

# Calcoli con Fogli Elettronici

## *Lezione 3*

**Corso di Laurea in Biotecnologie  
AA. 2010/2011**

**Docente del laboratorio:  
Maria Silvia Pini  
[mpini@math.unipd.it](mailto:mpini@math.unipd.it)**

# Frequenze

- Campione  $X$ : un insieme di  $N$  **osservazioni**  $\{x_1, x_2 \dots x_N\}$  misurati con
- Scala di misura  $Y$  con  $k$  **categorie**  $[y_1 \dots y_k]$ .  
Esempio: il sesso di 30 persone; categorie [M,F]
- Frequenza  $f_i$  di una categoria  $y_i \in Y$  nel campione  $X$ :  
il numero di osservazioni di  $y_i$  nel  $X$ .
- Proporzione  $p_i$  di una categoria  $y_i \in Y$  nel campione  $X$ :  
la *frequenza*  $f_i$  divise per il numero totale  $N$  di osservazioni nel campione  $X$ .
- Distribuzione di frequenze  $F(y)$  nel campione  $X$ :  
l'insieme di frequenze  $f_i$  di ciascuna categoria  $y_i \in Y$  nel campione  $X$ .
- Distribuzione di probabilità empirica  $P(y)$  nel campione  $X$ :  
l'insieme di proporzioni  $p_i$  di ciascuna categoria  $y_i \in Y$  nel campione  $X$ .

Codice\_studente    Voto

1	27
2	30
3	26
4	29
5	26
6	27
7	25
8	25
9	25
10	29



Voto	f
25	3
26	2
27	2
28	0
29	2
30	1

# Categorie ab**bin**ate

Se abbiamo dei **dati numerici con tanti livelli**, rischiamo di avere un'osservazione per livello ... Che cosa fare ?

Soluzione: definire una **scala S derivata** dalla scala originale

$Y=[y_1 \dots y_k]$ , con un limitato numero di livelli (bin)  $\{s_1, s_2, \dots, s_M\}$ .

Procedura:

- nella nuova scala S, ciascuno livello  $s_i$  raggruppa livelli  $\{y_{i1}, y_{i2}, \dots, y_{ik}\}$
- **Il numero di livelli raggruppati** in ciascun nuovo livello  $s_i$  deve essere **uguale**.
- Il numero di livelli in S dovrebbe essere **limitato** (per es., 10).

# Categorie ab**bin**ate

Codice impiegato	Data di nascita
1	01/06/57
2	02/08/72
3	03/05/62
4	07/08/76
5	23/02/75
6	01/03/72
7	06/06/51
8	19/08/72
9	18/02/69
10	04/05/75



Decennio	
50	2
60	2
70	6
Total result	10

# Distribuzioni di frequenze di dati numerici nei Fogli Elettronici

Se abbiamo dati  $X$  di tipo numerico (**scala ad intervalli**), possiamo utilizzare la funzione ***frequenza()*** per calcolare la distribuzione dei dati  $X$ .

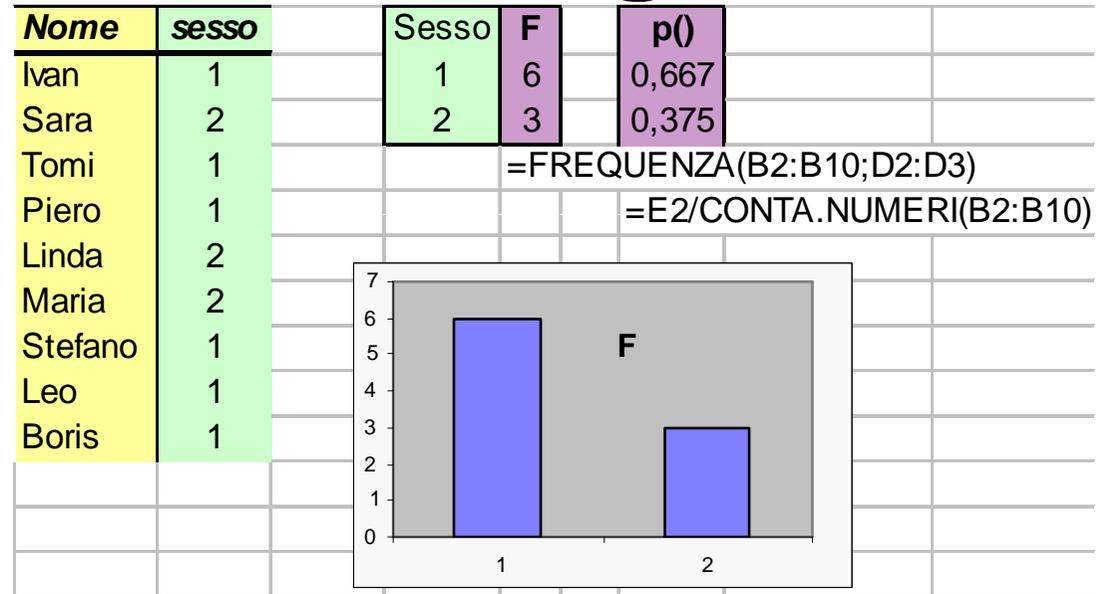
1. Avendo le **osservazioni**  $X$  in una colonna (es.: B2:B21)
2. Inserire i **livelli**  $S$  in un'altra colonna (es: D2:D11)
3. Applicare la funzione ***frequenza(vettore\_osserv.;*** ***vettore\_categorie)***
  - **Selezionare** una colonna per il **risultato** (es., E2:E11)
  - Scrivere '=frequenza('
  - Selezionare / riferire le **osservazioni**; scrivere ';'
  - Selezionare / riferire le **categorie**; scrivere ')'
  - Premere '***Ctrl-Maiusc-Invio***' (solo '***Invio***' calcola una frequenza solo!!)

**Si nota:** I **livelli indicati** definiscono una **nuova scala**  $S$  per calcolare la distribuzione dei dati numerici raggruppati. **La frequenza  $f_i$  di ciascuna categoria  $s_i$**  corrisponde al **numero delle osservazioni** con valori  **$(s_{i-1} \dots s_i]$** .

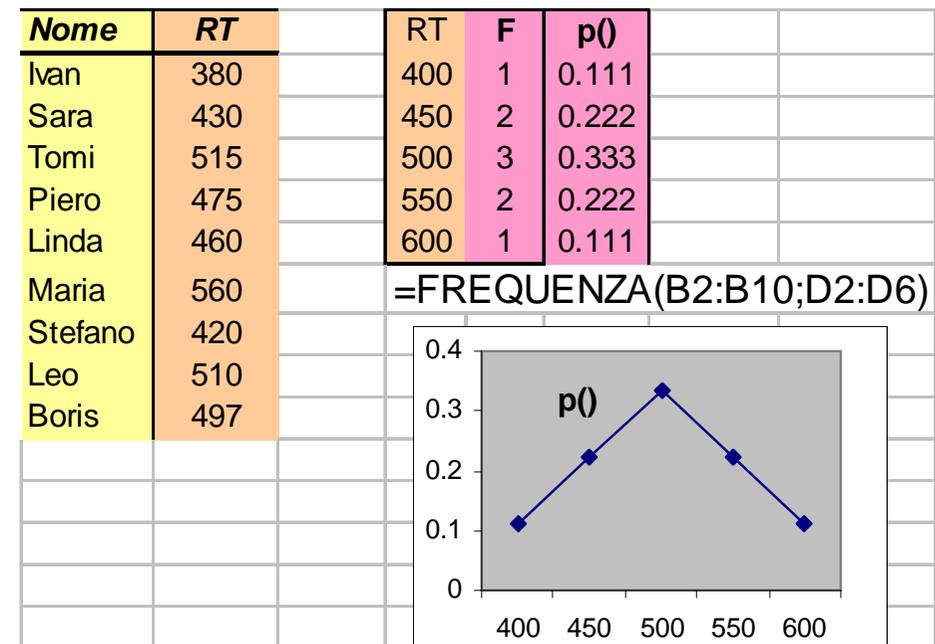
# Distribuzioni: calcolo e rappresentazione grafico

Grafici:

1. Pochi livelli: **Istogramma**  
(diagramma a barre)



2. Tanti livelli: **Grafico a linea**



# Esercizio: probabilità empirica

File dati: **S3\_RT.ods**

- Calcolare le **frequenze** degli RT suddivisi nei seguenti livelli [400, 450, 500, 550,...]
- Calcolare la **distribuzione delle probabilità**
- Fare il **grafico** della distribuzione di probabilità

# Tendenze centrali

1. Se abbiamo dati in una colonna (es.: B2:B21)
2. **Media** = somma divisa per il numero dei punti.
  - selezionare una cella per il **risultato** (di solito sotto i dati) (es: B23)
  - scrivere '=somma(';
  - selezionare le **osservazioni**; ')'; Premere '**Invio**'
  - in un'altra cella, dividere la somma per il numero dei dati (**conta.numeri()**)
3. **Mediana o moda**: ordinare i dati e trovare:
  - la categoria centrale (mediana)
  - più-frequente (moda)

# Variabilità

- **Range / Intervallo** (distanza tra i valori estremi):  $\max(X) - \min(X)$
  - **Somma dei quadrati delle distanze dalla media** (scarti)  $SS_x = \sum (x_i - \bar{X})^2$
  - **Varianza della popolazione** (la media degli scarti<sup>2</sup>):  $\sigma_x^2 = \frac{\sum (x_i - \bar{X})^2}{n}$
- $$\sigma_x^2 = \frac{\sum (x_i - \bar{X})^2}{n} = \frac{\sum (x_i^2 - 2x_i \bar{X} + \bar{X}^2)}{n} = \frac{\sum x_i^2}{n} - \bar{X}^2 = \overline{X^2} - \bar{X}^2$$
- **Deviazione standard** (Scarto quadratico medio):  $\sigma_x = \sqrt{\sigma_x^2}$

# Calcolo della varianza / dev.st

I. La **varianza** (della popolazione):

Algoritmo **A**: la media del quadrato delle scarti

1. Calcolare la media M
2. Calcolare gli scarti dalla media  $D_i=(X_i-M)$
3. Calcolare il quadrato degli scarti  $D_i^2=D_i*D_i$ .
4. Calcolare la media del quadrato delle distanze **var** = media( $D_i^2$ ).

$$Var_x = \frac{\sum (x_i - \bar{X})^2}{n}$$

Algoritmo **B**: la differenza tra la media dei quadrati e il quadrato della media

1. Calcolare la media M1
2. Calcolare i quadrati:  $X_i^2$
3. Calcolare la media dei quadrati M2
4. Calcolare il quadrato della media M1\*M1
5. Calcolare la differenza **var**=M2-M1\*M1.

$$Var_x = \overline{X^2} - \bar{X}^2$$

II. La **dev. standard**:

**dst**=radq(**var**)

$$\sigma_x = \sqrt{Var_x}$$

# Covarianza

- Due variabili aleatorie (stocastiche) X e Y possono co-variare. (per.es.: studenti con voto alto in Matematica hanno alto voto in Informatica.)
- La **covarianza** – la media del prodotto di X e Y normalizzati con loro medie – esprime il grado di dipendenza lineare tra X e Y:

$$\begin{aligned}Cov(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E(XY) - E(Y)E(X)\end{aligned}$$

- se **positivo**: al **crescere** di x in media **cresce** anche y,
  - se **negativo** al **crescere** di x in media **decrescere** y.
- Si nota: **la covarianza di una variabile con se stessa = la varianza !**

Algoritmo **A**:

1. Calcolare le medie  $M_x$  e  $M_y$
2. Calcolare il prodotto degli scarti  $D_i = (X_i - M_x)(Y_i - M_y)$
3. Calcolare la media  $Cov = \text{media}(D_i)$ .

Algoritmo **B**:

1. Calcolare i prodotti  $X_i Y_i$
2. Calcolare le medie  $M_x$ ,  $M_y$ ,  $M_{xy}$
3. Calcolare la differenza **Cov** =  $M_{xy} - M_x M_y$ .

# Funzioni per la statistica descrittiva

## Funzione

## Sintassi

(ad esempio, per le osservazioni in A2:A150)

- **MEDIA** *MEDIA*(A2:A150)
- **MEDIANA** *MEDIANA*(A2:A150)
- **Media distanza dalla Media** *MEDIA.DEV*(A2:A150)
- **VARIANZA della popolazione** *VAR.POP*(A2:A150)
- **SCARTO quadr.med. (pop.)** *DEV.ST.POP*(A2:A150)
- **COVARIANZA tra 2 campioni** *COVARIANZA*(A2:A150; B2:B150)

# Esercizio: medie e variabilità

- File dati: **S3\_RT.ods**  
calcolare **con e senza** funzioni statistiche:  
**la media, la varianza, e la dev.st.** dei valori RT
- File dati: **S2\_Math.ods**  
calcolare **con e senza** funzioni statistiche la **covarianza** tra  
*add\_time* e *sub\_time*

# Tabelle di contingenza (tabelle pivot)

- Sintetizzano una caratteristica dei dati rispetto ad altre loro caratteristiche
- I dati: una serie di casi con varie caratteristiche (tabella), di cui alcuni sono **causali** (indipendenti) ed altri sono **dipendenti**
- **Tabella pivot** ad una-, due-, o più- entrate (i fattori causali A, B, ... ) in cui:
  - gli elementi della **riga (colonna)** codificano i **livelli** delle **categorie causali** A (B),
  - le **celle** contengono una **misura di sintesi della variabile dipendente X** per ciascuna combinazione  $A_i B_j$  dei livelli dei fattori A e B
  - **l'ultima riga/colonna** codifica la misura di sintesi della variabile dipendente per ciascun livello del fattore A (B).
- **Misure di sintesi:**
  - *numero di casi*
  - *somma / media / dev.st.* dei valori  $X_i$  per ciascuna cella.
  - *altri*
  - i valori possono essere riportati come % rispetto  $A_i / B_j / \text{Totale}$

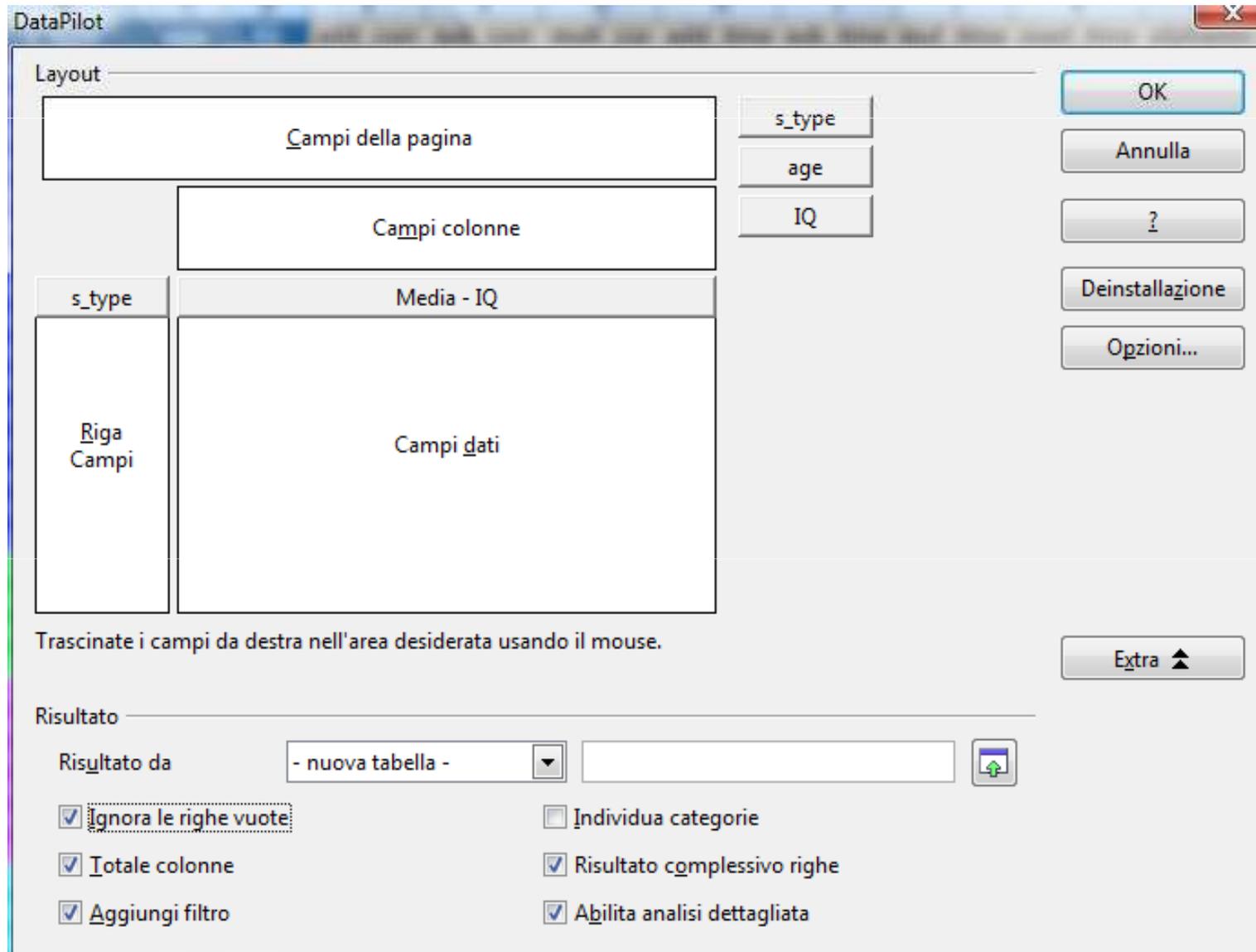
**Esempio:** i dati: {età, condizione, accuratezza};  
la tabella: accuratezza (condizione, età)

# Esempio pivot

oggetti	s_type	age	IQ		
C1	control	126	98		
C2	control	132	124	<b>Conteggio di s_type</b>	
C3	control	136	111	s_type	Totale
C4	control	131	91	control	16
C5	control	125	107	number fact	4
C6	control	129	105	Turners Syn.	6
C7	control	132	94	Williams Syn.	5
C8	control	129	107	Totale complessivo	31
C9	control	127	105		
C10	control	138	115		
C11	control	127	109	<b>Conteggio di s_type</b>	
C12	control	126	96	s_type	Totale
C13	control	130	98	control	51,61%
C14	control	134	118	number fact	12,90%
C15	control	135	91	Turners Syn.	19,35%
C16	control	141	127	Williams Syn.	16,13%
NF1	number fact	130	98	Totale complessivo	100,00%
NF2	number fact	132	101		

# Tabelle pivot – procedura (1)

1. Selezionare un intervallo compatto di dati.
2. Menu: **Dati > DataPilot > Avvia**



# Tabelle pivot – procedura (2)

## 3) Disegnare la tabella pivot

- trascinare le categorie **causali** in: colonna, riga o pagina.
- trascinare le variabili **dipendenti** nel centro della tabella;
- selezionare la **misura di sintesi**: conteggio, media, ... per ciascuna variabile (con doppio click; si possono scegliere più di una sintesi)

The screenshot shows a pivot table configuration window. On the right, a list of fields includes 's\_type', 'age', and 'IQ'. The main area is divided into four zones: 'Campi della pagina' (top), 'Campi colonne' (middle), 'Riga Campi' (left), and 'Campi dati' (center). A blue arrow points from the 'Media - IQ' field in the 'Campi dati' zone to the 'Media - IQ' field in the 'Campi colonne' zone. Below the main area, a dropdown menu is set to '- nuova tabella -'. At the bottom, there are several checkboxes: 'Ignora le righe vuote' (checked), 'Individua categorie' (unchecked), 'Totale colonne' (checked), 'Risultato complessivo righe' (checked), 'Aggiungi filtro' (checked), and 'Abilita analisi dettagliata' (checked).

- ## 4) Scegliere l'allocazione della tabella (Risultato da ..)
- (Ad esempio: *nuova tabella*)

# Esercizio

- File dati: **S3\_Dan.ods**  
Creare in fogli diversi di S3\_Dan.ods
  - una **tabella pivot** risp\_verbale(luogo, livello) dove si esamina la **media**
  - una **tabella pivot** risp\_verbale(luogo, livello) dove si esamina il **conteggio**
  - una **tabella pivot** risp\_verbale(luogo, livello) dove si esamina la **dev. st.**
  - una **tabella pivot** risp\_verbale(luogo, livello) e risp\_sem(luogo, livello) dove si esaminano la **media di risp\_verbale** e il **massimo di risp\_sem**