

Metodi statistici per l'analisi dei dati

Massimiliano Grosso

Dipartimento di Ingegneria Meccanica, Chimica e dei
Materiali

E-mail: massimiliano.grosso@dimcm.unica.it

Web: <http://people.unica.it/massimilianogrosso>

1

Metodi Statistici per l'Analisi dei Dati

INTRODUZIONE

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

2

Motivazioni

Richiami di statistica – Esperimenti replicati

- Obiettivo del corso:
 - **Pianificazione** degli esperimenti in maniera tale che i risultati della campagna sperimentali possano essere analizzati con **metodi statistici**, per giungere a conclusioni **oggettive** del processo in esame

Due fasi distinte:

{

1. Pianificazione della campagna sperimentale (Design Of Experiments: **DOE**)

2. Analisi **statistica** dei risultati

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

3

Motivazioni

Richiami di statistica – Esperimenti replicati

Fattori controllabili:

x_1 x_2 x_3 ... x_n

Inputs: Processo Outputs y

z_1 z_2 z_3 ... z_n

Fattori incontrollabili:

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

4

Metodi statistici per l'Analisi dei Dati
10 – 14 febbraio 2020

2

Motivazioni

Richiami di statistica – Esperimenti replicati

• Lo studio di un processo è una procedura iterativa

Congettura su un processo

→

Esperimenti sul processo

↖

Conoscenza del processo

↘

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

5

Progettazione campagna sperimentale – Concetti di base

Richiami di statistica – Esperimenti replicati

• I principi di base della progettazione della campagna sperimentale sono:

1. Replicazione
2. Randomizzazione
3. Blocking

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

6

Progettazione campagna sperimentale – Concetti di base

Richiami di
statistica –
Esperimenti
replicati

- **Replicazione**
 - Ripetere gli esperimenti nelle stesse condizioni più volte
 1. Permette di ottenere una stima «**genuina**» dell'errore sperimentale
 2. Permette una stima più **precisa** della variabile di output
 - **Replica della misura sperimentale** è un concetto diverso (e più potente) della **misura ripetuta**
 - Nell'ultimo caso si può valutare al più la variabilità intrinseca del sistema di misura

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

7

Progettazione campagna sperimentale – Concetti di base

Richiami di
statistica –
Esperimenti
replicati

- **Randomizzazione**
 - **Ordine** con cui sono eseguite le misure sperimentali deve essere del tutto **casuale**
 - Randomizzando l'ordine delle esperienze si possono *compensare* eventuali effetti di ulteriori fattori (non considerati nel modello) che possono essere presenti e che risultano essere soggetti a variazioni nel tempo

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

8

Progettazione campagna sperimentale – Concetti di base

**Richiami di
statistica –
Esperimenti
replicati**

- **Blocking**
- Tecnica di progettazione della campagna sperimentale usata per aumentare la precisione con cui sono effettuati i confronti tra i fattori di interesse.
- Il Blocking è usato per ridurre la variabilità relativa a fattori di disturbo
 - fattori che possono influenzare la risposta ma a cui non siamo interessati
- **Blocco – Definizione**
- Un insieme di condizioni sperimentali relativamente omogenee

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

9

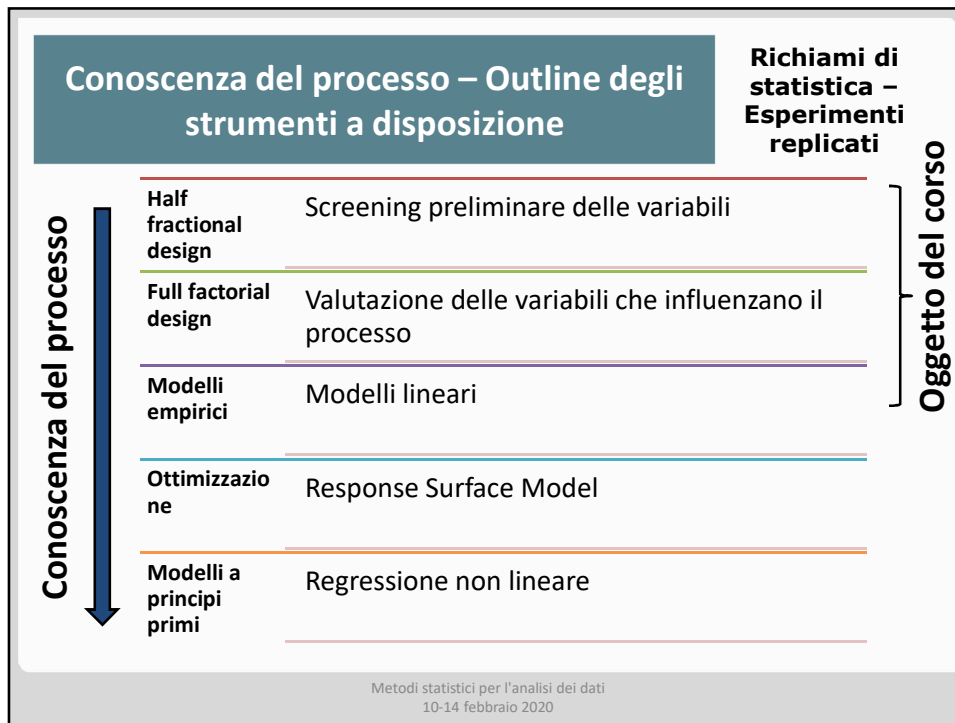
Linee guida per una campagna sperimentale

**Richiami di
statistica –
Esperimenti
replicati**

1. Definizione del problema
 2. Scelta dei fattori, livelli e intervalli
 3. Selezione delle variabili da misurare
 4. Pianificazione della campagna sperimentale
 5. Esperimenti
 6. Analisi statistica dei dati
 7. Conclusioni
- Le linee guida riportate sono valide qualunque sia il livello di conoscenza del processo

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

10



11

Metodi statistici per l'analisi dei dati

**RICHIAMI DI STATISTICA –
ESPERIMENTI REPLICATI**

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

12

Introduzione alla sezione

Richiami di statistica – Esperimenti replicati

- La discussione permetterà di rivedere diversi concetti di base di statistica
 - Variabili aleatorie
 - Distribuzioni di probabilità
 - Campioni aleatori
 - Distribuzioni di campionamento
 - Test delle ipotesi – Intervalli di fiducia
- Per il momento esperimenti effettuati **sempre nelle stesse condizioni**.
- N.B. Da non confondere **esperimenti replicati nelle stesse condizioni** con **misure ripetute**

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

13

Esempio introduttivo

Richiami di statistica – Esperimenti replicati

- Si intende monitorare la qualità di una crema destinata ad uso alimentare.
- A tal proposito sono considerati 10 diversi campioni della crema e, per ciascuno di essi è misurata la viscosità
 - 10 misure sperimentali di viscosità riportate in tabella
- L'insieme di misure di viscosità è un **campione** sperimentale.

j	Controllo (cp)
1	70.00
2	70.52
3	73.00
4	72.00
5	71.44
6	71.00
7	72.88
8	71.60
9	71.84
10	72.60
\bar{y}	71.69

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

14

Concetti di statistica di base

Richiami di statistica – Esperimenti replicati

- Le **prove sperimentali** (etichettate con il pedice j) differiscono tra loro per effetto delle fluttuazioni dovute all'**errore sperimentale**.
- La presenza dell'errore sperimentale implica che la singola misura sia l'esito di una **variabile aleatoria** (ovvero, non è possibile a priori la sua previsione).

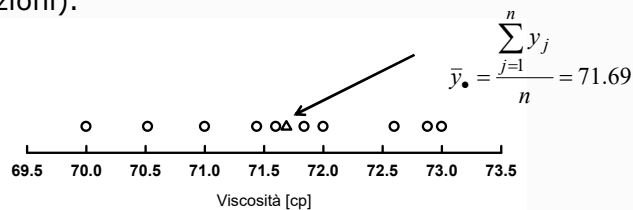
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

15

Concetti di statistica di base – Descrizioni grafiche della variabilità

Richiami di statistica – Esperimenti replicati

- **Diagramma per punti**
- Utile per campioni di piccole dimensioni (sino a 20 osservazioni).



- Il diagramma permette di riconoscere il **trend centrale** e la **dispersione** dei dati.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

16

Concetti di statistica di base – Indici di posizione e dispersione del campione

Richiami di statistica – Esperimenti replicati

- Scalari per identificare il trend centrale:
- **Media aritmetica**

$$\bar{y}_{\bullet} = \frac{\sum_{j=1}^n y_j}{n} = 71.69$$

- **Mediana**: rappresenta il valore centrale che divide il campione in due parti uguali costituiti rispettivamente dai valori inferiori e superiori ad esso

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

17

Concetti di statistica di base – Indici di posizione e dispersione del campione

Richiami di statistica – Esperimenti replicati

- Misure della **dispersione** dei dati:
- **Varianza**:

$$s^2 = \frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y}_{\bullet})^2$$

La somma dei quadrati è divisa per $(n-1)$ anziché n

- **Deviazione standard**
- È la radice quadrata della varianza

$$s = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y}_{\bullet})^2}$$

- Utile perché ha le stesse dimensioni della variabile y

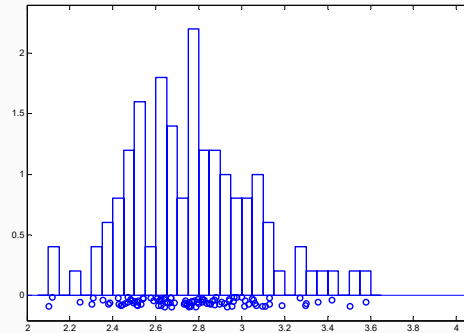
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

18

Concetti di statistica di base – Descrizioni grafiche della variabilità – Frequenze

Richiami di statistica – Esperimenti replicati

- In presenza di campioni di dimensioni maggiori è possibile riportare i dati negli istogrammi delle **frequenze assolute** (o **relative**) del campione di dati.
- L'istogramma è costruito dividendo l'asse orizzontale in intervalli (in genere di uguale lunghezza) e disegnando un rettangolo sul j -esimo intervallo la cui area sia proporzionale a n_j , numero di osservazioni che cadono nell'intervallo.



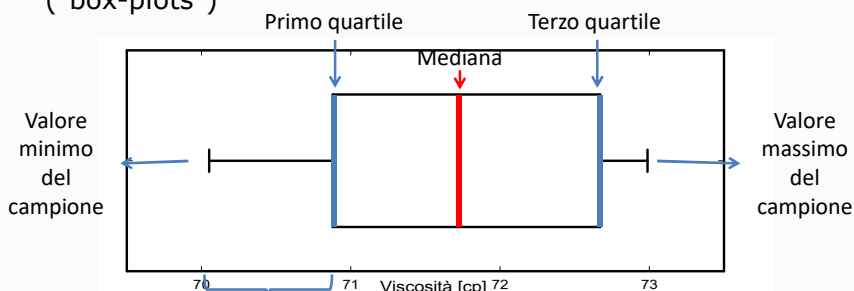
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

19

Concetti di statistica di base – Descrizioni grafiche della variabilità

Richiami di statistica – Esperimenti replicati

- Rappresentazione dei campioni tramite "diagrammi a scatola" ("box-plots")



Il 25% delle osservazioni cade in questo intervallo

Il 50% delle osservazioni cade in questo intervallo

Il 75% delle osservazioni cade in questo intervallo

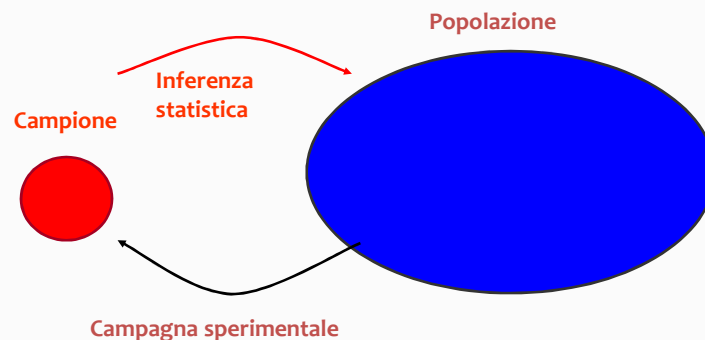
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

20

Campioni e distribuzioni campionarie

Richiami di statistica – Esperimenti replicati

- L'obiettivo dell'inferenza statistica è trarre delle conclusioni su una popolazione a partire da un suo campione



Dal campione si intende ottenere informazioni sulla popolazione generatrice non nota

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

21

Caratterizzazione della Popolazione – Distribuzioni di probabilità

Richiami di statistica – Esperimenti replicati

- La struttura di probabilità di una variabile aleatoria (**VA**) Y è descritta dalla sua **funzione densità di probabilità** (probability density function: **pdf**) $f(y)$.
- Proprietà fondamentali della pdf di una VA:

$$1. \quad P(a \leq y \leq b) = \int_a^b f(y) dy$$

$$2. \quad \int_{-\infty}^{+\infty} f(y) dy = 1$$

$$3. \quad f(y) \geq 0$$

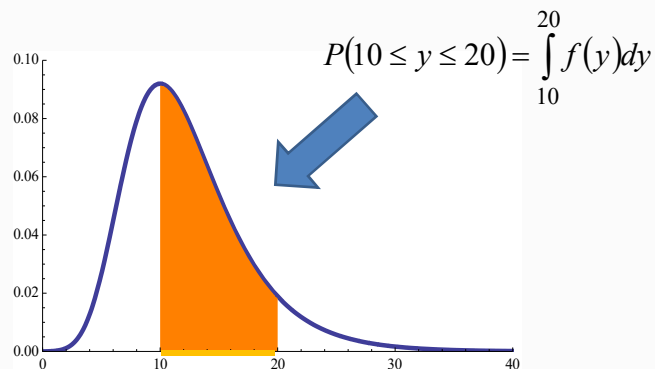
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

22

Caratterizzazione della Popolazione – Distribuzioni di probabilità

Richiami di
statistica –
Esperimenti
replicati

- Esempio di funzione densità di probabilità



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

23

Distribuzioni di probabilità – Scarsi associati

Richiami di
statistica –
Esperimenti
replicati

- **Media** di una variabile aleatoria Y (anche definito **valore atteso**)

$$\mu = \int_{-\infty}^{+\infty} y f(y) dy = E[Y]$$

- **Definizione**
- L'operatore **Valore Atteso** $E[X]$ restituisce il risultato medio che si osserverebbe per in presenza di infinite osservazioni della Variabile Aleatoria X

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

24

Caratterizzazione della Popolazione – Scalari associati ad una VA

**Richiami di
statistica –
Esperimenti
replicati**

- **Varianza** di una variabile aleatoria Y

$$\sigma^2 = V(Y) = \int_{-\infty}^{+\infty} (y - \mu)^2 f(y) dy = E[(Y - \mu)^2]$$

- Varianze piccole sono associate ad incertezze piccole.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

25

Caratterizzazione della Popolazione – Scalari associati ad una VA

**Richiami di
statistica –
Esperimenti
replicati**

- Alcune proprietà di interesse delle VA. 1/2

1. $E[c] = c$
2. $E[Y] = \mu$
3. $E[cY] = cE[Y] = c\mu$
4. $V[c] = 0$
5. $V[Y] = \sigma^2$
6. $V[cY] = c^2V[Y] = c^2\sigma^2$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

26

Caratterizzazione della Popolazione – Scalari associati ad una VA

**Richiami di
statistica –
Esperimenti
replicati**

- Alcune proprietà di interesse delle VA. 2/2
 - In presenza di più variabili aleatorie:
6. $E[Y_1 + Y_2] = E[Y_1] + E[Y_2] = \mu_1 + \mu_2$
7. $V[Y_1 + Y_2] = V[Y_1] + V[Y_2] + 2\text{cov}(Y_1, Y_2)$
- Dove è definita la covarianza delle VA Y_1 e Y_2 :

$$\text{cov}(Y_1, Y_2) = E[(Y_1 - \mu_1)(Y_2 - \mu_2)]$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

27

Caratterizzazione della Popolazione – Scalari associati ad una VA

**Richiami di
statistica –
Esperimenti
replicati**

- **Statistica – Definizione:**
- Una **statistica** è una **funzione** delle **osservazioni** di un **campione** che non contiene parametri **ignoti** della **popolazione** che ha generato il campione (es: media e varianza).
- Esempi di statistiche:

- **Media aritmetica**

$$\bar{Y} = \frac{\sum_{i=1}^n Y_i}{n}$$

- **Varianza campionaria**

$$S^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

28

Campioni e distribuzioni campionarie - Stimatori

Richiami di
statistica –
Esperimenti
replicati

- **Stimatore – definizione:**
- Uno **stimatore** di un **parametro** ignoto è una statistica che mira a valutare il parametro stesso.
- La media aritmetica e la varianza campionaria sono esempi di stimatori **puntuali**.
- Lo stimatore puntuale del generico parametro θ è in genere indicato con il simbolo del cappuccio:

$$\hat{\theta}$$

- Esempio media aritmetica:

$$\bar{Y} = \sum Y_i / n = \hat{\mu}$$

- Un valore numerico puntuale calcolato da un campione di dati, prende il nome di **stima**.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

29

Campioni e distribuzioni campionarie - Stimatori

Richiami di
statistica –
Esperimenti
replicati

- Proprietà stimatori
- **Imparzialità:** Uno stimatore si dice **imparziale** (*unbiased*) se il suo valore atteso coincide con il valore vero del parametro

$$E[\hat{\theta}] = \theta$$

- NB sebbene il valore vero non sarà mai noto è possibile valutare il verificarsi della imparzialità.
- **Efficienza:** È una misura della **varianza dello stimatore**. Se dispongo di più stimatori devo scegliere quello con varianza minima ovvero quello con la massima efficienza.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

30

Campioni e distribuzioni campionarie - Stimatori

Richiami di
statistica -
Esperimenti
replicati

- Lo stimatore **media aritmetica** è **imparziale**:

$$\begin{aligned} E[\bar{Y}] &= E\left[\frac{\sum_{i=1}^n Y_i}{n}\right] = \frac{1}{n} E\left[\sum_{i=1}^n Y_i\right] = \frac{1}{n} \sum_{i=1}^n E[Y_i] = \\ &= \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} n\mu = \mu \end{aligned}$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

31

Campioni e distribuzioni campionarie - Stimatori

Richiami di
statistica -
Esperimenti
replicati

- Lo stimatore media aritmetica è **efficiente**:

$$\begin{aligned} V[\bar{Y}] &= V\left[\frac{\sum_{i=1}^n Y_i}{n}\right] = V\left[\sum_{i=1}^n \frac{1}{n} Y_i\right] = \frac{1}{n^2} \sum_{i=1}^n V[Y_i] = \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n} \end{aligned}$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

32

Campioni e distribuzioni campionarie - Stimatori

Richiami di
statistica -
Esperimenti
replicati

- In maniera analoga si può dimostrare che la **varianza campionaria** S^2 è **imparziale**

$$\begin{aligned} E[S^2] &= E\left[\frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}\right] = \frac{1}{n-1} \sum_{i=1}^n E[(Y_i - \bar{Y})^2] = \\ &= \frac{1}{n-1} E[SS] \end{aligned}$$

- dove SS è la **somma corretta dei quadrati** delle osservazioni y_i

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

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

33

Campioni e distribuzioni campionarie - Stimatori

Richiami di
statistica -
Esperimenti
replicati

- Dimostrazione imparzialità varianza - Continua

$$\begin{aligned} E[SS] &= E\left[\sum_{i=1}^n (Y_i - \bar{Y})^2\right] = E\left[\sum_{i=1}^n Y_i^2 - n\bar{Y}^2\right] = \\ &= \sum_{i=1}^n (\mu^2 + \sigma^2) - n(\mu^2 + \sigma^2/n) = (n-1)\sigma^2 \end{aligned}$$

- da cui:

$$E[S^2] = \frac{1}{n-1} E[SS] = \sigma^2$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

34

Campioni e distribuzioni campionarie – Definizione gradi di libertà

**Richiami di
statistica –
Esperimenti
replicati**

- Il **numero di gradi di libertà** di una **somma di quadrati** è data dal numero di elementi **indipendenti** presenti nella somma.
- Esempio: SS ha n-1 g.d.l.

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

- In SS **non tutti gli elementi sono indipendenti**.
- Il valore della media \bar{y} deve essere tale da soddisfare il vincolo:

$$\sum_{i=1}^n (y_i - \bar{y}) = 0$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

35

Campioni e distribuzioni campionarie – Definizione gradi di libertà

**Richiami di
statistica –
Esperimenti
replicati**

- **Risultato generale:**
- Se y è una variabile aleatoria di varianza σ^2 e una somma degli scarti quadratici ricavata da essa ha v g.d.l., allora

$$E\left[\frac{SS}{v}\right] = \sigma^2$$

- **Proprietà importante** per le applicazioni successive

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

36

Caratterizzazione della Popolazione – Esempi di distribuzione

**Richiami di
statistica –
Esperimenti
replicati**

- **Distribuzione di tipo normale o Gaussiana**

- La densità di probabilità è data da:

$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(y-\mu)^2}{\sigma^2}\right) \quad -\infty < y < +\infty$$

- La funzione è definita lungo tutto l'asse reale (ovvero un qualunque numero reale può essere un esito di una VA di tipo normale)
- Il grafico di tale funzione è una curva a campana simmetrica rispetto a $y=\mu$
- La distribuzione dipende da due parametri, μ e σ^2 .

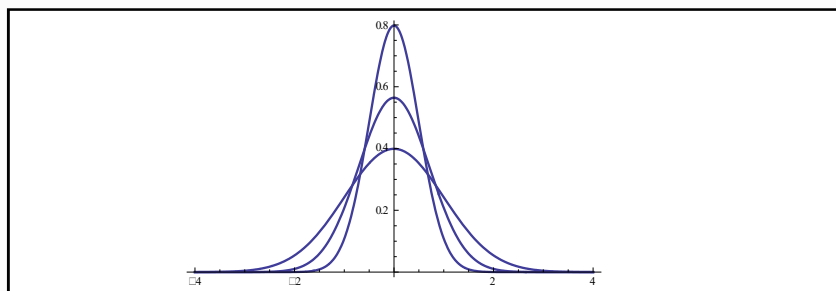
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

37

Distribuzione normale

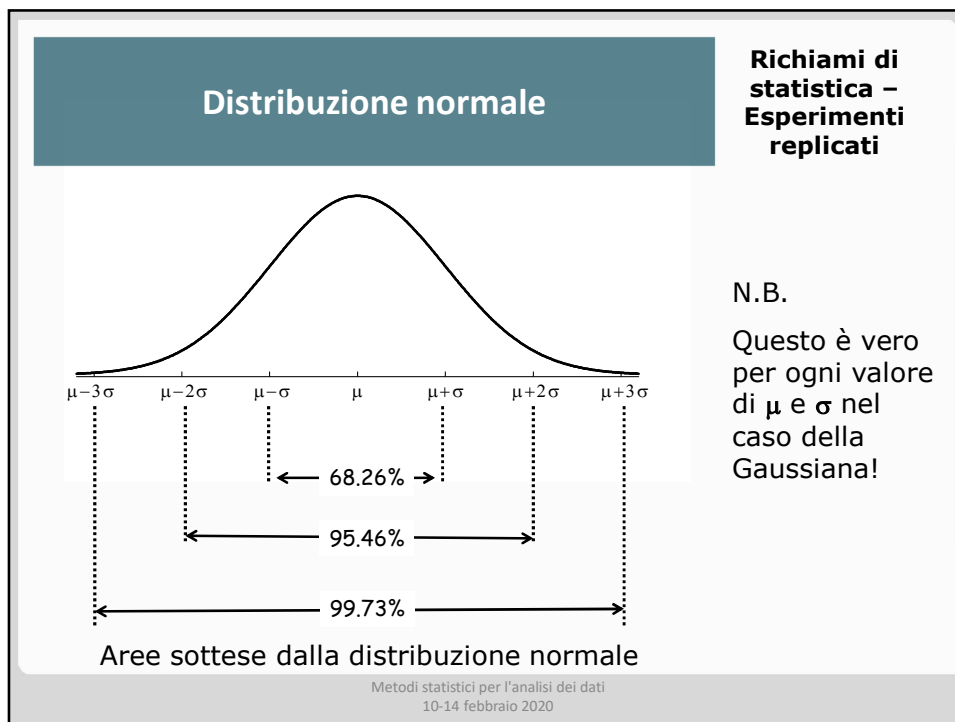
**Richiami di
statistica –
Esperimenti
replicati**

In figura sono riportate tre gaussiane con egual media e
varianza 0.25, 0.5, 1



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

38



39

Distribuzione normale di tipo standard – Definizione

Richiami di statistica – Esperimenti replicati

- Data una variabile aleatoria Y (di tipo gaussiano) di media μ e varianza σ^2

$$Y \sim N(\mu, \sigma^2)$$
- Si consideri la seguente trasformazione lineare:

$$Z = \frac{Y - \mu}{\sigma}$$
- È facile verificare che la nuova VA Z ha media 0 e varianza unitaria:

$Z \sim N(0,1)$ \longrightarrow **Gaussiana di tipo standard**

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

40

Funzioni di VA Gaussiane Trasformazioni lineari

**Richiami di
statistica –
Esperimenti
replicati**

- Nota la funzione di distribuzione standard è possibile ricavare le proprietà di una qualsiasi distribuzione gaussiana
- In particolare, è possibile calcolare la probabilità che si verifichi un dato evento per un generico processo, con media e varianza note.
- Questo è possibile sapendo solo i valori della distribuzione di tipo standard.

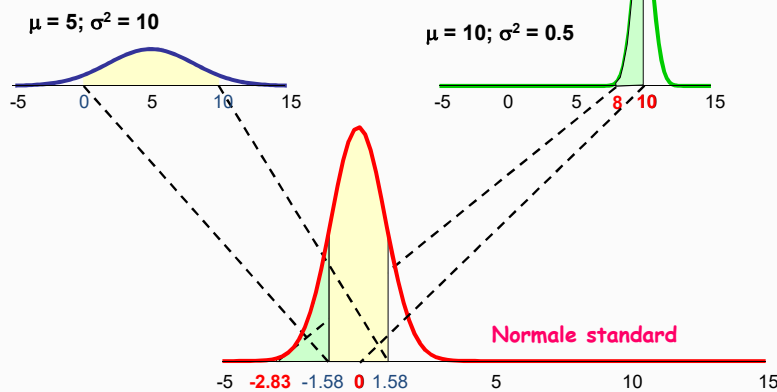
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

41

Calcolo probabilità per una Gaussiana generica

**Richiami di
statistica –
Esperimenti
replicati**

$$z = \frac{(y - \mu)}{\sigma}$$



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

42

Calcolo probabilità per una Gaussiana generica

Richiami di statistica – Esperimenti replicati

- Esempio: calcolare quale è la probabilità che si verifichi un evento appartenente all'intervallo $[0,5]$ per la variabile aleatoria di media 3 e deviazione standard 2:
- Si deve calcolare quale è la probabilità che la variabile aleatoria di tipo standard assuma un valore nell'intervallo corrispondente.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

43

Calcolo probabilità per una Gaussiana generica

Richiami di statistica – Esperimenti replicati

- Dobbiamo calcolare la probabilità:

$$P(0 < X < 5)$$

- Gli estremi dell'intervallo corrispondente per la distribuzione di tipo standard possono essere facilmente calcolati

$$z_1 = \frac{x_1 - \mu_X}{\sigma_X} = \frac{0 - 3}{2}$$
$$z_2 = \frac{x_2 - \mu_X}{\sigma_X} = \frac{5 - 3}{2} = 1$$



$$P(0 < X < 5) =$$
$$P(-1.5 < Z < 1) =$$
$$0.8413 - 0.0668 = 77.4\%$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

44

Calcolo probabilità per una Gaussiana generica

Richiami di statistica – Esperimenti replicati

- Esercizi
- Sia Y una variabile aleatoria di tipo normale, di media $\mu = 16$ e varianza $\sigma^2 = 25$
- Calcolare:
 - $P(Y > 20)$
 - $P(20 < Y < 25)$
 - $P(Y < 10)$
 - $P(12 < Y < 24)$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

45

Teorema del limite centrale

Richiami di statistica – Esperimenti replicati

- **Teorema del limite centrale**
- Sia y_1, y_2, \dots, y_n una successione di n VA indipendenti ed identicamente distribuite tali che $E[y_i] = \mu$ e $V(y_i) = \sigma^2$.
- Sia inoltre $x_n = y_1 + y_2 + \dots + y_n$
- Allora:

$$Z_n = \frac{X_n - n\mu}{\sqrt{n\sigma^2}}$$

- tende ad una **VA Gaussiana di tipo standard** per $n \rightarrow \infty$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

46

Teorema del limite centrale

**Richiami di
statistica –
Esperimenti
replicati**

- **VA Gaussiana è ideale per descrivere l'errore sperimentale**
- La VA di tipo normale è un valido modello matematico per descrivere le incertezze presenti nella misura sperimentale
 - È ragionevole assumere che le deviazioni dal valore vero provengano da diverse fonti indipendenti

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

47

Variabili Aleatorie derivate dalla gaussiana - Variabile χ^2

**Richiami di
statistica –
Esperimenti
replicati**

- Si considerino k VA di tipo Standard indipendenti z_1, z_2, \dots, z_k
- La variabile aleatoria scalare

$$X = Z_1^2 + Z_2^2 + \dots + Z_k^2$$

- prende il nome di variabile aleatoria χ^2 a k gradi di libertà.
- Tale variabile aleatoria è caratterizzata completamente da un solo parametro, il numero di gradi di libertà k .
- La pdf ha espressione:

$$f(x) = \frac{1}{2^{k/2} \Gamma\left(\frac{k}{2}\right)} x^{\frac{k}{2}-1} \exp\left(-\frac{x}{2}\right) \quad x > 0$$

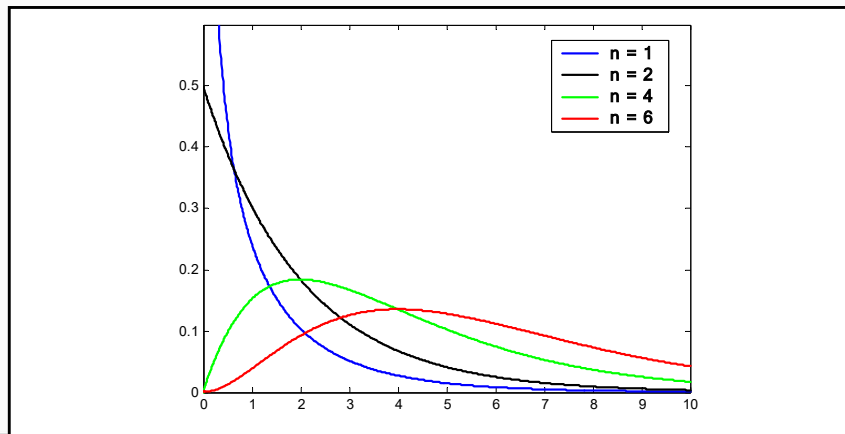
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

48

Variabile aleatoria χ^2

Richiami di
statistica –
Esperimenti
replicati

- Funzione densità di probabilità



49

Variabile aleatoria χ^2

Richiami di
statistica –
Esperimenti
replicati

- Proprietà di una variabile aleatoria χ^2 a k gradi di libertà

$$\mu = k$$

$$\sigma^2 = 2k$$

- Il massimo si osserva in corrispondenza di $y = n-2$.
- Per $n \rightarrow \infty$ la distribuzione χ^2 tende ad una gaussiana.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

50

Variabile aleatoria χ^2 – Esempio

**Richiami di
statistica –
Esperimenti
replicati**

- Esempio di VA che segue la distribuzione di tipo χ^2 :
- Siano y_1, y_2, \dots, y_n un campione di dati generati da una VA di tipo Gaussiano $N(\mu, \sigma^2)$. Allora:

$$\frac{SS}{\sigma^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sigma^2} \sim \chi_{n-1}^2$$

- Da cui, con semplici passaggi, si può ricavare la seguente relazione per la stima S^2 della varianza:

$$S^2 = \frac{SS}{n-1} \Rightarrow (n-1) \frac{S^2}{\sigma^2} \approx \chi_{n-1}^2$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

51

VA derivate dalla gaussiana Distribuzione T-student

**Richiami di
statistica –
Esperimenti
replicati**

- Siano dati una variabile aleatoria Z Gaussiana di tipo standard (ovvero $Z \sim \mathcal{N}(0,1)$), ed una χ^2 ad r gradi di libertà
- La variabile aleatoria :

$$T_r = \frac{Z}{\sqrt{\frac{\chi_r^2}{r}}}$$

è una distribuzione T di student ad r gradi di libertà.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

52

VA derivate dalla Gaussiana Distribuzione T di student

**Richiami di
statistica –
Esperimenti
replicati**

- Espressione analitica della T di student

$$f_r(y) = \frac{\Gamma\left(\frac{r+1}{2}\right)}{\sqrt{r\pi}\Gamma(r/2)} \frac{1}{\left(\frac{y^2}{r} + 1\right)^{\frac{r+1}{2}}} \quad -\infty < y < \infty$$

- Proprietà:
- Dipende da un solo parametro il numero intero r

Media: $\mu_{t,r} = 0$

Varianza: $\sigma_{t,r}^2 = \frac{r}{r-2} \quad (r > 2)$

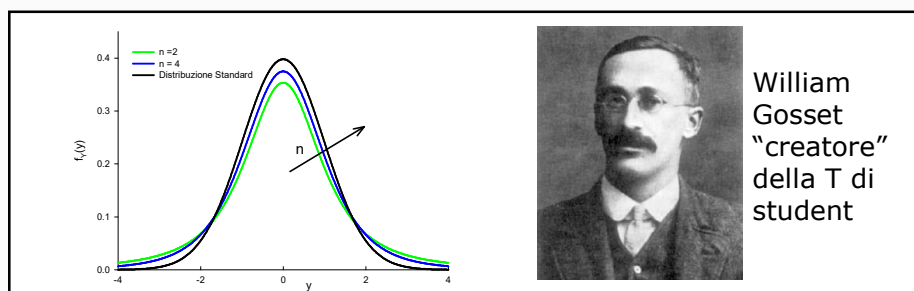
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

53

VA derivate dalla Gaussiana Distribuzione T di student

**Richiami di
statistica –
Esperimenti
replicati**

- In figura sono mostrate le funzioni densità per 1,3,6 gradi di libertà.



- La T è simmetrica rispetto a $y=0$
- Per $r \rightarrow +\infty$ la T di student tende ad una gaussiana di tipo standard.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

54

Variabile aleatoria di tipo t di student – Esempio

**Richiami di
statistica –
Esperimenti
replicati**

- Esempio di VA che segue la distribuzione di tipo t di student:
- Siano y_1, y_2, \dots, y_n un campione di dati generati da una VA di tipo Gaussiano $N(\mu, \sigma^2)$. Allora, la quantità:

$$t = \frac{\bar{y} - \mu}{\sqrt{S^2/n}}$$

- Segue una distribuzione di tipo t di student a $(n-1)$ g.d.l.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

55

VA derivate dalla Gaussiana La distribuzione F di Fisher

**Richiami di
statistica –
Esperimenti
replicati**

- Siano Y e W due VA di tipo χ^2 rispettivamente ad u e v gradi di libertà.
- Il rapporto

$$F_{u,v} = \frac{\chi_u^2 / u}{\chi_v^2 / v}$$

è una VA di tipo **F di Fisher** ad (u, v) gradi di libertà.

- La VA ha due parametri, u e v .

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

56

VA derivate dalla Gaussiana La distribuzione F di Fisher

**Richiami di
statistica –
Esperimenti
replicati**

- Espressione analitica della F di Fisher

$$f(y; u, v) = \frac{\Gamma\left(\frac{u+v}{2}\right)}{\Gamma\left(\frac{u}{2}\right)\Gamma\left(\frac{v}{2}\right)} \left(\frac{u}{v}\right)^{u/2} \frac{y^{\frac{u-2}{2}}}{\left(1 + \left(\frac{u}{v}\right)y\right)^{\frac{u+v}{2}}} \quad 0 < y < \infty$$

Media: $\mu_F = \frac{v}{v-2}, \quad (v > 2)$

Varianza: $\sigma_F^2 = \frac{2 v^2 (u+v-2)}{u(v-4)(v-2)^2} \quad v > 4$

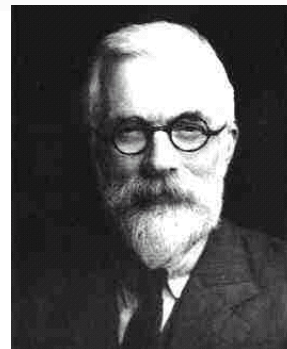
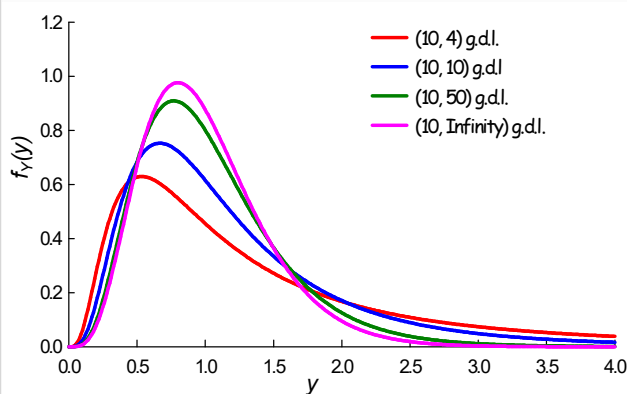
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

57

VA derivate dalla Gaussiana La distribuzione F di Fisher

**Richiami di
statistica –
Esperimenti
replicati**

- Grafici della F di Fisher al variare dei gradi di libertà



Sir Ronald Aylmer Fisher
1890 - 1962

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

58

Variabile aleatoria di tipo F di Fisher – Esempio

**Richiami di
statistica –
Esperimenti
replicati**

- Esempio di VA che segue la distribuzione di tipo F di Fisher:
- Siano:
 - $Y_{1,1}, Y_{1,2}, \dots, Y_{1,n_1}$ un campione di n_1 osservazioni provenienti da una data popolazione
 - $Y_{2,1}, Y_{2,2}, \dots, Y_{2,n_2}$ un campione di n_2 osservazioni provenienti da una altra popolazione
- Si suppone inoltre che la varianza σ^2 sia la stessa per entrambe le popolazioni. Allora:

$$\frac{S_1^2}{S_2^2} \approx F_{n_1-1, n_2-1}$$

- Dove S_1^2 e S_2^2 sono le due varianze campionarie calcolate per i due campioni

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

59

Analisi del campione di dati con strumenti statistici – Ulteriori sviluppi

**Richiami di
statistica –
Esperimenti
replicati**

- Nei prossimi lucidi si illustreranno delle tecniche utili per approfondire ulteriormente la conoscenza del campione di dati a disposizione:
 - **Test statistici delle ipotesi**
 - **Intervalli di fiducia**

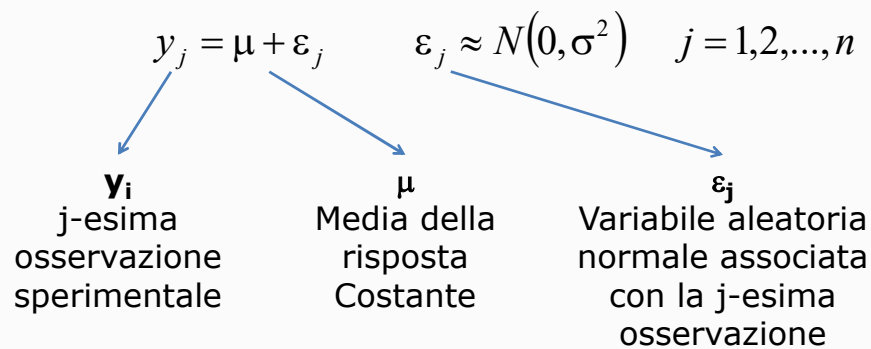
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

60

Analisi del campione di dati con strumenti statistici – Ulteriori sviluppi

Richiami di statistica – Esperimenti replicati

- **Modello statistico per il campione di dati:**



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

61

Test delle ipotesi – Introduzione

Richiami di statistica – Esperimenti replicati

- Torniamo al campione sperimentale di prodotti alimentari dell'esempio introduttivo.
- Da pregressi studi sull'impianto si sa che nella linea produttiva non sono graditi materiali troppo viscosi (motivi: perdite di carico, costi elevati di esercizio etc.).
- Da pregresse analisi si è stabilito un **valore di soglia** per la viscosità:

$$v=72.5$$

- al di sopra del quale risulta difficile la lavorazione del prodotto.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

62

Test delle ipotesi - Definizione

Richiami di statistica – Esperimenti replicati

- Un'**ipotesi statistica** è un'assunzione che noi facciamo sui parametri di una distribuzione o, equivalentemente, di un modello.
- L'ipotesi riflette qualche **congettura** sul problema in esame.
- Nel caso dell'esempio introduttivo, si vuole stabilire se
 - la viscosità della crema possa essere **almeno pari** al valore critico oppure
 - vi è una differenza **significativa** rispetto al valore $v=72.5$.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

63

Test statistici – Definizione del problema

Richiami di statistica – Esperimenti replicati

- Un **test statistico** di un'ipotesi è una procedura in cui si conclude se è possibile *non rigettare l'ipotesi* (cioè non si può escludere che essa sia vera) oppure *rigettare l'ipotesi*.
 - Si usa un campione e si cerca di concludere se tale campione è compatibile o meno con l'ipotesi nulla di partenza.
- Nell'esempio preso in considerazione, si vuole testare se il campione sperimentale possa derivare da una variabile aleatoria di media $v = 72.5$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

64

Test delle Ipotesi - Ipotesi nulla

Richiami di
statistica –
Esperimenti
replicati

- Il test delle ipotesi richiede l'introduzione di una **ipotesi nulla H_0** :

$$H_0: \mu = \mu_0 = 72.5$$

- In alternativa è possibile che la viscosità sia **realmente** minore del valore di soglia. Questa ipotesi, in **contrasto con l'ipotesi nulla**, è l'**ipotesi alternativa H_1** :

$$H_1: \mu < \mu_0 = 72.5$$

- Tutti i test delle ipotesi statistici richiedono la formulazione di un'ipotesi nulla e di un'ipotesi alternativa
- L'ipotesi nulla e l'ipotesi alternativa sono **esaustive e mutuamente esclusive**.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

65

Test statistici – Errori che si possono commettere nella procedura

Richiami di
statistica –
Esperimenti
replicati

- **Errore di tipo I (o errore α)**
- Probabilità di rigettare l'ipotesi nulla nonostante essa fosse vera

$$\alpha = P(\text{errore di tipo I}) = P(\text{rigetto } H_0 | H_0 \text{ è vera})$$

- è anche il **livello di significatività** del test.

- **Errore di tipo II (o errore β)**
- Probabilità di *non* rigettare l'ipotesi nulla nonostante essa fosse falsa

$$\beta = P(\text{errore di tipo II}) = P(\text{non rigetto } H_0 | H_0 \text{ è falsa})$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

66

Test statistici – Sviluppo della procedura

Richiami di statistica – Esperimenti replicati

- Parte della procedura consiste nel calcolo dell'insieme di valori che portano al rigetto di H_0 .
- Tale insieme di valori prende il nome di **regione critica** o **regione di rigetto** del test.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

67

Test statistici – Caso varianza σ^2 nota – Ricetta 1/4

Richiami di statistica – Esperimenti replicati

- N.B. Tale eventualità non è solo di interesse didattico: l'incertezza presente nelle misure sperimentali può essere nota a priori, per esempio da pregresse misure.
 - Per l'esempio si assume $\sigma^2=1$
1. Scegliere un **livello di significatività** α del test (in genere $\alpha=0.05$)
 2. Calcolare il **valore critico** z_α tale che:

$$P(Z \leq z_\alpha) = \alpha$$

- Nel caso in esame, per $\alpha=0.05$ si può leggere dalle tabelle $z_\alpha=-1.64485$

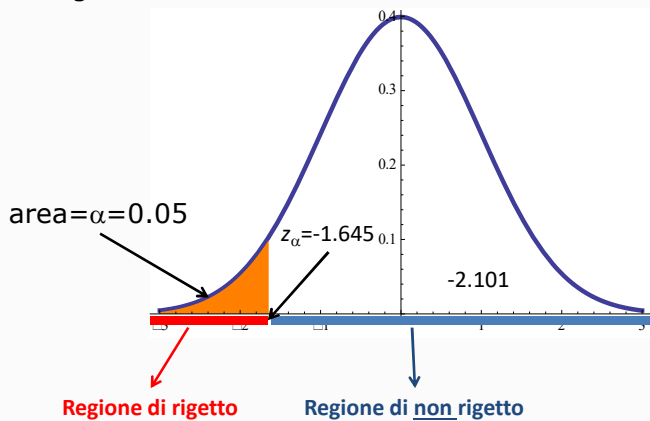
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

68

Test statistici – Esempio: Caso varianza nota – Ricetta 2/4

Richiami di statistica – Esperimenti replicati

- Distribuzione normale di tipo standard con l'evidenza delle regioni critiche



69

Test statistici – Esempio: Caso varianza nota – Ricetta 3/4

Richiami di statistica – Esperimenti replicati

- Calcolare

$$z_0 = \frac{\bar{y} - \mu_0}{\sqrt{\frac{\sigma^2}{n}}}$$

- Dove:
 - \bar{y} è la media campionaria
 - σ^2 è la varianza dell'errore sperimentale
 - n è la dimensione del campione

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

70

Test statistici – Esempio: Caso varianza nota – Ricetta 3/4

Richiami di statistica – Esperimenti replicati

- Si confronta il valore di z_0 osservato con il valore critico z_α

$$z_0 > z_\alpha$$

- **non** rigettiamo l'ipotesi nulla H_0 : non si hanno evidenze sperimentali tali da affermare che la media sia **significativamente minore del valore di riferimento**

$$z_0 < z_\alpha$$

- Si **rigetta** l'ipotesi nulla: la media è **significativamente minore** di μ_0 .
- Il «rischio» di affermare la conclusione sbagliata è pari al livello di significatività α del test

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

71

Test delle Ipotesi sulla media - Teoria

Richiami di statistica – Esperimenti replicati

- Caso varianza σ^2 nota
- Se l'ipotesi nulla

$$H_0 : \mu = \mu_0$$

- **fosse vera**, la variabile aleatoria media campionaria

$$\bar{Y} = \frac{\sum Y_i}{n}$$

- si comporterebbe come una distribuzione gaussiana di media μ_0 e varianza σ^2/n

$$\bar{Y} \approx N\left(\mu_0, \frac{\sigma^2}{n}\right)$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

72

Test delle Ipotesi sulla media - Teoria

Richiami di statistica – Esperimenti replicati

- Pertanto, se H_0 fosse vera, la variabile aleatoria

$$Z = \frac{\bar{Y} - \mu_0}{\sqrt{\frac{\sigma^2}{n}}}$$

- sarebbe **una distribuzione normale di tipo standard** e il valore osservato z_0 sarebbe un esito che rispetta tale VA.

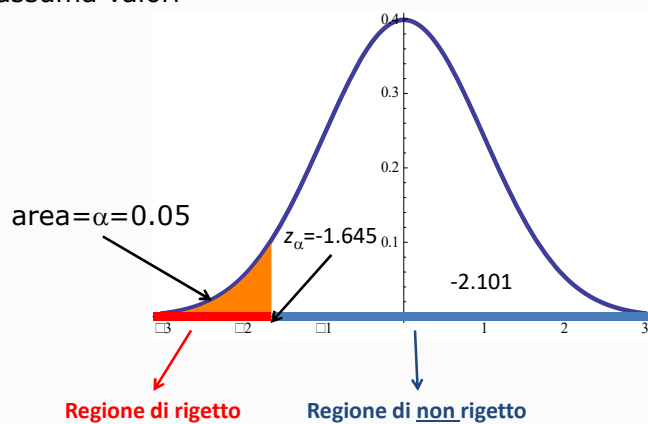
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

73

Test delle Ipotesi sulla media - Teoria

Richiami di statistica – Esperimenti replicati

- Al di sopra di z_α è poco plausibile che la variabile aleatoria Z assuma valori



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

74

Test delle ipotesi sulla media - Esempio

Richiami di
statistica –
Esperimenti
replicati

- Si consideri di nuovo l'esempio.
- Il test delle ipotesi è sul valore medio:

$$H_0 : \mu = \mu_0$$

$$H_1 : \mu < \mu_0$$

- Con un livello di significatività $\alpha = 5\%$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

75

Test delle ipotesi sulla media - Esempio

Richiami di
statistica –
Esperimenti
replicati

- Si valuta innanzitutto il valore z_α tale che $P(Z < z_\alpha) = \alpha = 0.05$.

$$P(Z < z_\alpha) = \alpha \Rightarrow z_\alpha = -1.645$$

- **Se** l'ipotesi nulla **fosse vera**, il risultato

$$z_0 = \frac{\bar{y} - \mu_0}{\sigma} \sqrt{n} = \frac{71.69 - 72.5}{1} \sqrt{10} = -2.568$$

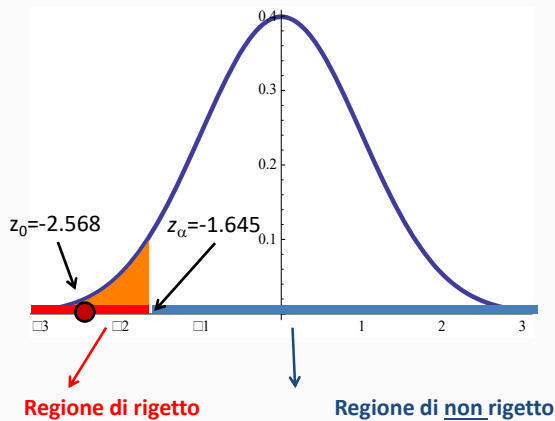
- **sarebbe** un valore osservato di una variabile aleatoria **normale di tipo standard**.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

76

Test delle ipotesi sulla media - Esempio

Richiami di statistica – Esperimenti replicati



Il valore osservato z_0 rientra nella regione in cui la variabile aleatoria ha poche probabilità di cadere

C'è un 5% di probabilità che il valore osservato appartenga alla VA supposta nell'ipotesi nulla H_0 e sia comunque rigettata

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

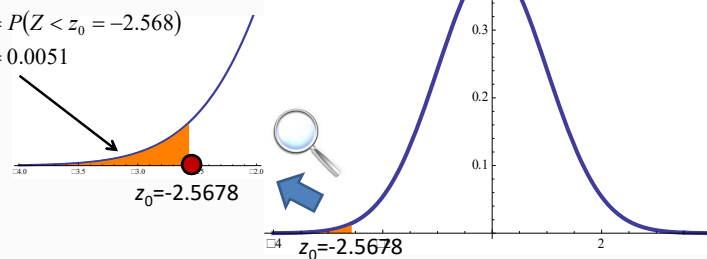
77

Test statistici – Uso del p-value

Richiami di statistica – Esperimenti replicati

- Approccio **alternativo** a quello classico dell'individuazione delle zone di rigetto.
- Il **p-value** rappresenta la probabilità che la statistica test stimata assuma un valore almeno uguale al valore osservato della statistica nel caso in cui l'ipotesi nulla fosse vera.
- Nel caso dell'esempio:

$$p\text{-value} = P(Z < z_0 = -2.568) = 0.0051$$



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

78

Test statistici – Uso del p-value

**Richiami di
statistica –
Esperimenti
replicati**

- **Pro**
- Informazione più quantitativa
- **Contro:**
- Necessita di calcolatori con programmi specifici (o comunque competenze di programmazione avanzata)

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

79

Test delle Ipotesi - Ipotesi alternative 1/4

**Richiami di
statistica –
Esperimenti
replicati**

- Nel problema in esame si assume che il nostro campione di dati sperimentali sia caratterizzato da una variabile aleatoria che abbia una funzione densità di probabilità che coinvolge un parametro ignoto θ e si assume l'ipotesi nulla che

$$H_0 : \theta = \theta_0$$

- L'ipotesi alternativa era:

$$H_1 : \theta < \theta_0$$

- Ma non è l'unica alternativa che possiamo considerare.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

80

Test delle Ipotesi - Ipotesi alternative 2/4

Richiami di statistica – Esperimenti replicati

- In altri casi la natura può suggerire altri tipi di alternative:
 $H_1 : \theta > \theta_0$
- Oppure
 $H_1 : \theta \neq \theta_0$
- Le prime 2 alternative si chiamano **one-sided**. L'ultima **two-sided**

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

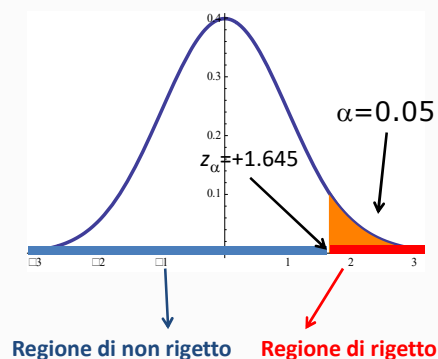
81

Test delle Ipotesi – Ipotesi alternative 3/4

Richiami di statistica – Esperimenti replicati

- Nel caso di ipotesi alternativa
 $H_1 : \theta > \theta_0$
- Si deve determinare il valore critico z_α tale che tutti i valori superiori ad esso abbiano una probabilità di verificarsi pari a α
- Dobbiamo escludere i valori per cui la distribuzione gaussiana standard assume valori tali che

$$P(Z > z_\alpha) = \alpha$$



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

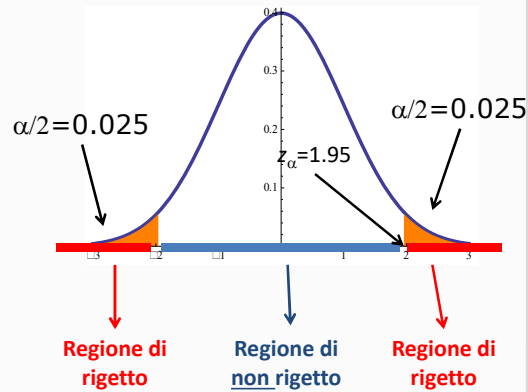
82

Test delle Ipotesi – Ipotesi alternative 4/4

Richiami di statistica – Esperimenti replicati

- Nel caso di ipotesi alternativa
 $H_1: \theta \neq \theta_0$
- Ricordiamo che è una ipotesi alternativa «two-sided»
- Si deve determinare il valore critico z_α tale che

$$P(|Z| > z_\alpha) = \alpha$$



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

83

Test delle ipotesi sulla media - Varianza ignota

Richiami di statistica – Esperimenti replicati

- Nel caso in cui non fosse nota la varianza s^2 non è possibile sfruttare la statistica per determinare i valori critici dei test statistici

$$z = \sqrt{n} \frac{\bar{Y} - \mu_0}{\sigma}$$

- È possibile ricorrere alla stima S^2 della varianza campionaria

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

- Se l'ipotesi nulla fosse vera, allora la variabile aleatoria

$$t = \frac{\bar{Y} - \mu_0}{\sqrt{S^2/n}}$$

- Sarebbe una distribuzione t di student ad $(n-1)$ gdl.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

84

Test delle ipotesi sulla media – Varianza ignota

Richiami di statistica – Esperimenti replicati

- Ricetta
- Fissare un livello di significatività del test (es: $\alpha = 5\%$)
- Calcolare il valore t_α per cui:

$$P(t \leq t_\alpha) = \alpha$$

- dove t è la distribuzione di student ad $r = n - 1$ gradi di libertà.
- Calcolare S^2 :

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

- Calcolare

$$t_0 = \sqrt{n} \frac{\bar{y} - \mu_0}{\sqrt{S^2}}$$

- $t_0 < t_\alpha$: rigettare H_0
- $t_0 > t_\alpha$: non rigettare H_0 .

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

85

Test delle ipotesi sulla media – Varianza ignota – Esercizio

Richiami di statistica – Esperimenti replicati

- Ritorniamo al campione in esame
- Si fissa un livello di significatività $\alpha = 0.05$ per il test
- Dalle tabelle si determina il valore t_α :

$$P(t \leq t_\alpha) = 0.05 \Rightarrow t_{\alpha, 19} = -1.833$$

- Si calcola il valore stimato per la varianza:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = 0.983$$

- Da cui è possibile calcolare la statistica t_0 :

$$t_0 = \sqrt{n} \frac{\bar{y} - \mu_0}{\sqrt{S^2}} = -2.589$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

86

Test delle ipotesi sulla media – Varianza ignota – Esercizio

Richiami di statistica – Esperimenti replicati

- Quindi

$$t_0 < t_{\alpha,19}$$

- Si **rigetta** l'ipotesi nulla.
- Alternativamente, è possibile calcolare il **p-value**

$$P(t_r \leq t_0 = -2.59) = 0.0146$$

- Da notare come il p-value sia più elevato rispetto a quello stimato nel caso della varianza nota
 - La mancanza di informazioni sul processo si riflette in delle conclusioni più incerte.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

87

Test delle ipotesi sulla media – Altre ipotesi alternative

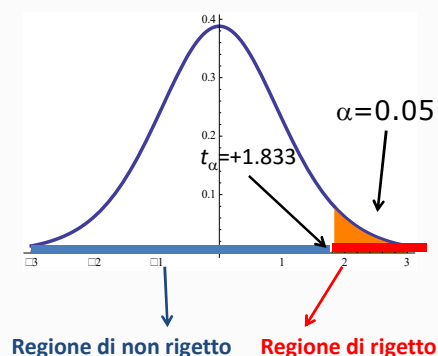
Richiami di statistica – Esperimenti replicati

- Nel caso di ipotesi alternativa

$$H_1: \mu > \mu_0$$

- Si deve determinare il valore critico t_α tale che tutti i valori superiori ad esso abbiano una probabilità di verificarsi pari a α
- Dobbiamo escludere i valori per cui la t di student assuma valori tali che

$$P(t_r > t_\alpha) = \alpha \quad (r = 9 \text{ gdl})$$



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

88

Test delle ipotesi sulla media – Altre ipotesi alternative

Richiami di statistica – Esperimenti replicati

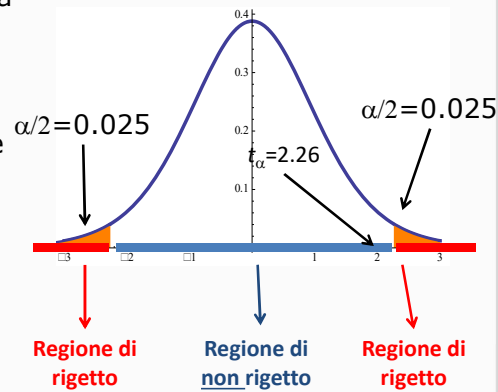
- Nel caso di ipotesi alternativa «two sided»

$$H_1: \mu \neq \mu_0$$

- Si deve determinare il valore critico z_α per cui

$$P(|T| > t_{r,\alpha}) = \alpha$$

$$P(t_r > t_\alpha) = \alpha \quad (r = 9 \text{ gdl})$$



Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

89

Intervalli di fiducia - Introduzione

Richiami di statistica – Esperimenti replicati

- Nell'esaminare un campione di dati sperimentali, si può essere interessati ad un'informazione più qualitativa di una semplice stima puntuale di parametri.
- Ad esempio, si può essere interessati a determinare un **intervallo** di valori in cui è molto probabile cada il valore vero del parametro.
- Tale tipo di inferenza prende il nome di **inferenza di intervallo** e il risultato della procedura è un **intervallo di fiducia** (anche denominato **intervallo di confidenza**)
- Per esempio, si può essere interessati ad un intervallo di fiducia per la media μ della viscosità.

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

90

Intervalli di fiducia Procedura

**Richiami di
statistica –
Esperimenti
replicati**

- Si suppone che θ sia il parametro incognito da stimare
- Si sceglie una probabilità γ vicina a 1 (in genere $\gamma=0.95$). Tale probabilità prende il nome di **livello di fiducia**.
- In seguito si determinano due quantità L e U tali che

$$P(L \leq \theta \leq U) = \gamma$$

- L'intervallo di estremi L e U prende il nome di **intervallo di fiducia** e si indica con il simbolo:

$$CONF\{L \leq \theta \leq U\}$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

91

Intervalli di fiducia della Media – Caso varianza non nota.

**Richiami di
statistica –
Esperimenti
replicati**

Determinazione intervallo di fiducia:

1. Scegliere un livello di fiducia $\gamma=1-\alpha$
2. Ricavare (per esempio da tabelle) il valore $t_{\alpha/2}$ tale che:

$$P(-t_{\alpha/2} \leq T_r \leq t_{\alpha/2}) = \gamma = 1 - \alpha$$

essendo T_r la T di student a $r=n-1$ gdl

3. Calcolare media e varianza del campione dei dati sperimentali.
3. L'intervallo di fiducia per la media sarà:

$$CONF\left\{\bar{y} - t_{\alpha/2} \sqrt{\frac{S^2}{n}} \leq \mu \leq \bar{y} + t_{\alpha/2} \sqrt{\frac{S^2}{n}}\right\}$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

92

Intervalli di fiducia della Media nel caso di varianza non nota.

Richiami di statistica – Esperimenti replicati

- La variabile aleatoria:

$$Z = \sqrt{n} \frac{\bar{Y} - \mu}{\sigma}$$

- È una variabile normale di tipo standard
- Si può ulteriormente dimostrare che la variabile aleatoria:

$$W = \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \bar{Y})^2 = (n-1) \frac{S^2}{\sigma^2} \approx \chi_{n-1}^2$$

- È una variabile aleatoria χ^2 a $n-1$ gradi di libertà

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

93

Intervalli di fiducia della Media nel caso di varianza non nota.

Richiami di statistica – Esperimenti replicati

- In conclusione la variabile aleatoria:

$$T = \frac{Z}{\sqrt{W/(n-1)}} = \sqrt{n} \frac{\frac{\bar{Y} - \mu}{\sigma}}{\sqrt{\frac{\sum (Y_i - \bar{Y})^2}{\sigma^2} / (n-1)}} = \frac{\bar{Y} - \mu}{\sqrt{\frac{S^2}{n}}}$$

- È una variabile aleatoria di tipo t di student ad $n-1$ gradi di libertà

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

94

Intervalli di fiducia della Media nel caso di varianza non nota.

**Richiami di
statistica –
Esperimenti
replicati**

- Dalla definizione di probabilità è possibile ricavare la relazione:

$$P(-t_{\alpha/2} \leq T_r \leq t_{\alpha/2}) = P\left(-t_{\alpha/2} \leq \frac{\bar{y} - \mu}{\sqrt{\frac{S^2}{n_1}}} \leq t_{\alpha/2}\right) = \gamma$$

- da cui con qualche passaggio è possibile ricavare l'intervallo di fiducia desiderato:

$$P\left(\bar{y} - t_{\alpha/2} S \sqrt{\frac{1}{n}} \leq \mu \leq \bar{y} + t_{\alpha/2} S \sqrt{\frac{1}{n}}\right) = \gamma \quad \text{CONF} \left\{ \underbrace{\bar{y} - t_{\alpha/2} S \sqrt{\frac{1}{n}}}_{L} \leq \mu \leq \bar{y} + t_{\alpha/2} S \sqrt{\frac{1}{n}} \right\}$$

L **θ** **U**

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

95

Intervalli di fiducia della Media – Esercizio

**Richiami di
statistica –
Esperimenti
replicati**

Determinazione intervallo di fiducia:

- Si sceglie un livello di fiducia $\gamma=95\%$
- Ricavare il valore $t_{\alpha/2}$ tale che:

$$P(-t_{\alpha/2} \leq T_r \leq t_{\alpha/2}) = 95\% \Rightarrow t_{\alpha/2} = 2.262$$

essendo T_r la T di student a $r=9$ gdl

- Calcolare media e varianza del campione dei dati sperimentali.

$$\bar{y} = 71.69, \quad S^2 = 0.9834$$

- L'intervallo di fiducia per la media sarà:

$$\text{CONF} \left\{ 71.69 - 2.262 \sqrt{\frac{0.9834}{10}} \leq \mu \leq 71.69 + 2.262 \sqrt{\frac{0.9834}{10}} \right\} = \{70.98 \leq \mu \leq 72.4\}$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

96

Intervalli di fiducia della Media – Esercizio

Richiami di statistica – Esperimenti replicati

- Da notare che nell'intervallo di fiducia calcolato non ricade il valore 72.5, confermando che tale valore è molto improbabile per la media della popolazione.
- In generale, si deve ricordare che, per le proprietà di simmetria della t di student:

$$P(T_r \leq -t_{\alpha/2}) = P(T_r \geq +t_{\alpha/2})$$

- Il valore di $t_{\alpha/2}$ può essere calcolato anche dalla relazione:

$$P(T_r \leq t_{\alpha/2}) = \frac{1}{2}(1 + \gamma)$$

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

97

Diagramma in scala probabilistica

Richiami di statistica – Esperimenti replicati

- Da notare che il modello statistico preso in considerazione parte dall'assunzione che i dati sperimentali seguano una distribuzione di tipo Gaussiano.
- Tale assunzione può essere verificata costruendo un **diagramma in scala probabilistica**.
- La procedura è abbastanza semplice e consiste in un'analisi di tipo **grafico**.
- Per costruire il diagramma si deve:
 - ordinare i dati dal più piccolo al più grande
 - le osservazioni così ordinate sono rappresentate rispetto la loro frequenza cumulativa osservata
 - la scala in ordinata non è lineare ma è tale che, se i dati rispettassero una dispersione di tipo **Gaussiano**, essi si disporrebbero approssimativamente **lungo una retta**

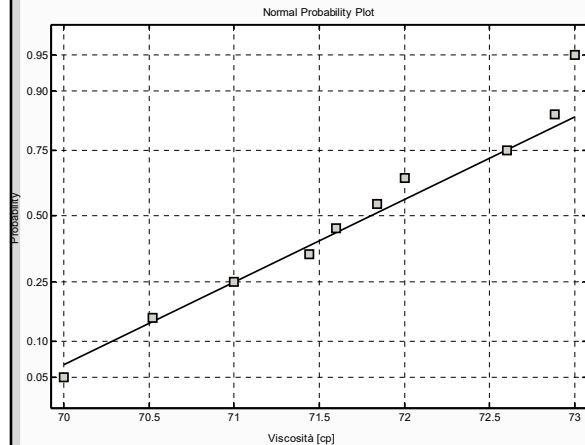
Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

98

Diagramma in scala probabilistica

Richiami di statistica – Esperimenti replicati

- Esempio dati crema



- In linea di principio, è possibile implementare il metodo a mano, ma risulta molto pesante.
- La maggior parte dei software di uso comune supportano la rappresentazione su carta probabilistica.

Metodi statistici per l'analisi dei dati
19-23 settembre 2017

99

Conclusioni – Concetti importanti

Richiami di statistica – Esperimenti replicati

- Esperimento come esito di una variabile aleatoria
 - VA di tipo Gaussiano
- Campagna sperimentale esito di una variabile aleatoria
 - VA di tipo student (o, in casi fortunati, di tipo Gaussiano)
- Con gli strumenti della statistica è possibile inferire conclusioni rigorose sul processo.
- Sono stati introdotti i concetti (verranno ampiamente ripresi nel seguito):
 - Test statistici
 - Intervalli di fiducia
- Diagrammi in scala probabilistica

Metodi statistici per l'analisi dei dati
10-14 febbraio 2020

100