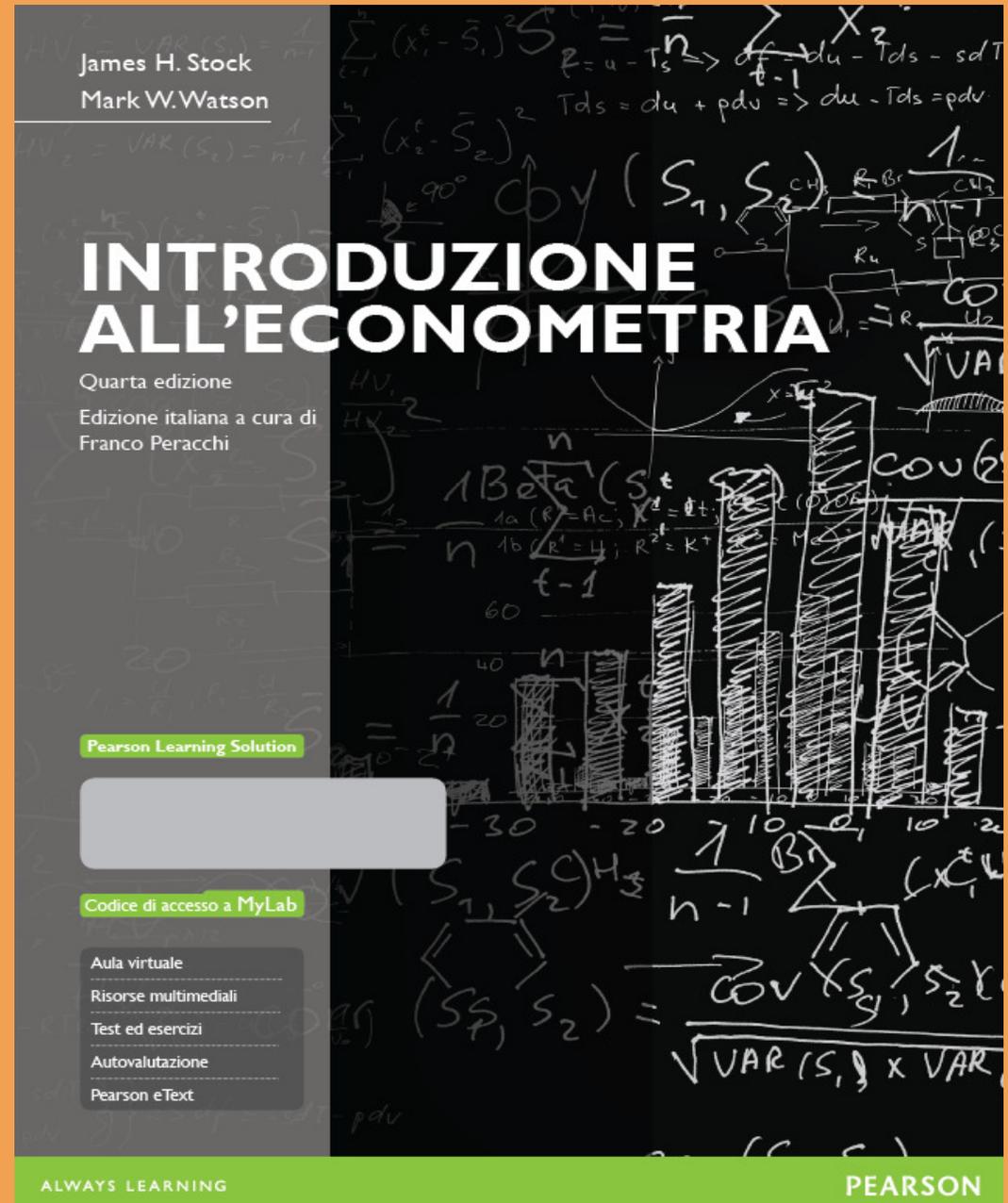


# Introduzione all' econometria

## Capitoli 1, 2 e 3

L' analisi statistica di dati  
economici (e correlati)



# Breve panoramica del corso

- L' economia suggerisce importanti relazioni, spesso con implicazioni politiche, ma praticamente mai fornisce dimensioni quantitative di effetti causali.
  - Qual è l' effetto *quantitativo* della riduzione delle dimensioni delle classi sui risultati degli studenti?
  - In che modo un anno in più di istruzione può influire sul reddito?
  - Qual è l' elasticità al prezzo delle sigarette?
  - Qual è l' effetto sulla crescita del prodotto interno lordo di un aumento di 1 punto percentuale nei tassi di interesse stabilito dalla Fed?
  - Qual è l' effetto sui prezzi delle case dei miglioramenti di tipo ambientale?

# Questo corso tratta dell' uso dei dati per misurare effetti causali

- Idealmente vorremmo un esperimento
  - Quale sarebbe un esperimento per stimare l' effetto della dimensione delle classi sui punteggi nei test standardizzati?
- Ma quasi sempre abbiamo a disposizione soltanto dati osservazionali (non sperimentali).
  - rendimenti dell' istruzione
  - prezzi delle sigarette
  - politica monetaria
- La maggior parte del corso affronta le difficoltà che derivano dall' uso di dati non sperimentali per stimare effetti causali
  - effetti perturbativi (fattori omessi)
  - causalità simultanea
  - “la correlazione non implica causalità”

## In questo corso:

- apprenderete metodi per stimare effetti causali usando dati non sperimentali;
- apprenderete l'uso di alcuni strumenti che possono essere impiegati per altri scopi, per esempio la previsione utilizzando serie di dati temporali;
- vi focalizzerete sulle applicazioni – si ricorre alla teoria solo ove necessario per comprendere i motivi alla base dei metodi;
- imparerete a valutare l'analisi di regressione effettuata da altri – questo significa che sarete in grado di leggere e comprendere articoli economici di carattere empirico in altri corsi di tipo economico;
- farete un po' di esperienza pratica con l'analisi di regressione nelle serie di esercizi.

# Richiami di probabilità e statistica (Capitoli 2, 3)

- **Problema empirico:** Dimensione della classe e risultato dell'istruzione
  - Domanda: qual è l'effetto sui punteggi nei test (o su un'altra misura di risultato) della riduzione della dimensione delle classi di uno studente per classe? E di 8 studenti per classe?
  - Dobbiamo utilizzare i dati per rispondere (esiste un modo per rispondere a questa domanda *senza* dati?)

# I dati dei punteggi nei test della California

Tutti i distretti scolastici K-6 e K-8 della California  
( $n = 420$ )

Variabili:

- Punteggi nei test del quinto anno (*Stanford-9 achievement test*, combinazione di matematica e lettura), media del distretto
- Rapporto studenti/insegnanti (STR) = numero di studenti nel distretto diviso per numero di insegnanti a tempo pieno equivalente

# Primo sguardo ai dati:

(dovreste già sapere come interpretare questa tabella)

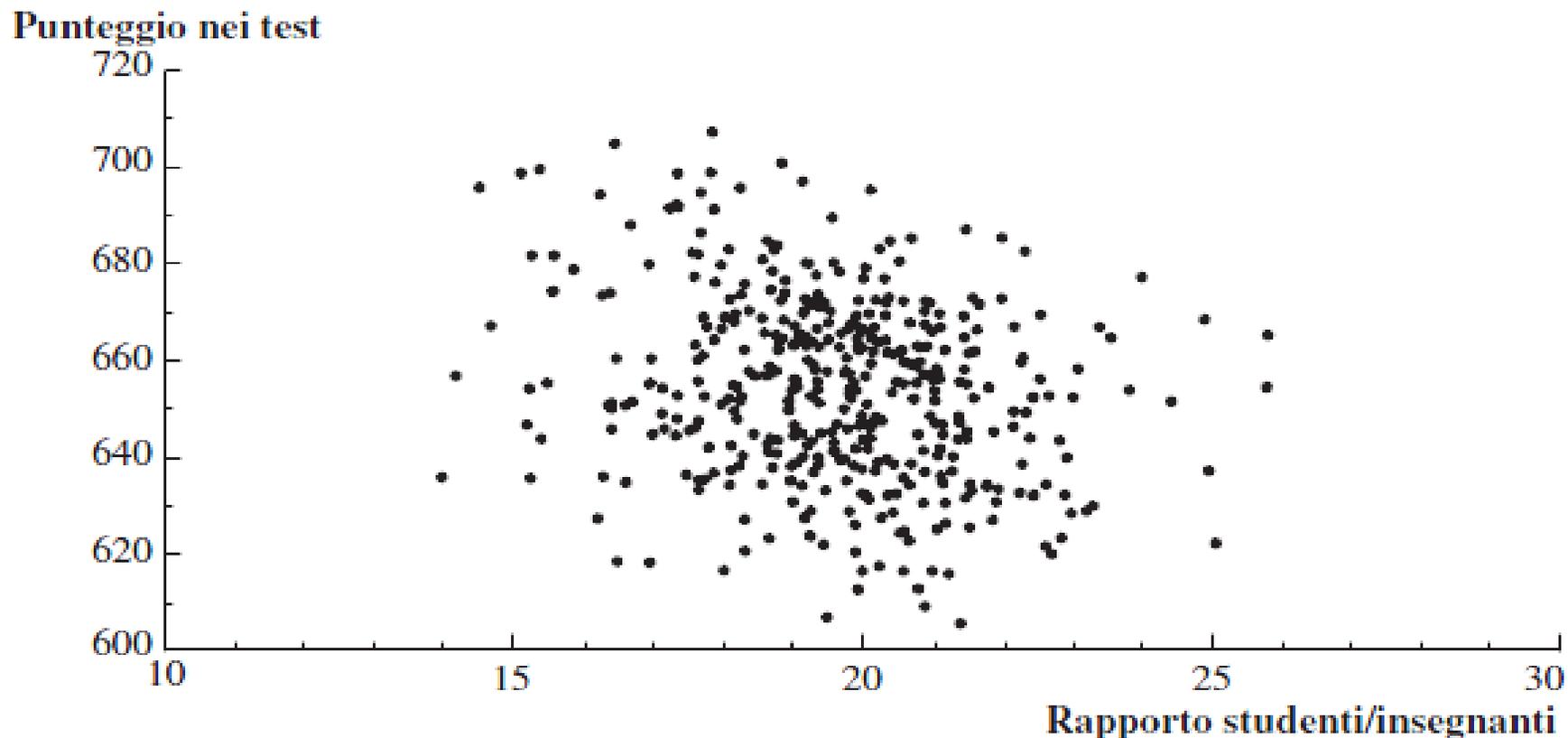
**Tabella 4.1** Sintesi della distribuzione del rapporto studenti/insegnanti e del punteggio nei test relativa al quinto grado d'istruzione (quinta elementare) per 420 distretti K-8 in California nel 1998.

|                              | Media | Deviazione standard | Percentile |       |       |               |       |       |       |
|------------------------------|-------|---------------------|------------|-------|-------|---------------|-------|-------|-------|
|                              |       |                     | 10%        | 25%   | 40%   | 50% (mediana) | 60%   | 75%   | 90%   |
| Rapporto studenti/insegnanti | 19,6  | 1,9                 | 17,3       | 18,6  | 19,3  | 19,7          | 20,1  | 20,9  | 21,9  |
| Punteggio nei test           | 654,2 | 19,1                | 630,4      | 640,0 | 649,1 | 654,5         | 659,4 | 666,7 | 679,1 |

Questa tabella non ci dice nulla sulla relazione tra punteggi nei test e *STR*.

# I distretti con classi più piccole ottengono punteggi più elevati nei test?

**Diagramma a nuvola** di punteggio nei test e STR



*Che cosa mostra questa figura?*

## **Dobbiamo ottenere evidenza numerica che indichi se i distretti con basso STR hanno punteggi nei test più alti – ma come?**

1. Confrontare i punteggi nei test nei distretti con basso STR con quelli con alto STR (“**stima**”)
2. Sottoporre a verifica l’ipotesi “nulla” che i punteggi medi nei test nei due tipi di distretti siano gli stessi, contro l’ipotesi “alternativa” che siano diversi (“**test di ipotesi**”)
3. Stimare un intervallo per la differenza nei punteggi medi nei test, nei distretti con alto vs basso STR (“**intervallo di confidenza**”)

**Analisi dei dati iniziali:** confrontare i distretti con dimensioni delle classi “piccole” ( $STR < 20$ ) e “grandi” ( $STR \geq 20$ ) :

| Dimensione classe | Punteggio medio( $\bar{Y}$ ) | Deviazione standard ( $s_Y$ ) | $n$ |
|-------------------|------------------------------|-------------------------------|-----|
| Piccola           | 657,4                        | 19,4                          | 238 |
| Grande            | 650,0                        | 17,9                          | 182 |

- 1. Stima** di  $\Delta$  = differenza tra medie dei gruppi
- 2. Verifica dell'ipotesi che**  $\Delta = 0$
- 3. Costruire un intervallo di confidenza** per  $\Delta$

# 1. Stima

$$\begin{aligned}\bar{Y}_{piccola} - \bar{Y}_{grande} &= \frac{1}{n_{piccola}} \sum_{i=1}^{n_{piccola}} Y_i - \frac{1}{n_{grande}} \sum_{j=1}^{n_{grande}} Y_j \\ &= 657,4 - 650,0 \\ &= 7,4\end{aligned}$$

È una differenza da considerare grande nel mondo reale?

- Deviazione standard dei distretti = 19,1
- Differenza tra 60-esimo and 75-esimo percentile della distribuzione dei punteggi nei test = 667,6 - 659,4 = 8,2
- È una differenza sufficientemente grande da risultare importante per discussioni sulla riforma della scuola, per i genitori o per un comitato scolastico?

## 2. Verifica di ipotesi

Test di differenza tra medie: calcolare la *statistica-t*,

$$t = \frac{\bar{Y}_s - \bar{Y}_l}{\sqrt{\frac{s_s^2}{n_s} + \frac{s_l^2}{n_l}}} = \frac{\bar{Y}_s - \bar{Y}_l}{SE(\bar{Y}_s - \bar{Y}_l)} \quad (\text{ricordate?})$$

- dove  $SE(\bar{Y}_s - \bar{Y}_l)$  è l' "errore standard " di  $(\bar{Y}_s - \bar{Y}_l)$ , i pedici **s** e **l** indicano distretti con STR "small" (piccolo) e "large" (grande), e

$$s_s^2 = \frac{1}{n_s - 1} \sum_{i=1}^{n_s} (Y_i - \bar{Y}_s)^2$$

(ecc.)

Calcolare la statistica- $t$  per la differenza tra medie

| Dim     | $\bar{Y}$ | $S_Y$ | $n$ |
|---------|-----------|-------|-----|
| piccola | 657,4     | 19,4  | 238 |
| grande  | 650,0     | 17,9  | 182 |

$$t = \frac{\bar{Y}_s - \bar{Y}_l}{\sqrt{\frac{s_s^2}{n_s} + \frac{s_l^2}{n_l}}} = \frac{657,4 - 650,0}{\sqrt{\frac{19,4^2}{238} + \frac{17,9^2}{182}}} = \frac{7,4}{1,83} = 4,05$$

$|t\text{-act}| > 1,96$ , perciò si rifiuta (al livello di significatività del 5%) l'ipotesi nulla che le due medie coincidano.

### 3. Intervallo di confidenza

Un intervallo di confidenza al 95% per la differenza tra medie  $\Delta = \mu_s - \mu_l$  è

$$\begin{aligned} (\bar{Y}_s - \bar{Y}_l) \pm 1,96 \times SE(\bar{Y}_s - \bar{Y}_l) \\ = 7,4 \pm 1,96 \times 1,83 = (3,8, 11,0) \end{aligned}$$

*Due affermazioni equivalenti:*

1. L'intervallo di confidenza al 95% per  $\Delta$  non include 0;
2. L'ipotesi nulla  $H_0: \Delta = 0$ , è rifiutata al livello di significatività del 5%.

## E ora...

- I meccanismi di stima, verifica di ipotesi e intervalli di confidenza dovrebbero risultare familiari
- Questi concetti si estendono direttamente a regressione e relative varianti
- Prima di passare alla regressione, tuttavia, rivedremo alcuni elementi della teoria alla base di stima, verifica di ipotesi e intervalli di confidenza:
  - Perché queste procedure funzionano, e perché utilizzare proprio queste invece di altre?
  - Rivedremo i fondamenti teorici di statistica ed econometria

# Richiami di teoria statistica

- 1. Quadro di riferimento probabilistico per l'inferenza statistica**
2. Stima
3. Verifica
4. Intervalli di confidenza

## **Quadro di riferimento probabilistico per l'inferenza statistica**

- a) Popolazione, variabile casuale e distribuzione
- b) Momenti di una distribuzione (media, varianza, deviazione standard, covarianza, correlazione)
- c) Distribuzione condizionata e media condizionata
- d) Distribuzione di un campione di dati estratto a caso da una popolazione:  $Y_1, \dots, Y_n$

# (a) Popolazione, variabile casuale e distribuzione

## ***Popolazione***

- Il gruppo o l'insieme di tutte le possibili unità di interesse (distretti scolastici)
- Considereremo le popolazioni infinitamente grandi ( $\infty$  è un'approssimazione di "molto grande")

## ***Variabile casuale Y***

- Rappresentazione numerica di un risultato casuale (punteggio medio nei test del distretto, STR del distretto)

## ***Distribuzione di Y***

- Le probabilità di diversi valori di  $Y$  che si verificano nella popolazione, per esempio  $\Pr[Y = 650]$  (quando  $Y$  è discreta)
- oppure: le probabilità di insiemi di questi valori, per esempio  $\Pr[640 \leq Y \leq 660]$  (quando  $Y$  è continua).
- In particolare per v.c. discrete e continue è detta **funzione di ripartizione o di probabilità cumulata** la funzione:

$$F(y) = \text{prob}\{Y \leq y\}, y \in R$$

## (b) Momenti di una distribuzione: media, varianza, deviazione standard, covarianza, correlazione

**media** = valore atteso (aspettativa) di  $Y$

$$= E(Y)$$

$$= \mu_Y$$

**varianza** =  $E(Y - \mu_Y)^2$

$$= \sigma_Y^2$$

= misura della dispersione quadratica della distribuzione

**deviazione standard** =  $\sqrt{\text{varianza}}$  =  $\sigma_Y$

= misura della dispersione della distribuzione nell'unità di misura della v.c.  $Y$

# Momenti (continua)

$$\mathbf{asimmetria} = \frac{E \left[ (Y - \mu_Y)^3 \right]}{\sigma_Y^3}$$

= misura di asimmetria di una distribuzione

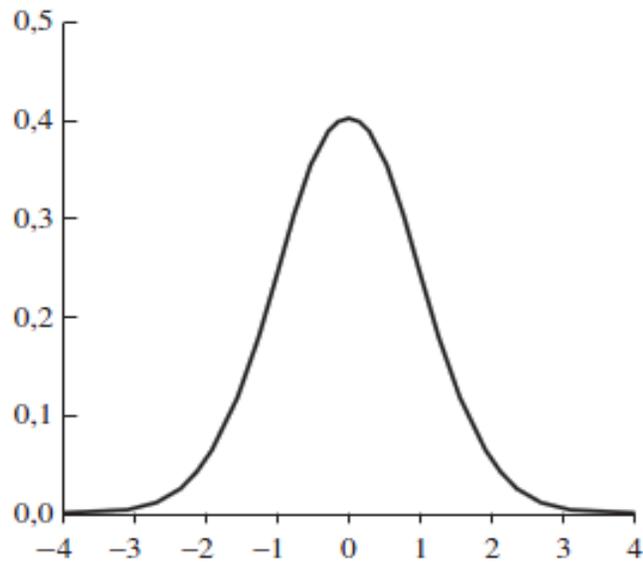
- *asimmetria* = 0: la distribuzione è simmetrica
- *asimmetria* > (<) 0: la distribuzione ha una coda lunga destra (sinistra)

$$\mathbf{curtosi} = \frac{E \left[ (Y - \mu_Y)^4 \right]}{\sigma_Y^4}$$

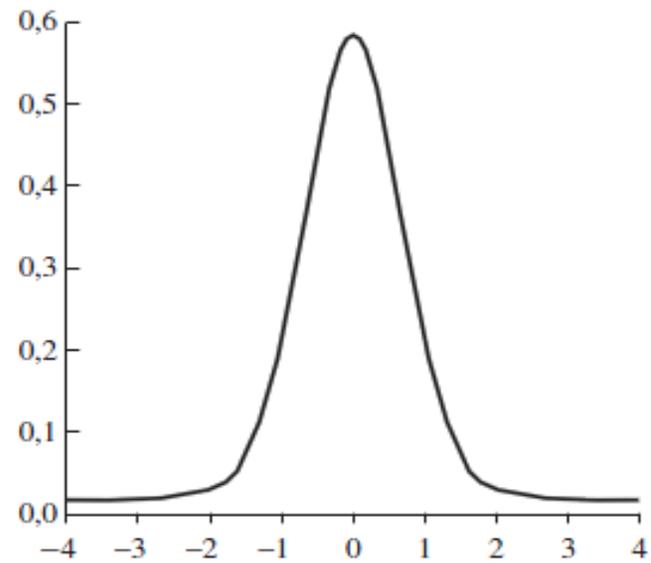
= misura di massa nelle code

= misura di probabilità di valori grandi

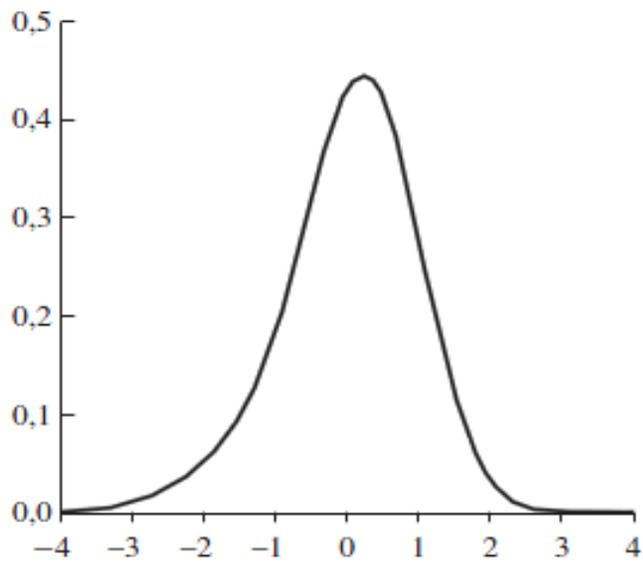
- *curtosi* = 3: distribuzione normale
- *curtosi* > 3: code pesanti (“**leptocurtica**”)



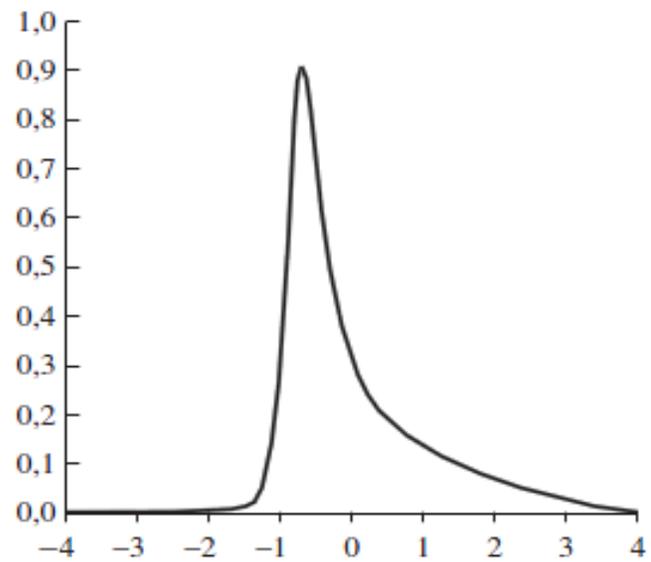
(a) Asimmetria = 0, Curtosi = 3



(b) Asimmetria = 0, Curtosi = 20



(c) Asimmetria = -0,1, Curtosi = 5



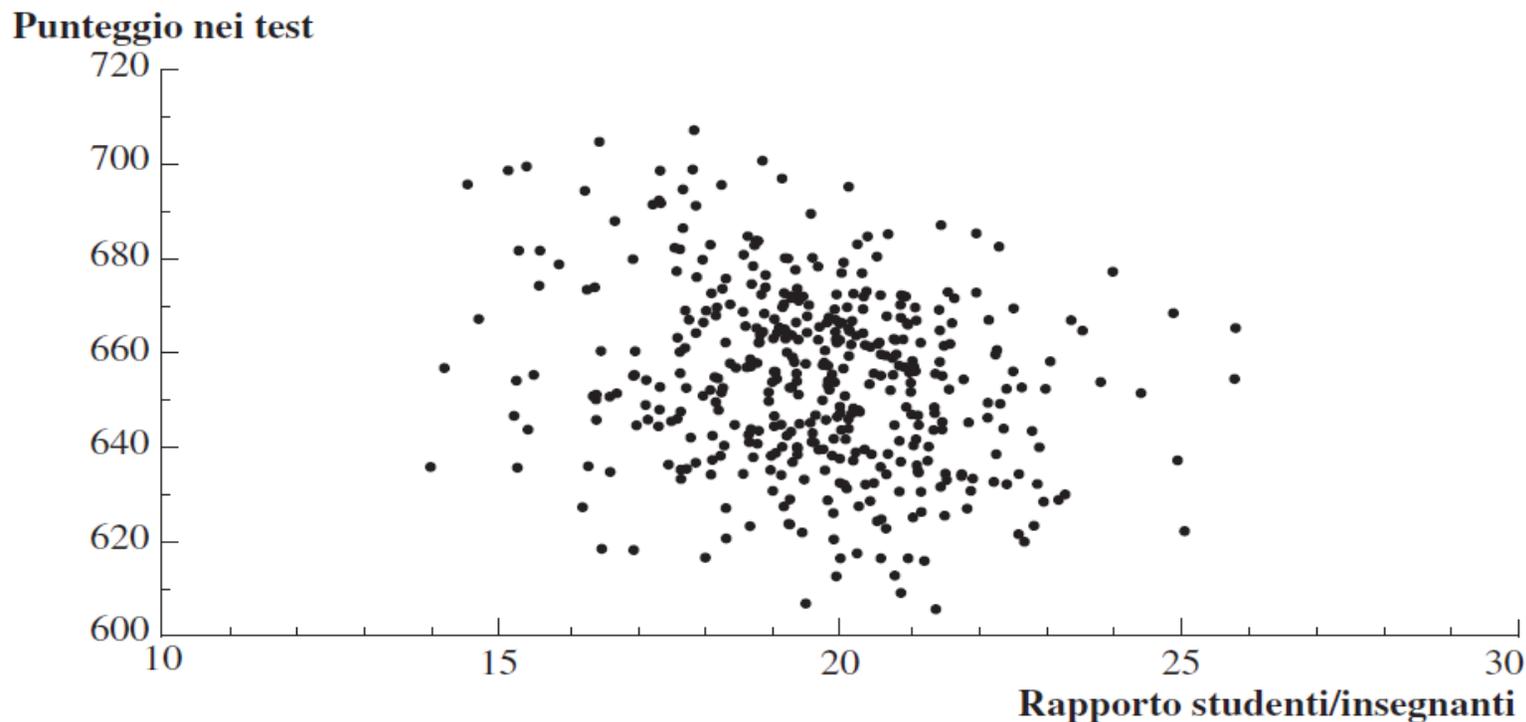
(d) Asimmetria = 0,6, Curtosi = 5

## 2 variabili casuali: distribuzioni congiunte e covarianza

- Le variabili casuali  $X$  e  $Z$  hanno una **distribuzione congiunta**
- La **covarianza** tra  $X$  e  $Z$  è
$$\text{cov}(X,Z) = E[(X - \mu_X)(Z - \mu_Z)] = \sigma_{XZ}$$
- La covarianza è una misura dell'associazione lineare tra  $X$  e  $Z$ ; le sue unità sono unità di  $X \times$  unità di  $Z$
- $\text{cov}(X,Z) > 0$  significa una relazione positiva tra  $X$  e  $Z$
- Se  $X$  e  $Z$  sono indipendentemente distribuite, allora  $\text{cov}(X,Z) = 0$  (ma non vale il vice versa!!)
- La covarianza di una variabile casuale con se stessa è la sua varianza:

$$\text{cov}(X,X) = E[(X - \mu_X)(X - \mu_X)] = E[(X - \mu_X)^2] = \sigma_X^2$$

La covarianza tra *Punteggio nei test* e *Rapporto studenti/insegnanti* è negativa:



**Figura 4.2**

**Diagramma a nuvola del punteggio nei test e del rapporto studenti/insegnanti (dati relativi ai distretti scolastici della California).**

Dati per i 420 distretti scolastici della California. C'è una debole relazione negativa tra il rapporto studenti/insegnanti e il punteggio nei test: la correlazione campionaria è pari a  $-0,23$ .

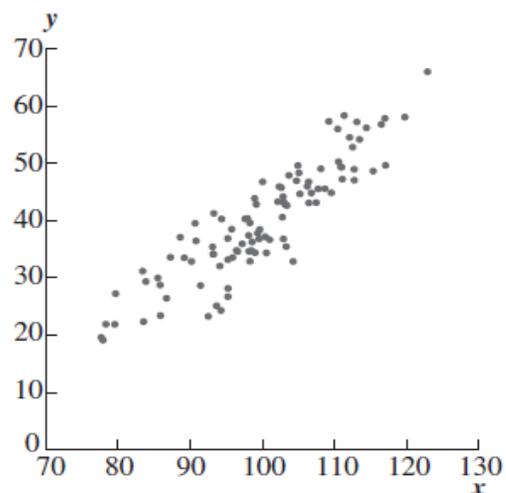
E così la ***correlazione***...

Il **coefficiente di correlazione** è definito in termini di covarianza:

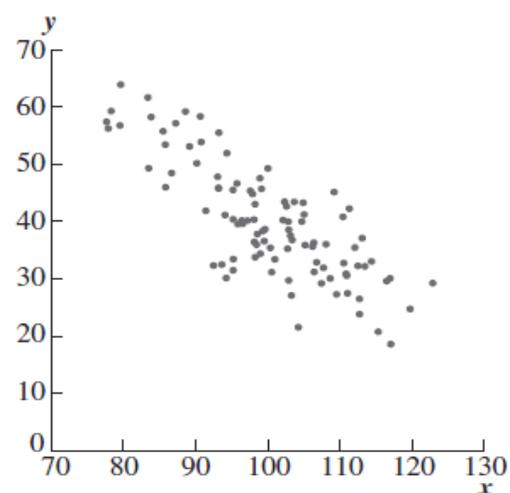
$$\text{corr}(X,Z) = \frac{\text{cov}(X,Z)}{\sqrt{\text{var}(X)\text{var}(Z)}} = \frac{\sigma_{XZ}}{\sigma_X \sigma_Z} = r_{XZ}$$

- $-1 \leq \text{corr}(X,Z) \leq 1$
- $\text{corr}(X,Z) = 1$  significa associazione lineare positiva perfetta
- $\text{corr}(X,Z) = -1$  significa associazione lineare negativa perfetta
- $\text{corr}(X,Z) = 0$  significa che non c'è associazione lineare

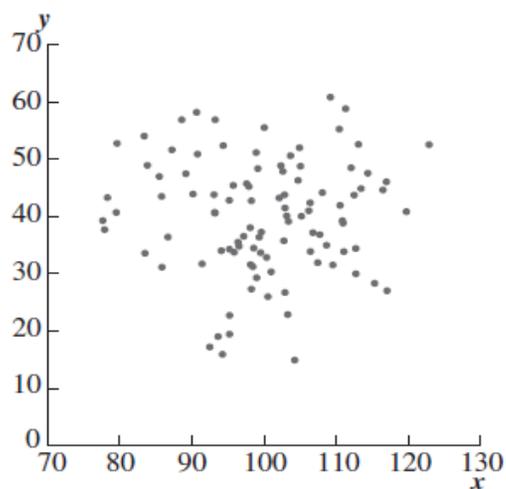
# ***Il coefficiente di correlazione misura l'associazione lineare***



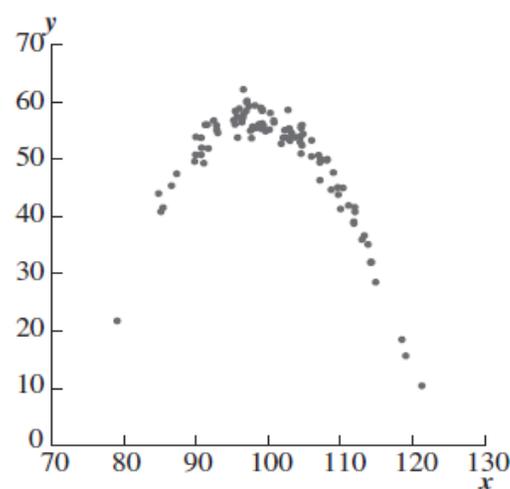
(a) Correlazione = +0,9



(b) Correlazione = -0,8



(c) Correlazione = 0,0



(d) Correlazione = 0,0 (quadratica)

## (c) Distribuzione condizionata e media condizionata

### ***Distribuzione condizionata***

- La distribuzione di  $Y$  dato il valore (o i valori) di un'altra variabile casuale  $X$
- Es: la distribuzione dei punteggi nei test dato  $STR < 20$

### ***Valore atteso condizionato e momento condizionato***

- *media condizionata* = media della distribuzione condizionata  
=  $E(Y|X = x)$  (**concetto e notazione importanti**)
- *varianza condizionata* = varianza della distribuzione condizionata
- *Esempio*:  $E(\text{Punteggio test} | STR < 20)$  = media dei punteggi nei test tra i distretti con dimensioni delle classi piccole

***La differenza in media è la differenza tra le medie di due distribuzioni condizionate:***

# ***Media condizionata (continua)***

$$\Delta = E(\text{Punteggio test} | STR < 20) - E(\text{Punteggio test} | STR \geq 20)$$

Altri esempi di media condizionata:

- Salari medi delle lavoratrici femmine ( $Y = \text{salari}$ ,  $X = \text{genere con femmina se } X=1$ )...  $E[Y|X=1]$
- Probabilità di morte di pazienti che ricevono una cura sperimentale ( $Y = \text{morto/vivo } (1/0)$ ;  $X = \text{trattato/non trattato}(1/0)$ )...  $E[Y|X=1] = \text{prob}(Y=1 | X=1)$
- Se  $E(X|Z=z) = \text{costante}$  al variare di  $z$ , allora  $\text{corr}(X,Z) = 0$  (tuttavia non vale necessariamente il vice versa)

***La media condizionata è un termine (forse nuovo) utilizzato per il concetto familiare di media di gruppo***

## **(d) Distribuzione di un campione di dati estratto a caso da una popolazione: $Y_1, \dots, Y_n$**

### ***Assumeremo un campionamento casuale semplice***

- Scegliere a caso un individuo (distretto, unità) dalla popolazione

### ***Casualità e dati***

- Prima della selezione del campione, il valore di  $Y$  è casuale perché l'individuo selezionato è casuale
- Una volta selezionato l'individuo e osservato il valore di  $Y$ ,  $Y$  è soltanto un numero – non casuale
- Il data set è  $(Y_1, Y_2, \dots, Y_n)$ , dove  $Y_i$  = valore di  $Y$  per l' $i$ -esimo individuo (distretto, unità) del campione

## ***Distribuzione di $Y_1, \dots, Y_n$ sotto campionamento casuale semplice***

- Poiché gli individui n. 1 e 2 sono selezionati a caso, il valore di  $Y_1$  non modifica la distribuzione di probabilità di  $Y_2$ . Quindi:
  - $Y_1$  e  $Y_2$  sono ***indipendentemente distribuiti***
  - $Y_1$  e  $Y_2$  provengono dalla stessa distribuzione, perchè estratti dalla stessa popolazione, cioè  $Y_1, Y_2$  sono ***identicamente distribuiti***
  - Ovvero, sotto campionamento casuale semplice,  $Y_1$  e  $Y_2$  sono indipendentemente e identicamente distribuiti (***i.i.d.***).
  - Più in generale, sotto campionamento casuale semplice,  $\{Y_i\}$ ,  $i = 1, \dots, n$ , sono i.i.d.

**Questo quadro consente rigorose inferenze statistiche sui momenti della distribuzione di  $Y$  utilizzando un campione di dati tratto dalla stessa popolazione ...**

1. Quadro probabilistico per inferenza statistica
2. **Stima**
3. Verifica
4. Intervalli di confidenza

## **Stima**

$\bar{Y}$  è lo stimatore naturale del valore atteso  $E(Y) = \mu_Y$ . Ma:

- a) quali sono le proprietà di  $\bar{Y}$  ?
- b) Perché dovremmo usare  $\bar{Y}$  anziché un altro stimatore?
  - $Y_1$  (prima osservazione)
  - forse pesi non uniformi – non media semplice
  - mediana( $Y_1, \dots, Y_n$ )

Il punto di partenza è la distribuzione campionaria di  $\bar{Y}$  ...

## (a) La distribuzione campionaria di $\bar{Y}$

$\bar{Y}$  è una variabile casuale e le sue proprietà sono determinate dalla **distribuzione campionaria** di  $\bar{Y}$

- Gli individui nel campione sono estratti a caso.
- Quindi i valori di  $(Y_1, \dots, Y_n)$  sono casuali
- Quindi funzioni di  $(Y_1, \dots, Y_n)$ , come  $\bar{Y}$ , sono casuali: se si fosse estratto un campione diverso, esso avrebbe assunto valori differenti
- La distribuzione di  $\bar{Y}$  su diversi possibili campioni di dimensione  $n$  si chiama **distribuzione campionaria** di  $\bar{Y}$ .
- La media e la varianza di  $\bar{Y}$  sono la media e la varianza della sua distribuzione campionaria,  $E(\bar{Y})$  e  $\text{var}(\bar{Y})$ .
- Il concetto di distribuzione campionaria è alla base di tutta l'econometria.

# La distribuzione campionaria di $\bar{Y}$ (continua)

**Esempio:** Si supponga che  $Y$  assuma il valore 0 o 1 (variabile casuale di **Bernoulli**) con la distribuzione di probabilità

$$\Pr[Y = 0] = 0,22, \Pr(Y = 1) = 0,78$$

Allora

$$E(Y) = p \times 1 + (1 - p) \times 0 = p = 0,78$$

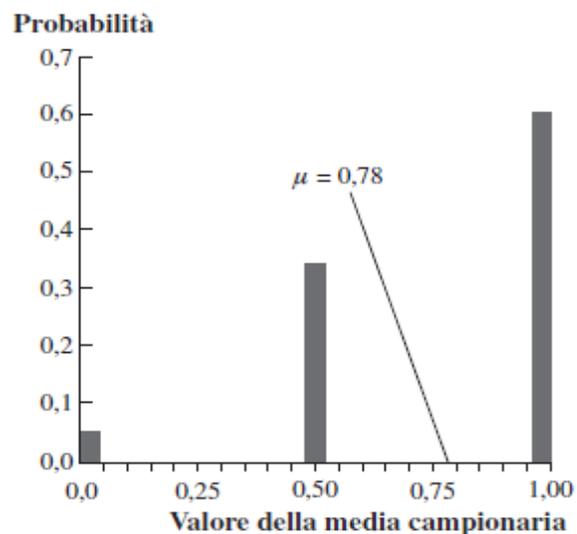
$$\begin{aligned}\sigma_Y^2 &= E[Y - E(Y)]^2 = p(1 - p) \text{ [ricordate?]} \\ &= 0,78 \times (1 - 0,78) = 0,1716\end{aligned}$$

La distribuzione campionaria di  $\bar{Y}$  dipende da  $n$ .

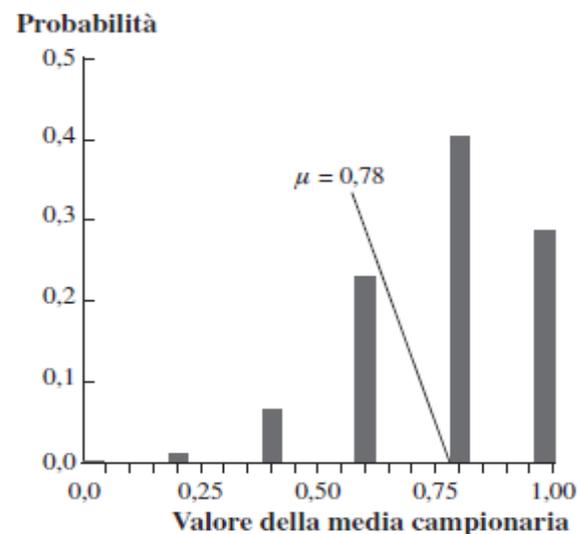
Si consideri  $n = 2$ . La distribuzione campionaria di  $\bar{Y}$  è

- $\Pr(\bar{Y} = 0) = 0,22^2 = 0,0484$
- $\Pr(\bar{Y} = \frac{1}{2}) = 2 \times 0,22 \times 0,78 = 0,3432$
- $\Pr(\bar{Y} = 1) = 0,78^2 = 0,6084$

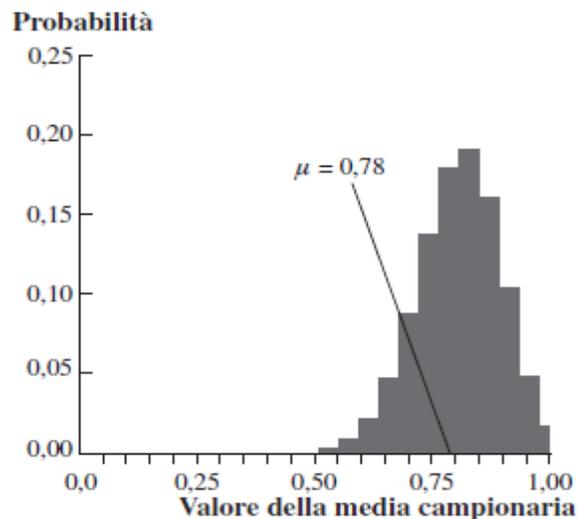
Distribuzione campionaria di  $\bar{Y}$  quando  $Y$  è di Bernoulli ( $p = 0,78$ ):



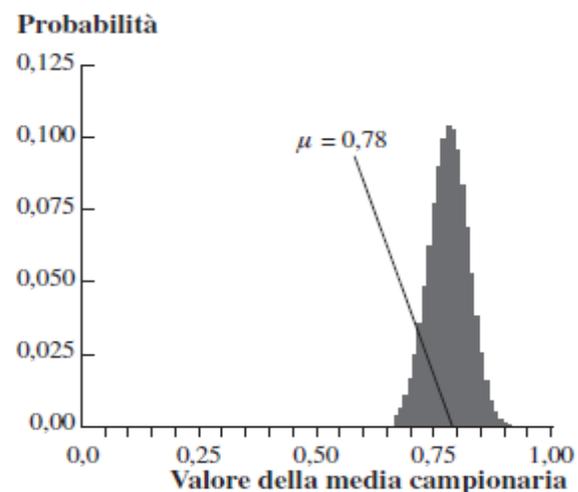
(a)  $n = 2$



(b)  $n = 5$



(c)  $n = 25$



(d)  $n = 100$

# Cose che vogliamo sapere sulla distribuzione campionaria:

- Qual è la media di  $\bar{Y}$ ?
  - Se  $E(\bar{Y}) = \mu = 0,78$ , allora  $\bar{Y}$  è uno stimatore **non distorto** di  $\mu$
- Qual è la varianza di  $\bar{Y}$ ?
  - In che modo  $\text{var}(\bar{Y})$  dipende da  $n$  (famosa formula  $1/n$ )
- Si avvicina a  $\mu$  quando  $n$  è grande?
  - Legge dei grandi numeri:  $\bar{Y}$  è uno stimatore **consistente** di  $\mu$
- $\bar{Y} - \mu$  assume forma a campana per  $n$  grande... questo è vero in generale?
  - In effetti,  $\bar{Y} - \mu$  è approssimato da una distribuzione normale per  $n$  grande (teorema limite centrale)

# Media e varianza della distribuzione campionaria di $\bar{Y}$

- Caso generale – cioè, per  $Y_i$  i.i.d. da qualsiasi distribuzione, non solo di Bernoulli:

- media:  $E(\bar{Y}) = E\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n} \sum_{i=1}^n E(Y_i) = \frac{1}{n} \sum_{i=1}^n \mu_Y = \mu_Y$

- Varianza:  $\text{var}(\bar{Y}) = E[\bar{Y} - E(\bar{Y})]^2$   
 $= E[\bar{Y} - \mu_Y]^2$

$$= E \left[ \left( \frac{1}{n} \sum_{i=1}^n Y_i \right) - \mu_Y \right]^2$$
$$= E \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_Y) \right]^2$$

quindi

$$\begin{aligned}
 \text{var}(\bar{Y}) &= E \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_Y) \right]^2 \\
 &= E \left\{ \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_Y) \right] \times \left[ \frac{1}{n} \sum_{j=1}^n (Y_j - \mu_Y) \right] \right\} \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E \left[ (Y_i - \mu_Y)(Y_j - \mu_Y) \right] \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{cov}(Y_i, Y_j) \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sigma_Y^2 \\
 &= \frac{\sigma_Y^2}{n}
 \end{aligned}$$

# Media e varianza della distribuzione campionaria di $\bar{Y}$ (continua)

$$E(\bar{Y}) = \mu_Y$$

$$\text{var}(\bar{Y}) = \frac{\sigma_Y^2}{n}$$

*Implicazioni:*

1.  $\bar{Y}$  è uno stimatore non distorto di  $\mu_Y$  (cioè  $E(\bar{Y}) = \mu_Y$ )
2.  $\text{var}(\bar{Y})$  è inversamente proporzionale a  $n$ 
  1. la dispersione della distribuzione campionaria è proporzionale a  $1/\sqrt{n}$
  2. Quindi l'incertezza campionaria associata con  $\bar{Y}$  è proporzionale a  $1/\sqrt{n}$  (grandi campioni, meno incertezza, ma legge con radice quadrata)

# Distribuzione campionaria di $\bar{Y}$ quando $n$ è grande

Per piccoli campioni, la distribuzione di  $\bar{Y}$  è complicata, ma se  $n$  è grande, la distribuzione campionaria è semplice!

1. All'aumentare di  $n$ , la distribuzione di  $\bar{Y}$  diventa più strettamente centrata su  $\mu_Y$  (*legge dei grandi numeri*)
1. Inoltre, la distribuzione di  $\bar{Y} - \mu_Y$  diventa normale (*teorema limite centrale*)

# Legge dei grandi numeri:

Uno stimatore è **consistente** se la probabilità che ricada entro un intervallo del vero valore della popolazione tende a uno all'aumentare della dimensione del campione.

Se  $(Y_1, \dots, Y_n)$  sono i.i.d. e  $\sigma_Y^2 < \infty$ , allora  $\bar{Y}$  è uno stimatore consistente di  $\mu_Y$ , cioè

$$\Pr[|\bar{Y} - \mu_Y| < \mu] \xrightarrow{p} 1 \text{ per } n \rightarrow \infty$$

che si può scrivere  $\bar{Y} \xrightarrow{p} \mu_Y$

(“ $\bar{Y} \xrightarrow{p} \mu_Y$ ” significa “ $\bar{Y}$  converge in probabilità a  $\mu_Y$ ”).

(*matematica*: per  $n \rightarrow \infty$ ,  $\text{var}(\bar{Y}) = \frac{\sigma_Y^2}{n} \rightarrow 0$ , il che implica che  $\Pr[|\bar{Y} - \mu_Y| < \varepsilon] \rightarrow 1$ .)

## **Teorema limite centrale (TLC):**

Se  $(Y_1, \dots, Y_n)$  sono i.i.d. e  $0 < \sigma_Y^2 < \infty$ , allora quando  $n$  è grande la distribuzione di  $\bar{Y}$  è bene approssimata da una distribuzione normale.

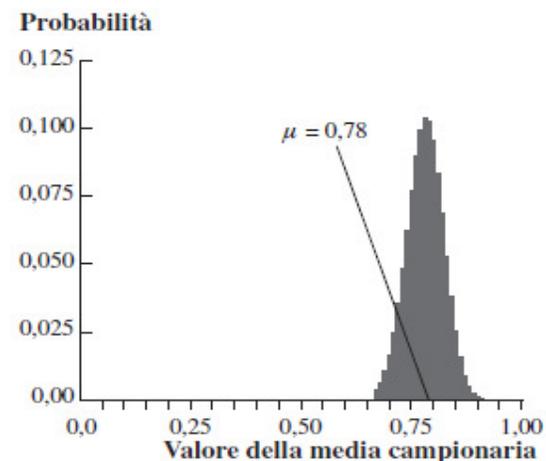
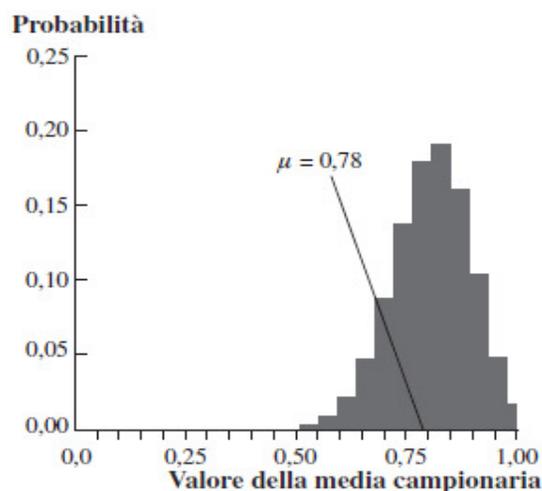
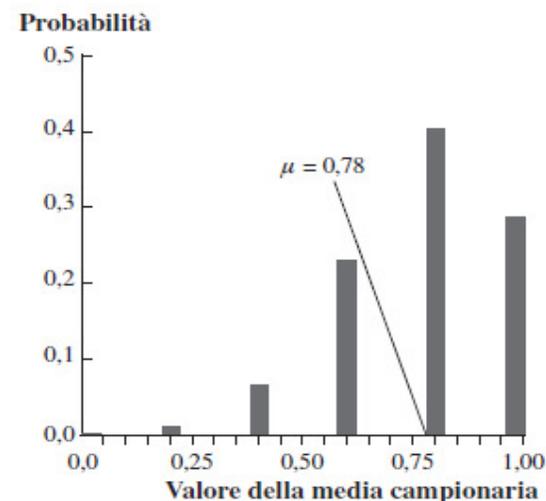
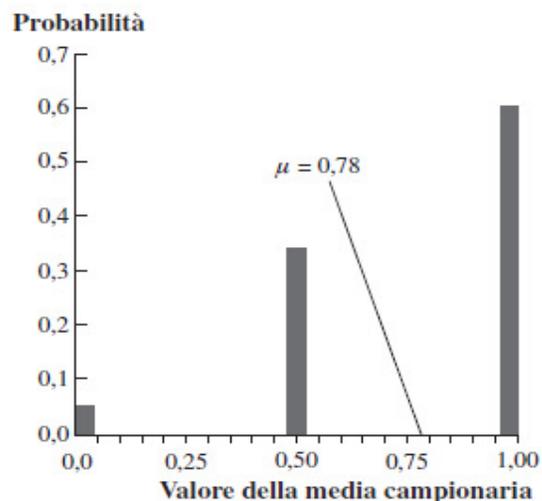
-  $\bar{Y}$  è distribuita approssimativamente come  $N(\mu_Y, \frac{\sigma_Y^2}{n})$   
 (“distribuzione normale con media  $\mu_Y$  e varianza  $\sigma_Y^2/n$ ”)

-  $\sqrt{n} (\bar{Y} - \mu_Y)/\sigma_Y$  è distribuita approssimativamente come  $N(0,1)$  (normale standard)

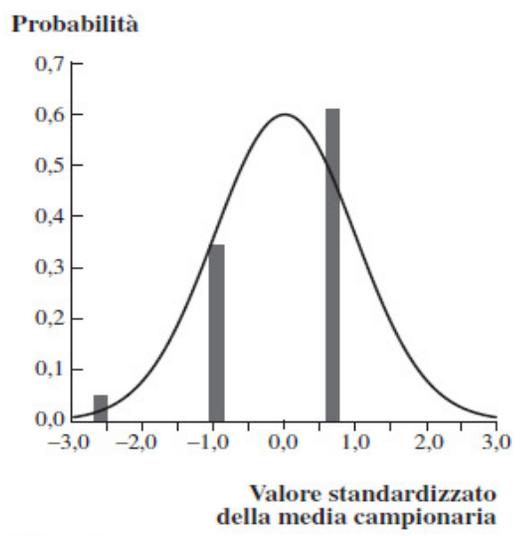
- **Cioè,  $\bar{Y}$  “standardizzata”**  $= \frac{\bar{Y} - E(\bar{Y})}{\sqrt{\text{var}(\bar{Y})}} = \frac{\bar{Y} - \mu_Y}{\sigma_Y / \sqrt{n}}$  **è distribuita approssimativamente come  $N(0,1)$**

- **Più grande è  $n$ , migliore è l’ approssimazione.**

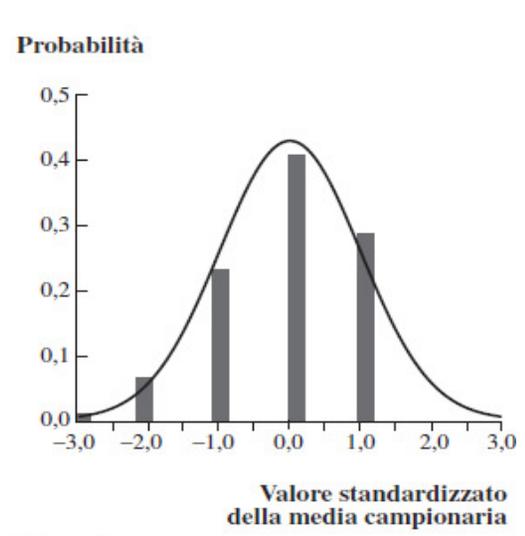
# Distribuzione campionaria di $\bar{Y}$ quando $Y$ è di Bernoulli, $p = 0,78$ :



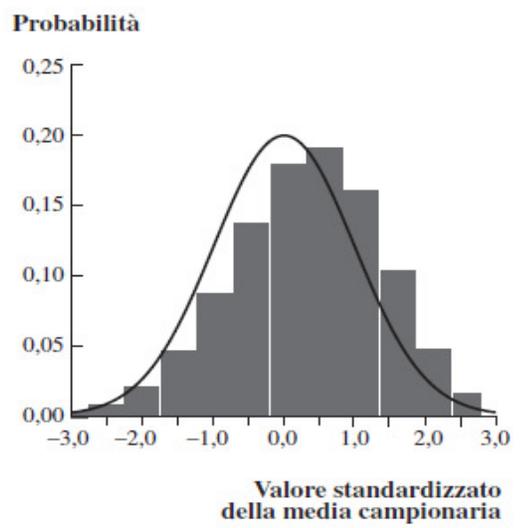
Stesso esempio: distribuzione campionaria di  $\frac{\bar{Y} - E(\bar{Y})}{\sqrt{\text{var}(\bar{Y})}}$  :



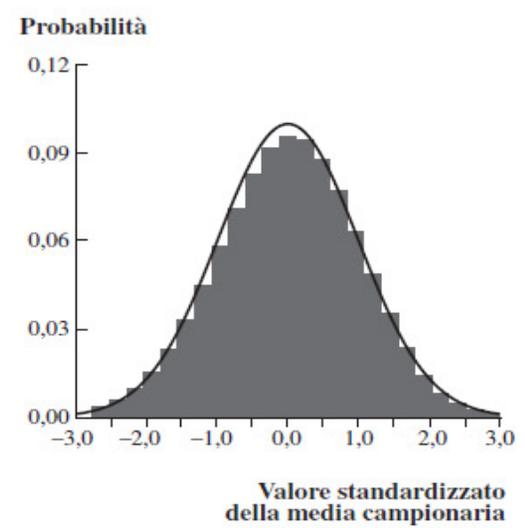
(a)  $n = 2$



(b)  $n = 5$



(c)  $n = 25$



(d)  $n = 100$

## Riepilogo: distribuzione campionaria di $\bar{Y}$

Per  $Y_1, \dots, Y_n$  i.i.d. con  $0 < \sigma_Y^2 < \infty$ ,

- La distribuzione campionaria esatta (campione finito) di  $\bar{Y}$  ha media  $\mu_Y$  (“ $\bar{Y}$  è uno stimatore non distorto di  $\mu_Y$ ”) e varianza  $\sigma_Y^2/n$
- Al di là di media e varianza, la distribuzione esatta di  $\bar{Y}$  è complessa e dipende dalla distribuzione di  $Y$  (la distribuzione della popolazione)
- Quando  $n$  è grande, la distribuzione campionaria si semplifica:

–  $\bar{Y} \xrightarrow{p} \mu_Y$  (Legge dei grandi numeri)

–  $\frac{\bar{Y} - E(\bar{Y})}{\sqrt{\text{var}(\bar{Y})}}$  è approssimata da  $N(0,1)$  (TLC)

## (b) Perché usare $\bar{Y}$ per stimare $\mu_Y$ ?

- $\bar{Y}$  è non distorto:  $E(\bar{Y}) = \mu_Y$
- $\bar{Y}$  è consistente:  $\bar{Y} \xrightarrow{P} \mu_Y$
- $\bar{Y}$  è lo stimatore “dei minimi quadrati” di  $\mu_Y$ ;  $\bar{Y}$  risolve

$$\min_m \sum_{i=1}^n (Y_i - m)^2$$

quindi  $\bar{Y}$  minimizza la somma dei quadrati dei “residui”: *derivazione facoltativa (cfr. anche Appendice 3.2)*

$$\frac{d}{dm} \sum_{i=1}^n (Y_i - m)^2 = \sum_{i=1}^n \frac{d}{dm} (Y_i - m)^2 = 2 \sum_{i=1}^n (Y_i - m)$$

Si pone la derivata a zero e si denota il valore ottimale di  $m$  con:  $\hat{m}$

$$\sum_{i=1}^n Y_i = \sum_{i=1}^n \hat{m} = n\hat{m} \quad \text{o} \quad \hat{m} = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}$$

# Perché usare $\bar{Y}$ per stimare $\mu_Y$ (*continua*)

- $\bar{Y}$  ha una varianza minore di tutti gli altri *stimatori lineari non distorti*: si consideri lo stimatore  $\hat{\mu}_Y = \frac{1}{n} \sum_{i=1}^n a_i Y_i$ , dove gli  $\{a_i\}$  sono tali che  $\hat{\mu}_Y$  è non distorto; allora  $\text{var}(\bar{Y}) \leq \text{var}(\hat{\mu}_Y)$

(dimostrazione: Capitolo 17 del volume stampato)

- $\bar{Y}$  non è l'unico stimatore di  $\mu_Y$  – vi viene in mente un caso in cui potrebbe essere preferibile utilizzare la mediana?
  1. Quadro di riferimento probabilistico per l'inferenza statistica
  2. Stima
  - 3. Verifica di ipotesi**
  4. Intervalli di confidenza

# Verifica di ipotesi

Il problema della **verifica di ipotesi** (per la media): prendere una decisione preliminare in base all'evidenza disponibile che un'ipotesi nulla è vera, o che è vera, invece, un'ipotesi alternativa. Cioè verificare

- $H_0: E(Y) = \mu_{Y,0}$  vs.  $H_1: E(Y) > \mu_{Y,0}$  (monodirezionale, >)
- $H_0: E(Y) = \mu_{Y,0}$  vs.  $H_1: E(Y) < \mu_{Y,0}$  (monodirezionale, <)
- $H_0: E(Y) = \mu_{Y,0}$  vs.  $H_1: E(Y) \neq \mu_{Y,0}$  (bidirezionale)

# Terminologia per la verifica di ipotesi statistiche

**valore-p** = probabilità di ricavare una statistica (per es.  $\bar{Y}$ ) sfavorevole all'ipotesi nulla almeno quanto il valore effettivamente calcolato con i dati, supponendo che l'ipotesi nulla sia corretta.

Il **livello di significatività** di un test è una probabilità predeterminata di rifiutare in modo errato l'ipotesi nulla, quando invece è corretta.

**Calcolo del valore-p** in base a  $\bar{Y}$  :

$$\text{valore-p} = \Pr_{H_0} [|\bar{Y} - \mu_{Y,0}| > |\bar{Y}^{act} - \mu_{Y,0}|]$$

Dove  $\bar{Y}^{act}$  è il valore di  $\bar{Y}$  effettivamente osservato (non casuale)

## Calcolo del valore- $p$ (continua)

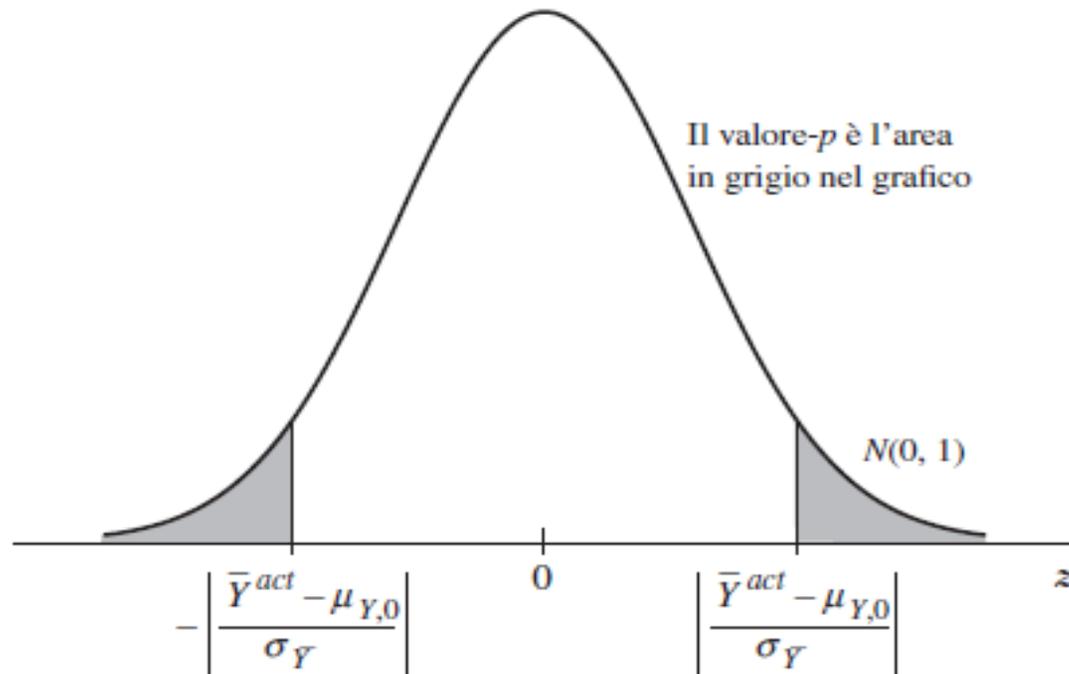
- Per calcolare il valore- $p$  è necessario conoscere la distribuzione campionaria di  $\bar{Y}$ , che è complessa se  $n$  è piccolo.
- Se  $n$  è grande, si può usare l' approssimazione normale (TLC):

$$\begin{aligned}\text{valore-}p &= \Pr_{H_0} \left[ \left| \frac{\bar{Y} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}} \right| \right] \\ &= \Pr_{H_0} \left[ \left| \bar{Y} - \mu_{Y,0} \right| > \left| \bar{Y}^{act} - \mu_{Y,0} \right| \right] \\ &= \Pr_{H_0} \left[ \left| \frac{\bar{Y} - \mu_{Y,0}}{\sigma_{\bar{Y}}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_{\bar{Y}}} \right| \right]\end{aligned}$$

$\cong$  probabilità sotto code  $N(0,1)$  sin+destra

dove  $\sigma_{\bar{Y}} = \text{dev. std della distribuzione di } \bar{Y} = \sigma_Y / \sqrt{n}$ .

## Calcolo del valore- $p$ con $\sigma_Y$ nota:



- Per  $n$  grande, valore- $p$  = probabilità che una variabile casuale  $N(0,1)$  ricada al di fuori dell'intervallo  $- \left| (\bar{Y}^{act} - \mu_{Y,0}) / \sigma_{\bar{Y}} \right|, + \left| (\bar{Y}^{act} - \mu_{Y,0}) / \sigma_{\bar{Y}} \right|$
- In pratica,  $\sigma_{\bar{Y}}$  è ignota – deve essere stimata

## ***Stimatore della varianza di Y:***

$$s_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \text{“varianza campionaria di Y”}$$

Se  $(Y_1, \dots, Y_n)$  sono i.i.d. e  $E(Y^4) < \infty$ , allora

$$s_Y^2 \xrightarrow{p} \sigma_Y^2$$

Perché si applica la legge dei grandi numeri?

- Perché  $s_Y^2$  è una media campionaria; cfr. Appendice 3.3
- Nota tecnica: si assume  $E(Y^4) < \infty$  perché la media non è di  $Y_i$ , ma del suo quadrato; cfr. Appendice 3.3.

## Calcolo del valore-p con $\sigma_Y^2$ stimato

$$\begin{aligned}\text{valore-p} &= \Pr_{H_0} \left[ \left| \frac{\bar{Y} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_Y / \sqrt{n}} \right| \right] \\ &= \Pr_{H_0} \left[ \left| \bar{Y} - \mu_{Y,0} \right| > \left| \bar{Y}^{act} - \mu_{Y,0} \right| \right] \\ &\cong \Pr_{H_0} \left[ \left| \frac{\bar{Y} - \mu_{Y,0}}{s_Y / \sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,0}}{s_Y / \sqrt{n}} \right| \right] \quad (n \text{ grande})\end{aligned}$$

quindi

$$\text{valore-p} = \Pr_{H_0} \left[ \left| t \right| > \left| t^{act} \right| \right] \quad \left( \sigma_Y^2 \text{ stimato} \right)$$

$\cong$  probabilità sotto code normali al di fuori di  $-|t^{act}|$ ,  $+|t^{act}|$

$$\text{dove } t = \frac{\bar{Y} - \mu_{Y,0}}{s_Y / \sqrt{n}} \quad (\text{la consueta statistica } t)$$

# Che collegamento c'è tra il valore- $p$ e il livello di significatività?

- Il livello di significatività è specificato in anticipo. Per esempio, se tale livello è del 5%,
  - si rifiuta l'ipotesi nulla se  $|t| \geq 1,96$ .
  - in modo equivalente, la si rifiuta se  $p \leq 0,05$ .
  - Il valore- $p$  è detto talvolta **livello di significatività marginale**.
  - Spesso è meglio comunicare il valore- $p$  che limitarsi a indicare se un test rifiuta o no – il valore- $p$  contiene più informazioni di un semplice risultato “sì/no” in riferimento a un test.

A questo punto potreste chiedervi...

**Che ne è della tabella- $t$  e dei gradi di libertà?**

## **Digressione: la distribuzione $t$ di Student**

Se  $Y_i, i = 1, \dots, n$  sono i.i.d.  $N(\mu_Y, \sigma_Y^2)$ , allora la statistica  $t$  ha la distribuzione  $t$  di Student *con*  $n - 1$  gradi di libertà.

I valori critici della distribuzione  $t$  di Student sono elencati in tutti i libri di statistica. Ricordate la procedura?

1. Calcolare la statistica  $t$
2. Calcolare i gradi di libertà,  $n - 1$
3. Cercare il valore critico al 5%
4. Se la statistica  $t$  supera (in valore assoluto) questo valore critico, rifiutare l'ipotesi nulla.

# Commenti su questa procedura e sulla distribuzione $t$ di Student

1. La teoria della distribuzione  $t$  è stata uno dei primi trionfi della statistica matematica. È davvero sorprendente: se  $Y$  è i.i.d. e normale, allora è possibile conoscere la distribuzione *esatta*, a campione finito della statistica  $t$  – è la  $t$  di Student. Perciò si possono costruire intervalli di confidenza (usando il valore critico  $t$  di Student) che hanno *esattamente* lo stesso tasso di copertura, indipendentemente dalla dimensione del campione. Questo risultato è stato molto utile in tempi in cui “calcolatore” era una posizione lavorativa, la raccolta di dati era costosa e il numero di osservazioni si aggirava attorno alla decina. È anche un risultato concettualmente splendido, e anche la matematica è molto elegante – il che probabilmente spiega perché i docenti amano insegnare la distribuzione  $t$ . Tuttavia...

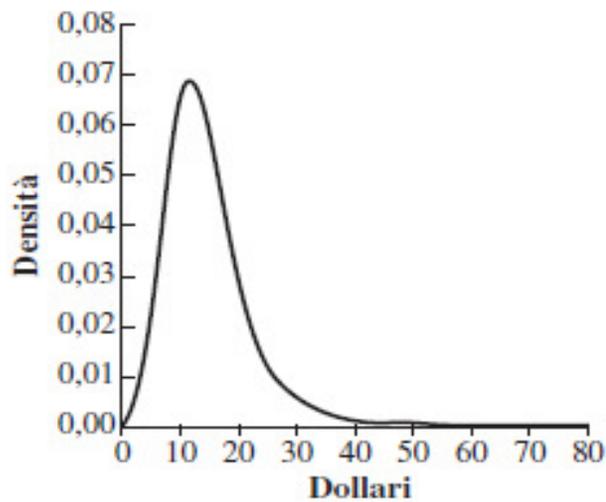
# Commenti sulla distribuzione $t$ di Student (*continua*)

2. Se la dimensione del campione è moderata (varie dozzine) o grande (centinaia o più), la differenza tra la distribuzione  $t$  e i valori critici  $N(0,1)$  è trascurabile. Riportiamo di seguito alcuni valori critici al 5% per test bidirezionali:

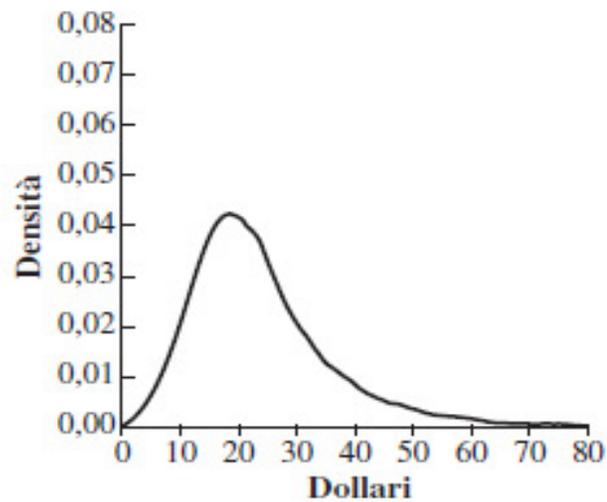
| gradi di libertà<br>( $n - 1$ ) | valore critico distribuzione<br>$t$ al 5% |
|---------------------------------|---|
| 10                              | 2,23                                      |
| 20                              | 2,09                                      |
| 30                              | 2,04                                      |
| 60                              | 2,00                                      |
| $\infty$                        | 1,96                                      |

# Commenti sulla distribuzione $t$ di Student (*continua*)

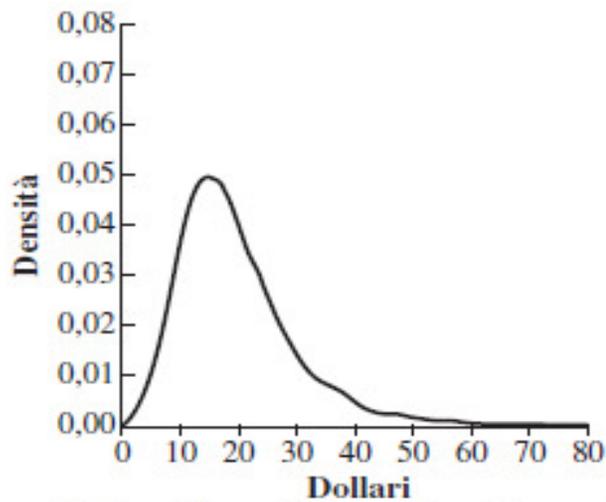
3. Perciò la distribuzione  $t$  di Student è di interesse soltanto quando la dimensione del campione è molto piccola; ma in quel caso, affinché sia corretta, è necessario assicurarsi che la distribuzione di  $Y$  sia normale. Per dati economici, l'assunzione di normalità è raramente credibile. Ecco le distribuzioni di alcuni dati economici.
- Pensate che i guadagni abbiano distribuzione normale?
  - Supponete di avere un campione di  $n = 10$  osservazioni da una di queste distribuzioni – vi sentireste di usare la distribuzione  $t$  di Student?



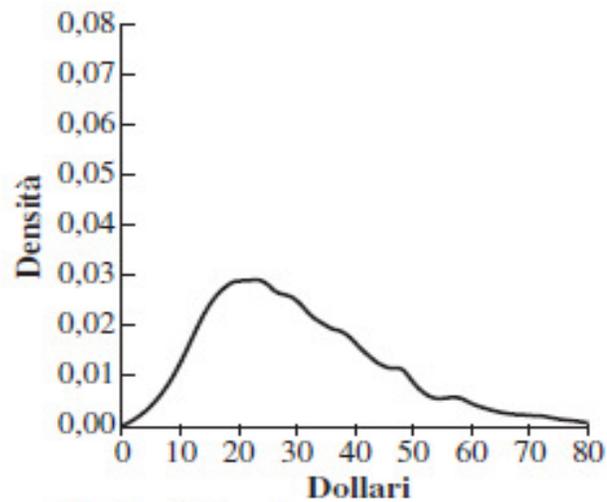
(a) Donne con diploma di scuola superiore



(b) Donne con laurea



(c) Uomini con diploma di scuola superiore



(d) Uomini con laurea

**Figura 2.4**

**Distribuzione condizionata delle retribuzioni orarie medie dei lavoratori statunitensi a tempo pieno nel 2004, dati il livello d'istruzione e il sesso.**

Le quattro distribuzioni delle retribuzioni sono per uomini e donne, per coloro che hanno solo un diploma di scuola superiore (a e c) e coloro che hanno una laurea (b e d).

# Commenti sulla distribuzione $t$ di Student (*continua*)

4. Forse non lo sapete. Considerate la statistica  $t$  che verifica l'ipotesi che due medie (gruppi  $s, l$ ) siano uguali:

$$t = \frac{\bar{Y}_s - \bar{Y}_l}{\sqrt{\frac{s_s^2}{n_s} + \frac{s_l^2}{n_l}}} = \frac{\bar{Y}_s - \bar{Y}_l}{SE(\bar{Y}_s - \bar{Y}_l)}$$

Anche se la distribuzione di  $Y$  nei due gruppi è normale, questa statistica non ha una distribuzione  $t$  di Student!

Esiste una statistica che verifica questa ipotesi e ha distribuzione normale, la statistica  $t$  “a varianza aggregata” – cfr. il Paragrafo 3.6 del volume stampato – tuttavia essa è valida soltanto se le varianze delle distribuzioni normali sono le stesse nei due gruppi. Pensate che questo sia vero, per esempio, per i salari di uomini vs donne?

# La distribuzione $t$ di Student – Riepilogo

- L'ipotesi che  $Y$  abbia distribuzione  $N(\mu_Y, \sigma_Y^2)$  è raramente plausibile nella pratica (reddito? numero di figli?)
- per  $n > 30$ , la distribuzione  $t$  e  $N(0,1)$  sono molto vicine (al crescere di  $n$ , la distribuzione  $t_{n-1}$  converge a  $N(0,1)$ )
- La distribuzione  $t$  è un artefatto che risale ai tempi in cui le dimensioni dei campioni erano piccole e i “calcolatori” erano persone
- Per motivi storici, il software statistico utilizza generalmente la distribuzione  $t$  per calcolare valori- $p$  ma questo è irrilevante quando la dimensione del campione è moderata o grande.
- Per questi motivi, in questo corso ci concentreremo sull'approssimazione con  $n$  grande data dal TLC
  1. Quadro probabilistico per l'inferenza statistica
  2. Stima
  3. Verifica
  4. **Intervalli di confidenza**

# Intervalli di confidenza

- Un ***intervallo di confidenza al 95%*** per  $\mu_Y$  è un intervallo che contiene il valore vero di  $\mu_Y$  nel 95% dei campioni ripetuti.
- *Digressione*: qual è la casualità qui? I valori di  $Y_1, \dots, Y_n$  e quindi qualsiasi funzione degli stessi – incluso l'intervallo di confidenza, che differirà da un campione all'altro. Il parametro della popolazione,  $\mu_Y$ , non è casuale; semplicemente, non lo conosciamo.

## ***Intervalli di confidenza (continua)***

Un intervallo di confidenza al 95% può sempre essere costruito come insieme di valori dei  $\mu_Y$  non rifiutati da un test di ipotesi con un livello di significatività del 5%.

$$\begin{aligned} \{ \mu_Y : \left| \frac{\bar{Y} - \mu_Y}{s_Y / \sqrt{n}} \right| \leq 1,96 \} &= \{ \mu_Y : -1,96 \leq \frac{\bar{Y} - \mu_Y}{s_Y / \sqrt{n}} \leq 1,96 \} \\ &= \{ \mu_Y : -1,96 \leq -\mu_Y \leq 1,96 \frac{s_Y}{\sqrt{n}} \} \\ &= \{ \mu_Y \in ( \bar{Y} - 1,96 \frac{s_Y}{\sqrt{n}} , \bar{Y} + 1,96 \frac{s_Y}{\sqrt{n}} ) \} \end{aligned}$$

*Questo intervallo di confidenza si basa sugli  $n$ -grande risultati che  $\bar{Y}$  è approssimata da una distribuzione normale e  $s_Y^2 \xrightarrow{p} \sigma_Y^2$*

# Riepilogo:

Dalle due ipotesi di:

1. campionamento casuale semplice di una popolazione, cioè  $\{Y_i, i = 1, \dots, n\}$  sono i.i.d.
2.  $0 < E(Y^4) < \infty$

abbiamo sviluppato, per grandi campioni ( $n$  grande):

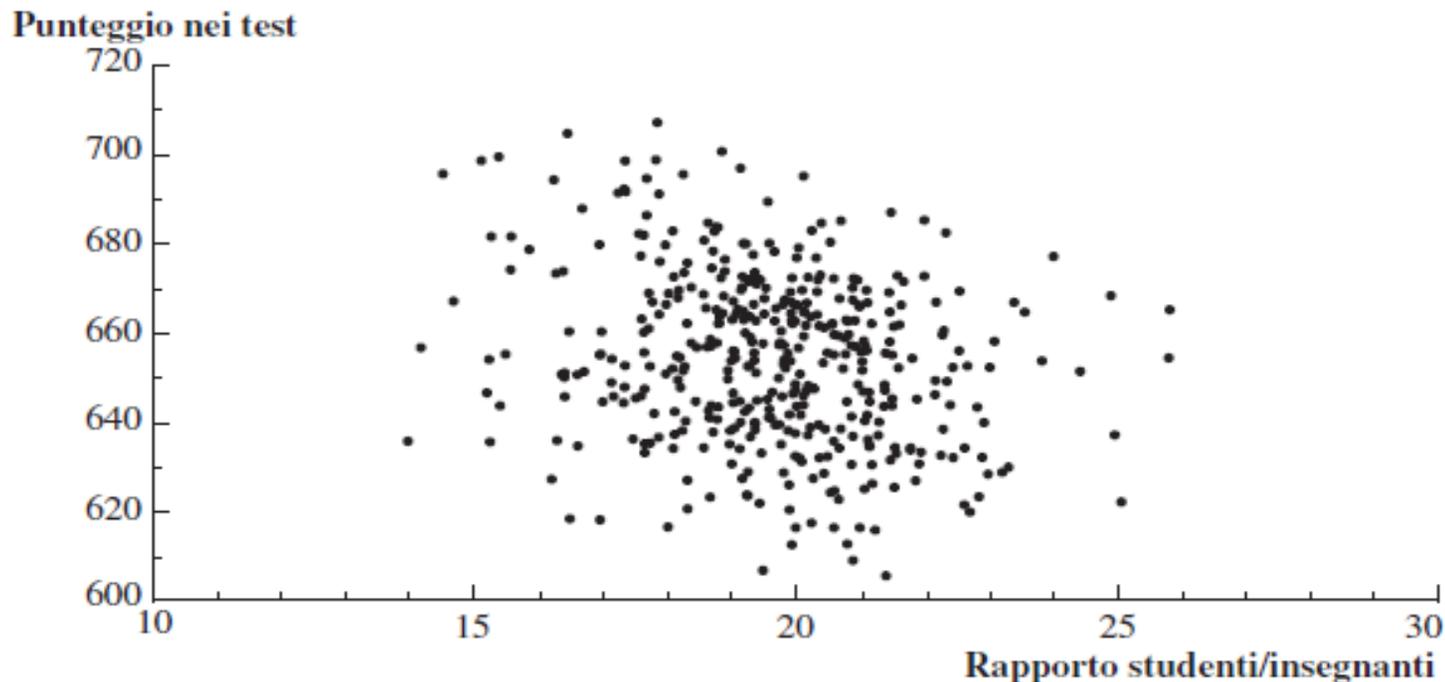
- Teoria della stima (distribuzione campionaria di  $\bar{Y}$ )
- Teoria della verifica di ipotesi (distribuzione con  $n$  grande della statistica  $t$  e calcolo del valore- $p$ )
- Teoria degli intervalli di confidenza (costruita invertendo la statistica test)

Le ipotesi (1) e (2) sono plausibili nella pratica? **Sì**

# Torniamo alla domanda politica di partenza:

Qual è l'effetto sui punteggi nei test della riduzione della dimensione delle classi di uno studente per classe?

*Abbiamo risposto a questa domanda?*



**Figura 4.2**

**Diagramma a nuvola del punteggio nei test e del rapporto studenti/insegnanti (dati relativi ai distretti scolastici della California).**

Dati per i 420 distretti scolastici della California. C'è una debole relazione negativa tra il rapporto studenti/insegnanti e il punteggio nei test: la correlazione campionaria è pari a  $-0,23$ .