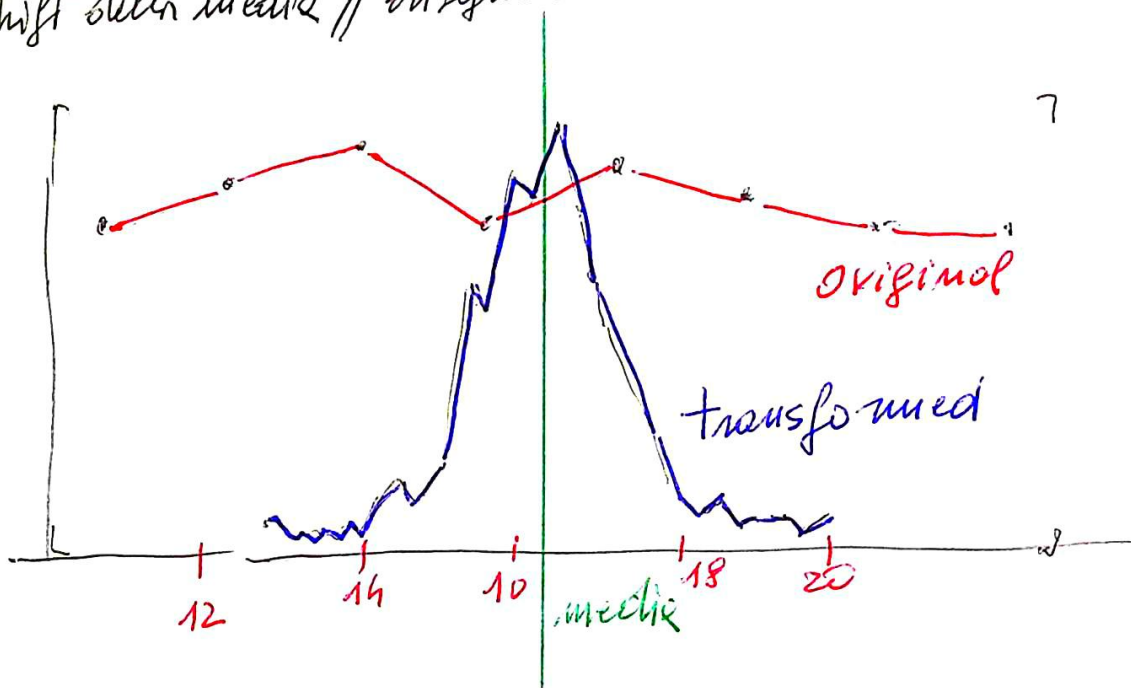
• inverremo ora i QQ plots



infine spostare i dati trasformati in T in modo che abbiamo $media's = \Phi$

- 3) la normalizzazione per le dati? Risolvere la bit per x_i , se questo è possibile \rightarrow la trasformazione è invertibile. Fare questo anche con codice, mediante 25 interi casuali. Tra H e H_0 z-transformiamo, anti-transformiamo. Verificare che il risultato coincida con originale.
- 4) applichiamo molte trasformazioni / numeri & sempre presentando la media dei dati.
- 313 numeri casuali da distribuzione uniforme tra 3π e π e incerti che tutti i campioni siano reali.
- hist di raw data, hist transformed data
- description procedure:

1) genero i dati // 2) normal (-PPP, PPP) // Fisher-z //
 shift della media // distribuzione normale // *sivento 6000*



5) \mathcal{Z} mp function x z-score \rightarrow scipy.stats - intervalli 3 e 9

scipy.stats.zscore oppure scale

calcolare entrambi e confrontarli che siano uguali

- ora proviamo con una matrice - calcolare \bar{x} e s da ciascuna colonna separatamente, o da intervalli matriciale?

- creare una matrice con valori di x^2 con $0 < x < 11$

efficienza scipy.stats.zscore

$\begin{bmatrix} 0 & 1 & 4 \end{bmatrix}$

$\begin{bmatrix} 9 & 16 & 25 \end{bmatrix}$

$\begin{bmatrix} 36 & 49 & 64 \end{bmatrix}$

$\begin{bmatrix} 81 & 100 & 121 \end{bmatrix}$

column-wise z-scoring

collo

Matrix-wise z-scoring

pono calcolare std, medie x ogni colonna, ma e' sufficiente ispezione visiva

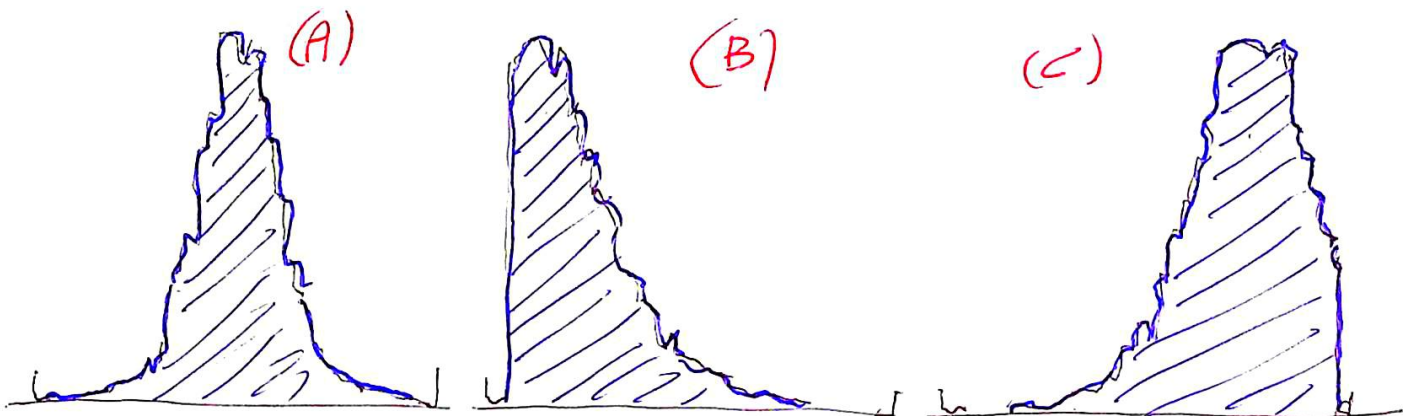
z-scoring x Matrice puo' essere appropriato quando tutte le colonne contengono caratteristiche nella stessa scala e intervallo numerico

6) trovare modo di usare Fisher-Z e generare una distribuzione di numeri casuali (A) *nessuna asimmetria*
 (B) *forte asimmetria positiva* (C) *leggera asimmetria negativa*
 - calcolare l'asimmetria empirica, graficare

X1 np. ortanti (np. random. uniform (-0.999, 0.999, nfe = 500) (A)
 X2 _____ (-0.2, 0.999) (B)
 X3 _____ (-0.999, 0.8) (C)

$$\text{skew}[\phi] = \text{stats.skew}(X1)$$

_____ 1 _____ X2
 _____ 2 _____ X3



skew = 0,02

skew = 1,25

skew = -0,74

7) esplorare relazione tra std e MAD in distribuzioni che si allontanano da Gauss

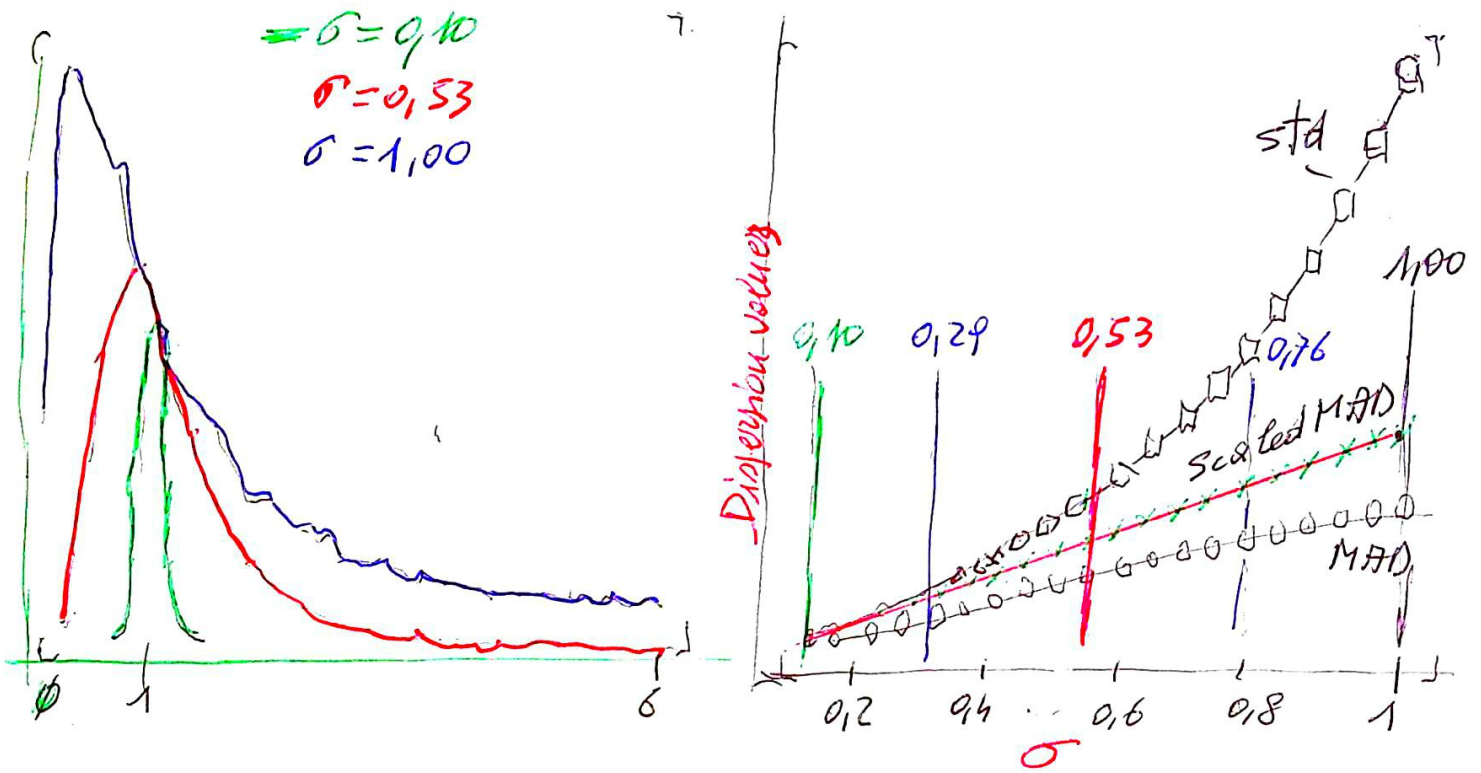
- dati casuali esponentioli tra 0,2 e 1 - calcolare tre statistiche descrittive (1) std (2) MAD (3) Scaled MAD, Graficare

$$\text{data}[i] = \text{np.exp}(X * s) \quad \text{casuali}$$

$$M[i, \phi] = \text{mp.std}(\text{data}[i], \text{ddof} = 1)$$

$$M[i, 1] = \text{mp.median}(\text{np.abs}(\text{data}[i] - \text{mp.median}(\text{data}[i])))$$

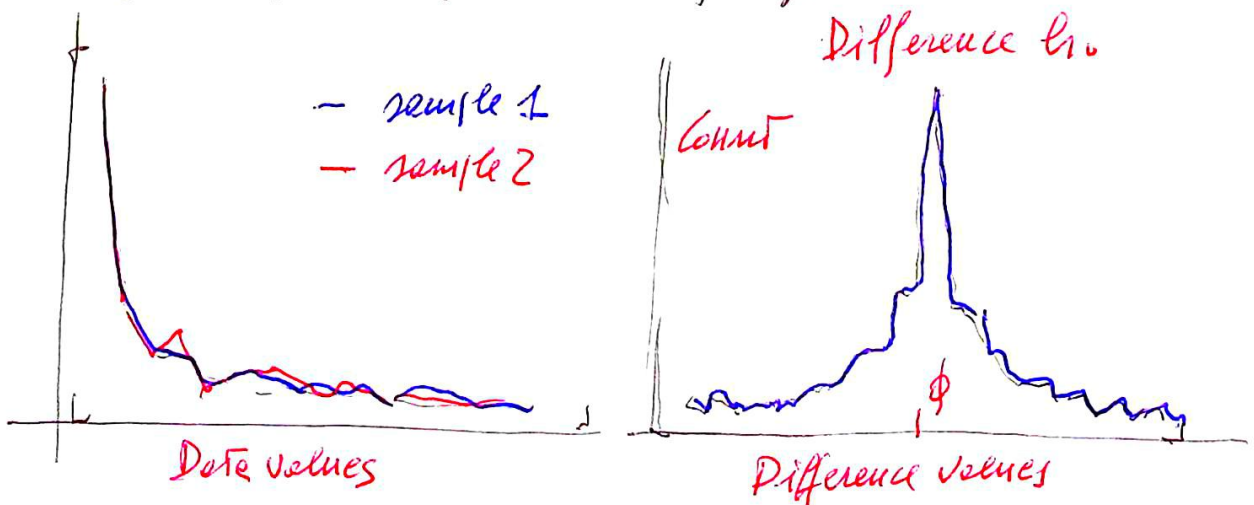
$$\text{sigmas} = \text{np.linspace}(0, 1, 20) \quad \text{vari valori di } \sigma$$



8) La differenza tra due variabili distribuite in modo simile, non normale. \rightarrow può essere \sim Geom

2 set di dati $N=300$ casuali da legge di potenza

loro integramenti, e integramenti differenza



— FINE CAPITOLO 6 —