

# Metodi statistici per l'analisi dei dati

Confronto tra due trattamenti

1

## Motivazioni

**Confronto  
tra  
trattamenti**

- Obiettivo:
  - Confrontare due diverse condizioni (anche definiti trattamenti) per cui sono stati condotti gli esperimenti.

2

## Esempio introduttivo

### Confronto tra trattamenti

- L'impasto dell'esempio introduttivo è modificato aggiungendo un additivo che dovrebbe ridurre la viscosità.
- A tal proposito si prendono in considerazione
  - 10 misure sperimentali di viscosità del trattamento originale (**controllo**)
  - 10 misure sperimentali di viscosità del prodotto alimentare modificato.

j	Crema modificata (cp)	Controllo (cp)
1	67.40	70.00
2	65.60	70.52
3	68.84	73.00
4	65.40	72.00
5	66.08	71.44
6	68.16	71.00
7	67.84	72.88
8	68.60	71.60
9	66.36	71.84
10	66.28	72.60
$\bar{y}$	<b>67.06</b>	<b>71.69</b>

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

3

3

## Introduzione – Esempio

### Confronto tra trattamenti

- Il valore medio delle misure del prodotto modificato è minore del valore medio del prodotto originale
- Due possibili spiegazioni:
  - Le due formulazioni sono realmente differenti
  - La differenza osservata nei trattamenti è dovuta alle fluttuazioni inevitabilmente presenti nelle misure sperimentali, e le due formulazioni sono di fatto equivalenti
- **Obiettivo:**
- Stabilire in maniera oggettiva se i due campioni sono **realmente** differenti ricorrendo a strumenti statistici rigorosi

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

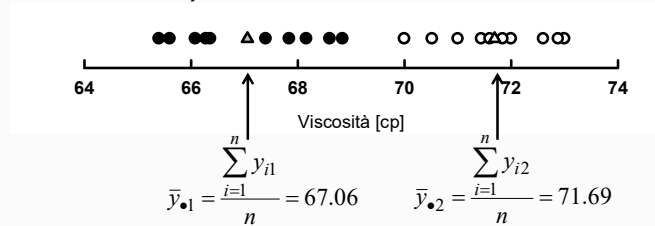
4

4

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

### Confronto tra trattamenti

- **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.
- Un'ispezione visiva (**qualitativa**) del diagramma rivela che i due trattamenti sono verosimilmente differenti.

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

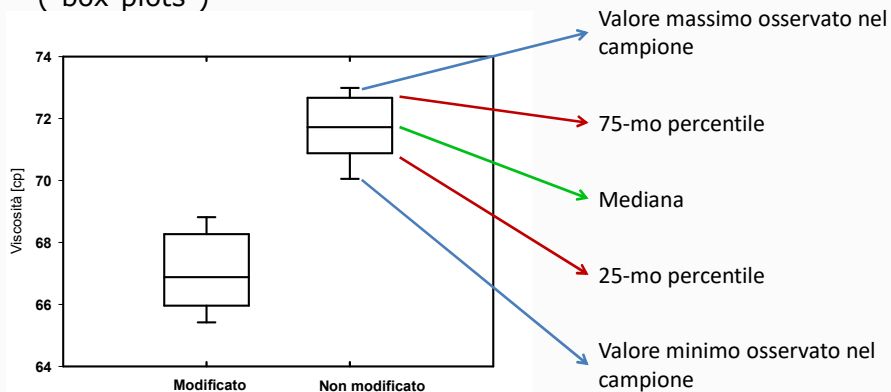
5

5

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

### Confronto tra trattamenti

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



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

6

6

## Confronto tra campioni

**Confronto  
tra  
trattamenti**

- I due diversi trattamenti possono essere visti come due **livelli del fattore** "formulazioni"
- **Modello statistico per i dati:**

$$y_{ij} = \mu_i + \varepsilon_{ij} \quad \left\{ \begin{array}{l} i = 1,2 \\ j = 1,2,\dots,n_i \end{array} \right.$$

**$y_{ij}$**   
j-esima  
osservazione dal  
livello i-esimo

**$\mu_i$**   
Media della  
risposta al livello  
i-esimo

**$\varepsilon_{ij}$**   
Variabile aleatoria  
normale associata  
con la j-esima  
osservazione

Metodi statistici per l'analisi dei dati  
10 - 14 febbraio 2020
7

7

## Confronto tra campioni

**Confronto  
tra  
trattamenti**

- **Assunzione:**
- Gli errori  $\varepsilon$  sono normali, indipendenti e identicamente distribuiti (**i.i.d.**):

$$\varepsilon \approx N(0, \sigma^2)$$

Metodi statistici per l'analisi dei dati  
10 - 14 febbraio 2020
8

8

## Confronto tra campioni – Test statistici – Definizione del problema

## Confronto tra trattamenti

- Dal punto di vista formale, si possono definire le seguenti **ipotesi nulla  $H_0$**  e l'**ipotesi alternativa  $H_1$** :

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

9

## Confronto tra campioni – Test statistici – Definizione del problema

## Confronto tra trattamenti

- **Tipologia di errori** che possono essere commessi durante un test statistico:
- **Errore di tipo I: rigettare** l'ipotesi  $H_0$  di partenza nonostante essa fosse **vera**  
$$\alpha = P(\text{errore di tipo I}) = P(\text{rigetto } H_0 \mid H_0 \text{ è vera})$$
- **Errore di tipo II: non rigettare**  $H_0$  nonostante essa fosse **falsa**  
$$\beta = P(\text{errore di tipo II}) = P(\text{non rigetto } H_0 \mid H_0 \text{ è falsa})$$
- Si deve ridurre al minimo il rischio (le probabilità  $\alpha$  e  $\beta$ ) di incorrere in questi due tipi di errore

10

## Test statistici – Esempio: t-Test per il confronto di due campioni – Ricetta 1/4

### Confronto tra trattamenti

- Scegliere un **livello di significatività**  $\alpha$  del test (in genere  $\alpha=0.05$ ).
- Rappresenta la probabilità di sbagliare, nel caso in cui arrivassimo alla conclusione di rigetto dell'ipotesi nulla.
- Calcolare il **valore critico**  $t_{\alpha/2}$  tale che:

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

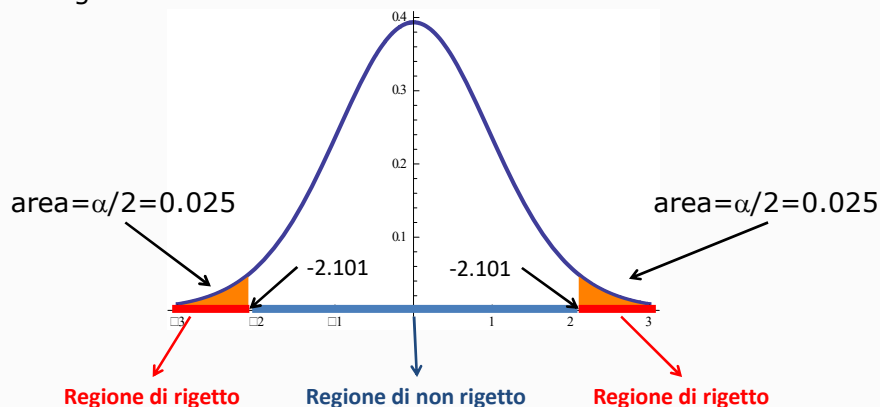
dove  $T$  è la distribuzione  $t$  di Student a  $(n_1+n_2-2)$  gradi di libertà

11

## Test statistici – Esempio: t-Test per il confronto di due campioni – Ricetta 2/4

### Confronto tra trattamenti

- Distribuzione  $T$  di student a 18 gdl con l'evidenza delle regioni critiche



12

## Test statistici – Esempio: t-Test per il confronto di due campioni – Ricetta 3/4

Confronto tra trattamenti

- Calcolare

$$t_0 = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{S_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

- dove  $S_p^2$  è la stima della varianza campionaria comune  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ :

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

- Nota: nel caso di  $n = n_1 = n_2$ , l'espressione si riduce a

$$t_0 = \sqrt{n} \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{S_1^2 + S_2^2}}$$

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

13

13

## Test statistici – Esempio: t-Test per il confronto di due campioni – Ricetta 4/4

Confronto tra trattamenti

- Si confronta il  $t_0$  osservato con il valore critico  $t_{\alpha/2}$

$$|t_0| < t_{\alpha/2}$$

- **non** rigettiamo l'ipotesi nulla  $H_0$ .
- Non si hanno evidenze sperimentali tali da affermare che i due trattamenti siano **significativamente diversi**

$$|t_0| > t_{\alpha/2}$$

- Si **rigetta** l'ipotesi nulla: i due trattamenti sono significativamente diversi
- Si corre un «rischio» di affermare la conclusione sbagliata pari al livello di significatività  $\alpha$  del test

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

14

14

## Test statistici – Esempio: t-Test per il confronto di due campioni – Teoria

### Confronto tra trattamenti

- La differenza delle medie dei due trattamenti può essere vista come una VA di tipo normale:

$$\bar{y}_1 - \bar{y}_2 \sim N[\mu_1 - \mu_2, \sigma^2(1/n_1 + 1/n_2)]$$

- Se** i due trattamenti provenissero dalla **stessa popolazione** (ovvero  $\mu_1 = \mu_2$ ,  $\sigma_1^2 = \sigma_2^2$ ), e la **varianza  $\sigma^2$  fosse nota**, la statistica:

$$z_0 = \frac{\bar{y}_1 - \bar{y}_2}{\sigma \sqrt{1/n_1 + 1/n_2}}$$

- sarebbe un valore osservato di una Gaussiana di tipo standard ( $Z \sim N(0,1)$ )
- $t_0$  è un'osservazione di una  $T$  di student a  $(n_1 + n_2 - 2)$  gdl:

15

## Test statistici – Esempio: t-Test per il confronto di due campioni – Teoria

### Confronto tra trattamenti

- Dato che la **varianza è ignota**, si ricorre alla sua stima  $S_p^2$ , per cui la statistica  $t_0$ :

$$t_0 = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{S_p^2(1/n_1 + 1/n_2)}}$$

- è relativa ad una **T di student a  $(n_1 + n_2 - 2)$  g.d.l.**
- Se  $H_0$  fosse vera**, ci si aspetta quindi che la probabilità per  $t_0$  di cadere nell'intervallo  $[-t_{\alpha/2}, t_{\alpha/2}]$  sia pari a  $100(1-\alpha)$ .
- Un campione che produce dei risultati al di fuori di questo intervallo non sarebbe usuale
  - è quindi più plausibile che l'ipotesi nulla di partenza sia sbagliata

16



## Test statistici – Esempio: t-Test per il confronto di due campioni – Applicazione

Confronto tra trattamenti

- Per le emulsioni alimentari dell'esempio
- Calcolo valore critico  $t_{\alpha/2}$ :
  - (da tabelle disponibili in letteratura o, equivalentemente, con l'ausilio di software)

$$t_{\alpha/2} = 2.101$$

- Calcolo statistica  $t_0$ :

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} = \frac{9 \cdot 1.602 + 9 \cdot 0.983}{18} = 1.2928$$

$$t_0 = \frac{\bar{y}_1 - \bar{y}_2}{S_p \sqrt{1/n_1 + 1/n_2}} = \frac{67.06 - 71.69}{1.137 \sqrt{1/10 + 1/10}} = -9.109$$

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

17

17

## Test statistici – Esempio: t-Test per il confronto di due campioni – Applicazione

Confronto tra trattamenti

- Conclusioni:

$$|t_0| \gg t_{\alpha/2}$$

- Il valore osservato  $t_0$  è maggiore (in valore assoluto) del valore critico  $t_{\alpha/2}$



- I **due trattamenti sono significativamente differenti** con un **rischio (errore di tipo I) del 5%** che tale conclusione sia sbagliata

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

18

18

## Test statistici – Definizione ipotesi alternative

### Confronto tra trattamenti

- L'ipotesi alternativa presa in considerazione

$$H_1: \mu_1 \neq \mu_2$$

- contempla l'eventualità che i due trattamenti presi in considerazione siano solo **differenti**

$$H_1: \mu_1 < \mu_2 \text{ oppure } \mu_1 > \mu_2$$

- Tale ipotesi prende il nome di **ipotesi alternativa "two-sided"** (valori sia maggiori che minori portano al rigetto di  $H_0$ ).
- In altri casi, la natura del problema può suggerire altre espressioni per le ipotesi alternative.

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

19

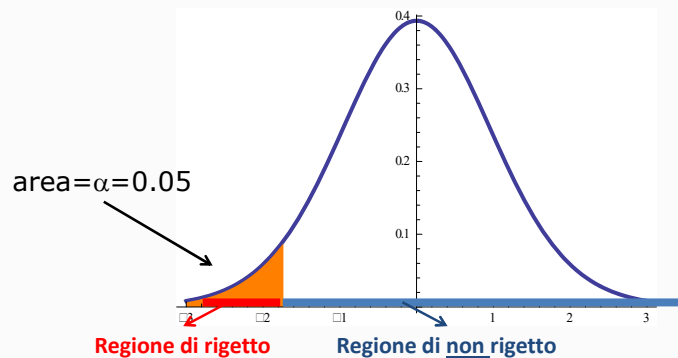
19

## Test statistici – Definizione ipotesi alternative «one sided»

### Confronto tra trattamenti

$$H_1: \mu_1 < \mu_2$$

- Si arriva al rigetto dell'ipotesi nulla solo se  $\mu_1$  è significativamente **minore** di  $\mu_2$



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

20

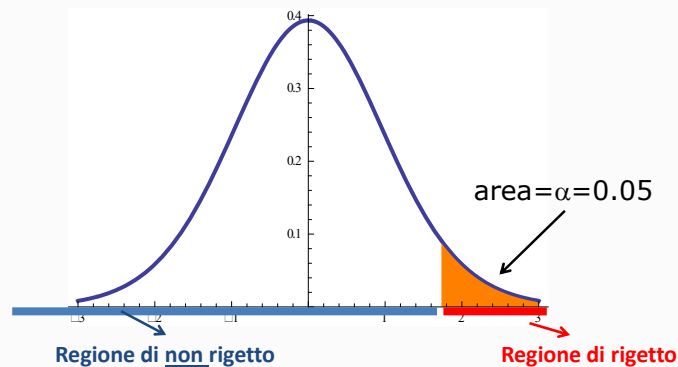
20

## Test statistici – Definizione ipotesi alternative «one sided»

Confronto tra trattamenti

$$H_1: \mu_1 > \mu_2$$

- Si arriva al rigetto dell'ipotesi nulla solo se  $\mu_1$  è significativamente **maggiore** di  $\mu_2$



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

21

21

## Intervalli di fiducia – Differenza di medie

Confronto tra trattamenti

- Supponiamo di essere interessati a determinare un intervallo di fiducia al 95% per la differenza  $\delta_1 = \mu_1 - \mu_2$  delle medie.
- A tale scopo, si consideri la statistica:

$$t = \frac{d_1 - \delta_1}{S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{(\bar{y}_1 - \bar{y}_2) - (\mu_1 - \mu_2)}{S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

- essa è distribuita secondo una  $t_{n_1+n_2-2}$ .

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

22

22

## Intervalli di fiducia – Differenza di medie

Confronto  
tra  
trattamenti

- Da cui:

$$P\left\{-t_{\alpha/2} \leq \frac{(\bar{y}_1 - \bar{y}_2) - (\mu_1 - \mu_2)}{S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \leq t_{\alpha/2}\right\} = \gamma$$

- ovvero:

$$P\left\{(\bar{y}_1 - \bar{y}_2) - t_{\alpha/2} S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \leq (\mu_1 - \mu_2) \leq (\bar{y}_1 - \bar{y}_2) + t_{\alpha/2} S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}\right\} = \gamma$$

- Per cui

$$CONF\left\{\underbrace{(\bar{y}_1 - \bar{y}_2) - t_{\alpha/2} S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}_L \leq (\mu_1 - \mu_2) \leq \underbrace{(\bar{y}_1 - \bar{y}_2) + t_{\alpha/2} S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}_U\right\}$$

$\theta$

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

23

23

## Intervalli di fiducia differenza di due medie - Applicazione

Confronto  
tra  
trattamenti

- Per l'esempio introduttivo

$$CONF\left\{\underbrace{(67.06 - 71.69)}_{(\bar{y}_1 - \bar{y}_2)} - \underbrace{(2.101)}_{t_{\alpha/2}} \underbrace{1.137}_{S_P} \underbrace{\sqrt{\frac{1}{10} + \frac{1}{10}}}_{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \leq (\mu_1 - \mu_2) \leq \underbrace{(67.06 - 71.69)}_{(\bar{y}_1 - \bar{y}_2)} + \underbrace{(2.101)}_{t_{\alpha/2}} \underbrace{1.137}_{S_P} \underbrace{\sqrt{\frac{1}{10} + \frac{1}{10}}}_{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}\right\}$$

- ovvero:

$$CONF\{-5.70 \leq \mu_1 - \mu_2 \leq -3.56\}$$

- Da notare che il valore  $\mu_1 - \mu_2 = 0$  non è incluso nell'intervallo, confermando ulteriormente che l'ipotesi  $\mu_1 = \mu_2$  non è plausibile per il livello di significatività  $\alpha = 1 - \gamma = 5\%$ .

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

24

24

## Test statistici – Esempio: t-Test per il confronto di due campioni – Estensioni

Confronto tra trattamenti

- $\sigma_1^2 \neq \sigma_2^2$  – **Parte 1/2**
- In presenza di un test delle ipotesi:

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

- in cui non è possibile assumere che le varianze dei campioni siano uguali, si può ricorrere alla statistica test

$$t_0 = \frac{(\bar{y}_1 - \bar{y}_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

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

25

25

## Test statistici – Esempio: t-Test per il confronto di due campioni – Estensioni

Confronto tra trattamenti

- **Caso in cui  $\sigma_1^2 \neq \sigma_2^2$  – Parte 2/2**
- essendo la statistica  $t$  ben approssimata da una  $t$  di student a  $\nu$  gradi di libertà:

$$\nu = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{(S_1^2/n_1)^2}{n_1 - 1} + \frac{(S_2^2/n_2)^2}{n_2 - 1}}$$

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

26

26

## Test statistici – Esempio: Test per il confronto di due campioni – Estensioni

Confronto  
tra  
trattamenti

- $\sigma_1^2$  e  $\sigma_2^2$  **note – 1/2**
- In presenza di un test delle ipotesi:

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

- in cui le varianze sono a priori note, si può usare la statistica:

$$z_0 = \frac{(\bar{y}_1 - \bar{y}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

- Se entrambi le popolazioni sono normali (o le dimensioni dei campioni sono grandi) la statistica  $z_0 \sim N(0,1)$

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

27

27

## Test statistici – Esempio: Test per il confronto di due campioni – Estensioni

Confronto  
tra  
trattamenti

- $\sigma_1^2$  e  $\sigma_2^2$  **note – 2/2**
- La regione critica per i test delle ipotesi può essere calcolata usando la distribuzione **normale** di tipo **standard** anziché la t di student. Sarà necessario calcolare il valore critico  $z_\alpha$  tale che, per esempio (caso test ipotesi two-sided):

$$P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$$

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

28

28

## Esempio introduttivo – Confronto medie con “blocco”

### Confronto tra trattamenti

- Si consideri di nuovo il caso dell’esempio introduttivo.
- La randomizzazione dei campioni ha portato alla selezione casuale di:
  - 10 campioni da modificare
  - 10 campioni non trattati che saranno usati come controllo
- Tale politica implica una inevitabile sorgente di variazione:
  - la **casualità** nella **scelta** dei campioni da destinare ai due trattamenti si aggiunge alla eventuale variazione indotta dalla differenza dei trattamenti

Metodi statistici per l’analisi dei dati  
10 – 14 febbraio 2020

29

29

## Esempio introduttivo – Confronto medie con “blocco”

### Confronto tra trattamenti

- Strategia possibile per eliminare tale sorgente di incertezza:
  1. Si selezionano 10 campioni (anziché 20) non modificati
  2. Se ne misura la viscosità
  3. In seguito, gli **stessi** campioni sono modificati con l’aggiunta dell’additivo.
  4. Si ripete la misura della viscosità sui campioni modificati

Metodi statistici per l’analisi dei dati  
10 – 14 febbraio 2020

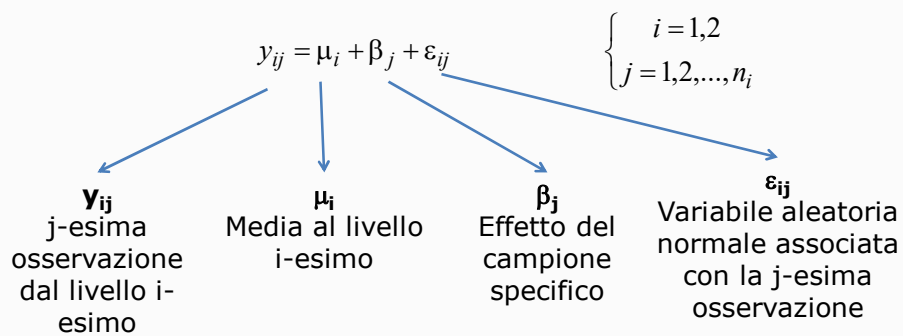
30

30

## Esempio introduttivo – Confronto medie con “blocco”

**Confronto tra trattamenti**

- Il modello statistico per descrivere la procedura è il seguente:



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

31

31

## Esempio introduttivo – Confronto medie con “blocco”

**Confronto tra trattamenti**

- In questo modo è possibile valutare la differenza per la coppia *j*-ma:

$$d_j = y_{1j} - y_{2j} \quad j = 1, \dots, n$$

- Il valore atteso di questa differenza è

$$\begin{aligned} \mu_d &= E[d_j] \\ &= E[y_{1j} - y_{2j}] \\ &= E[y_{1j}] - E[y_{2j}] \\ &= \mu_1 + \beta_j - (\mu_2 + \beta_j) \\ &= \mu_1 - \mu_2 \end{aligned}$$



È possibile fare inferenza sulla differenza di medie lavorando semplicemente con **le differenze delle coppie**

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

32

32



## Esempio introduttivo – Confronto medie con “blocco”

### Confronto tra trattamenti

- Testare  $H_0: \mu_1 = \mu_2$  è equivalente a testare

$$H_0: \mu_d = 0$$

$$H_1: \mu_d \neq 0$$

- La statistica test per questa ipotesi è ancora una  $t$  a  $(n-1)$  gdl:

$$t_0 = \frac{\bar{d}}{\sqrt{\frac{S_d^2}{n}}}$$

- dove

$$\bar{d} = \frac{\sum_{j=1}^n d_j}{n} \quad \text{e} \quad S_d^2 = \frac{\sum_{j=1}^n (d_j - \bar{d})^2}{n-1}$$

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

33

33

## Esempio introduttivo – Confronto medie con “blocco”

### Confronto tra trattamenti

- La possibilità di una campagna sperimentale con dati accoppiati, risulta da preferire alla campagna sperimentale completamente randomizzata.
- Si elimina una sorgente di incertezza dovuta alle differenze eventualmente presenti tra i campioni (e non dovute ai trattamenti).
- La filosofia del “blocco” permette di eliminare sorgenti di incertezza presenti nel campione di dati.

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

34

34