



# Esercitazione 4

## Statistica

Alfonso Iodice D'Enza  
iodicede@gmail.com

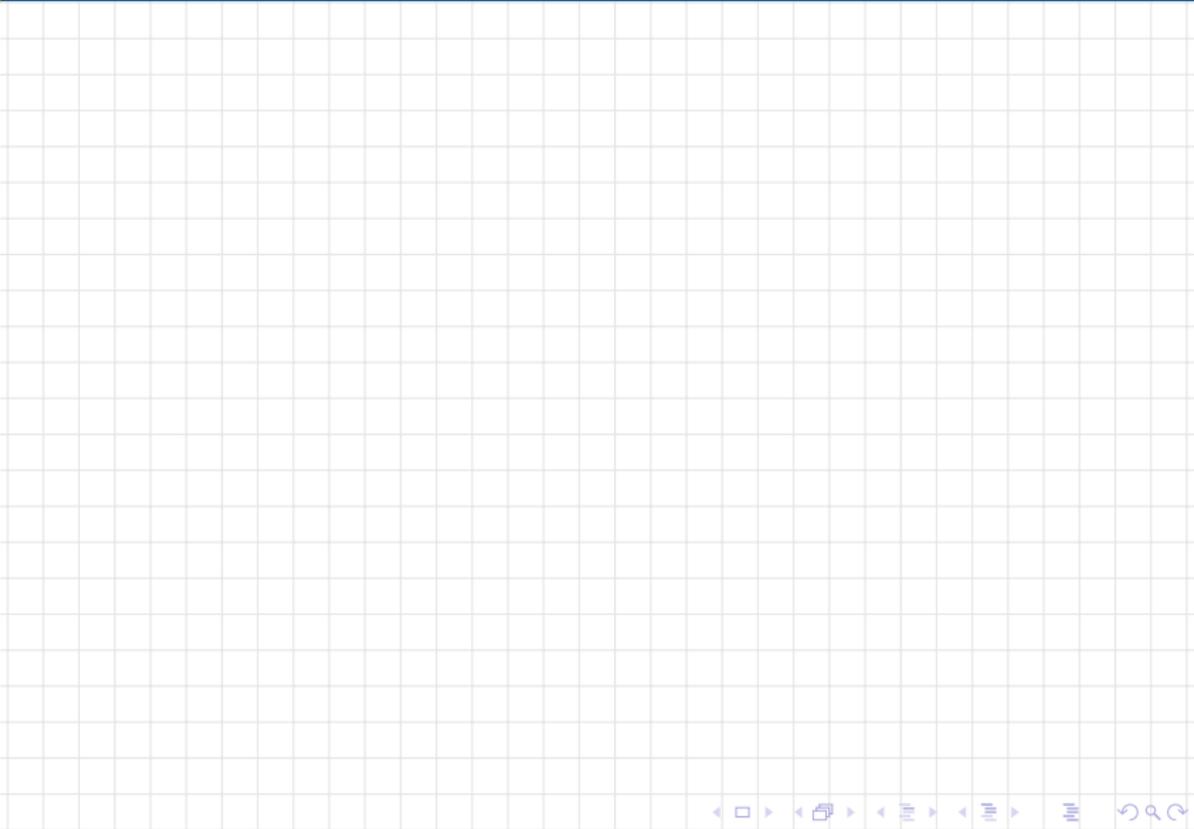
Università degli studi di Cassino



# Outline

Esercitazione  
4

A. Iodice





Si consideri una variabile  $Y$  quantitativa ed una variabile  $X$  qualitativa con ad esempio modalità  $(A, B, C)$ . Siano  $n_a$ ,  $n_b$  e  $n_c$  il numero di unità che presentano ciascuna delle modalità della variabile  $X$ , quindi  $n = n_a + n_b + n_c$ .

La media di  $Y$  si ottiene considerando la distribuzione di  $Y$  è

$$\mu_y = \bar{y} = \frac{1}{n} \sum_{i=1}^h y_i$$

Le medie di  $Y$  condizionate a ciascuna delle modalità della variabile  $X$  è

$$\bar{y}_a = \bar{y}|A = \frac{1}{n_a} \sum_{i=1}^{n_a} y_i$$

$$\bar{y}_b = \bar{y}|B = \frac{1}{n_b} \sum_{i=1}^{n_b} y_i$$

$$\bar{y}_c = \bar{y}|C = \frac{1}{n_c} \sum_{i=1}^{n_c} y_i$$



Si consideri una variabile  $Y$  quantitativa ed una variabile  $X$  qualitativa con ad esempio modalità  $(A, B, C)$ . Siano  $n_a$ ,  $n_b$  e  $n_c$  il numero di unità che presentano ciascuna delle modalità della variabile  $X$ , quindi  $n = n_a + n_b + n_c$ .

La media di  $Y$  si ottiene considerando la distribuzione di  $Y$  è

$$\mu_y = \bar{y} = \frac{1}{n} \sum_{i=1}^h y_i$$

Le medie di  $Y$  condizionate a ciascuna delle modalità della variabile  $X$  è

$$\bar{y}_a = \bar{y}|A = \frac{1}{n_a} \sum_{i=1}^{n_a} y_i$$

$$\bar{y}_b = \bar{y}|B = \frac{1}{n_b} \sum_{i=1}^{n_b} y_i$$

$$\bar{y}_c = \bar{y}|C = \frac{1}{n_c} \sum_{i=1}^{n_c} y_i$$



Ricordando che la devianza il numeratore della varianza... sia  $j = 1, 2, 3$  se si fa riferimento alle modalità  $A, B, C$  rispettivamente.

$$\begin{aligned} Dev_y &= \sum_{j=1}^3 \sum_{i=1}^{n_j} (y_i - \bar{y})^2 = \\ &= \sum_{j=1}^3 \sum_{i=1}^{n_j} (y_i - \bar{y}_j + \bar{y}_j - \bar{y})^2 = \\ &= \sum_{j=1}^3 \sum_{i=1}^{n_j} (y_i - \bar{y}_j)^2 + \sum_{j=1}^3 \sum_{i=1}^{n_j} (\bar{y}_j - \bar{y})^2 + \\ &+ 2 \sum_{j=1}^3 \sum_{i=1}^{n_j} (y_j - \bar{y}_j)(\bar{y}_j - \bar{y}) \end{aligned}$$

$$\begin{aligned} &= \sum_{j=1}^3 \left[ \sum_{i=1}^{n_j} (y_i - \bar{y}_j)^2 \right] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j + \\ &+ 2 \sum_{j=1}^3 (y_i - \bar{y}_j) \sum_{i=1}^{n_j} (\bar{y}_j - \bar{y}) = \\ &= \sum_{j=1}^3 [Dev(Y | X = x_j)] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j = \\ &= \text{Devianza(interna)} + \text{Devianza(esterna)} \end{aligned}$$

$$\begin{aligned} &= \sum_{j=1}^3 \left[ \sum_{i=1}^{n_j} (y_i - \bar{y}_j)^2 \right] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j + \\ &+ 2 \sum_{j=1}^3 (y_i - \bar{y}_j) \sum_{i=1}^{n_j} (\bar{y}_j - \bar{y}) = \\ &= \sum_{j=1}^3 [Dev(Y | X = x_j)] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j = \\ &= Devianza(interna) + Devianza(esterna) \end{aligned}$$

$$\begin{aligned}
&= \sum_{j=1}^3 \left[ \sum_{i=1}^{n_j} (y_i - \bar{y}_j)^2 \right] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j + \\
&+ 2 \sum_{j=1}^3 (y_i - \bar{y}_j) \sum_{i=1}^{n_j} (\bar{y}_j - \bar{y}) = \\
&= \sum_{j=1}^3 [Dev(Y | X = x_j)] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j = \\
&= Devianza(interna) + Devianza(esterna)
\end{aligned}$$

$$\begin{aligned} &= \sum_{j=1}^3 \left[ \sum_{i=1}^{n_j} (y_i - \bar{y}_j)^2 \right] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j + \\ &+ 2 \sum_{j=1}^3 (y_i - \bar{y}_j) \sum_{i=1}^{n_j} (\bar{y}_j - \bar{y}) = \\ &= \sum_{j=1}^3 [Dev(Y | X = x_j)] + \sum_{j=1}^3 (\bar{y}_j - \bar{y})^2 n_j = \\ &= \textit{Devianza(interna)} + \textit{Devianza(esterna)} \end{aligned}$$



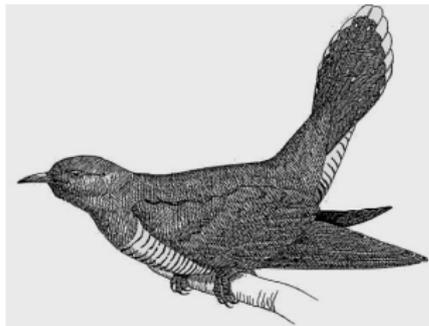
## Esempio di calcolo della decomposizione della devianza

### Esercitazione 4

#### A. Iodice

#### Il nido del cuculo

Il **cuculo** è un uccello caratterizzato da una particolare abitudine: depone le uova nei nidi di altri uccelli, e lascia dunque che siano altre specie a covarle. Ovviamente, il tutto funziona se la dimensione delle uova nel nido ospite sono compatibili con quelle del nido ospitante. In alcuni territori, il cuculo depone le uova in nidi di **scricciolo**, in altri sceglie nidi di **pettirosso**.



Si consideri di aver osservato la lunghezza di  $n_1 = 15$  uova di cuculo ritrovate in nidi di scricciolo e  $n_2 = 16$  uova di cuculo ritrovate in nidi di pettirosso. Si vuole **verificare se la lunghezza delle uova deposte cambia in media a seconda del tipo di nido in cui vengono deposte**.

## Scricciolo

Sia  $S$  la lunghezza delle uova di cuculo nei nidi di scricciolo

```
> S
      [,1]
 [1,] 19.85
 [2,] 20.05
 [3,] 20.25
 [4,] 20.85
 [5,] 20.85
 [6,] 20.85
 [7,] 21.05
 [8,] 21.05
 [9,] 21.05
[10,] 21.25
[11,] 21.45
[12,] 22.05
[13,] 22.05
[14,] 22.05
[15,] 22.25
```

```
> summary(scriccio)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 19.85  20.85  21.05  21.13  21.75  22.25
```



## Pettirosso

Sia  $P$  la lunghezza delle uova di cuculo nei nidi di pettirosso

```
> P
      [,1]
 [1,] 21.05
 [2,] 21.85
 [3,] 22.05
 [4,] 22.05
 [5,] 22.05
 [6,] 22.25
 [7,] 22.45
 [8,] 22.45
 [9,] 22.65
[10,] 23.05
[11,] 23.05
[12,] 23.05
[13,] 23.05
[14,] 23.05
[15,] 23.25
[16,] 23.85
```

```
> summary(pettirosso)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 21.05  22.05  22.55  22.57  23.05  23.85
```





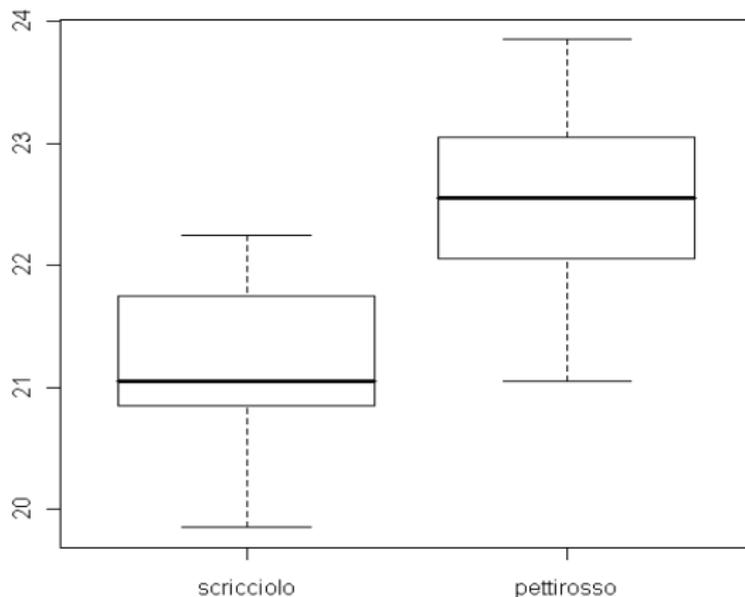
## Esempio di calcolo della decomposizione della devianza

### Esercitazione 4

#### A. Iodice

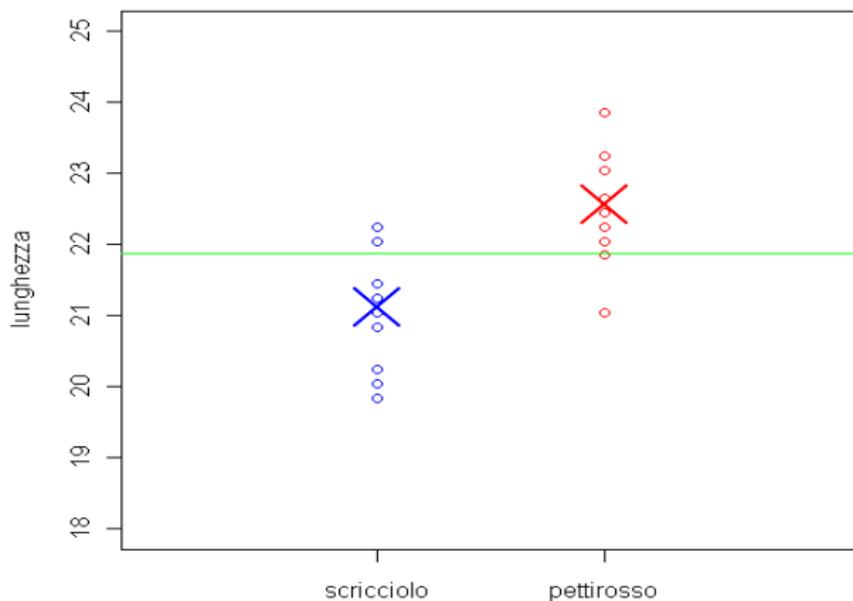
#### Confronto tra le distribuzioni

Un primo confronto grafico via box plot tra le due distribuzioni mostra che le uova deposte in nidi di pettirosso hanno una lunghezza maggiore di quelle deposte in nidi di scricciolo.



## Confronto tra le distribuzioni

Un ulteriore confronto grafico tra le due distribuzioni consiste in un diagramma per punti: sono riportate graficamente le medie condizionate, mentre la media generale  $\bar{y}_{\cdot}$  è rappresentata dalla linea orizzontale.





Si indica con  $\mu_X = 21.875$  la lunghezza media delle  $n = n_1 + n_2$  uova complessivamente considerate. Le medie condizionate al nido in cui le uova sono state deposte sono rispettivamente  $\mu_{X|S} = 21.13$  e  $\mu_{X|P} = 22.57$ . La devianza delle medie condizionate rispetto alla media generale è dunque

$$dev_b = (21.13 - 21.875)^2 \times 15 + (22.57 - 21.875)^2 \times 16 = 16.165$$

mentre la devianza complessiva è data da

$$dev_{tot} = (19.85 - 21.875)^2 + (20.05 - 21.875)^2 + \dots + (23.25 - 21.875)^2 + (23.85 - 21.875)^2 = 30.94$$