

Lezione : Varianza totale e generalizzata

Docente: Aldo Solari

Nel caso $p = 1$, la variabilità (o dispersione) presente nelle misurazioni della variabile considerata è descritta da un singolo numero, la varianza $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$. Nel caso $p > 1$, la variabilità presente nelle misurazioni delle p variabili considerate è descritta da p varianze $s_{jj}, j = 1, \dots, p$ e $\frac{1}{2}p(p-1)$ covarianze $s_{jk}, j \neq k = 1, \dots, p$, ovvero

$$p + \frac{1}{2}p(p-1)$$

numerici, contenuti nella matrice di varianze/covarianze S $p \times p$. Possiamo riassumere la variabilità descritta da S in un singolo numero (senza perdere troppa informazione)?

1 Varianza totale

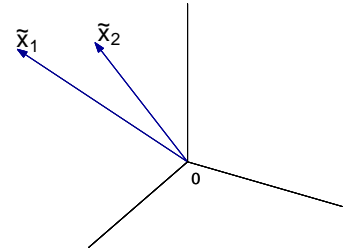
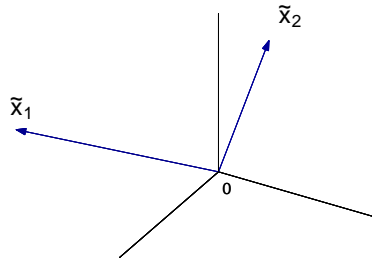
$$\text{Varianza totale} = \text{tr}(S) = \sum_{j=1}^p s_{jj}$$

1.1 $\text{tr}(S)$ nello spazio delle osservazioni

Nello spazio delle osservazioni, la varianza totale può essere interpretata come la somma delle lunghezze al quadrato dei p vettori scarto dalla media $\tilde{x}_j, j = 1, \dots, p$, divisa per n .

$$\text{tr}(S) = \frac{1}{n} \sum_{j=1}^p n s_{jj} = \frac{1}{n} \sum_{j=1}^p \|\tilde{x}_j\|_{p \times 1}^2$$

Sintetizzando la matrice di varianze/covarianze con un singolo numero dato dalla varianza totale, perdiamo tutta l'informazione sulla struttura di correlazione (di covarianza) tra le p variabili.



$$\tilde{X}_{3 \times 2} = \begin{bmatrix} 2 & -2 \\ -3 & 0 \\ 1 & 2 \end{bmatrix} \quad \text{tr}({}_2S) = \frac{14}{3} + \frac{8}{3} \quad r_{12} = -0.19$$

$$\tilde{X}_{3 \times 2} = \begin{bmatrix} 1 & 0 \\ -3 & -2 \\ 2 & 2 \end{bmatrix} \quad \text{tr}({}_2S) = \frac{14}{3} + \frac{8}{3} \quad r_{12} = 0.95$$

1.2 $\text{tr}(S)$ nello spazio delle variabili

La distanza Euclidea

- tra due unità statistiche u_i' e u_l' :
 $1 \times p \quad 1 \times p$

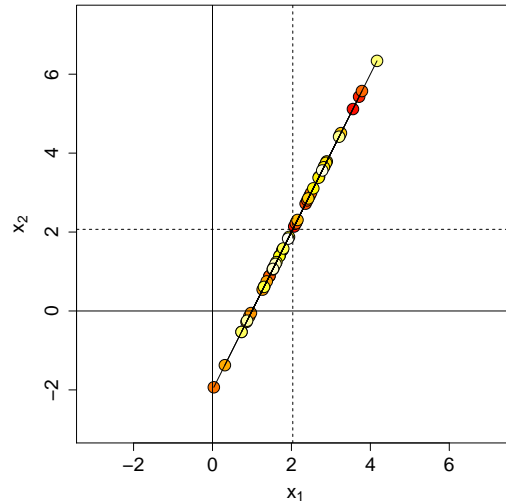
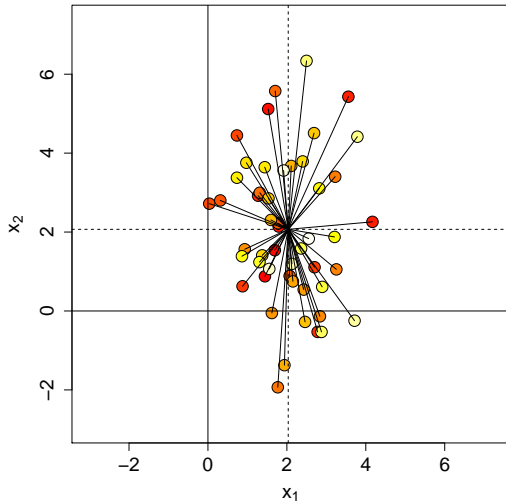
$$d(u_i, u_l) = \sqrt{(u_i - u_l)'(u_i - u_l)} = \sqrt{\sum_{j=1}^p (x_{ij} - x_{lj})^2}$$

- tra l' i -sima unità statistica u_i' e il baricentro \bar{x}' :
 $1 \times p \quad 1 \times p$

$$d(u_i, \bar{x}) = \sqrt{(u_i - \bar{x})'(u_i - \bar{x})} = \sqrt{\sum_{j=1}^p (x_{ij} - \bar{x}_j)^2}$$

Nello spazio delle variabili, la varianza totale può essere interpretata come la media aritmetica delle distanze Euclidee al quadrato delle n unità statistiche $u'_i, i = 1, \dots, n$, dal baricentro \bar{x}'

$$\text{tr}(S) = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - \bar{x}_j)^2 = \frac{1}{n} \sum_{i=1}^n d^2(u_i, \bar{x})$$



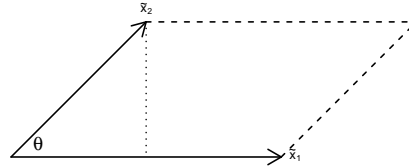
2 Varianza generalizzata

Varianza generalizzata = $\det(S)$
 $p \times p$

2.1 $\det(S)$ nello spazio delle osservazioni

Consideriamo geometricamente l'area generata da $p = 2$ vettori scarto dalla media \tilde{x}_1 e \tilde{x}_2 nello

$n \times 1$ $n \times 1$



spazio n -dimensionale

$$\begin{aligned} \text{Area parallelogramma} &= \text{base paral.} \cdot \text{altezza paral.} \\ &= \|\tilde{x}_1\| \cdot \|\tilde{x}_2\| \sqrt{1 - \cos^2(\theta)} \\ &= n \sqrt{s_{11}s_{22}(1 - r_{12}^2)} \end{aligned}$$

$$\begin{aligned} \det(S_{2 \times 2}) &= \det \left(\begin{bmatrix} s_{11} & \sqrt{s_{11}}\sqrt{s_{22}}r_{12} \\ \sqrt{s_{11}}\sqrt{s_{22}}r_{12} & s_{22} \end{bmatrix} \right) \\ &= s_{11}s_{22} - s_{11}s_{22}r_{12}^2 \\ &= s_{11}s_{22}(1 - r_{12}^2) \end{aligned}$$

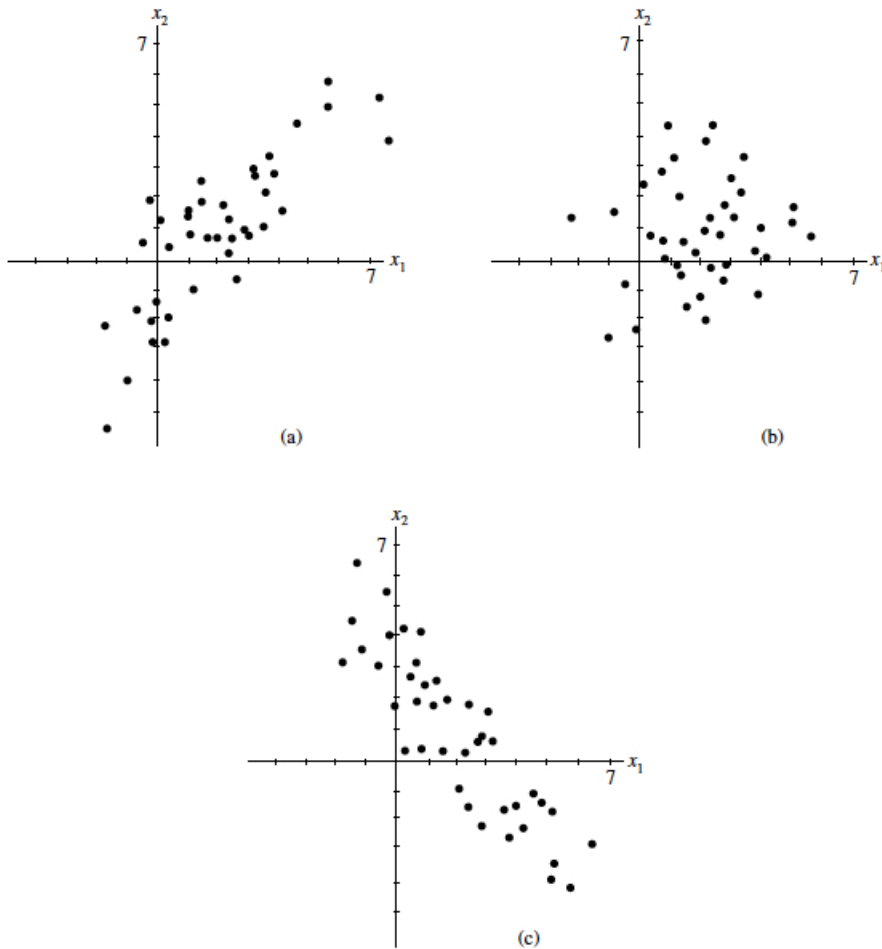
Quindi

$$\det(S_{2 \times 2}) = \frac{(\text{Area parallelogramma})^2}{n^2}$$

In generale, per p vettori n -dimensionali $\tilde{x}_j, j = 1, \dots, p$:

$$\det(S_{p \times p}) = \frac{(\text{Volume parallelepipedo } p\text{-dimensionale})^2}{n^p}$$

2.2 $\det(S)$ nello spazio delle variabili



Alla matrice di varianze/covarianze $S_{p \times p}$ possiamo associare p coppie di autovalori e autovettori

$$(\lambda_1, v_1), (\lambda_2, v_2), \dots, (\lambda_p, v_p)$$

$\begin{matrix} p \times 1 & & p \times 1 & & p \times 1 & & p \times 1 \end{matrix}$

dove gli autovalori (*eigenvalues*) sono ordinati in maniera decrescente, ovvero

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

e dove gli autovettori (*eigenvectors*) v_1, \dots, v_p sono tali che

- hanno lunghezza unitaria $\|v_1\| = \dots = \|v_p\| = 1$
- sono mutualmente perpendicolari: $\langle v_k, v_j \rangle = v_j' v_k = 0$ per $j \neq k$

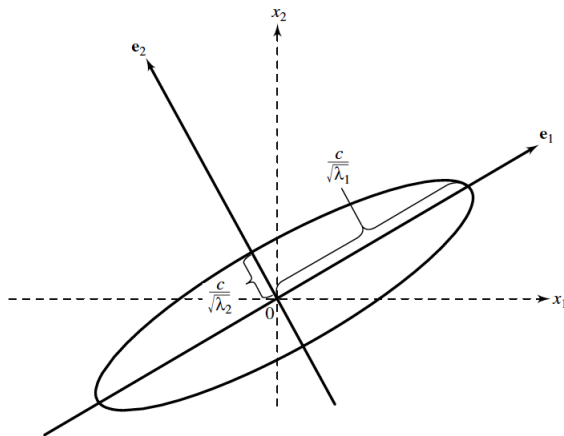


Figure 2.6 Points a constant distance c from the origin ($p = 2, 1 \leq \lambda_1 < \lambda_2$).

L'equazione

$$(x - \bar{x})' \underset{1 \times p}{S^{-1}} \underset{p \times p}{(x - \bar{x})} = c^2$$

definisce l'iper-ellissoide

- centrato sul baricentro \bar{x}'
 $1 \times p$
- con il j -simo asse orientato secondo il j -simo autovettore v_j di S
 $p \times 1$ $p \times p$
- di lunghezza $c\sqrt{\lambda_j}$, proporzionale al j -simo autovalore λ_j di S ,
 $p \times p$

dove stiamo assumendo che S è una matrice definita positiva in modo da garantire l'esistenza di S^{-1} . Il volume dell'iper-ellissoide è funzione della varianza generalizzata:

$$\text{Volume di } \left\{ x' : (x - \bar{x})' \underset{1 \times p}{S^{-1}} \underset{p \times p}{(x - \bar{x})} \leq c^2 \right\} = k_p c^p \sqrt{\det(S)}$$

dove $k_p = \frac{2\pi^{p/2}}{p\Gamma(p/2)}$ e $\Gamma(\cdot)$ è la funzione Gamma. Quindi

$$(\text{Volume iperellissoide})^2 = (\text{costante})(\text{varianza generalizzata})$$

Varianza generalizzata: cosa perdiamo: Sintetizzando la matrice di varianze/covarianze con un singolo numero dato dalla varianza generalizzata, perdiamo l'informazione riguardante l'orientamento della nuvola di punti p -dimensionale formata dalle n unità statistiche

2.3 Quando la varianza generalizzata è zero?

Proposition 2.1. *La varianza generalizzata è 0 se e solo se le colonne di \tilde{X} sono linearmente dipendenti.*

Dimostrazione. Si ricordi che le colonne di \tilde{X} , ovvero i vettori \tilde{x}_j , $j = 1, \dots, p$, sono linearmente dipendenti se esiste un vettore non nullo $c \neq 0$ tale che

$$\tilde{X} c = c_1 \tilde{x}_1 + \dots + c_p \tilde{x}_p = 0$$

⇐

Se le colonne di \tilde{X} sono linearmente dipendenti, esiste $c \neq 0$ tale che

$$0 = \tilde{X} c$$

quindi

$$n S c = \tilde{X}' \tilde{X} c = \tilde{X}' 0 = 0.$$

Segue che esiste $c \neq 0$ tale che $S c = 0$, ovvero che S è una matrice singolare, e quindi

$$\det(S) = 0$$

⇒

Se $\det(S) = 0$, allora S è singolare ed esiste $c \neq 0$ tale che $S c = 0$, ovvero

$$\begin{aligned} 0 &= n S c = \tilde{X}' \tilde{X} c \\ c' 0 &= c' \tilde{X}' \tilde{X} c \\ 0 &= \|\tilde{X} c\|^2 \end{aligned}$$

e quindi per avere lunghezza 0 dobbiamo avere $\tilde{X} c = 0$, ovvero le colonne di \tilde{X} sono linearmente dipendenti. □

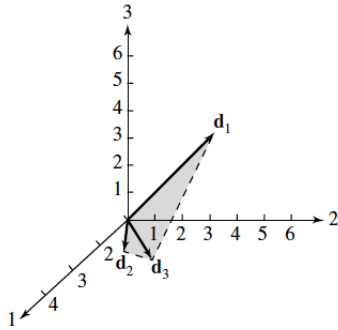
Example 2.2. $X = \begin{bmatrix} 1 & 2 & 5 \\ 4 & 1 & 6 \\ 4 & 0 & 4 \end{bmatrix}$, $\tilde{X} = \begin{bmatrix} -2 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & -1 & -1 \end{bmatrix}$, quindi poichè

$$\tilde{x}_3 = \tilde{x}_1 + 2\tilde{x}_2$$

le colonne \tilde{X} sono linearmente dipendenti, ovvero $\tilde{X} c = 0$ con $c = \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \neq 0$.

Geometricamente questo significa che uno dei vettori scarto dalla media, ad esempio \tilde{x}_3 , giace

nel piano generato da \tilde{x}_1 e \tilde{x}_2 . Di conseguenza, il volume del parallelepipedo tridimensionale è 0.



Proposition 2.3. Se $n \leq p$, allora $\det(S) = 0$

Dimostrazione. Sia $\tilde{u}_i = [\tilde{x}_{i1} \ \cdots \ \tilde{x}_{ip}]$ l' i -sima riga di \tilde{X} . Abbiamo

$$\tilde{X}' \mathbf{1} = 1 \cdot \tilde{u}_1 + \dots + 1 \cdot \tilde{u}_n = \begin{bmatrix} \sum_{i=1}^n \tilde{x}_{i1} \\ \vdots \\ \sum_{i=1}^n \tilde{x}_{ip} \end{bmatrix} = \mathbf{0}_{p \times 1}$$

quindi le righe di \tilde{X} sono linearmente dipendenti. Allora $\text{rango}(\tilde{X}) < n \leq p$. Segue che

$$\text{rango}(\tilde{X}) = \text{rango}(\tilde{X}' \tilde{X}) = \text{rango}(n S) = \text{rango}(S) < p$$

e quindi S è singolare, e risulta $\det(S) = 0$ □

2.4 Varianza generalizzata per dati standardizzati

Varianza generalizzata per
dati standardizzati Z $= \det(S^Z) = \det(R)$

$$\begin{aligned} \det(S) &= \det(D^{1/2} R D^{1/2}) \\ &= \det(D^{1/2}) \det(R) \det(D^{1/2}) \\ &= (s_{11} s_{22} \cdots s_{pp}) \det(R) \\ &= \left(\prod_{j=1}^p s_{jj} \right) \det(R) \end{aligned}$$

dove $D^{1/2} = \text{diag}(\sqrt{s_{11}}, \dots, \sqrt{s_{pp}})$

Example 2.4. *Se cambiamo l'unità di misura per la prima variabile x_1 , ad esempio da Kg a gr, e quindi moltiplicando x_1 per 1000, abbiamo che la varianza s_{11} aumenta di un fattore moltiplicativo pari a 1000^2 . Questo cambio di unità di misura da Kg a gr influenza la varianza generalizzata:*

$$\det(S^{gr}) = ((1000^2 s_{11}) s_{22} \cdots s_{pp}) \det(R) = 1000^2 \det(S^{Kg})$$

Per questo motivo, spesso è conveniente calcolare la varianza generalizzata considerando i dati standardizzati Z

2.5 Indice relativo di variabilità

$$0 \leq \frac{\text{Indice relativo di variabilità}}{\text{di variabilità}} = \det(R) = \frac{\det(S)}{\prod_{j=1}^p s_{jj}} \leq 1$$