

Esempiocap. investite
 L_i prezzo titoli
 P_i quantità
 q_i

Obiettivo: determinare prezzo medio che lascia invariata la disponibilità monetaria globale.

$$\Rightarrow \bar{p} = ?$$

$$\sum_{i=1}^r L_i = \sum_{i=1}^r P_i \cdot q_i \Rightarrow \sum_{i=1}^r P_i \cdot q_i = \sum_{i=1}^r \bar{p} \cdot q_i \Rightarrow \bar{p} = \frac{\sum_{i=1}^r P_i \cdot q_i}{\sum_{i=1}^r q_i}$$

q_i è il sistema di pesi

NB: \bar{p} non è una media in senso stretto.

risolvo per \bar{p}

Esempio

Si considerino n unità statistiche variabili. D = disp. finanziaria;
 S = spese fisse t.c. $D \geq S$. X = propensione all'investimento $\in (0, 1)$.
 k = ammontare investimenti di ciascun soggetto.

Relazione tra le variabili: $k_i = (d_i - s_i) \cdot n_i$

Si determini la propensione media all'investimento che lascia invariato il tot. spese fisse.

$$\bar{n} = ? \quad \sum_{i=1}^r s_i \stackrel{\text{tot spese}}{=} \text{costante} = \varphi(\cdot)$$

Ma devo vedere la x nella funzione invariante φ

$$\Rightarrow k_i = (d_i - s_i) n_i \Rightarrow s_i = d_i - \frac{k_i}{n_i} \Rightarrow \sum_{i=1}^r s_i = \sum_{i=1}^r (d_i - \frac{k_i}{n_i})$$

$$\Rightarrow \sum_{i=1}^r (d_i - \frac{k_i}{n_i}) = \sum_{i=1}^r (d_i - \frac{k_i}{\bar{n}}) \Rightarrow \sum_{i=1}^r d_i - \sum_{i=1}^r \frac{k_i}{\bar{n}} = \sum_{i=1}^r d_i - \sum_{i=1}^r \frac{k_i}{\bar{n}}$$

$$\Rightarrow \frac{1}{\bar{n}} \sum_{i=1}^r k_i - \sum_{i=1}^r \frac{k_i}{\bar{n}} = 0 \Rightarrow \frac{1}{\bar{n}} \sum_{i=1}^r k_i = \sum_{i=1}^r \frac{k_i}{\bar{n}} \Rightarrow \frac{1}{\bar{n}} = \frac{\sum_{i=1}^r \frac{k_i}{n_i}}{\sum_{i=1}^r k_i}$$

$$\Rightarrow \bar{n} = \frac{\sum_{i=1}^r k_i}{\sum_{i=1}^r \frac{k_i}{n_i}}$$

\rightarrow non è una media in senso stretto, tuttavia è riconoscibile come media armonica con sistema di pesi k_i .

Lezione 9 - 21/10/21

② Metodo della minimizzazione del danno

Danno = perdita di informazione

$X = n_1, n_2, \dots, n_n$: osservazioni

Sintesi : $\alpha = \alpha(X) = \alpha(n_1, n_2, \dots, n_r; n_1, n_2, \dots, n_r)$

2 α sintetizza la distribuzione di frequenza di X

Obiettivo : minimizzare il danno globale

ERRORE = $e_i = n_i - \alpha$ dove α è una media, ma bisogna definire quale.

es. Due studenti : A e B ; entrambi media voti = 8

ma : - A studia tutti i giorni e prende sempre (e solo) 8

- B studia poco : prende 6 ma poi recupera con un 10

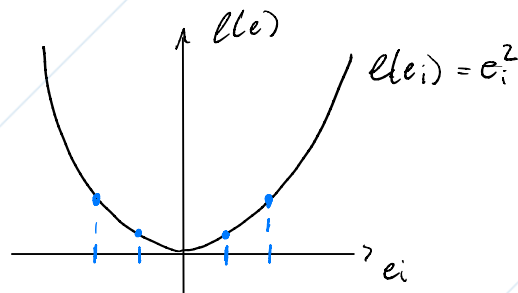
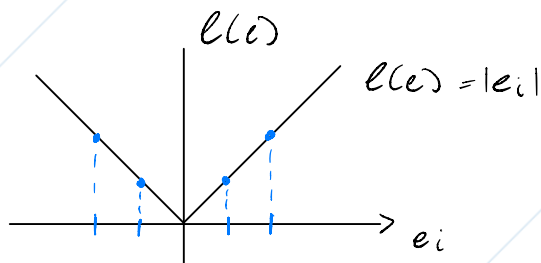
Media = 8 rappresenta (sintetizza) A con nessuna perdita di informazione, cosa che avviene invece con B.

Per quantificare l'errore globale : somma di tutti gli e_i

Ma gli errori e_i possono compensarsi

\Rightarrow annulliamo la compensazione applicando : - potenze pari
- valore assoluto

$l(e)$ = funzione l degli errori e_i = ?



$l(e)$ deve essere :

- simmetrica
- continua
- minima (o meglio ancora nulla) nell'origine
- monotona decrescente per $e_i < 0$
- monotona crescente per $e_i > 0$

$$l(e) = |e|^s \quad s = 1, 2, \dots \quad (\text{vanno già bene } s=1 \text{ e } s=2)$$

Danno = Perdita globale = D

$$D = \sum_{i=1}^r l(e_i) n_i \quad \text{dove } e_i = n_i - \alpha$$

α è una media, ancora non precisata
 cerchiamo α che minimizzi D
 funzione applicata agli errori

$$\Rightarrow D = \sum_{i=1}^r |n_i - \alpha|^s \cdot n_i \quad \text{NB: } s=0 \text{ può servire per studiare D al limite } s \rightarrow 0$$

$$s=0 \quad D = \sum_{i=1}^r |n_i - \alpha|^0 \cdot n_i \quad (\text{messa così è inutile, } \bar{e} = n)$$

$$l(e) = l(n_i - \alpha) \begin{cases} = 0 & \Leftrightarrow n_i = \alpha \\ \text{costante} & \Leftrightarrow n_i \neq \alpha \end{cases}$$

Passaggio al limite:

caso ① : $n_i \neq \alpha$

$$\Rightarrow |n_i - \alpha| > 0 \quad \forall i \quad \lim_{s \rightarrow 0} |n_i - \alpha|^s \rightarrow 1$$

$$\Rightarrow D = \lim_{s \rightarrow 0} \sum_{i=1}^r |n_i - \alpha|^s \cdot n_i = \sum_{i=1}^r n_i = n$$

caso ② : $n_j = \alpha$, $n_i \neq \alpha \quad i=1, 2, \dots, j-1, j+1, \dots, r$

$$\Rightarrow \begin{cases} |n_i - \alpha| > 0 & i \neq j \\ |n_i - \alpha| = 0 & i = j \end{cases}$$

Quindi $\sum_{i=1}^r |n_i - \alpha|^s \cdot n_i = \sum_{i \neq j} |n_i - \alpha|^s n_i + 0$

$$D = \sum_{i=j} \lim_{s \rightarrow 0} |n_i - \alpha|^s \cdot n_i = n_1 + n_2 + \dots + n_{j-1} + n_{j+1} + \dots + n_r = n - n_j < n$$

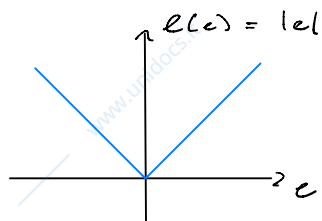
$$\Rightarrow D = n - n_j < n \quad \leftarrow \text{vogliamo minimizzarla}$$

Per minimizzare D bisogna massimizzare n_j

$$\Rightarrow \alpha = n_i = \text{MODA}$$

$S = 1$

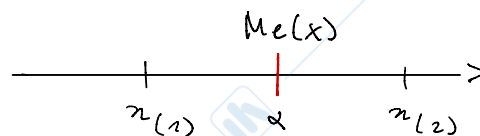
$l(c) = |c|$



$D = \sum_{i=1}^k |x_i - \alpha| n_i$

$D_{min} \Rightarrow \alpha = \text{MEDIANA}$

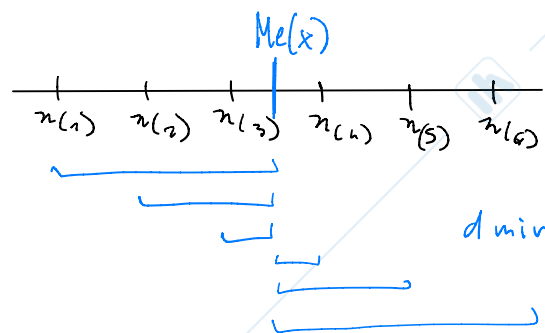
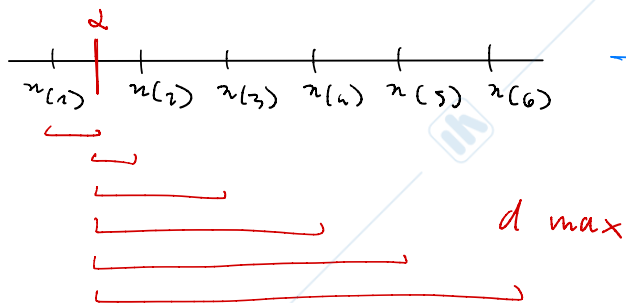
Ipotezziamo di avere solo 2 n_i ($k=2$):



$\Rightarrow n=2 \quad \left\{ \begin{array}{l} n_{(1)} \\ n_{(2)} \end{array} \right\} \quad \left\{ \begin{array}{l} n_1 < n_2 \\ d = n_{(2)} - n_{(1)} \end{array} \right.$

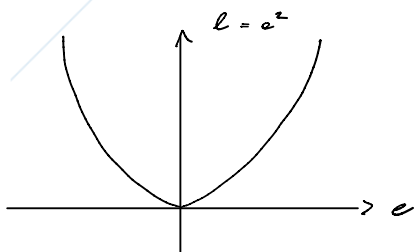
$$\sum_k |n_k - \alpha| = \begin{cases} 2|n_{(1)} - \alpha| + d & \alpha < n_{(1)} \\ 2|n_{(2)} - \alpha| + d & \alpha > n_{(2)} \\ d & n_{(1)} \leq \alpha \leq n_{(2)} \end{cases} \quad \text{: condizione di minimo}$$

Ipotezziamo $n=6$



$S = 2$

$l(c) = c^2$



$D = \sum_{i=1}^r (x_i - \alpha)^2 n_i$

\downarrow
 $x_i^2 + \alpha^2 - 2\alpha x_i$: funzione di II grado di argomento α

$\frac{dD}{d\alpha} = -2 \sum_{i=1}^r (x_i - \alpha) n_i = 0 \quad \sum_{i=1}^r (x_i - \alpha) n_i = 0$

$\alpha = \text{Media aritmetica } \mu$ (in quanto μ annulla la somma degli scarti)

$\frac{d^2 D}{d\alpha^2} = 2 \sum_{i=1}^r n_i = 2n > 0$

VARIABILITÀ

: oscillazioni intorno a un indice di posizione

Mediana: non risente di valori estremi (non considera il valore delle osservazioni bensì il loro ordine)

Media: risente della presenza di valori estremi

Indici di variabilità (o di dispersione)

Variabilità: attitudine di un fenomeno ad assumere intensità di valore differenti.

Nel caso di fenomeni qualitativi si parla di **mutabilità** o eterogeneità.

[**NOTA**: non tratteremo indici di variabilità globale]

Proprietà (sia per variabili quantitative che qualitative)

$$\textcircled{1} V(x) \geq 0 \quad [V(x) = \text{variabilità di } x \text{ (per ora)}]$$

$$\textcircled{2} V(x) = 0 \quad \text{se } x_i \text{ sono costanti}$$

$$\textcircled{3} V(x + d) = V(x) \quad : \text{invarianza per traslazione}$$

\Rightarrow se aggiungo o sottraggo una stessa costante ad ogni x_i non cambia la sua variabilità.

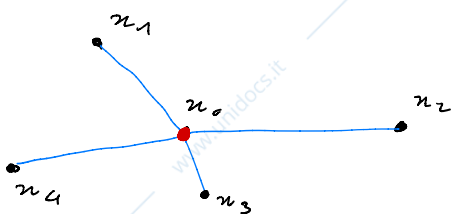
$$\textcircled{4} |x_i - x_j| \leq |y_i - y_j| \quad \forall i, j \Rightarrow V(x) \leq V(y)$$

\Rightarrow monotona crescente rispetto agli scarti in modulo

Lezione 10 - 25/10/21

$$d = |x_i - x_0| \quad i = 1, 2, \dots, n$$

x_0 = CENTRO (media, moda, mediana) o POLO



Il modulo (o la potenza pari) serve per evitare la compensazione degli scarti positivi e negativi dal centro x_0 .

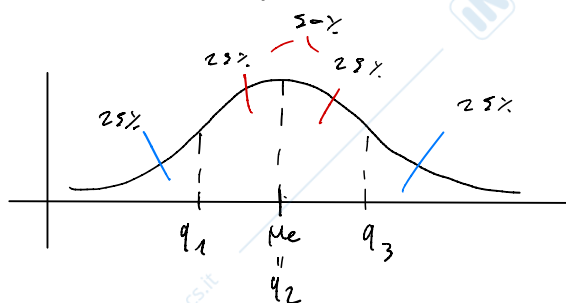
Campo di variazione . Range

$$\text{Range} = x_{\max} - x_{\min}$$

È una misura "embrionale" dello scostamento da un centro (che infatti qui non è espresso)

Differenza interquartile

$DQ = q_3 - q_1 \Rightarrow$ considero il 50% centrale della distribuzione



Anche con la DQ non si considera la dispersione da un centro

Vediamo ora gli indici di variabilità rispetto a un centro.

$$d_i = |n_i - n_0| \quad i = 1, 2, \dots, n$$

Scostamenti medi assoluti

$$n_0 S_r = \left(\frac{\sum_{i=1}^m |n_i - n_0|^r n_i}{n} \right)^{1/r} \quad r > 0$$

$n_0 S_r$ è una media potenziata di ordine r di d_i .

Come scegliere la media potenziata e l'indice di posizione n_0 ?

NB: oltre a rispettare le proprietà degli indici di variabilità, $n_0 S_r$ eredita anche le proprietà delle medie potenziate e dell'indice di posizione n_0 scelto.

- $n_0 = \text{Mediana}$ $r = 1 \Rightarrow$ SMA dalla Mediana

$$Me S_1 = \frac{1}{n} \sum_{i=1}^m |n_i - Me| n_i \quad (\text{che è minima per proprietà della Mediana})$$

- $n_0 = \text{Media aritmetica}$ $r = 1 \Rightarrow$ SMA dalla Media

$$\mu S_1 = \frac{1}{n} \sum_{i=1}^m |n_i - \mu| n_i$$

- $n_0 = \text{Media aritmetica}$ $r = 2 \Rightarrow$ SCARTO QUADRATICO MIPLO

$$\mu S_2 = \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^m (n_i - \mu)^2 n_i}$$

$$(\mu S_2)^2 = \sigma^2 = \text{VARIANZA}$$

Proprietà

① $Me S_1 \leq n_0 S_1$

\Rightarrow la mediana minimizza lo scostamento rispetto a qualsiasi altro indice di posizione per l'ordine 1 (per proprietà della Mediana).

② $\mu S_2 \leq n_0 S_2$

\hookrightarrow Media aritmetica

③ $Me S_1 \leq Me S_2$, $\mu S_1 \leq \mu S_2$, ...

\Rightarrow a parità di indici di posizione, ad un ordine maggiore di media potenziata (potenzia) corrisponde una SMA maggiore.

① + ③ $Me S_1 \leq \mu S_2$

Particolare rilevanza ha la **VARIANZA**

$$\sigma^2 = (SQM)^2 = \frac{1}{n} \sum_{i=1}^m (x_i - \mu)^2 n_i = M[(x - \mu)^2]$$

operatore media

Formula Operativa per il calcolo della varianza

$$\sigma^2 = M(x^2) - \mu^2 = \frac{\sum_{i=1}^m x_i^2 n_i}{n} - \mu^2$$

D.M.

$$\begin{aligned} \sigma^2 &= M[(x - \mu)^2] = M[x^2 + \mu^2 - 2x\mu] = M(x^2) + \mu^2 - 2\mu M(x) = \\ &= M(x^2) + \mu^2 - 2\mu^2 = M(x^2) - \mu^2 \quad \underline{CVD} \end{aligned}$$

Esempio : Formula operativa

x_i	h_i	$x_i h_i$	$x_i^2 h_i$
5	5	25	125
6	10	60	360
7	20	140	980
8	10	80	640
9	5	45	405
	50	350	2510

$$\mu = 7 \quad M(x) = \frac{350}{50} = 7 \quad M(x^2) = \frac{2510}{50} = 50,2$$

$$\sigma^2 = 50,2 - 7^2 = 1,2$$

$$\sigma = \sqrt{1,2} = 1,095$$

Metodo "classico" (con $\sigma^2 = \frac{1}{n} \sum_{i=1}^m (x_i - \mu)^2 n_i$)

x_i	n_i	$x_i - \mu$	$(x_i - \mu)^2 n_i$
5	5	-2	20
6	10	-1	10
4	20	0	0
8	10	1	10
9	5	2	20
	50	0	60

$$\sigma^2 = \frac{60}{50} = 1,2$$

NB: La varianza σ^2 restituisce una quantità espressa nell'unità di misura di partenza al quadrato, non utile nella pratica; per questo si usa lo s.d.m.

Proprietà della varianza

Consideriamo l'operatore $\text{Var}(x)$ (o $V(x)$) e la costante λ .

① $V(\lambda) = 0 \Rightarrow \sigma^2$ di una costante è ovviamente nulla

$$\text{DIM. } M(\lambda^2) - [M(\lambda)]^2 = \lambda^2 - \lambda^2 = 0$$

② $V(\lambda x) = \lambda^2 V(x)$

NB: $V(\lambda x)$ non è una traslazione, bensì un cambiamento di scala.

$$\begin{aligned} \text{DIM. } V(\lambda x) &= \overset{\text{DEF.}}{M[\lambda x - M(\lambda x)]^2} = M[\lambda x - \lambda M(x)]^2 = \\ &= \lambda^2 M[x - M(x)]^2 \overset{\text{DEF.}}{=} \lambda^2 V(x) \end{aligned}$$

③ $V(\lambda + x) = V(x) \Rightarrow$ invariante per traslazione

$$\begin{aligned} \text{DIM. } V(\lambda + x) &= M[(\lambda + x) - M(\lambda + x)]^2 = M[(\lambda + x) - M(\lambda) - M(x)]^2 = \\ &= M[\lambda + x - \lambda - M(x)]^2 = M[x - M(x)]^2 = V(x) \end{aligned}$$

④ $V(x + y) = \text{VAR}(x) + \text{VAR}(y) + 2 M[(x - M(x))(y - M(y))]$

$$\begin{aligned} \Rightarrow \text{VAR}(x + y) &= M[(x + y) - M(x + y)]^2 = M[(x + y) - M(x) - M(y)]^2 = \\ &= M[(x - M(x)) + (y - M(y))]^2 = M[(x - M(x))^2] + M[(y - M(y))^2] + \\ &+ 2 M[(x - M(x))(y - M(y))] = \text{VAR}(x) + \text{VAR}(y) + 2 \text{COVARIANZA}(x, y) \end{aligned}$$

COVARIANZA

④ bis : se $\text{COV}(x, y) = 0 \Rightarrow \text{VAR}(x + y) = \text{VAR}(x) + \text{VAR}(y)$

Indici di variabilità relativa

COEFFICIENTE DI VARIATIONE

$CV = \frac{\sigma}{\mu} \Rightarrow$ Numero puro \Rightarrow adimensionale, indipendente da unità di misura

↳ rapporto tra indice di dispersione e indice di posizione

NB: $CV = \frac{\sigma}{\mu}$ deve essere ≥ 0 ; se $\mu < 0 \Rightarrow CV = \frac{\sigma}{|\mu|}$

Indici di variabilità normalizzati

Indici normalizzati:

- numeri puri \Rightarrow privi di unità di misura
- campo di variazione $[0, 1]$

NB: È buona norma normalizzare solo con caratteri trasferibili.

Normalizzare un indice:

$$I_{\min} \leq I \leq I_{\max}$$

$$I_N = \frac{I - I_{\min}}{I_{\max} - I_{\min}} \Rightarrow 0 \leq I_N \leq 1$$

$$\sigma_N^2 = \frac{\sigma^2 - 0}{\sigma_{\max}^2 - 0} = \frac{\sigma^2}{\sigma_{\max}^2} \Rightarrow 0 \leq \sigma_N^2 \leq 1$$

$\sigma_{\min}^2 = 0$

Ma $\sigma_{\max}^2 = ?$

Intervallo di dispersione $[A, B]$

n_i	n_i
A	n_A
B	n_B
	n

① $n = n_A + n_B \Rightarrow$ invarianza della popolazione

② $\mu = \frac{A n_A + B n_B}{n} \Rightarrow$ invarianza della media

$$\Rightarrow \begin{cases} n_A + n_B = n \\ \frac{A n_A + B n_B}{n} = \mu \end{cases} \quad \begin{matrix} n_A = ? \\ n_B = ? \end{matrix} \quad \begin{cases} n_A = n - n_B \\ n\mu = A(n - n_B) + B n_B \end{cases}$$

$$\begin{cases} n_A = n - n_B \\ n\mu = nA - A n_B + B n_B \end{cases} \quad \begin{cases} - \\ n\mu = n_B(B - A) + nA \end{cases}$$

$$\begin{cases} n_A = n - n_B \\ n_B = \frac{n\mu - nA}{B-A} \end{cases} \quad n_A = n - \frac{n\mu - nA}{B-A} = \frac{n(B-A) - n(\mu - A)}{B-A} = \frac{n(B-A - \mu + A)}{B-A} = n \frac{B-\mu}{B-A}$$

$$\Rightarrow n_A = n \frac{B-\mu}{B-A} \Rightarrow f_A = \frac{B-\mu}{B-A}$$

$$\Rightarrow n_B = n \frac{\mu - A}{B-A} \Rightarrow f_B = \frac{\mu - A}{B-A} = 1 - f_A$$

Quindi:

x_i	n_i
A	$n \frac{B-\mu}{B-A}$
B	$\frac{n - n_A}{n}$

$$\sigma_{\max}^2 = M(x^2) - \mu^2 = \frac{A^2}{n} \cdot \frac{B-\mu}{B-A} + \frac{B^2}{n} (n - n_A) - \mu^2$$

$$= A^2 \frac{B-\mu}{B-A} + \frac{B^2}{n} \left(n - n \frac{B-\mu}{B-A} \right) - \mu^2 =$$

$$= A^2 \frac{B-\mu}{B-A} + B^2 \left(1 - \frac{B-\mu}{B-A} \right) - \mu^2 =$$

$$= A^2 \frac{B-\mu}{B-A} + B^2 - B^2 \frac{B-\mu}{B-A} - \mu^2 = \frac{B-\mu}{B-A} (A^2 - B^2) + B^2 - \mu^2$$

$$= \frac{B-\mu}{B-A} (A+B)(A-B) + (B^2 - \mu^2) = -(B-\mu)(A+B) + (B+\mu)(B-\mu)$$

$$= (B-\mu)(-A-B+B+\mu) = (B-\mu)(\mu-A) \quad \underline{\text{CVD}}$$

Consideriamo ora l'insieme classe $[0, T]$ $T = \sum_{i=1}^r x_i n_i$

x_i	n_i
0	$n-1$
T	1
	<hr/> n

Infatti:

x_i	n_i
0	$n \frac{B-\mu}{B-A} = n \frac{(n\mu - \mu)}{n\mu - 0} = n \frac{\mu(n-1)}{n\mu}$
T	$n - n_A = 1 = n - 1$
	<hr/> n

$$\sigma_{\max}^2 = (B-\mu)(\mu-A) = (n\mu - \mu)(\mu - 0) = \mu(n-1)\mu$$

$$= \mu^2(n-1)$$

Esempio

x_i	n_i
6	3
7	8
8	6
9	4
10	19
	40

$$[A, B] = [6, 10]$$

$$n_A = n \frac{B-A}{B-A} = 13$$

$$n_B = 40 - 13 = 27$$

$$\mu = 8,7 \quad \sigma^2 = 2,01$$

$$\sigma_{\max}^2 = (B - \mu)(\mu - A) = (10 - 8,7)(8,7 - 6) = 3,51$$

$$\sigma_0^2 = \frac{\sigma^2}{\sigma_{\max}^2} = \frac{2,01}{3,51} = 0,5726$$

Lezione 12 - 28/10/21

Motabilità = variabilità di variabili statistiche qualitative (= motabili)

Indici di variabilità per fenomeni qualitativi

① Motabilità nulla = una sola modalità assume la totalità delle frequenze

$$\Rightarrow \exists n_i = n, n_j = 0 \forall j \neq i$$

② Motabilità massima = frequenze equidistribuite tra le modalità

$$k = \text{numero delle modalità} \quad n_1 = n_2 = \dots = n_k = \frac{n}{k}$$

NB: è importante definire i due "paletti" di min e max di un indice in quanto sono essenziali per normalizzarlo.

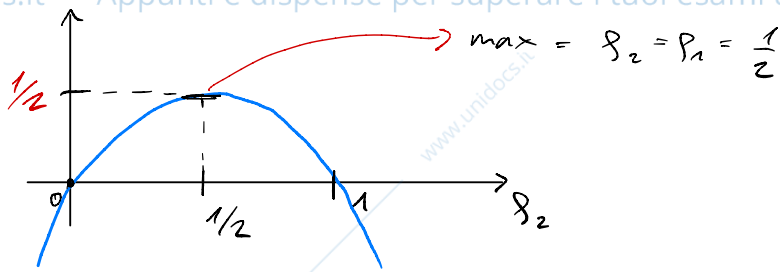
Indice di Gini $E(X)$ (o di eterogeneità)

Ipotizziamo X dicotomica $[0, 1]$

x_i	n_i	p_i
0	n_1	p_1
1	n_2	$p_2 = 1 - p_1$
	n	1

$$\begin{aligned} \sigma_x^2 &= M(x^2) - \mu^2 = (0^2 \cdot p_1 + 1^2 \cdot p_2) - (0 \cdot p_1 + 1 \cdot p_2)^2 \\ &= \underbrace{p_2 - p_2^2}_{\text{parabola}} = p_2(1 - p_2) \end{aligned}$$

$$y = p_2 - p_2^2$$



Ipotezziamo k modalità

$$E(x) = \sum_{i=1}^k p_i (1 - p_i) = 1 - \sum_{i=1}^k p_i^2$$

$$\begin{aligned} \text{D.M.} \quad E(x) &= \sum_{i=1}^k p_i (1 - p_i) = \sum_{i=1}^k p_i - \sum_{i=1}^k p_i^2 = \sum_{i=1}^k p_i - \sum_{i=1}^k p_i^2 = \\ &= 1 - \sum_{i=1}^k p_i^2 \end{aligned}$$

Proprietà del Gini

$$\textcircled{1} \quad \min E(x) = 0$$

$$\text{es.} \quad p_1 = 1 \quad p_2 = p_3 = \dots = p_k = 0 \quad (\text{minima eterogeneità})$$

$$E(x) = 1 - \left(\underbrace{p_1^2}_1 - \underbrace{\sum_{i=2}^k p_i^2}_0 \right) = 0 \quad \checkmark \text{ rispetta la condizione } \textcircled{1}$$

$$\textcircled{2} \quad \max E(x)$$

$$p_1 = p_2 = \dots = p_k = \frac{1}{k}$$

$$E(x) = 1 - \sum_{i=1}^k \left(\frac{1}{k}\right)^2 = 1 - k \left(\frac{1}{k}\right)^2 = 1 - \frac{1}{k}$$

Indice di Gini normalizzato

$$E_N(x) = \frac{E(x) - E_{\min}(x)}{E_{\max}(x) - E_{\min}(x)} = \frac{E(x)}{E_{\max}(x)}$$

$$\text{Ovviamente: } 0 \leq E_N(x) \leq 1$$

$$\Rightarrow E_N(x) = \frac{1 - \sum_{i=1}^k p_i^2}{1 - \frac{1}{k}}$$

Esempio n° addetti x settore

x_i	h_i	f_i	f_i^2
AGR	120	0,0311	0,0010
IND	1673	0,4334	0,1878
TER	2067	0,5355	0,2868
	3860	1	0,4756

$$E(x) = 1 - \sum_{i=1}^k f_i^2 = 1 - 0,4756 = 0,5244$$

$$E_{\max}(x) = 1 - \frac{1}{k} = 1 - \frac{1}{3} = 0,6666$$

$$E_N(x) = \frac{E(x)}{E_{\max}(x)} = \frac{0,5244}{0,6666} = 0,7866$$

⇒ il Gini è vicino al suo max
⇒ forte eterogeneità

⇒ forte disomogeneità tra i settori

Momenti di una distribuzione

Momento r-esimo di x da α = media aritmetica delle potenze r-esime degli scarti da α

$$\mu_r^\alpha = M[(x_i - \alpha)^r] \cdot h_i = 1 \quad \neq 1$$

$$L = \frac{1}{h} \sum_{i=1}^n (x_i - \alpha)^r \quad \text{se le frequenze sono unitarie}$$

$$L = \frac{1}{h} \sum_{i=1}^k (x_i - \alpha)^r h_i \quad \text{se le frequenze non sono unitarie ma si ripropongono su più unità statistiche}$$

$$\alpha = 0$$

Momenti dall'origine

$${}_0\mu_r \equiv \mu_r = M[(x-0)^r] = M[x^r] = \frac{1}{h} \sum_{i=1}^r x_i^r h_i$$

Notare che: $\left\{ \begin{array}{l} r=1 \Rightarrow \mu_1 = \text{Media aritmetica} \\ r=2 \Rightarrow \mu_2 = M(x^2) \end{array} \right.$

$$\alpha = \mu$$

Momenti dalla media • Momenti CENTRALI

$$\bar{\mu}_r = \frac{1}{h} \sum_{i=1}^k (x_i - \mu)^r h_i$$

Notare che: $\left\{ \begin{array}{l} r=1 \Rightarrow \bar{\mu}_1 = 0 \quad (\text{Proprietà media aritmetica}) \\ r=2 \Rightarrow \bar{\mu}_2 = \text{Varianza} \end{array} \right.$

Quindi : $\mu_1 = M(x) = \mu$

$$\mu_2 = M(x^2) = \frac{1}{n} \sum_{i=1}^k x_i^2 h_i$$

$$\bar{\mu}_1 = M[(x - \mu)] = 0$$

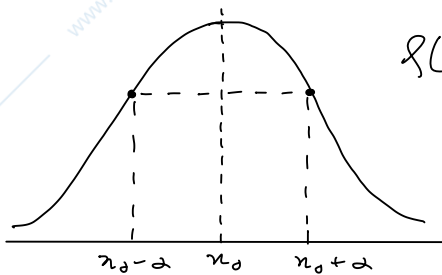
$$\bar{\mu}_2 = M[(x - \mu)^2] = \sigma^2 = M(x^2) - [M(x)]^2$$

NB: I momenti sono tra loro legati (infatti, ad es., $\bar{\mu}_1$ è funzione di μ_1)

Indicatori di forma

Asimmetria

$$x = x_0, \alpha > 0$$



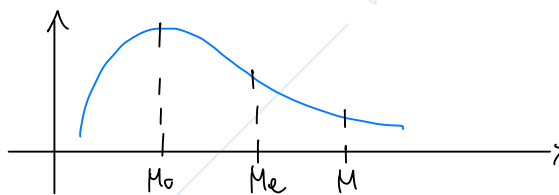
$$f(x_0 - \alpha) = f(x_0 + \alpha) \Rightarrow f \text{ simmetrica rispetto a } x_0$$

- Vogliamo:
- Funzione continua
 - Distr. unimodale (un solo punto di "picco")

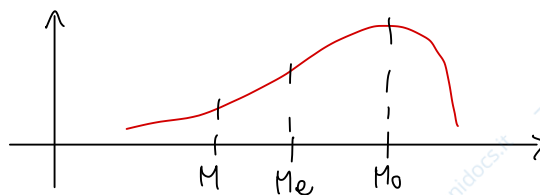
Se la distribuzione è simmetrica $\Rightarrow M(x) = Me(x) = Mo(x)$

Positiva $\Rightarrow Mo(x) < Me(x) < M(x)$

Asimmetria



Negativa $\Rightarrow Mo(x) > Me(x) > M(x)$



Indice di (a)simmetria di Fisher γ_1

In una distribuzione simmetrica, i momenti centrali di ordine dispari sono nulli.

$$\bar{\mu}_{2r+1} = 0 \quad r = 0, 1, 2, \dots$$

$$\hookrightarrow \bar{\mu}_3 = \bar{\mu}_5 = \dots = 0$$

Per verificarlo, dovremmo verificare tutti i momenti centrali di ordine dispari, che sono ∞ .

$\Rightarrow \gamma_1$ dice solo se una distribuzione è asimmetrica (se trova un m.c. di ordine dispari $\neq 0$)

NB: Simmetria $\Rightarrow Me = Mo = M$

$$\underline{MA} \quad Me = Mo = M \quad \not\Rightarrow \text{Simmetria}$$

L' unico modo per verificare la simmetria è per via grafica.

Tornando all'asimmetria:

- Asimmetria positiva: prevalgono gli scarti positivi $\Rightarrow \bar{\mu}_{2r+1} > 0$
- Asimmetria negativa: prevalgono gli scarti negativi $\Rightarrow \bar{\mu}_{2r+1} < 0$

$$\gamma_1 = \frac{\bar{\mu}_3}{\sigma^3} \quad \gamma_1 \text{ è adimensionale}$$

$$\gamma_1 > 0 \Rightarrow \bar{\mu}_3 > 0 \Rightarrow \text{asimmetria positiva}$$

$$\gamma_1 < 0 \Rightarrow \bar{\mu}_3 < 0 \Rightarrow \text{asimmetria negativa}$$

$$\gamma_1 = 0 \quad \text{Condizione necessaria ma non sufficiente per la simmetria}$$

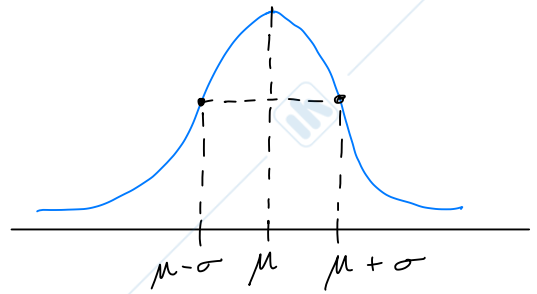
Quindi γ_1 è un SINTOMO di simmetria.

Indice di Curtosi

CURTOSI = grado di "appiattimento" di una curva (rispetto alla gaussiana)

Gaussiana:

- variabile statistica continua
- unimodale (un solo punto di massimo, in $\mu = Me = Mo$)
- Area sottesa dalla curva:



$$f(x) = \frac{1}{\sqrt{2\pi b}} \exp\left\{-\frac{1}{2} \left(\frac{x-a}{b}\right)^2\right\}$$

dove $a = M(x) = \mu$, $b = V(x) = \sigma^2$

$$\int_{-\infty}^{+\infty} f(x) dx = 1$$

- Simmetrica, con asse di simmetria $x = \mu$
- $\bar{\mu}_3 = \bar{\mu}_5 = \dots = 0$ e $\bar{\mu}_2 = \sigma^2$, $\bar{\mu}_4 = 3\sigma^4$
- $\gamma_1 = 0$
- $\gamma_2 = \frac{\bar{\mu}_4}{\sigma^4} = \frac{3\sigma^4}{\sigma^4} = 3$

\Rightarrow Se $\gamma_2 = 3$, la curva in questione ha la forma di una Normale.

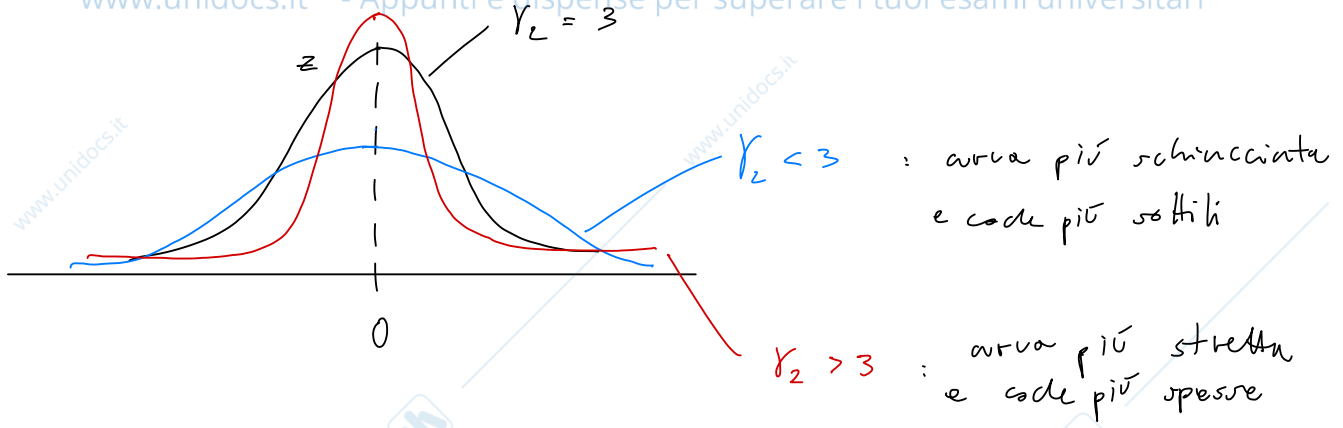
Ci interessa studiare γ_2

- < 3 Platicurtica
- $= 3$ Normocurtica
- > 3 Leptocurtica

Standardizzare una variabile statistica:

$$z = \frac{x - \mu}{\sigma}$$

z ha $\mu = 0$, $\sigma^2 = 1$



STATISTICA DESCRITTIVA BIVARIATA

: rilevazione di due caratteri per ogni unità statistica (tipicamente mettendoli in relazione)

Distribuzione **CONGIUNTA** di frequenze assolute dei caratteri X, Y in una tabella a doppia entrata :

$y_j \backslash x_i$	x_1	x_2	...	x_i	...	x_r	$h_{\cdot j}$
y_1	n_{11}	n_{21}		n_{i1}		n_{r1}	$h_{\cdot 1}$
y_2	n_{12}	n_{22}		n_{i2}		n_{r2}	$h_{\cdot 2}$
\vdots	\vdots	\vdots		\vdots		\vdots	\vdots
y_j	n_{1j}	n_{2j}		n_{ij}		n_{rj}	$h_{\cdot j}$
\vdots	\vdots	\vdots		\vdots		\vdots	\vdots
y_m	n_{1m}	n_{2m}		n_{im}		n_{rm}	$h_{\cdot m}$
$h_{i \cdot}$	$h_{1 \cdot}$	$h_{2 \cdot}$		$h_{i \cdot}$		$h_{r \cdot}$	h

numerosità totale

Le frequenze relative si ottengono : $f_{ij} = \frac{n_{ij}}{h}$
 (al posto di h c'è 1)

Dalla tabella si possono ricavare diverse distribuzioni univariate :

- Distribuzioni marginali :
 - rispetto a X
 - rispetto a Y
- Distribuzioni condizionale :
 - Y condizionato x_1
 - X condizionato y_2

oss. le condizionate permettono di evidenziare come varia X dato un certo valore di Y .

Distribuzioni marginali

X		Y	
x_i	$n_{i\cdot}$	y_j	$n_{\cdot j}$
x_1	$n_{1\cdot}$	y_1	$n_{\cdot 1}$
x_2	$n_{2\cdot}$	y_2	$n_{\cdot 2}$
\vdots	\vdots	\vdots	\vdots
x_t	$n_{t\cdot}$	y_m	$n_{\cdot m}$
	n		n

Quante unità statistiche presentano una certa modalità di Y indipendentemente da come si comporta X

Variabili statistiche condizionate

$X|y_j$ o $X|Y=y_j$ in matematica : $x = f(y)$

\Rightarrow studio come si comporta X (= come varia X) dato un a certa modalità di Y fissata.

$x_i y_j$	n_{ij}
x_1	n_{1j}
x_2	n_{2j}
\vdots	\vdots
x_t	n_{tj}
	$n_{\cdot j}$

Fisso la y_j e scotto tra le colonne x_i

$j = 1, 2, \dots, m \Rightarrow$ Ci sono m possibili distribuzioni univariate di questo tipo ($X|Y=y_j$)

$n_{\cdot j} \rightarrow$ totale di riga

$Y|x_i$

$y_j x_i$	n_{ij}
y_1	n_{i1}
y_2	n_{i2}
\vdots	\vdots
y_m	n_{im}

Fisso la colonna x_i e scotto tra le righe y_j

$i = 1, 2, \dots, t$

$n_{i\cdot} \rightarrow$ totale di colonna

Quindi:

① Distribuzione marginale di X

$$n_{i\cdot} = \sum_{j=1}^m n_{ij} \qquad \sum_{i=1}^r n_{i\cdot} = n$$

↳ somma di tutte le + freq. assolute marginali

② Distribuzione marginale di Y

$$n_{\cdot j} = \sum_{i=1}^r n_{ij} \qquad \sum_{j=1}^m n_{\cdot j} = n$$

③ Distribuzione condizionata $X|Y=y_j$

$(x_i ; n_{ij})$ con j fissato

④ Distribuzione condizionata $Y|X=x_i$

$(y_j ; n_{ij})$ con i fissato

In termini di frequenze relative:

$$f_{ij} = \frac{n_{ij}}{n} \qquad f_{i\cdot} = \frac{n_{i\cdot}}{n} \qquad f_{\cdot j} = \frac{n_{\cdot j}}{n}$$

$$\sum_{i=1}^r \sum_{j=1}^m f_{ij} = \sum_{i=1}^r f_{i\cdot} = \sum_{j=1}^m f_{\cdot j} = 1$$

$$f_{y_j | x_i} = \frac{n_{ij}}{n_{i\cdot}} \qquad f_{x_i | y_j} = \frac{n_{ij}}{n_{\cdot j}}$$

$$\sum_{j=1}^m f_{y_j | x_i} = \sum_{j=1}^m \frac{n_{ij}}{n_{i\cdot}} = 1 \qquad \sum_{i=1}^r f_{x_i | y_j} = \sum_{i=1}^r \frac{n_{ij}}{n_{\cdot j}} = 1$$

Esempio

FATT. \ C.A.	DET.	G.D.	GROSS.	
0-10	3	0	1	4
10-20	3	0	5	8
20-30	0	2	2	4
30-40	0	2	1	3
40-50	0	0	1	1
	6	4	10	20

n_i	$n_{i\cdot}$	$f_{i\cdot}$	$n_{\cdot j}$	$f_{i\cdot} \cdot f_{\cdot j}$
DET	6	6/20 = 0,3	0	0
G.D.	4	4/20 = 0,2	2	2/4 = 0,5
GR.	10	10/20 = 0,5	2	2/4 = 0,5
	20	1	4	1

Connessione e Indipendenza stocastica

es. $Y =$ professione paterna $X =$ scolarità dei figli

$Y \backslash X$	L.M.	DIP.	LAUREA	$n_{.j}$
OP.	3	2	5	10
IMP.	9	6	15	30
DIP.	18	12	30	60
$n_{i.}$	30	20	50	100

$Y|X = x_i$ (professione paterna condizionata da istruzione dei figli?)

y_j	$n_{ij}/n_{.1}$
OP.	0,1
IMP.	0,3
DIP.	0,6
	1

y_j	$n_{ij}/n_{.2}$
OP.	0,1
IMP.	0,3
DIP.	0,6
	1

y_j	$n_{ij}/n_{.3}$
OP.	0,1
IMP.	0,3
DIP.	0,6
	1

le f. relative condizionate sono uguali indipendentemente dalla x_i che scelgo come condizionante.

$X|Y = y_j$

x_i	$n_{ij}/n_{.1}$
L.M.	0,3
DIP.	0,2
LAUREA	0,5
	1

x_i	$n_{ij}/n_{.2}$
L.M.	0,3
DIP.	0,2
LAUREA	0,5
	1

x_i	$n_{ij}/n_{.3}$
L.M.	0,3
DIP.	0,2
LAUREA	0,5
	1

Idem come prima.

TEC. (Indipendenza stocastica)

Sia (X, Y) una variabile (o mutabile) statistica doppia.¹

Si dice che le componenti X e Y sono stocasticamente indipendenti se le variabili condizionate hanno la stessa distribuzione percentuale.

¹ Lavorando con le frequenze e non con le modalità, l'indipendenza stocastica vale sia per le variabili che per le mutabili.

UB: Indipendenza stocastica \Rightarrow assenza di legame funzionale.

Quindi :

$$\textcircled{a} \quad X|Y = y_j \quad \Rightarrow \quad \begin{array}{l} n_{1j}/n_{0j} = H_1 \\ \vdots \\ n_{tj}/n_{0j} = H_t \end{array} \quad \begin{array}{l} \forall j, \text{ fissandolo ogni volta} \\ H \text{ \u00e9 indipendente da } j, \forall j \end{array}$$

$$\textcircled{b} \quad Y|X = x_i \quad \Rightarrow \quad \begin{array}{l} n_{i1}/n_{i0} = k_1 \\ \vdots \\ n_{im}/n_{i0} = k_m \end{array} \quad \begin{array}{l} \forall i \text{ fissato} \\ k \text{ \u00e9 indipendente da } i \end{array}$$

Teo.

i. Condizione necessaria e sufficiente per l'indipendenza stocastica \u00e9 la fattorizzazione delle frequenze congiunte.

$$n_{ij} = \frac{n_{i0} \cdot n_{0j}}{n} \quad \forall i, j$$

(\u2192 il prodotto delle marginali diviso per la numerosit\u00e0 totale deve dare la frequenza congiunta)

ii. L'indipendenza stocastica \u00e9 simmetrica.

$$X \text{ \u00e9 indep. stoc. da } Y \iff Y \text{ \u00e9 indep. stoc. da } X$$

iii. Le distribuzioni condizionate sono uguali alle marginali percentuali in caso di indipendenza stocastica.

$$H_i = \frac{n_{ij}}{n_{0j}} = \frac{n_{i0}}{n} \quad \leftarrow \textcircled{a} \quad k_j = \frac{n_{ij}}{n_{i0}} = \frac{n_{0j}}{n} \quad \leftarrow \textcircled{b}$$

DIM. (Y|X)

$$\frac{n_{ij}}{n_{i0}} = k_j \quad \forall i \quad \Rightarrow \quad n_{ij} = k_j \cdot n_{i0} \quad \Rightarrow \quad \sum_{i=1}^t n_{ij} = \sum_{i=1}^t k_j \cdot n_{i0}$$

$$\Rightarrow n_{0j} = \sum_{i=1}^t k_j \cdot n_{i0} \quad \Rightarrow \quad n_{0j} = k_j \sum_{i=1}^t n_{i0} \quad \Rightarrow \quad n_{0j} = k_j \cdot n$$

$$\Rightarrow k_j = \frac{n_{0j}}{n} \quad (\text{DIM iii})$$

$$\frac{n_{ij}}{n_{i0}} = \frac{n_{0j}}{n} \quad \Rightarrow \quad n_{ij} = \frac{n_{i0} \cdot n_{0j}}{n} \quad (\text{DIM i})$$