

Analisi Esplorativa (Analisi Statistica Multivariata)

28 Febbraio 2022

Tempo a disposizione: 120 minuti

Modalità di consegna: svolgere gli esercizi di teoria (parte A) riportando le soluzioni sul foglio protocollo, e consegnare il foglio protocollo assieme al testo della prova d'esame. Successivamente, accedere alla piattaforma esaminonline tramite computer e svolgere gli esercizi di analisi dei dati (parte B). In questo caso la consegna si svolge tramite piattaforma esaminonline. Il tempo da dedicare alla parte A e alla parte B è a discrezione dello studente.

Compilare con nome, cognome e numero di matricola.

NOME:

COGNOME:

MATRICOLA:

PARTE A: esercizi di teoria

Esercizi da svolgere sul foglio protocollo senza l'ausilio di R/Rstudio.

Problema 1

- Si supponga che la matrice dei dati standardizzati $Z_{n \times 2}$ consista di due colonne $z_1_{n \times 1}$ e $z_2_{n \times 1}$ con la seguente matrice di correlazione $R = \begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix}$ dove $-1 \leq r \leq 1$. Si consideri la matrice dei dati trasformati $Y_{n \times 2}$ dove $y_1_{n \times 1} = z_1_{n \times 1}$ e $y_2_{n \times 1} = b z_2_{n \times 1}$ per una costante $b > 0$.
 - Calcolare la varianza totale per i dati Y .
 - Calcolare la varianza generalizzata per i dati Y .
 - Calcolare l'indice relativo di variabilità per i dati Y .
- Si consideri il modello fattoriale con le seguenti matrici di pesi fattoriali e varianze specifiche:

$$\Lambda = \begin{bmatrix} 4 & 1 \\ 7 & 2 \\ -1 & 6 \\ 1 & 8 \end{bmatrix}, \quad \Psi = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix}$$

- Calcolare la covarianza tra la prima variabile e il primo fattore comune, i.e. $\text{Cov}(x_1, f_1)$
- Calcolare la varianza della prima variabile, i.e. $\text{Var}(x_1)$
- Calcolare la covarianza tra le prime due variabili, i.e. $\text{Cov}(x_1, x_2)$
- Calcolare la correlazione tra le prime due variabili, i.e. $\text{Corr}(x_1, x_2)$

3. Si consideri la seguente matrice dei dati $X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$.

- Calcolare la matrice $D = \{d^2(u_i, u_l)\}$ delle distanze Euclidee al quadrato.
 - Calcolare la distanza totale $T = \frac{1}{2n} \sum_{i=1}^n \sum_{l=1}^n d^2(u_i, u_l)$.
 - Si consideri l'algoritmo delle K -medie con $K = 2$. Calcolare la distanza entro i gruppi $W = \sum_{k=1}^K W(G_k)$ dove $W(G_k) = \frac{1}{2n_k} \sum_{i: u_i \in G_k} \sum_{l: u_l \in G_k} d^2(u_i, u_l)$ per la partizione $G_1 = \{u_1\}$ e $G_2 = \{u_2, u_3\}$.
-

PARTE B: esercizi di analisi dei dati

Esercizi da svolgere con il computer sulla piattaforma esamionline con l'ausilio di R/Rstudio.

Problema 2

Si consideri la matrice dei dati $X_{n \times 3}$, a cui è associata la matrice di varianze/covarianze S . Sapendo che i) la prima riga di X è pari a $u'_1 = (20.6, 87, 77)$ e il baricentro è pari a $\bar{x}' = (13.9, 76.2, 32.9)$; ii) la seconda colonna di V è $v_2 = (0.1, -1, 0.2)'$, dove V è la matrice degli autovettori relativa a S ; la varianza della prima variabile è $s_{11} = 9.8$, mentre il secondo autovalore di S è $\lambda_2 = 26.5$. (nel vostro esercizio il valori potrebbero essere diversi).

- Calcolare $d_\infty(u_1, \bar{x})$. Riportare il risultato arrotondando al **terzo decimale** (si ricordi l'uso della **virgola** per i decimali).
- Calcolare il valore del primo elemento di $y_2 = \tilde{X}v_2$, dove \tilde{X} è la matrice dei dati centrati e y_2 è il vettore dei punteggi della seconda componente principale. Riportare il risultato arrotondando al **terzo decimale** (si ricordi l'uso della **virgola** per i decimali).
- Calcolare la correlazione tra la prima colonna \tilde{x}_1 di \tilde{X} e i punteggi y_2 della seconda componente principale. Riportare il risultato arrotondando al **terzo decimale** (si ricordi l'uso della **virgola** per i decimali).

Problema 3

Si consideri il dataset **quakes** presente nella libreria **datasets**. Si tratta di 1000 osservazioni misurate su 5 variabili: **lat** Latitude of event; **long** Longitude; **depth** Depth (km); **mag** Richter Magnitude; **stations** Number of stations reporting. Sia X la matrice 999×4 corrispondente al dataset **quakes**, escludendo la riga 1 (nel vostro esercizio il numero di riga potrebbe essere diverso) e la variabile **stations**.

- Calcolare la matrice dei dati standardizzati Z e riportare il numero di osservazioni anomale verificando se la distanza di Mahalanobis al quadrato di ciascuna osservazione di Z dal baricentro 0 è superiore alla soglia s , dove s corrisponde al quantile 0.95 di una variabile casuale χ_p^2 (dove p è il numero di colonne di Z).
- Rimuovere da Z le osservazioni anomale individuate al punto precedente e su questi dati utilizzare l'algoritmo delle K -medie i) inizializzando i centroidi con le prime K unità statistiche; ii) impostando il numero di iterazioni (argomento **iter.max**) a 100; iii) utilizzando l'algoritmo di Lloyd (argomento **algorithm**). Decidere il numero di gruppi K ottimale secondo l'indice CH, per $K = 2, \dots, 12$. Riportare la distanza entro i gruppi W corrispondente alla scelta effettuata, arrotondando al **terzo decimale** (si ricordi l'uso della **virgola** per i decimali).
- Calcolare il valore medio della *silhouette* corrispondente alla raggruppamento selezionato al punto precedente, utilizzando la matrice delle distanze Euclidee. Riportare il risultato arrotondando al **terzo decimale** (si ricordi l'uso della virgola per i decimali).