



Gestione e Analisi dei Dati

Lezione 7

Introduzione all'analisi dei dati

Popolazioni, campioni, unità statistica, variabili

Studio e rappresentazione della frequenza

Obiettivi del modulo

- Scopo del corso è la presentazione di metodi per la estrazione di informazioni da insiemi di dati
 - Le informazioni possono essere
 - di sintesi (p.e. la media di un insieme di valori numerici)
 - di andamento (p.e. il grafico dei valori di una variabile nel tempo)
- La base concettuale è la *statistica (descrittiva e qualche elemento di inferenziale)*
- Lo strumento operativo è il foglio elettronico (spreadsheet), in particolare OpenOffice Calc

Tipiche domande

- Che incassi sono stati registrati l'anno passato per ciascuna regione e ciascuna categoria di prodotto?
- Che correlazione esiste tra l'andamento dei titoli azionari dei produttori di PC e i profitti trimestrali lungo gli ultimi 5 anni?
- Quali sono gli ordini che massimizzano gli incassi?
- Quale di due nuove terapie risulterà in una diminuzione della durata media di un ricovero?
- Che rapporto c'è tra i profitti realizzati con spedizioni di meno di dieci elementi e quelli realizzati con spedizioni di più di dieci elementi?

Outline

- Elementi di statistica descrittiva e uso di OpenOffice Calc per i calcoli:
 - Variabili e tabelle
 - Manipolazione di tabelle (sort ecc.)
 - Distribuzioni di frequenze e grafici
 - Statistica descrittiva per una singola variabile (media, varianza, ...)
 - Statistica descrittiva: misure di associazione.
- Data warehouses e OLAP (On Line Analytical Processing)
- OLAP in oCalc: pivot tables
- Elementi di data cleaning

Concetti di base

- Una *popolazione* (population) include tutti gli oggetti di interesse
- Esempi di popolazione:
 - tutti i potenziali votanti per l'elezione del rettore
 - tutti gli abbonati alla RAI
 - tutte le fatture ricevute dal Dipartimento di Informatica nel 2003
- Un *campione* (sample) è un sottoinsieme di una popolazione, spesso scelto in modo casuale e preferibilmente rappresentativo dell'intera popolazione.

Variabili e Osservazioni

- Una *unità statistica* {osservazione, caso} è una ennupla di valori che caratterizza un elemento di una popolazione o di un campione
- Una *variabile* {carattere, attributo, campo} è l'identificatore (nome) di uno dei valori dell'osservazione
- Se una popolazione (campione) è rappresentata in forma di *tabella* le righe della tabella sono le osservazioni, i nomi delle colonne sono le variabili e il contenuto di ciascuna riga è la lista dei valori delle variabili

Un esempio

Name	Gender	DomesticGross	ForeignGross	Salary
Angela Bassett	F	32	17	2,5
Jessica Lange	F	21	27	2,5
Winona Ryder	F	36	30	4
Michelle Pfeiffer	F	66	31	10
Whoopi Goldberg	F	32	33	10
Emma Thompson	F	26	44	3
Julia Roberts	F	57	47	12
Sharon Stone	F	32	47	6
Meryl Streep	F	34	47	4,5
Susan Sarandon	F	38	49	3
Nicole Kidman	F	55	51	4
Holly Hunter	F	51	53	2,5
Meg Ryan	F	43	55	8,5
Andie Macdowell	F	26	75	2
Jodie Foster	F	62	85	9

Tipi di Valori (I)

- Possiamo classificare le variabili in base alla tipologia dei valori che possono assumere;
- Distinguiamo tra:
 - Variabili numeriche (*quantitative*), se sui valori è possibile compiere un insieme significativo di operazioni aritmetiche;
 - Variabili categoriche (*qualitative*), altrimenti

Tipi di Valori (2)

- **Le variabili numeriche possono essere:**
 - discrete, se i valori possono essere contati
 - continue, se sono il risultato di una misura continua
- **Le variabili categoriche possono essere:**
 - ordinali, se esiste un ordine naturale sui possibili valori (es. giudizi scolastici: insufficiente, sufficiente ecc.)
 - nominali, altrimenti (es. colori)

Esempio

EnvironmentalPolicy.xls - OpenOffice.org Calc

File Modifica Visualizza Inserisci Formato Strumenti Dati Finestra ?

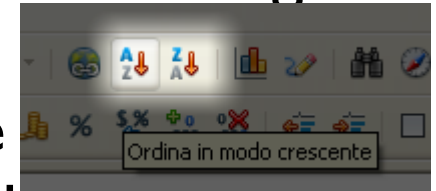
Arial 10 G C S

F7 = 5

	A	B	C	D	E	F	G	H
1	<u>Data from a questionnaire on environmental policy</u>							Visit the a
2								Or visit the
3	<u>Age</u>	<u>Gender</u>	<u>State</u>	<u>Children</u>	<u>Salary</u>	<u>Opinion</u>		
4	35	Male	Minnesota	1	\$65.400	5		
5	61	Female	Texas	2	\$62.000	1		
6	35	Male	Ohio	0	\$63.200	3		
7	37	Male	Florida	2	\$52.000	5		
8	32	Female	California	3	\$81.400	1		
9	33	Female	New York	3	\$46.300	5		
10	65	Female	Minnesota	2	\$49.600	1		
11	45	Male	New York	1	\$45.900	5		
12	40	Male	Texas	3	\$47.700	4		
13	32	Female	Texas	1	\$59.900	4		
14	57	Male	New York	1	\$48.100	4		
15	38	Female	Virginia	0	\$58.100	3		
16	37	Female	Illinois	2	\$56.000	1		
17	42	Female	Virginia	2	\$53.400	1		
18	38	Female	New York	2	\$39.000	2		
19	48	Male	Michigan	1	\$61.500	2		
20	40	Male	Ohio	0	\$37.700	1		
21	57	Female	Michigan	2	\$36.700	4		
22	44	Male	Florida	2	\$45.200	3		

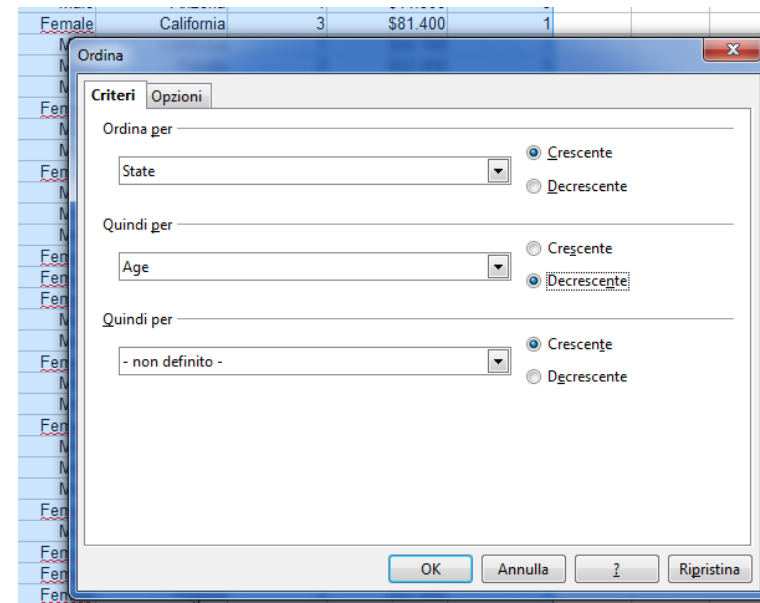
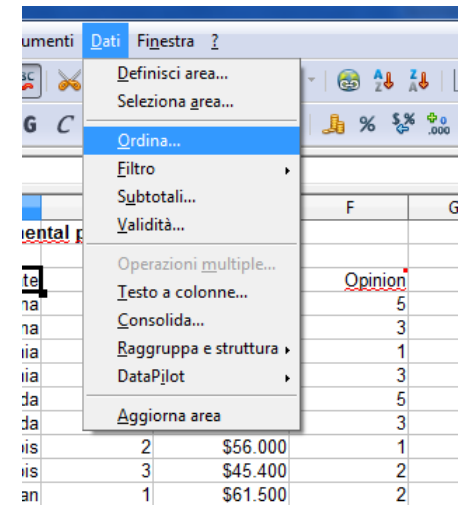
Interpretazione dei dati: ordinamento

- Il più potente e il più semplice strumento di elaborazione di una tabella è l'*ordinamento*
- In oCalc è possibile ordinare la tabella scegliendo:
 - la colonna (l'attributo)
 - la modalità: crescente, decrescente
- E' inoltre possibile effettuare ordinamenti annidati, ovvero ordinare rispetto ad un attributo e, tra gli elementi che hanno lo stesso valore dell'attributo, effettuare l'ordinamento rispetto a un secondo attributo, e così via;
- Esempio su environmental policy (slide precedente);



Ordinamento

- E' anche possibile effettuare l'ordinamento di una tabella di dati tramite il menu *Dati*→*Ordina...*
- Dalla finestra di dialogo è possibile impostare i criteri per l'ordinamento anche su più attributi



Interpretazione mediante distribuzione dei valori degli attributi

- Nello studio di un campione di osservazioni in cui alcune variabili sono di tipo categorico o categorizzabile, può essere molto informativo vedere come le osservazioni si distribuiscono sulle categorie;
- Una *tabella di frequenze* riporta il numero di osservazioni che ricadono in ciascuna delle categorie stabilite;
- Un *istogramma* è una tecnica di visualizzazione di una tabella di frequenza tramite un diagramma a barre.

Creazione di una tabella di frequenze (I)

- oCalc consente di creare una tabella di frequenze a partire da una tabella di dati mediante la funzione:
 - **FREQUENZA(dati,classi)**
- dove:
 - **dati** è una tabella monodimensionale (array) che contiene l'insieme di valori di cui vogliamo calcolare le frequenze.
 - **classi** è una tabella monodimensionale che contiene gli intervalli in cui vogliamo raggruppare i valori in **dati**.

Creazione di una tabella di frequenze (2)

- *passo 1*: in un'area dello stesso foglio che contiene la tabella dei dati o su un altro foglio costruire una tabella monodimensionale (un array) contenente il valore superiore per ciascuna categoria
- *passo 2*: per rendere la tabella delle frequenze leggibile, creare una colonna affiancata a quella dei limiti superiori contenente una descrizione della categoria.

<u>Female</u>	Texas	2	\$62.000	1			
Male	Ohio	0	\$63.200	3			
Male	Florida	2	\$52.000	5			
<u>Female</u>	California	3	\$81.400	1			
<u>Female</u>	New York	3	\$46.300	5		34 <34 anni	
<u>Female</u>	Minnesota	2	\$49.600	1		59 34-59	
Male	New York	1	\$45.900	5		>59	
Male	Texas	3	\$47.700	4			
<u>Female</u>	Texas	1	\$59.900	4			
Male	New York	1	\$48.100	4			
<u>Female</u>	Virginia	0	\$58.100	3			
<u>Female</u>	Illinois	2	\$56.000	1			
Female	Virginia	2	\$53.400	1			

Creazione di una tabella di frequenze (3)

- Selezionare una colonna contenente un numero di celle pari alle categorie create al passo precedente
- digitare la formula utilizzando lo strumento “Creazione guidata funzione”:
 - Per il campo “Dati” selezionare la colonna delle età; per il campo “Classi” selezionare i valori nella tabella di riferimento creata prima {34, 59, }
- Completare i passi della creazione guidata
- *Esercizio*: generare le distribuzioni di Salary e ForeignGross a partire dal file ACTORS

H	I	J	K
	34 <34 anni		
	59 34-59		
	>59		

EnvironmentalPolicy.xls - OpenOffice.org Calc

File Modifica Visualizza Inserisci Formato Strumenti Dati

Arial 10 G C S

FREQUENZA

A B [Creazione guidata funzione]

1 Data from a questionnaire on environmental po

2

3 Age Gender State

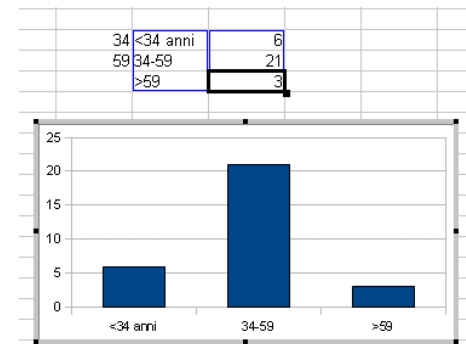
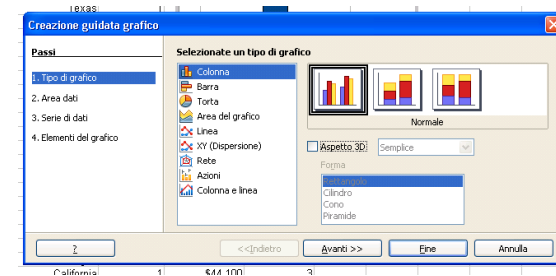
4 35 Male Minnesota

	34 <34 anni	6	
	59 34-59	21	
	>59	3	

Creazione di un Istogramma

- Utilizzare il menu Inserisci → Grafico a partire dalle colonne:
 - delle spiegazioni delle categorie
 - delle frequenze (calcolate al passo precedente)
- e scegliendo come output il diagramma a barre
- *Esercizio:* costruire gli istogrammi relativi alle tabelle di frequenze ottenute per il file ACTORS

	F	G	H	I	J
	5				
	1				
	3				
	5				
	1				
	5		34	<34 anni	6
	1		59	34-59	21
	5			>59	3
	4				
	4				
	4				
	3				
	1				
	1				
	2				
	2				



Esercitazione (I)

- Generare le distribuzioni di Salary a partire dal file ACTORS
- Procediamo a definire le soglie per i valori di Salary, in modo tale da avere dei range uniformi per ogni soglia
 - Selezioniamo il MIN e il MAX per i valori utilizzando le funzioni del foglio di calcolo
 - =MIN(E5:E71)
 - =MAX(E5:E71)

Or visit the [Duxbury site for this book at http://v](http://v)

Salary	Salary		
2,5	MIN		2
2,5	MAX		20
4			
10			
10			
3			

Esercitazione (2)

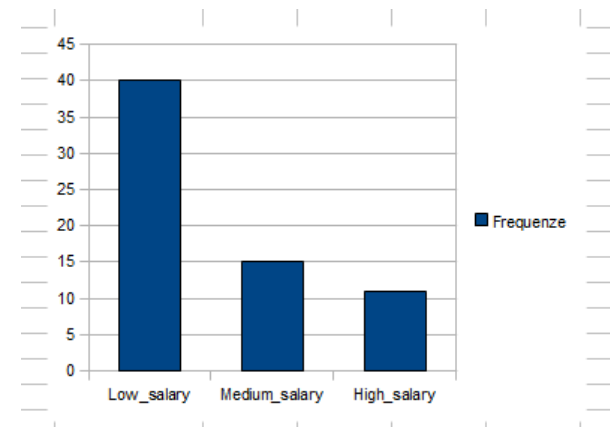
- I valori di Salary sono compresi nell'intervallo [2,20]
- La dimensione dell'intervallo è $20-2=18$
- Fissiamo tre soglie: Low_salary, medium_salary e high_salary
- La dimensione di ogni sottointervallo sarà: $\text{dim_intervallo}/\text{num_intervalli}$, ovvero $18/3=6$
- I valori soglia per ogni intervallo saranno
 - $\text{Min}+\text{dim_intervallo} = 2+6 = 8$
 - $\text{soglia_precedente}+\text{dim_intervallo} = 8+6=14$
 - $\text{soglia_precedente}+\text{dim_intervallo} = 14+6=20$

Salary		
MIN	2	
MAX	20	
Soglie		
Low_salary	8	
Medium_salary	14	
High_salary	20	

Esercitazione (3)

- La distribuzione delle frequenze si ottiene dalle tre soglie definite prima utilizzando la funzione:
 - =FREQUENZA(E6:E71;H10:H12)
- Dai valori di frequenza ottenuti si deriva l'istogramma tramite il diagramma a barre

	E	F	G	H	I
4					
5	<u>Salary</u>		<u>Salary</u>		
6	2,5		MIN	2	
7	2,5		MAX	20	
8	4				
9	10		Soglie		Frequenze
10	10		<u>Low_salary</u>	8	40
11	3		<u>Medium_salary</u>	14	15
12	12		<u>High_salary</u>	20	11
13	6				



Esercitazione (I)

- La fabbrica di ascensori OTIS ha misurato il diametro dei cavi da ascensore prodotti (file OTIS1). Generare la distribuzione (scegliendo opportunamente le categorie) e l'istogramma relativo e discuterne la forma
- I valori contenuti nel file sono 399
 - In questo caso, sceglieremo tre soglie in modo che ogni sottointervallo contenga lo stesso numero di valori

Esercitazione (2)

- Iniziamo ordinando i valori in senso crescente
- Dato che vogliamo tre intervalli della stessa dimensione e che i valori totali sono 399, la prima soglia sarà il valore corrispondente alla 133esima riga, la seconda soglia sarà la 266esima riga e l'ultimo valore sarà l'ultima riga

Esercitazione (3)

- Possiamo evitare di contare le righe sullo schermo utilizzando una serie numerica da affiancare ai valori ottenuti
- Inseriamo un 1 accanto alla prima riga e poi trasciniamo l'angolo in basso a destra della cella appena inserita, fino all'ultima cella della prima colonna

	A	B
1	Diameters of elevato	
2		
3	Note: All diameters an	
4		
5	Diameter	
6	0,450	1
7	0,451	
8	0,453	
9	0,459	
10	0,461	
11	0,461	
12	0,461	
13	0,462	
14	0,463	
15	0,463	

4			
5	Diameter		
6	0,450	1	
7	0,451	2	
8	0,453	3	
9	0,459	4	
10	0,461	5	
11	0,461	6	
12	0,461	7	
13	0,462	8	
14	0,463	9	

Esercitazione (4)

- Cerchiamo le righe dove vogliamo effettuare il taglio e annotiamo i valori corrispondenti come soglie
 - 0.492
 - 0.510
 - 0.548
- La distribuzione della frequenza si ottiene come prima

137	0,492	132
138	0,492	133
139	0,493	134
140	0,493	135

270	0,510	265
271	0,510	266
272	0,510	267
273	0,510	268

403	0,547	398
404	0,548	399
405		
406		

Σ = {=FREQUENZA(A6:A404;E6:E8)}			
C	D	E	F
are measured in fractions of an inch			
	Soglie		Frequenze
1	Bassa	0,49	129
2	Media	0,51	138
3	Alta	0,55	131
4			