

## La v.a. normale standardizzata

## La distribuzione normale standardizzata

- La distribuzione normale è difficilmente trattabile dal punto di vista calcolatorio, a causa dei suoi due parametri,  $\mu$  e  $\sigma^2$ .
- Il ricorso alla “**distribuzione normale standardizzata**” permette invece di individuare facilmente le probabilità relative agli intervalli di valori, utilizzando opportune *tavole statistiche*.

## La distribuzione normale standardizzata

- La **distribuzione normale standardizzata** (detta "Z") si ottiene mediante una trasformazione della variabile X, di questo tipo ("punteggi z"):

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

e pertanto

$$z^2 = \frac{(X - \mu)^2}{\sigma^2}$$

21 e 22 novembre 2011

Statistica sociale

3

## La distribuzione normale standardizzata

La **standardizzazione** è una **trasformazione dei dati** che consiste nel:

- rendere la media **nulla** ( $\mu = 0$ ), dato che ad ogni valore della variabile originaria viene sottratta la media della variabile stessa;
- assumere la deviazione standard  $\sigma$  quale **unità di misura** ( $\sigma = 1$ ) della nuova variabile, dato che ogni valore viene diviso per  $\sigma$ .

La distribuzione **normale standardizzata** viene indicata con **N(0,1)**.

I valori della Z sono **tabulati**: tra qualche diapositiva vedremo la tavola della Z.

21 e 22 novembre 2011

Statistica sociale

4

## La distribuzione Normale standardizzata

- La funzione di densità di probabilità della distribuzione normale standardizzata,  $f(z)$ , assume la forma:

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$$
$$(-\infty < z < +\infty)$$

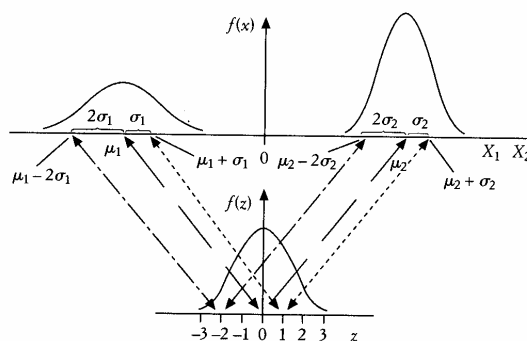
Osservazione: questa funzione non contiene più i parametri.

21 e 22 novembre 2011

Statistica sociale

5

## La distribuzione normale standardizzata



La v.a. normale standardizzata ha

MEDIA=0 e

DEVIATION STANDARD=1, per cui è rappresentata da UNA SOLA CURVA, mentre la distribuzione normale generale è rappresentata da infinite curve, che variano a seconda dei valori di  $\mu$  e  $\sigma$ .

21 e 22 novembre 2011

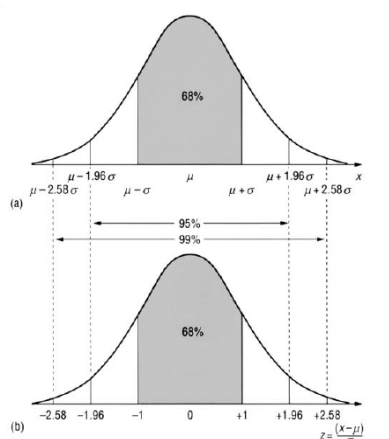
Statistica sociale

6

## Aree sottese dalla curva normale generale

- La probabilità che un valore estratto casualmente da una v.a.  $N(\mu, \sigma^2)$  sia compreso nell'intervallo  $(\mu - \sigma, \mu + \sigma)$  è pari al **68%**;
- Il 95% dei valori assunti da una distribuzione Normale cadono nell'intervallo  $(\mu - 1,96\sigma, \mu + 1,96\sigma)$ ;
- Il 99%, invece, nell'intervallo  $(\mu - 2,58\sigma, \mu + 2,58\sigma)$

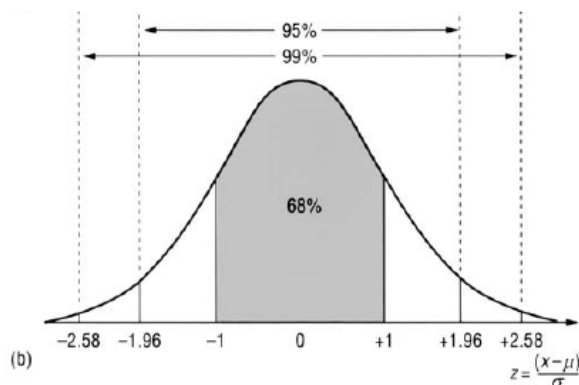
21 e 22 novembre 2011



Statist

7

## Aree sottese dalla curva normale standardizzata



21 e 22 novembre 2011

Statistica sociale

8

## Aree sottese dalla curva normale standardizzata

- La distribuzione normale standardizzata è importante perché le probabilità corrispondenti alle aree sottese dalla curva normale possono essere calcolate. Queste probabilità vengono riportate in apposite tavole.
  - In questo modo è possibile evitare il ricorso a complessi calcoli integrali per trovare le probabilità che una v.a.  $X$  assuma valori compresi all'interno di determinati intervalli.

21 e 22 novembre 2011

Statistica sociale

9

## Aree sottese dalla curva normale standardizzata

- È noto che il 68,26% dell'area totale è compreso tra  $\pm 1$  *deviazioni standard* attorno alla media, cioè a  $\pm 1$  *punti z* dalla media; mentre il 95,44% è racchiuso tra  $\pm 2$  *deviazioni standard* attorno alla media: quindi a  $\pm 2$  *punti z* dalla media.

21 e 22 novembre 2011

Statistica sociale

10

## Aree sottese dalla curva normale standardizzata

In virtù della proprietà di simmetria della distribuzione normale, le tavole riportano soltanto i valori dell'area compresa fra lo zero e l'ascissa  $+X$ , poiché, per la simmetria, l'area sottesa dall'altra metà della curva è ovviamente uguale.

21 e 22 novembre 2011

Statistica sociale

11

## Aree sottese dalla curva normale standardizzata

- Osservando la tavola, si troveranno i punti  $z$  nella colonna di sinistra con una cifra decimale; la seconda cifra decimale è posta nella prima riga in alto della stessa tavola.

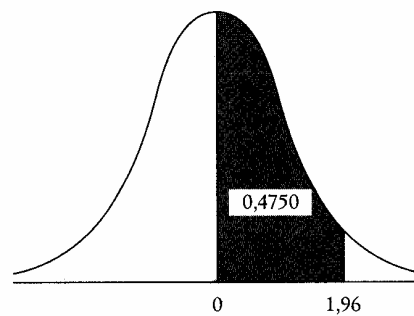
21 e 22 novembre 2011

Statistica sociale

12

## In termini pratici ...

Supponiamo di voler conoscere l'area compresa tra le ascisse pari, rispettivamente, a  $z=0$  e  $z=1,96$ .



21 e 22 novembre 2011

Statistica sociale

13

## In termini pratici ...

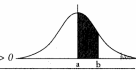
- Osservando la colonna dei punti  $z$ , si deve scendere fino a trovare  $z=1,9$ , e poi rimanere nella stessa riga fino a trovarsi in quella indicata con  $\delta$ .
- Il punteggio che si trova in quel punto indica la porzione di area compresa tra i due valori di  $z$ : 0,4750. Poiché l'area totale sottesa dalla curva nella sua parte positiva è pari a 0,500, l'area che si trova alla destra del valore  $z = 1,96$  sarà data da:
  - $0,5000 - 0,4750 = 0,0250$ .

21 e 22 novembre 2011

Statistica sociale

14

Tav. B. Area della distribuzione normale standard tra  $a = 0$  e  $b > 0$



z	0	1	2	3	4	5	6	7	8	9
0.0	.0000	.0040	.0080	.0120	.0160	.0199	.0239	.0279	.0319	.0359
0.1	.0398	.0438	.0478	.0517	.0557	.0596	.0636	.0675	.0714	.0754
0.2	.0793	.0832	.0871	.0910	.0948	.0987	.1026	.1064	.1103	.1141
0.3	.1179	.1217	.1255	.1293	.1331	.1368	.1406	.1443	.1480	.1517
0.4	.1554	.1591	.1628	.1664	.1700	.1736	.1772	.1808	.1844	.1879
0.5	.1915	.1950	.1985	.2019	.2054	.2088	.2123	.2157	.2190	.2224
0.6	.2258	.2291	.2324	.2357	.2389	.2422	.2454	.2486	.2518	.2549
0.7	.2580	.2612	.2642	.2673	.2704	.2734	.2764	.2794	.2823	.2852
0.8	.2881	.2910	.2939	.2967	.2996	.3023	.3051	.3078	.3106	.3133
0.9	.3159	.3186	.3212	.3238	.3264	.3289	.3315	.3340	.3365	.3389
1.0	.3413	.3438	.3461	.3485	.3508	.3531	.3554	.3577	.3599	.3621
1.1	.3643	.3665	.3686	.3708	.3729	.3749	.3770	.3790	.3810	.3830
1.2	.3849	.3869	.3888	.3907	.3925	.3944	.3962	.3980	.3997	.4015
1.3	.4032	.4049	.4066	.4082	.4099	.4115	.4131	.4147	.4162	.4177
1.4	.4192	.4207	.4222	.4236	.4251	.4265	.4279	.4292	.4306	.4319
1.5	.4332	.4345	.4357	.4370	.4382	.4394	.4406	.4418	.4429	.4441
1.6	.4452	.4463	.4474	.4484	.4495	.4505	.4515	.4525	.4535	.4545
1.7	.4554	.4564	.4573	.4582	.4591	.4599	.4608	.4616	.4625	.4633
1.8	.4641	.4649	.4656	.4664	.4671	.4678	.4686	.4693	.4699	.4706
1.9	.4713	.4719	.4726	.4732	.4738	.4744	.4750	.4756	.4761	.4767
2.0	.4772	.4778	.4783	.4788	.4793	.4798	.4803	.4808	.4812	.4817
2.1	.4821	.4826	.4830	.4834	.4838	.4842	.4846	.4850	.4854	.4857
2.2	.4861	.4864	.4868	.4871	.4875	.4878	.4881	.4884	.4887	.4890
2.3	.4893	.4896	.4898	.4901	.4904	.4906	.4909	.4911	.4913	.4916
2.4	.4918	.4920	.4922	.4925	.4927	.4929	.4931	.4932	.4934	.4936
2.5	.4938	.4940	.4941	.4943	.4945	.4946	.4948	.4949	.4951	.4952
2.6	.4953	.4955	.4956	.4957	.4959	.4960	.4961	.4962	.4963	.4964
2.7	.4965	.4966	.4967	.4968	.4969	.4970	.4971	.4972	.4973	.4974
2.8	.4974	.4975	.4976	.4977	.4977	.4978	.4979	.4979	.4980	.4981
2.9	.4981	.4982	.4982	.4983	.4984	.4984	.4985	.4985	.4986	.4986
3.0	.4987	.4987	.4987	.4988	.4988	.4989	.4989	.4989	.4990	.4990
3.1	.4990	.4991	.4991	.4991	.4992	.4992	.4992	.4992	.4993	.4993
3.2	.4993	.4993	.4994	.4994	.4994	.4994	.4994	.4995	.4995	.4995
3.3	.4995	.4995	.4995	.4996	.4996	.4996	.4996	.4996	.4996	.4997
3.4	.4997	.4997	.4997	.4997	.4997	.4997	.4997	.4997	.4997	.4998
3.5	.4998	.4998	.4998	.4998	.4998	.4998	.4998	.4998	.4998	.4998
3.6	.4998	.4998	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999
3.7	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999
3.8	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999	.4999
3.9	.5000	.5000	.5000	.5000	.5000	.5000	.5000	.5000	.5000	.5000

21 e 22 novembre

# Inferenza statistica



## Dalla popolazione al campione: l'inferenza statistica

- Il reperimento dei DATI STATISTICI attraverso una RILEVAZIONE è un'operazione che ha dei COSTI, sia in termini di TEMPO IMPIEGATO che in termini ECONOMICI.
- In molti ambiti scientifici, come in biologia e in medicina, si ha spesso a che fare con dati di origine sperimentale, per i quali quello di cui si dispone è **sempre un campione**, visto che la popolazione di riferimento è virtualmente "infinita" (es. un campione di animali da laboratorio rappresenta idealmente **tutti** gli esemplari di quella specie di animali da laboratorio).

21 e 22 novembre 2011

Statistica sociale

17

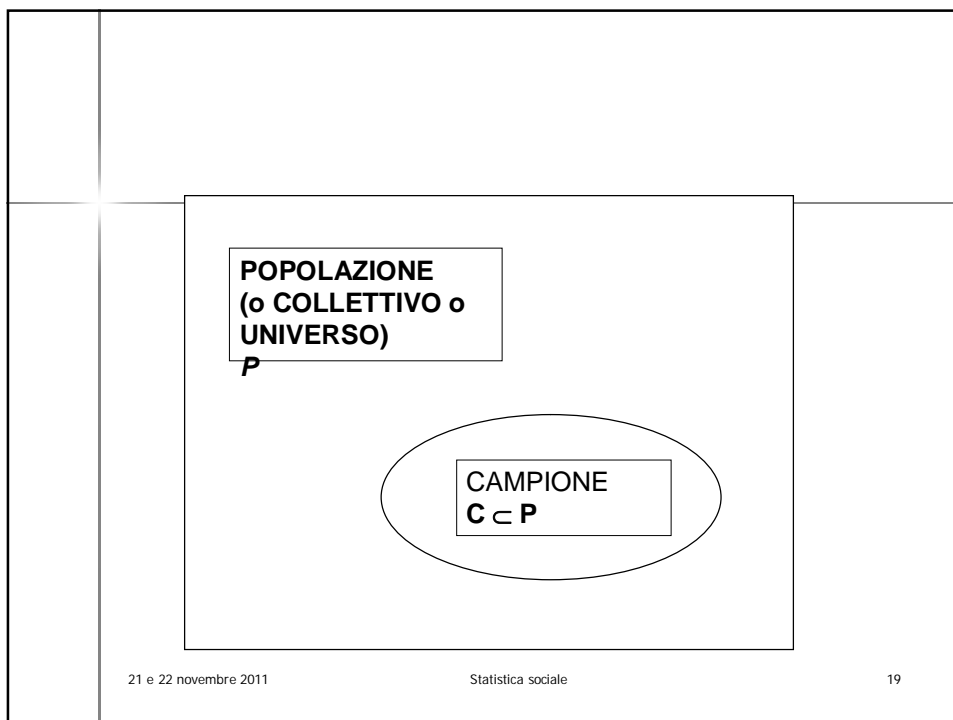
## Dalla popolazione al campione

- Nei primi anni del '900, integrandosi con alcuni risultati fondamentali del calcolo delle probabilità, la statistica ha cominciato ad interessarsi ai CAMPIONI.
- Un CAMPIONE è un SOTTOINSIEME del collettivo (popolazione) rispetto al quale si sta studiando un certo fenomeno: di solito, è un sottoinsieme di numerosità MOLTO INFERIORE a quella del collettivo di partenza.
- La TEORIA DELL'**INFERENZA STATISTICA** studia le tecniche per ricavare informazioni attendibili dai dati campionari.

21 e 22 novembre 2011

Statistica sociale

18



■ Le **TECNICHE CAMPIONARIE (TEORIA DEI CAMPIONI)** permettono di **OTTIMIZZARE** i criteri di **ESTRAZIONE DEL CAMPIONE** (il cosiddetto **DISEGNO DI CAMPIONAMENTO**), in maniera tale da **RICAVARE DAL CAMPIONE PRESSOCHÉ LE STESSA INFORMAZIONI CHE SI SAREBBERO RICAVATE DISPONENDO DELL'INTERO COLLETTIVO.**

21 e 22 novembre 2011 Statistica sociale 20

## IN STATISTICA:

- **RAPPRESENTATIVITÀ DEL CAMPIONE**

=

- **ESTRAZIONE CASUALE DELLE UNITA' DEL CAMPIONE**

21 e 22 novembre 2011

Statistica sociale

21

## Statistica e calcolo delle probabilità

- Le tecniche di inferenza statistica si basano tutte sulla "somiglianza" (che in termini tecnici si chiama "verosimiglianza") del campione rispetto alla popolazione da cui è stato estratto; estratto, non dobbiamo mai dimenticarlo, con criteri rigorosamente **casuali**.
- L'inferenza statistica è resa possibile dalla conoscenza delle leggi fondamentali del calcolo delle probabilità.

21 e 22 novembre 2011

Statistica sociale

22

## Variabili aleatorie e distribuzioni teoriche di probabilità

- In statistica descrittiva abbiamo visto cosa intendiamo per **variabile** statistica. Spostandoci in campo probabilistico, se una variabile può assumere valori esclusivamente dovuti al caso (o, per essere più precisi, a un **esperimento aleatorio**), essa prende il nome di **variabile aleatoria**.
- Una **v.a.** è un numero **X** che assume un valore in **R** (asse dei numeri reali), determinato sulla base di un evento **E**, che si riferisce a un certo esperimento aleatorio. A ciascun valore assunto da **X**, legato all'evento **E**, si associa una probabilità **P** (**P** varia tra 0 e 1).
- Una **distribuzione di probabilità** è una funzione che sintetizza la relazione tra i valori di una variabile casuale **X** e la probabilità ad essi associata. Una distribuzione di probabilità descrive il comportamento della v.a. a cui è associata.

21 e 22 novembre 2011

Statistica sociale

23

## Variabili aleatorie e distribuzioni teoriche di probabilità

- La conoscenza della distribuzione di probabilità di una variabile aleatoria fornisce ai ricercatori uno strumento potente per descrivere una **popolazione**, dalla quale verranno estratti i campioni che, successivamente, saranno studiati.
- Molti fenomeni naturali e biologici si distribuiscono secondo una **distribuzione normale**.
- Una distribuzione di probabilità è solitamente rappresentata da una formula (*funzione di probabilità* se la v.a. è discreta, *funzione di densità* se la v.a. è continua)

21 e 22 novembre 2011

Statistica sociale

24

## Variabili aleatorie e distribuzioni teoriche di probabilità

- La forma di una distribuzione può essere **simmetrica** rispetto al valore centrale, oppure vi può essere una coda più lunga da un lato piuttosto che dall'altro. Se la coda "lunga" è a destra (o, viceversa, a sinistra) la distribuzione avrà **asimmetrica positiva** (o, viceversa, **negativa**).
- Alcune distribuzioni teoriche di probabilità comunemente usate per descrivere i dati sono: la distribuzione Normale (o Gaussiana), la distribuzione Binomiale, la distribuzione di Poisson.
- Abbiamo già parlato ampiamente della v.a. **Normale** (o **Gaussiana**).

21 e 22 novembre 2011

Statistica sociale

25

## Dalla popolazione al campione: l'inferenza statistica

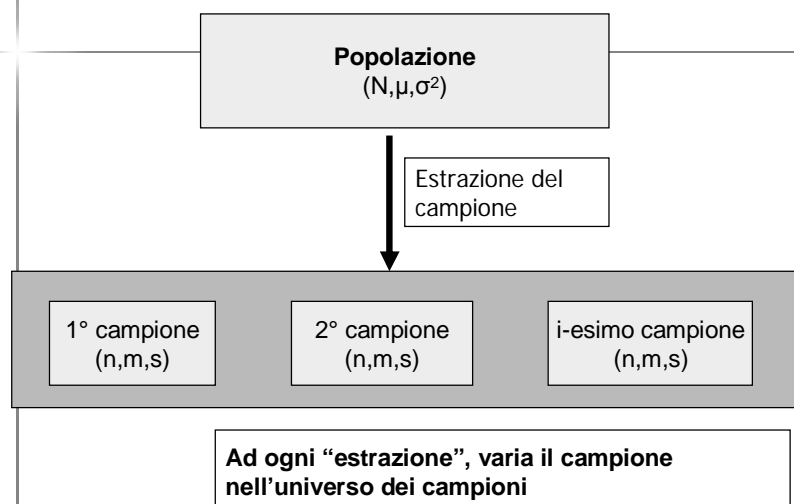
- Nelle indagini campionarie, l'obiettivo è quello di fare "**inferenza**" dal campione alla popolazione.
- In questo passaggio, poiché il campionamento è stato effettuato con criteri casuali, le "stime" che si ottengono sul campione, rispetto alla popolazione di partenza, sono per loro stessa natura affette da un errore, che si chiama errore di campionamento.
- Ad esempio, se io calcolo la media delle altezze di un campione di persone, per le quali so che la popolazione di partenza ha media 175 cm, non otterrò mai esattamente una media pari a 175 cm. Piuttosto, ripetendo infinite volte la stima su "infiniti" campioni, otterrò certamente una **DISTRIBUZIONE DI MEDIE**, che sarà dispersa attorno al vero valore della media, cioè 175 cm.

21 e 22 novembre 2011

Statistica sociale

26

## Popolazione e campione



21 e 22 novembre 2011

Statistica sociale

27

## La stima dei parametri di una popolazione

**I parametri** sono dei valori caratteristici della popolazione, come la media aritmetica, la probabilità del verificarsi di un evento, ecc.

Le stime sono effettuate in funzione delle osservazioni campionarie e, pertanto, dipendono dagli elementi del campione (media aritmetica del campione, frequenza di un certo evento nel campione, ecc.)

L'insieme di tutti i campioni estraibili casualmente da una popolazione è detto “spazio campionario”. Se la popolazione è finita, si parla di “universo dei campioni”.

Al variare del campione nell'universo campionario la stima assume valori diversi, valori dei quali sarà possibile costruire la distribuzione (distribuzione campionaria).

21 e 22 novembre 2011

Statistica sociale

28

## Stime e stimatori

- Per stimare un parametro (ad es. la media) della popolazione originaria si estrae un solo campione.
- Tutti i possibili campioni, virtualmente, sono estraibili, e sono pertanto possibili **diverse stime del parametro**, in numero corrispondente a quello dei possibili campioni.
- Si possono costruire pertanto le distribuzioni delle medie campionarie che, in termini probabilistici, costituiscono anch'esse delle variabili aleatorie, descrivibili da modelli discreti o continui.

21 e 22 novembre 2011

Statistica sociale

29

## Gli stimatori

- Lo **stimatore** è una variabile aleatoria definita nell'universo dei campioni; lo stimatore è una variabile aleatoria, che assume valori in ciascun campione compreso nell'universo dei campioni.
- Mentre una **stima** è una determinazione empirica (una "realizzazione") del corrispondente stimatore.
- Vediamo perché.

21 e 22 novembre 2011

Statistica sociale

30

## Stimatori e statistiche

- Supponiamo di voler stimare, nella popolazione, un certo parametro  $\theta$ .
- A partire dai dati osservati sul campione è possibile calcolare una **statistica t**, cioè una certa **funzione dei dati del campione** utilizzata allo scopo di stimare il parametro  $\theta$  incognito:
- $t = f(x_1, x_2, \dots, x_n)$

21 e 22 novembre 2011

Statistica sociale

31

## Stimatori e statistiche

- Una statistica è, pertanto, una funzione applicata all'insieme degli  $n$  dati del campione. Ognuno di questi dati, poiché i campioni variano nell'universo dei campioni, descrive a sua volta una variabile aleatoria.
- Queste variabili aleatorie, se il campionamento è perfettamente casuale, sono tra loro indipendenti e identicamente distribuite.

21 e 22 novembre 2011

Statistica sociale

32



## Stimatori e statistiche

- Da ciò consegue che la statistica si può vedere, in ultima analisi, come una particolare determinazione di una funzione di variabili aleatorie.
- Definiamo pertanto stimatore questa funzione,  $T_\theta$ , funzione di  $n$  variabili aleatorie, utilizzata per stimare il parametro  $\theta$ .
- $T_\theta = f(X_1, X_2, \dots, X_n)$

21 e 22 novembre 2011

Statistica sociale

33

## Proprietà degli estimatori

- Ad esempio, la **media campionaria** è uno **stimatore** del parametro "media" della popolazione;
- Come abbiamo visto, uno stimatore è una v.a., funzione dei valori del campione che, una volta calcolata, restituisce una **stima**;
- Uno stimatore può godere di alcune **proprietà** desiderabili.

21 e 22 novembre 2011

Statistica sociale

34

## Correttezza o non distorsione

Uno stimatore  $T$  si dice corretto, o non distorto, se il suo valore atteso coincide con il valore del parametro che intende stimare. Ad esempio la media campionaria è uno stimatore **corretto** della media della popolazione. Se, invece, il valore atteso non coincide con il parametro lo stimatore si dice **distorto**.

$$E(T_\theta) = \theta \quad \text{stimatore non distorto}$$

$$E(T_\theta) \neq \theta \quad \text{stimatore distorto}$$

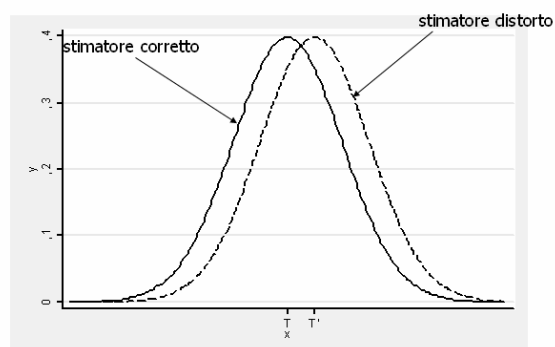
$$B = E(T_\theta) - \theta \quad \text{distorsione (bias)}$$

21 e 22 novembre 2011

Statistica sociale

35

## Distribuzione di uno stimatore



21 e 22 novembre 2011

Statistica sociale

36

## Esempi di stimatori corretti e distorti

- La media campionaria è uno stimatore corretto della media della popolazione: infatti, il suo valore atteso è proprio  $\mu$ ;
- La varianza campionaria è uno stimatore *distorto* (se  $n$  è piccolo) della varianza della popolazione

21 e 22 novembre 2011

Statistica sociale

37

## Varianza campionaria corretta

Si può dimostrare che la varianza campionaria **non** è uno stimatore corretto perché il suo valore atteso non coincide con  $\sigma^2$ . Si usa allora la sua "versione corretta" che si calcola nel seguente modo:

$$s^2 = \frac{\sum_{j=1}^n (x_j - \bar{x})^2}{n-1}$$

21 e 22 novembre 2011

Statistica sociale

38

# Efficienza

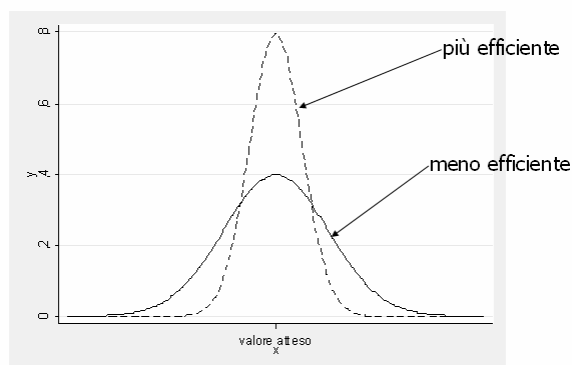
- Da uno stimatore ci aspettiamo non solo un valore medio uguale a quello della popolazione, ma anche che lo stimatore sia abbastanza "concentrato" attorno al valore medio.
- Se valutiamo due stimatori entrambi corretti, dobbiamo preferire quello con varianza minore perché è più **efficiente**, cioè riduce il margine di errore.

21 e 22 novembre 2011

Statistica sociale

39

## Confronto tra due stimatori entrambi corretti



21 e 22 novembre 2011

Statistica sociale

40

## Efficienza relativa

Quello di efficienza è un concetto sempre relativo. Dati due stimatori entrambi corretti,  $T_1$  e  $T_2$ , si definisce "efficienza relativa" del primo rispetto al secondo il rapporto tra la varianza del secondo e la varianza del primo:

$$\text{Eff. relativa } (T_1, T_2) = \frac{\text{Var}(T_2)}{\text{Var}(T_1)}$$

Ad esempio, la mediana campionaria è meno efficiente della media campionaria perché ha una varianza più elevata. Il rapporto tra varianza della mediana e della media campionaria è circa 1.5. Dunque la media campionaria è preferibile come stimatore della media della popolazione.

## Consistenza

- Se, all'aumentare della dimensione  $n$  del campione, il valore atteso di uno stimatore tende a concentrarsi intorno al vero valore del parametro, lo stimatore si dice **consistente**.

## Metodi per il reperimento dello stimatore "migliore"

- Abbiamo visto finora quali sono le proprietà "auspicabili" per uno stimatore.
- Gli stimatori finalizzati alla stima di un parametro, però, possono essere vari: certamente non c'è un solo stimatore possibile.
- In questo corso, non ci addentriamo sui metodi con i quali è possibile reperire gli stimatori "ottimali" per i rispettivi parametri.
- Il metodo di gran lunga più utilizzato è il cosiddetto "**metodo della massima verosimiglianza**".
- Si basa sul presupposto secondo cui il miglior stimatore possibile è quello che rende "più probabile" l'estrazione di un campione rispetto a una certa popolazione (la cui distribuzione è nota), dalla quale il campione viene estratto.
- Uno stimatore ottenuto con questo metodo si dice **stimatore di massima verosimiglianza** di un dato parametro,  $\theta$ .

21 e 22 novembre 2011

Statistica sociale

43

## La variabilità campionaria

- Come abbiamo visto, una certa statistica, applicata al campione, il cui scopo è quello di stimare il corrispondente parametro incognito nella popolazione si dice **stimatore**.
- L'aspetto che va sottolineato è che uno stimatore, che è funzione di  $n$  variabili aleatorie, è esso stesso una variabile aleatoria.

21 e 22 novembre 2011

Statistica sociale

44

# La variabilità campionaria

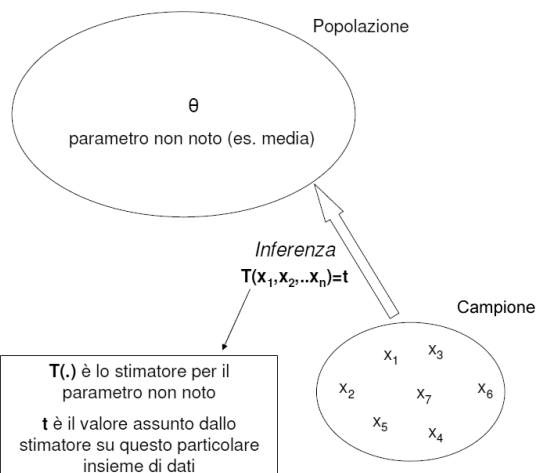
- Infatti, il valore assunto dallo stimatore varia al variare dei dati sui quali si applica (varia al variare del campione).
- Si può definire la probabilità che lo stimatore assuma dei particolari valori. La probabilità sui possibili valori assunti dallo stimatore è indotta dalla **legge di probabilità** definita sullo spazio dei possibili campioni.
- La distribuzione di probabilità di una statistica è chiamata **distribuzione campionaria** della statistica.

21 e 22 novembre 2011

Statistica sociale

45

# Schematicamente



21 e 22 novembre 2011

Statistica sociale

46

## Distribuzione della media campionaria

- Sia  $Y$  una variabile aleatoria con media  $\mu$  e varianza  $\sigma^2$
- Sia dato un campione di numerosità  $n$  su cui si osserva la variabile  $Y$ , che nel campione assumerà le determinazioni  $Y_1, Y_2, \dots, Y_n$
- La media campionaria è lo stimatore:

$$\bar{Y} = \frac{\sum_{j=1}^n Y_j}{n}$$

21 e 22 novembre 2011

Statistica sociale

47

## Media campionaria ed "errore standard"

- Si dimostra che la distribuzione dello stimatore "media campionaria" ha media  $\mu$  (la stessa della popolazione da cui il campione è stato estratto) e varianza  $\sigma^2/n$ .
- Lo stimatore media campionaria, pertanto, è non distorto.
- La deviazione standard della distribuzione campionaria è chiamata **errore standard**.
- Il termine **errore standard** è utilizzato per distinguere la deviazione standard di una statistica campionaria da quella della popolazione da cui il campione è stato estratto.

21 e 22 novembre 2011

Statistica sociale

48



## Teorema centrale del limite

- Se un campionamento viene ripetuto "infinite" volte, per il **teorema centrale del limite**, la media campionaria tende a distribuirsi in modo normale, anche quando non lo è la popolazione da cui i campioni sono stati estratti.
- Se  $n$  è "sufficientemente grande" ( $>30$ ), la forma della distribuzione campionaria delle medie è approssimativamente normale, indipendentemente dalla forma della distribuzione della popolazione di origine.

21 e 22 novembre 2011

Statistica sociale

49

## Teorema centrale del limite

- Il valore medio dell'insieme di tutte le possibili medie campionarie è uguale alla media  $\mu$  della popolazione di origine.
- La deviazione standard dell'insieme di tutte le possibili medie campionarie di campioni di numerosità  $n$ , detta **errore standard**, è funzione sia della deviazione standard della popolazione, sia della numerosità del campione:

$$ES(X) = \sigma_{\bar{X}} = \frac{\sigma_{POP}}{\sqrt{n}}$$

21 e 22 novembre 2011

Statistica sociale

50

## Errore standard stimato nel campione

Raramente conosceremo il valore di  $\sigma$  nella popolazione. Più spesso, dovremo **stimare** il valore di  $\sigma$  con il valore di  $s$  campionario (ottenuto con l'opportuno stimatore, che abbiamo visto in precedenza), e di conseguenza, calcolare il corrispondente valore dell'errore standard campionario.

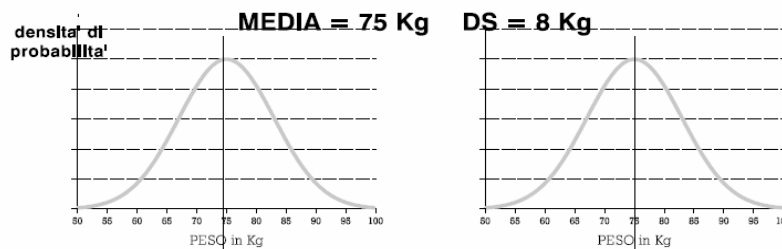
$$ES(X)_{stimato} = s_{\bar{X}} = \frac{s_{campione}}{\sqrt{n}}$$

21 e 22 novembre 2011

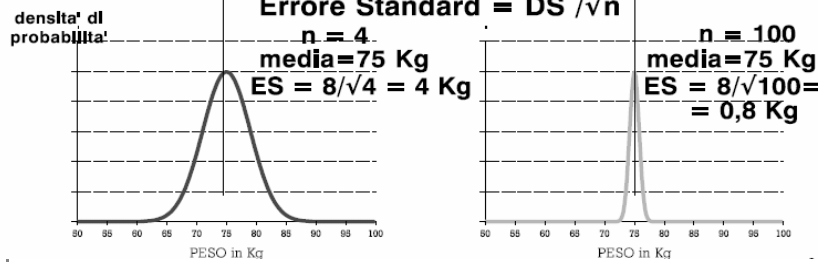
Statistica sociale

51

### DISTRIBUZIONE del PESO nella POPOLAZIONE ORIGINARIA



### DISTRIBUZIONE delle MEDIE CAMPIONARIE del PESO



21 e 22 novembre 2011

Statistica sociale

52

## Un semplice esempio

- Una v.a.,  $Y$ , assume valori discreti da 1 a 5, tutti con uguale probabilità ( $p=0,2=20\%$ ).
- La popolazione è formata da 5 unità.
- Si tratta di una v.a. uniforme discreta.

Y	Freq.	Percent	Cum.
1	1	20.00	20.00
2	1	20.00	40.00
3	1	20.00	60.00
4	1	20.00	80.00
5	1	20.00	100.00
Total	5	100.00	

21 e 22 novembre 2011

Statistica sociale

53

## Un semplice esempio

- Supponiamo di estrarre un campione di  $n=2$  unità dalla popolazione di  $N=5$  unità.
- I momenti della v.a. media campionaria saranno quindi, secondo la regola appena vista:
- $E(Y_{\text{medio}}) = E(Y) = (1+2+3+4+5)/5 = 15/5 = 3$
- $\text{Var}(Y_{\text{medio}}) = s^2/n = [[(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2]/5]/2 = [4 + 1 + 0 + 1 + 4]/5/2 = [10/5]/2 = 2/2 = 1$

21 e 22 novembre 2011

Statistica sociale

54

## Proviamo a verificarlo empiricamente

- Se estraiamo tutti i possibili campioni di 2 unità (che sono  $5^2 = 25$ ), otteniamo 25 possibili medie, qui schematizzate:

	1	2	3	4	5
1	1	1.5	2	2.5	3
2	1.5	2	2.5	3	3.5
3	2	2.5	3	3.5	4
4	2.5	3	3.5	4	4.5
5	3	3.5	4	4.5	5

21 e 22 novembre 2011

Statistica sociale

55

## Distribuzione empirica della media campionaria

- Se raccogliamo "in distribuzione" la tabella appena vista, otteniamo la distribuzione che segue:

$\bar{Y}$	Freq.	Percent	Cum.
1	1	4.00	4.00
1.5	2	8.00	12.00
2	3	12.00	24.00
2.5	4	16.00	40.00
3	5	20.00	60.00
3.5	4	16.00	76.00
4	3	12.00	88.00
4.5	2	8.00	96.00
5	1	4.00	100.00
Total	25	100.00	

21 e 22 novembre 2011

Statistica sociale

56

## Calcoliamo la media e la varianza della media campionaria

$$\begin{aligned} \blacksquare E(Y_{\text{medio}}) &= \\ & (1+1,5 \cdot 2+2 \cdot 3+2,5 \cdot 4+3 \cdot 5+3,5 \cdot 4+4 \cdot 3+4,5 \cdot 2+ \\ & +5)/25 = 75/25 = \mathbf{3} \end{aligned}$$

$$\begin{aligned} \blacksquare \text{Var}(Y_{\text{medio}}) &= [(1-3)^2+(1,5-3)^2 \cdot 2+(2-3)^2 \cdot 3+ \\ & +(2,5-3)^2 \cdot 4+(3-3)^2 \cdot 5+(3,5-3)^2 \cdot 4+(4-3)^2 \cdot 3+ \\ & +(4,5-3)^2 \cdot 2+(5-3)^2)]/25 = 25/25 = \mathbf{1} \end{aligned}$$

Abbiamo così verificato "empiricamente" i valori, rispettivamente, della media e della varianza campionaria.

21 e 22 novembre 2011

Statistica sociale

57

## Dalla popolazione al campione. Per riassumere:

- **INFERENZA STATISTICA** : Insieme delle operazioni compiute dal ricercatore per "stimare" alcune caratteristiche (i **parametri**) di una popolazione, non interamente esplorabile, attraverso la selezione da questa di un sottoinsieme casuale di unità, detto **campione**.
- **PARAMETRO**: "Vero" valore ( $\theta$ ) assunto da una caratteristica misurata a livello di popolazione (somma, media, varianza, proporzione, coefficiente di regressione, coefficiente di correlazione, ecc.). Il parametro è, quasi sempre, incognito.
- **STIMATORE**: si dice **stimatore** qualunque **statistica**  $T(X_1, X_2, \dots, X_n)$ , ovvero una funzione applicata alle unità statistiche comprese nel campione, le cui determinazioni vengono utilizzate per ottenere una misura (stima puntuale) del parametro incognito  $\theta$ . Pertanto, uno stimatore è esso stesso una variabile casuale e possiede una sua distribuzione, con i relativi momenti: valore atteso, varianza, ecc.

21 e 22 novembre 2011

Statistica sociale

58

## Errore di campionamento

- **ERRORE DI CAMPIONAMENTO:** Differenza tra il valore empirico dello stimatore (esempio: media) e il corrispondente valore che si sarebbe ottenuto analizzando la totalità delle unità statistiche della popolazione.
- L'errore di campionamento si verifica, come già accennato, perché quello osservato è solo un sottoinsieme (talvolta molto piccolo) delle unità della popolazione. L'errore di campionamento tende a diminuire all'aumentare della numerosità campionaria.

21 e 22 novembre 2011

Statistica sociale

59

## Stima puntuale, stima intervallare e test di ipotesi

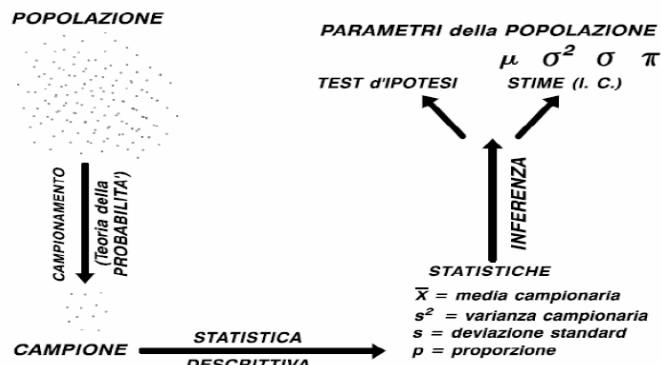
- **OPERAZIONE DI STIMA:** l'insieme delle regole e attraverso le quali è stato ottenuto un particolare valore dello stimatore.
- La stima può riguardare soltanto il parametro, e in tal caso si parla di "stima puntuale", oppure un intervallo attorno al parametro stesso, e allora si parla di "stima intervallare".
- **TEST DI IPOTESI:** Processo decisionale, basato sul controllo di ipotesi statistiche effettuate sulla realtà osservata. Tale processo porta a rifiutare, oppure non rifiutare, una certa ipotesi (statistica) formulata sulla popolazione.

21 e 22 novembre 2011

Statistica sociale

60

# Schematicamente



21 e 22 novembre 2011

Statistica sociale

61

# Stima puntuale

- Sui dati del campione di cui disponiamo, il valore calcolato, empirico dello stimatore (di un certo parametro) è la **stima puntuale del parametro**.
- Solitamente si usa:
  - La **media campionaria** per stimare la **media della popolazione**;
  - La **varianza campionaria** per stimare la **varianza della popolazione**;
  - La **differenza tra due medie campionarie** per stimare la **differenza tra due valori medi a livello di popolazione**;
  - Eccetera

21 e 22 novembre 2011

Statistica sociale

62

## Stima per intervallo

- Con questa procedura di stima si determina un set di valori a partire dal campione che con una certa **probabilità**,
- $(1-\alpha)\%$ , contiene il parametro incognito.
- $(1-\alpha)\%$  indica il **livello (o grado) di confidenza**; l'intervallo che si ottiene è detto **intervallo di confidenza**.
- Gli estremi dell'intervallo dipendono dal campione estratto, quindi sono sottoposti a **variazioni casuali**.

21 e 22 novembre 2011

Statistica sociale

63

## Intervalli di confidenza

- Un intervallo di confidenza è quindi un **insieme di valori plausibili** per il parametro incognito sulla base dell'**evidenza empirica**.
- **Attenzione**: il livello di confidenza rappresenta il **grado di affidabilità della procedura**, non il grado di plausibilità **del risultato**, dovuto al singolo campione. La plausibilità del risultato è invece espressa dalla ampiezza dell'intervallo di confidenza.
- Generalmente, si usa come livello di confidenza il 95% ( $\alpha=5\%$ )

21 e 22 novembre 2011

Statistica sociale

64



## Ampiezza dell'intervallo di confidenza

- L'**ampiezza** dell'intervallo è molto importante:
- Quanto **più** l'intervallo è **ridotto**, tanto **maggiore** è il grado di **precisione della stima**.

21 e 22 novembre 2011

Statistica sociale

65

## Ampiezza dell'intervallo di confidenza

- L'**ampiezza dell'intervallo** dipende quindi:
- dal grado di confidenza ( $1-\alpha$ ): **al diminuire di  $\alpha$**  [cioè al **crescere del grado di confidenza ( $1-\alpha$ )**], l'ampiezza dell'intervallo **augmenta**;
- dalla **variabilità** del fenomeno studiato: al **crescere della variabilità** dei dati che stiamo osservando cresce anche l'incertezza e quindi l'ampiezza dell'intervallo **augmenta**;
- dalla numerosità campionaria,  $n$ : al **crescere di  $n$**  aumenta la quantità di **informazione** disponibile e quindi l'ampiezza dell'intervallo **diminuisce**.

21 e 22 novembre 2011

Statistica sociale

66

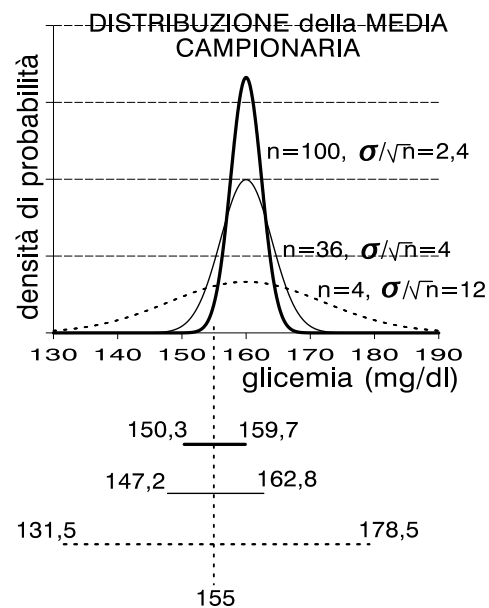
## Come si può far diminuire l'ampiezza di un intervallo di confidenza?

- L'ampiezza di un intervallo di confidenza diminuisce se:
  - 1) diminuisce il **livello di confidenza** ( $1-\alpha$ )  
(es. dal 99% al 95% al 90%)
  - 2) aumenta la **numerosità** del campione  
(es. da  $n=4$  a  $n=36$  a  $n=100$ )
  - 3) diminuisce la **variabilità nella popolazione**  
(es. da  $\sigma=100$  a  $\sigma=36$  a  $\sigma=4$ )

21 e 22 novembre 2011

Statistica sociale

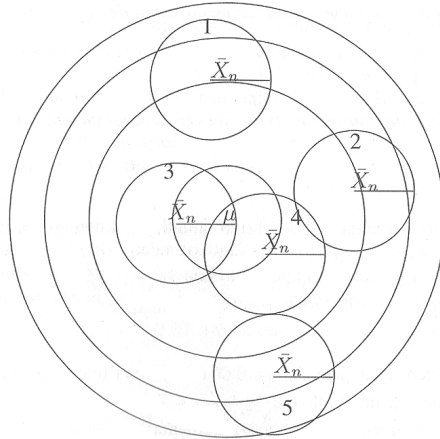
67



21 e 22 novembre

68

## Esempio: "tiro al bersaglio" con la media campionaria



- Al centro del bersaglio c'è la media incognita  $\mu$ . I cerchi con al centro i valori calcolati della media campionaria sono gli intervalli di confidenza di ampiezza costante, il cui raggio dipende dalla numerosità campionaria,  $n$ , e dal livello di confidenza  $\alpha$ .
- Come si nota, alcuni intervalli (cerchi) non contengono il valore  $\mu$  (si tratta degli intervalli 1, 2 e 5), mentre altri invece lo contengono (gli intervalli 3 e 4).
- Si può interpretare il livello di confidenza,  $(1-\alpha)$ , come la probabilità che i cerchi (intervalli) contengano il valore incognito  $\mu$ .

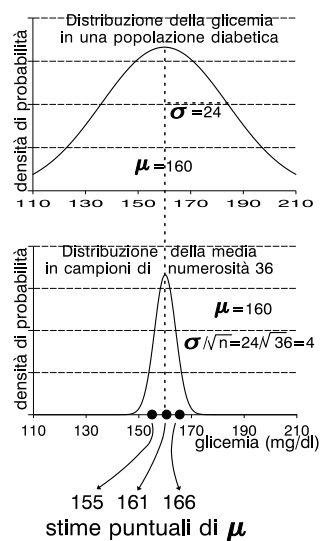
21 e 22 novembre 2011

Statistica sociale

69

## Stima puntuale e stima intervallare

- Una singola stima del valore medio è una **stima puntuale** (è un unico numero).
- Ma: il valore medio di un unico campione può essere una buona stima del valore medio di una popolazione, quando sappiamo che, prendendo diversi campioni della stessa popolazione, otterremo sempre un valore diverso?



21 e 22 novembre 2011

Statisti

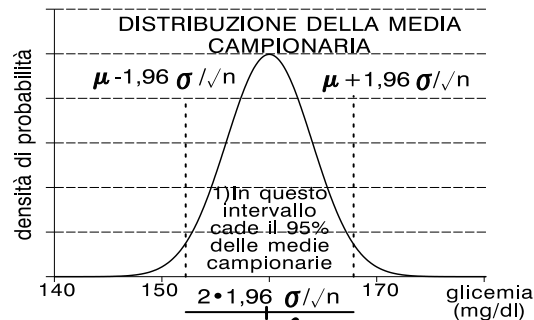
## Stima puntuale e stima intervallare

- Sappiamo che i campioni tratti da una stessa popolazione variano in modo sistematico.
- Oltre alla stima puntuale, possiamo conoscere anche una stima dell'intervallo in cui è "molto probabile" che cada il valore vero del parametro per la popolazione.
- Chiameremo questa valutazione è chiamata **stima intervallare**.

21 e 22 novembre 2011

Statistica sociale

71



2) Riportiamo l'intervallo intorno a ciascuna media campionaria

3) Il 95% di questi intervalli (di confidenza) contiene la media vera della popolazione

21 e 22 nov

72

## Stima puntuale e stima intervallare

- ❑ La **stima puntuale** fornisce un singolo valore. Da ciò consegue che:
  - ❑ 1) questo valore non coincide quasi mai con il valore vero (parametro) della popolazione;
  - ❑ 2) campioni diversi forniscono stime puntuali diverse.
  
- ❑ La **stima intervallare** fornisce un intervallo, che ha una predeterminata probabilità di contenere il valore vero del parametro. Pertanto, quest'intervallo ha una determinata probabilità prefissata (ad esempio, il 95%) di contenere il valore vero del parametro (della popolazione).

21 e 22 novembre 2011

Statistica sociale

73

## Esempio: intervallo di confidenza per la media, da una popolazione normale

- Intervallo di confidenza, con grado di confidenza  $(1-\alpha)\%$  =
  - $= m \pm z_{\alpha/2} * [\sigma(X)/\sqrt{n}]$
- Dove:
  - $m$  = valore calcolato dello stimatore "media campionaria";
  - $z_{\alpha/2}$  = valore critico della "normale standardizzata" in corrispondenza di  $\alpha$  ;
  - $[\sigma(X)/\sqrt{n}]$  = deviazione standard della v.c. "media campionaria (questa quantità è detta "errore standard" di X).

21 e 22 novembre 2011

Statistica sociale

74

## L'INTERVALLO DI CONFIDENZA, AL 95%, PER LA MEDIA CAMPIONARIA (da popolazione NORMALE)

- Se:
- $X \sim N(\mu, (\sigma^2/n))$  e quindi  $Z \sim N(0,1)$  [normale standardizzata]
- Sappiamo che, con probabilità pari al 95%, l'intervallo compreso tra gli estremi:
- $[ m - z_{\alpha/2} * [\sigma(X)/\sqrt{n}]; m + z_{\alpha/2} * [\sigma(X)/\sqrt{n}] ]$
- comprenderà il **vero valore del parametro  $\mu$** , non noto.

21 e 22 novembre 2011

Statistica sociale

75

## Un esempio

- Supponiamo di voler fare inferenza sul peso medio, alla nascita, di neonati di sesso maschile con età gestazionale di 39 settimane.
- **ESERCIZIO.**
- Sapendo che il peso alla nascita è una v.a. normale, con media incognita ( $\mu$ ) e deviazione standard ( $\sigma$ ) nota, pari a 535 grammi, si calcoli l'intervallo al 95% per  $\mu$  a partire da un campione casuale semplice, estratto dalla popolazione, di numerosità pari a 16.
- $n = 16$
- Media campionaria calcolata = 3434 g
- $\sigma$  (nota) = 535 g

21 e 22 novembre 2011

Statistica sociale

76

## Un esempio

- Limite superiore dell' I.C. 95% =  $3434 + 1,96 * (535/\sqrt{16}) = 3172$
- Limite inferiore dell' I.C. 95% =  $3434 - 1,96 * (535/\sqrt{16}) = 3696$
- Intervallo di confidenza al 95%: [3172 ; 3696]
- **QUINDI:**
- Con probabilità pari al 95%, il peso medio alla nascita dei neonati maschi, nati alla 39ma settimana di gestazione, è un valore compreso tra **3172 e 3696**.

21 e 22 novembre 2011

Statistica sociale

77

## Intervalli di confidenza: un esempio tratto dalla letteratura

- Mostriamo un esempio tratto dalla ricerca medica, anche se gli intervalli di confidenza riportati sono stati ricavati con un'altra tecnica statistica;
- Dati tratti da un articolo apparso sul New England Journal of Medicine nel 1996: *Bone mineral density in women with depression*;
- L'ipotesi è che il soffrire o l'aver sofferto di depressione in passato provochi un **calo** della densità ossea nelle donne (meccanismi endocrini).

21 e 22 novembre 2011

Statistica sociale

78

## Intervalli di confidenza: un esempio tratto dalla letteratura

**TABLE 3. BONE MINERAL DENSITY IN 24 DEPRESSED AND 24 NORMAL WOMEN.\***

BONE MEASURED†	DEPRESSED WOMEN	NORMAL WOMEN	MEAN DIFFERENCE (95% CI)	P VALUE
Lumbar spine (anteroposterior) Density (g/cm <sup>2</sup> )	1.00±0.15	1.07±0.09	0.08 (0.02 to 0.14)	0.02
Lumbar spine (lateral)‡ Density (g/cm <sup>2</sup> )	0.74±0.09	0.79±0.07	0.05 (0.00 to 0.09)	0.03
Femoral neck Density (g/cm <sup>2</sup> )	0.76±0.11	0.88±0.11	0.11 (0.06 to 0.17)	<0.001
Ward's triangle Density (g/cm <sup>2</sup> )	0.70±0.14	0.81±0.13	0.11 (0.06 to 0.17)	<0.001
Trochanter Density (g/cm <sup>2</sup> )	0.66±0.11	0.74±0.08	0.08 (0.04 to 0.13)	<0.001
Radius Density (g/cm <sup>2</sup> )	0.68±0.04	0.70±0.04	0.01 (-0.01 to 0.04)	0.25

\*Plus-minus values are means ±SD. CI denotes confidence interval.

79