

METODI STATISTICI PER LA BIOINGEGNERIA (B)

PARTE 18: PRINCIPAL COMPONENT ANALYSIS (PCA)

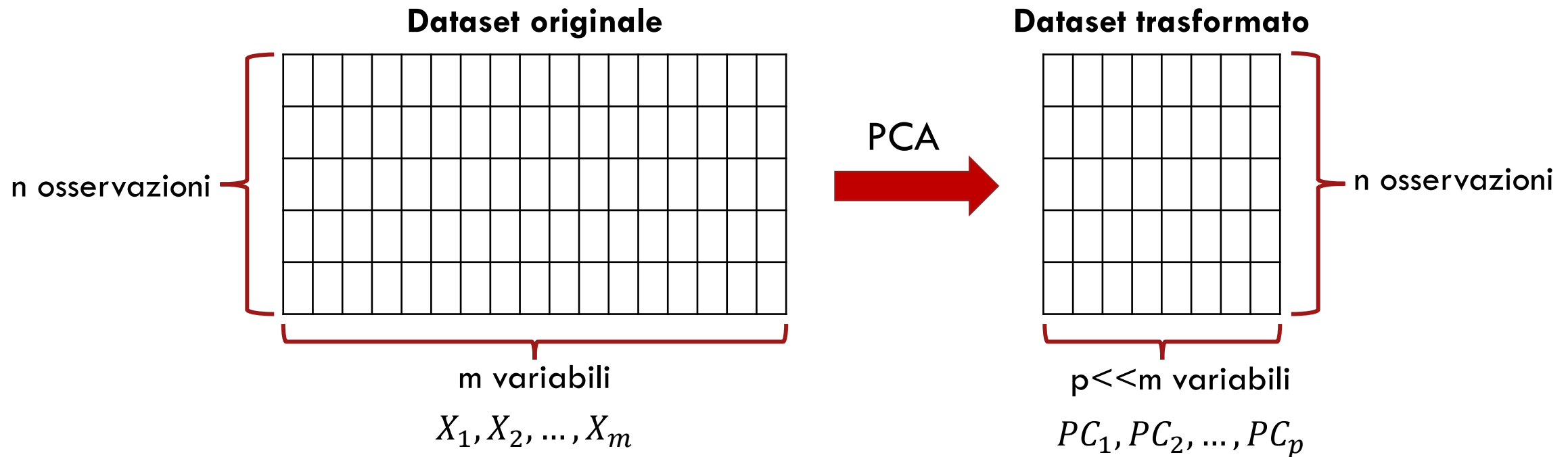
A.A. 2024-2025

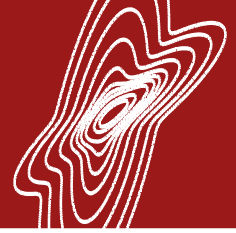
Prof. Martina Vettoretti

PRINCIPAL COMPONENT ANALYSIS (PCA)



- **Principal Component Analysis (PCA)**, o analisi delle componenti principali: tecnica di apprendimento non supervisionato per **ridurre la dimensionalità dei dati**.
- Obiettivo: rappresentare un set di dati di dimensione $n \times m$ (n osservazioni, m variabili) in uno spazio a dimensionalità ridotta con $p \ll m$ variabili.

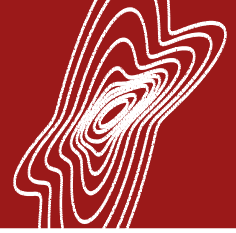




LE COMPONENTI PRINCIPALI



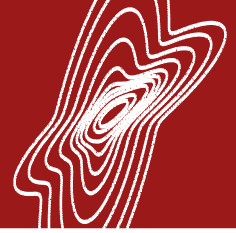
- Le p nuove variabili sono dette **componenti principali** (o *principal components*).
- Caratteristiche delle componenti principali:
 - Le componenti principali sono **nuove variabili**, create artificialmente, nessuna di loro coincide con alcuna delle m variabili di partenza.
 - Ogni componente principale è una **combinazione lineare** delle m variabili originali.
 - Le componenti principali sono tali da **riassumere quanta più informazione possibile** sulle m variabili originali.
 - Le componenti principali sono ordinate in base a quanta informazione del dataset originale racchiudono (PC_1 è la componente più informativa, PC_2 la seconda più informativa ecc.)
 - Le componenti principali per costruzione sono **tra loro scorrelate**.



PERCHE' PUO' ESSERE UTILE LA PCA?



- Visualizzazione dei dati
- Compressione dei dati
- Eliminare la correlazione tra le variabili di ingresso in un modello di regressione lineare multipla (o di altro tipo)

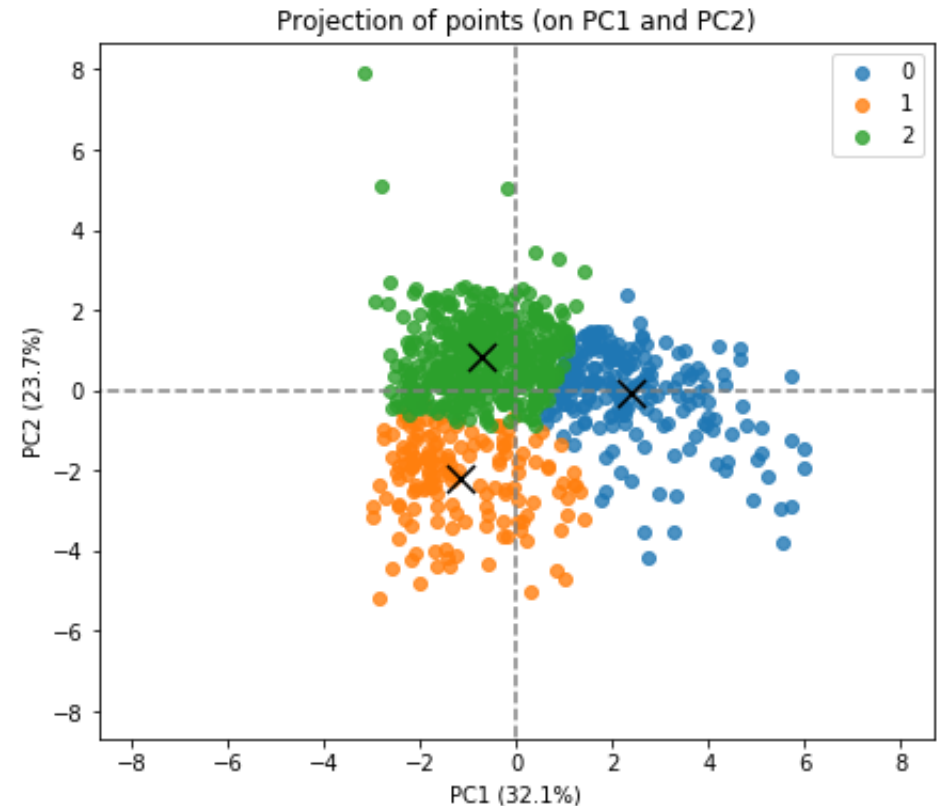


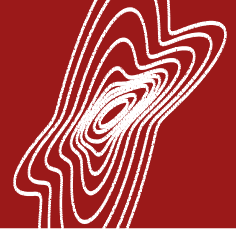
- Quando abbiamo dataset con tante variabili (m grande) diventa difficile rappresentare graficamente i dati per analizzarli dal punto di vista visivo.
 - Tipicamente si visualizzano le distribuzioni delle singole variabili o al più gli scatterplot di coppie di variabili.
 - Problemi:
 - Lo scatterplot di due variabili rappresenta solo una piccola quantità dell'informazione contenuta nei dati.
 - Con m variabili dovremmo fare $m \cdot (m - 1)/2$ scatterplot. Se $m = 10 \rightarrow 10 \times 9/2 = 45$ scatterplot!
- Possiamo sfruttare la PCA per riassumere l'informazione contenuta nei dati usando poche variabili, più semplici da rappresentare.
 - Potremmo realizzare lo scatterplot delle prime 2 componenti principali, PC_1 e PC_2 , ovvero quelle più informative!

ESEMPIO



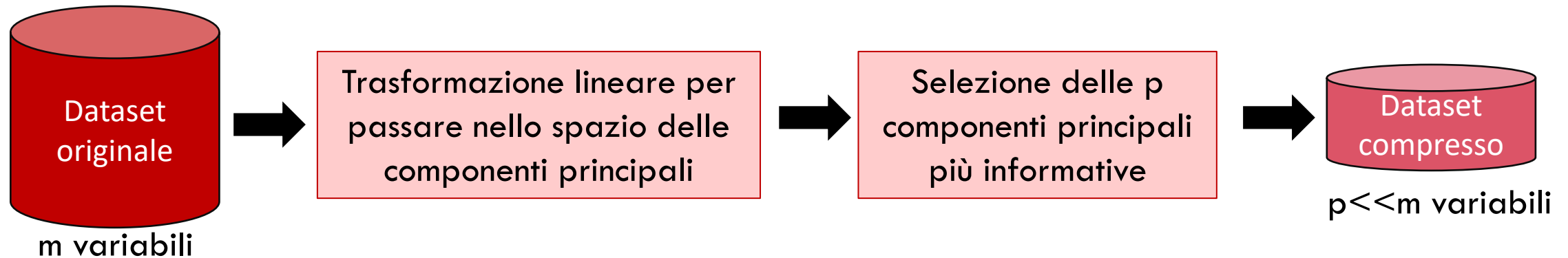
- Abbiamo un dataset con 50 variabili e realizziamo un clustering K-means per suddividere le osservazioni in 3 gruppi. Vogliamo verificare dal punto di vista visivo quanto sono separati i 3 cluster.
- Impossibile rappresentare i dati nello spazio a 50 dimensioni del dataset originale!
- Soluzione: applichiamo la PCA e rappresentiamo i cluster nello spazio bidimensionale delle prime 2 componenti principali, le più informative.



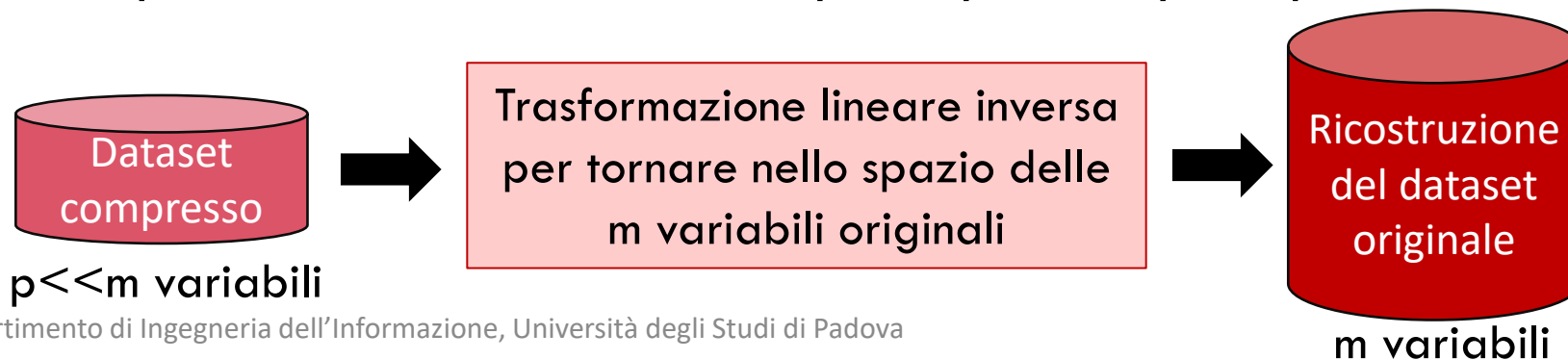


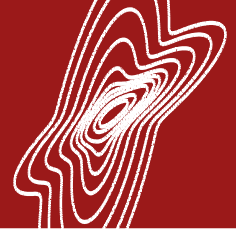
COMPRESSIONE DEI DATI

- **Compressione:** al posto di archiviare le m variabili di partenza, archivio le p componenti principali, con un risparmio in memoria.



- **Decompressione:** utilizzando le p componenti principali ricostruisco le m variabili di partenza. La ricostruzione non sarà perfetta, l'errore introdotto dalla compressione dipende da quanto informative erano le p componenti principali selezionate.





ESEMPIO



- Compressione di un'immagine di dimensione 256×256 (matrice di dati 256×256) mediante PCA.

INPUT IMAGE



RESTORED IMAGE-10 COMPONENTS



RESTORED IMAGE-50 COMPONENTS



RESTORED IMAGE-100 COMPONENTS



RESTORED IMAGE-150 COMPONENTS



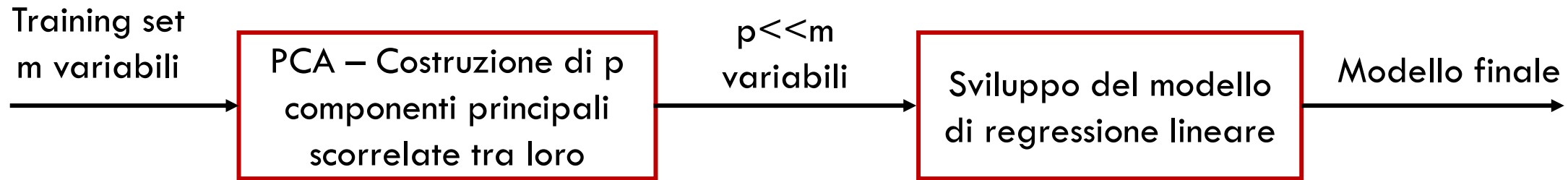
RESTORED IMAGE-256 COMPONENTS



ELIMINARE LA CORRELAZIONE TRA LE VARIABILI IN INGRESSO IN UN MODELLO DI REGRESSIONE LINEARE MULTIPLA



- La multicollinearità tra le variabili di ingresso è problematica per i modelli di regressione lineare multipla → possiamo eliminare la multicollinearità mediante PCA.
- Idea: allenare il modello sulle p componenti principali anziché sulle m variabili originali.

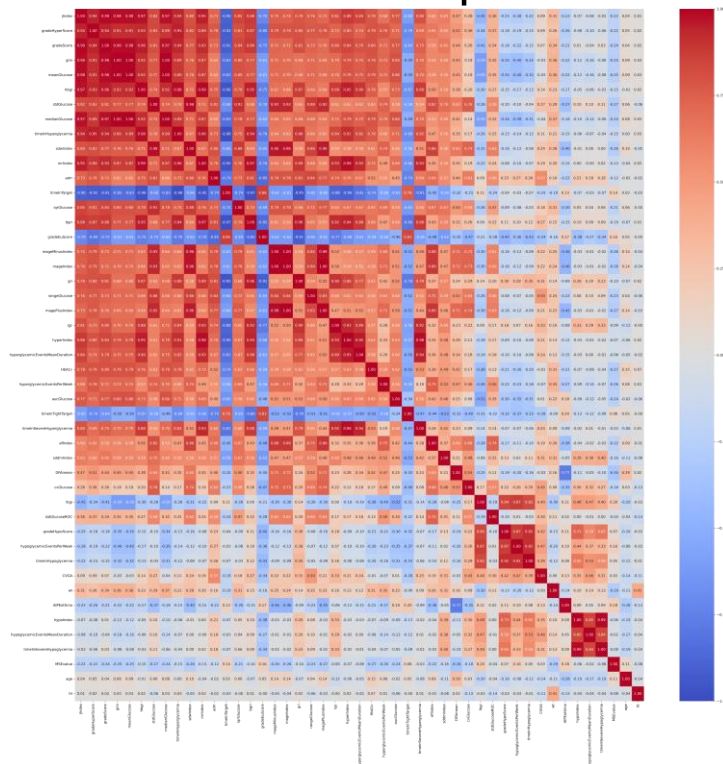


- Attenzione: le p componenti principali sono nuove variabili che non hanno un significato fisico riconducibile alle m variabili di partenza. → Stimando i coefficienti del modello di regressione lineare sulle componenti principali, perdiamo la possibilità di interpretare i coefficienti del modello.

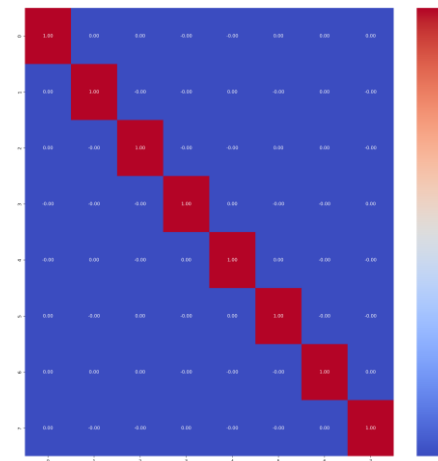
ESEMPIO

- Modello di regressione logistica per la classificazione dei pazienti diabetici vs non diabetici sulla base di indici di variabilità glicemica.

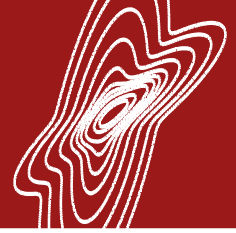
Matrice di correlazione delle
47 variabili di partenza



Matrice di correlazione delle 8
prime componenti principali



Modello di regressione logistica
allenato sulle 8 componenti
principali anziché sulle 47
variabili originali



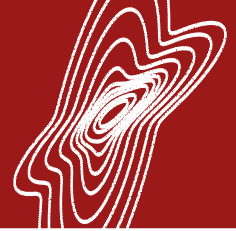
PCA, COME SI REALIZZA?



- Dataset originale: m variabili, X_1, X_2, \dots, X_m , a media nulla.
- Nota: se le variabili non sono a media nulla, prima di realizzare la PCA, i dati vanno centrati \rightarrow ad ogni variabile va sottratta la sua media.
- **Trasformazione lineare normalizzata** delle m variabili \rightarrow m nuove variabili, PC_1, PC_2, \dots, PC_m , dette componenti principali tra loro scorrelate

$$\begin{aligned} PC_1 &= v_{11}X_1 + v_{21}X_2 + \dots + v_{m1}X_m \\ PC_2 &= v_{12}X_1 + v_{22}X_2 + \dots + v_{m2}X_m \\ &\dots \\ PC_m &= v_{1m}X_1 + v_{2m}X_2 + \dots + v_{mm}X_m \end{aligned}$$

- I coefficienti v_{ik} consentono di trasformare le variabili di partenza nelle nuove variabili, le componenti principali.



I LOADINGS

- I coefficienti v_{ik} , $i = 1, \dots, m$ relativi alla componente principale k-esima, PC_k , si dicono **loadings** relativi alla componente k-esima:

$$\mathbf{v}_k = (v_{1k}, v_{2k}, \dots, v_{mk})^T$$

essi rappresentano il **contributo di ciascuna delle variabili originali** alla componente principale k-esima.

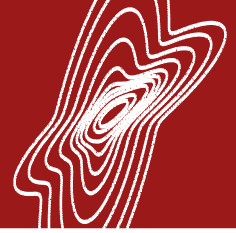
- I loadings di ciascuna componente principale hanno **norma 1**:

$$\sum_{i=1}^m v_{ik}^2 = 1 \quad \forall k$$

Per questo la trasformazione lineare effettuata dalla PCA si dice normalizzata.

- I vettori di loadings delle diverse componenti principali sono tra loro **ortogonali**:

$$\mathbf{v}_i^T \cdot \mathbf{v}_j = 0 \quad (\text{prodotto scalare})$$



VARIANZA DELLE COMPONENTI PRINCIPALI

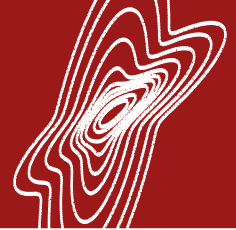


- I loadings devono essere tali che **le prime componenti principali riassumano quanta più varianza possibile** delle variabili di partenza:

$$\text{var}(PC_1) > \text{var}(PC_2) > \text{var}(PC_3) > \dots > \text{var}(PC_m)$$

- PC_1 da sola deve essere in grado di spiegare quanta più varianza possibile dei dati di partenza.
- PC_2 da sola deve essere in grado di spiegare quanta più varianza possibile della porzione di varianza non spiegata da PC_1 .
- ...

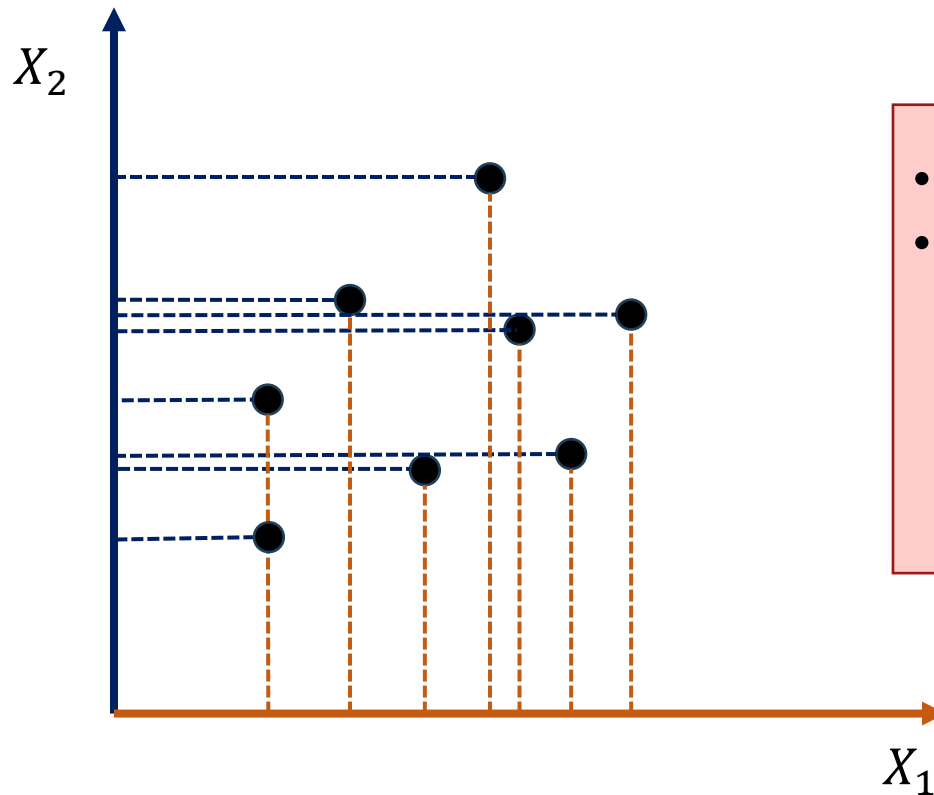
- **Considerando solo le prime p componenti principali** possiamo efficacemente ridurre la dimensionalità dei dati \rightarrow nuovo set di $p \ll m$ variabili che riassume la maggior parte della varianza delle m variabili di partenza.



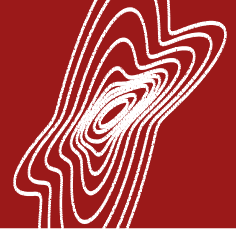
INTERPRETAZIONE GEOMETRICA



- In pratica la PCA effettua una **proiezione ortogonale** dei dati su uno spazio definito da nuove dimensioni dette componenti principali. Queste sono tali per cui la varianza delle coordinate dei dati proiettati sulle nuove dimensioni è massima per le prime dimensioni.



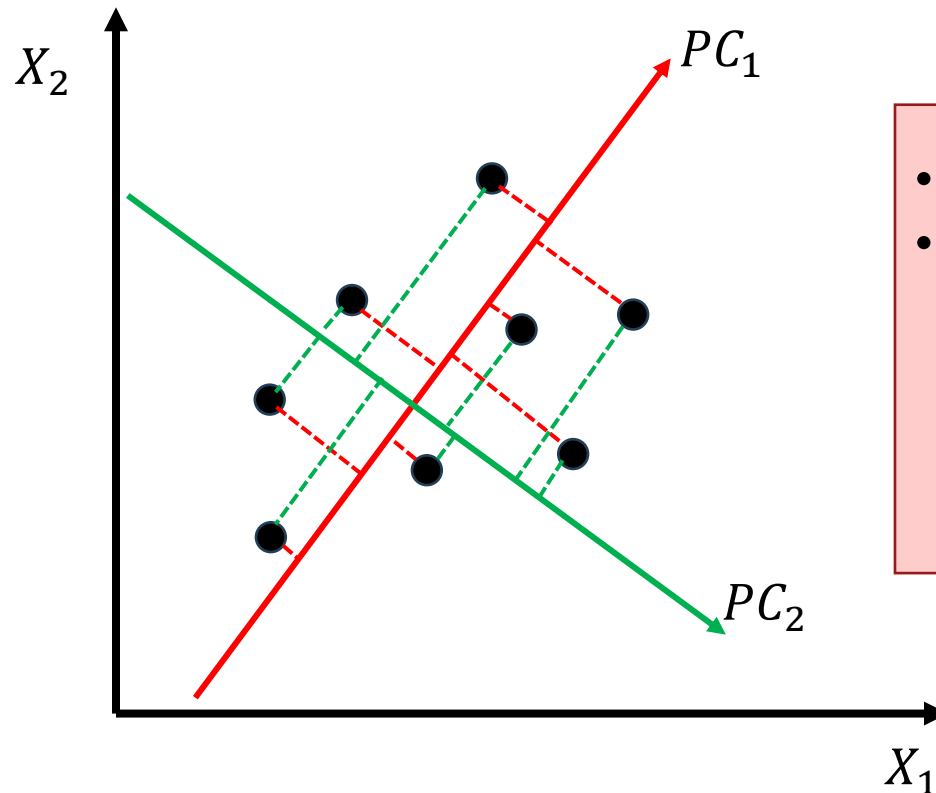
- Ogni osservazione è un punto.
- Le linee tratteggiate rappresentano le coordinate dei punti nello spazio definito dalle variabili di partenza X_1 e X_2 .



INTERPRETAZIONE GEOMETRICA



- In pratica la PCA effettua una **proiezione ortogonale** dei dati su uno spazio definito da nuove dimensioni dette componenti principali. Queste sono tali per cui la varianza delle coordinate dei dati proiettati sulle nuove dimensioni è massima per le prime dimensioni.



- Ogni osservazione è un punto.
- Le linee tratteggiate rappresentano le coordinate dei punti nello spazio definito dalle componenti principali PC_1 e PC_2 .

PASSARE DALLO SPAZIO ORIGINALE ALLO SPAZIO DELLE COMPONENTI PRINCIPALI



- Per passare dalle coordinate dei punti nello spazio originale a quelle nello spazio definito dalle componenti principali, basta applicare la trasformazione lineare definita dai loadings.
- Dataset originale: n osservazioni \times m variabili (a media nulla).
 - $x_i \rightarrow$ vettore colonna contenente le n osservazioni per la variabile X_i
 - X matrice dei dati originali (n righe, m colonne).

$$X = [x_1 \ x_2 \ \dots \ x_m] = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix}$$

- La media campionaria di ciascuna colonna di X deve essere 0.

COORDINATE DEI DATI NEL SISTEMA DELLE COMPONENTI PRINCIPALI



- Calcoliamo le nuove coordinate dei dati riferite al sistema di riferimento delle componenti principali, PC_1, PC_2, \dots, PC_m :

$$\underbrace{\begin{bmatrix} y_{11} & y_{12} & \dots & y_{1m} \\ y_{21} & y_{22} & \dots & y_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nm} \end{bmatrix}}_{\mathbf{Y} \text{ (n x m)}} = \underbrace{\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix}}_{\mathbf{X} \text{ (n x m)}} \cdot \underbrace{\begin{bmatrix} v_{11} & v_{12} & \dots & v_{1m} \\ v_{21} & v_{22} & \dots & v_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ v_{m1} & v_{m2} & \dots & v_{mm} \end{bmatrix}}_{\mathbf{V} \text{ (m x m)}}$$

Matrice dei dati trasformati, detti scores.
La riga i-esima contiene le coordinate dell'osservazione i-esima nello spazio delle componenti principali.

Matrice dei dati originali.
La riga i-esima contiene le coordinate dell'osservazione i-esima nello spazio delle variabili originali.

Matrice dei loadings, i coefficienti che definiscono la combinazione lineare. La colonna k-esima rappresenta i loadings della componente principale k-esima.

$$Y = X \cdot V$$

Dati trasformati
o scores

Dati originali

Coefficienti o loadings

- La matrice **V** è di fatto una **matrice di rotazione** che consente di passare dalle coordinate nel sistema di riferimento originale, alle coordinate nel sistema di riferimento delle componenti principali.



- Poiché le colonne di \mathbf{V} sono tra loro ortogonali e hanno norma 1, si ha che:

$$\mathbf{V} \cdot \mathbf{V}^T = \mathbf{V}^T \cdot \mathbf{V} = \mathbf{I}$$

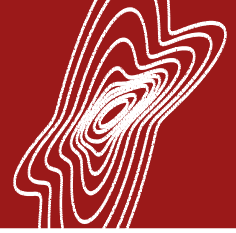
- Calcolo delle coordinate nello spazio delle variabili originali partendo dalle coordinate nello spazio delle componenti principali:

$$\mathbf{Y} = \mathbf{X} \cdot \mathbf{V}$$

$$\mathbf{Y} \cdot \mathbf{V}^T = \mathbf{X} \cdot \mathbf{V} \cdot \mathbf{V}^T$$



$$\mathbf{X} = \mathbf{Y} \cdot \mathbf{V}^T$$



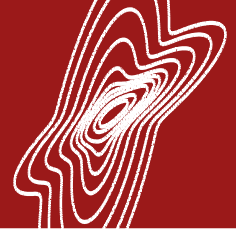
RIDUZIONE DELLA DIMENSIONALITA' (1 / 2)



- La PCA ci fornisce dunque una nuova rappresentazione dei dati nello spazio delle componenti principali, PC_1, PC_2, \dots, PC_m , che ha 2 caratteristiche fondamentali:
 - Le componenti principali sono scorrelate
 - La varianza dei dati proiettati lungo le componenti principali decresce all'aumentare delle componenti, la maggior parte della varianza complessiva è concentrata nelle prime componenti principali

$$\text{var}(PC_1) > \text{var}(PC_2) > \text{var}(PC_3) > \dots > \text{var}(PC_m)$$

- Per ridurre la dimensionalità possiamo considerare solo le prime p componenti principali.



RIDUZIONE DELLA DIMENSIONALITA' (2/2)

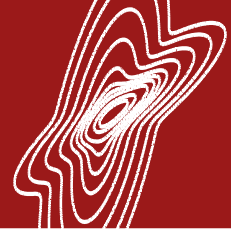


- Invece che memorizzare Y , memorizziamo solo le prime p colonne di Y .
- p viene scelto in modo da garantire che le prime p colonne di Y contengano la maggior parte della varianza dei dati originali.
- Se la PCA viene utilizzata per realizzare una compressione dei dati, utilizzando la matrice V^T possiamo ricostruire, con un certo errore dovuto alla compressione, i dati originali X .

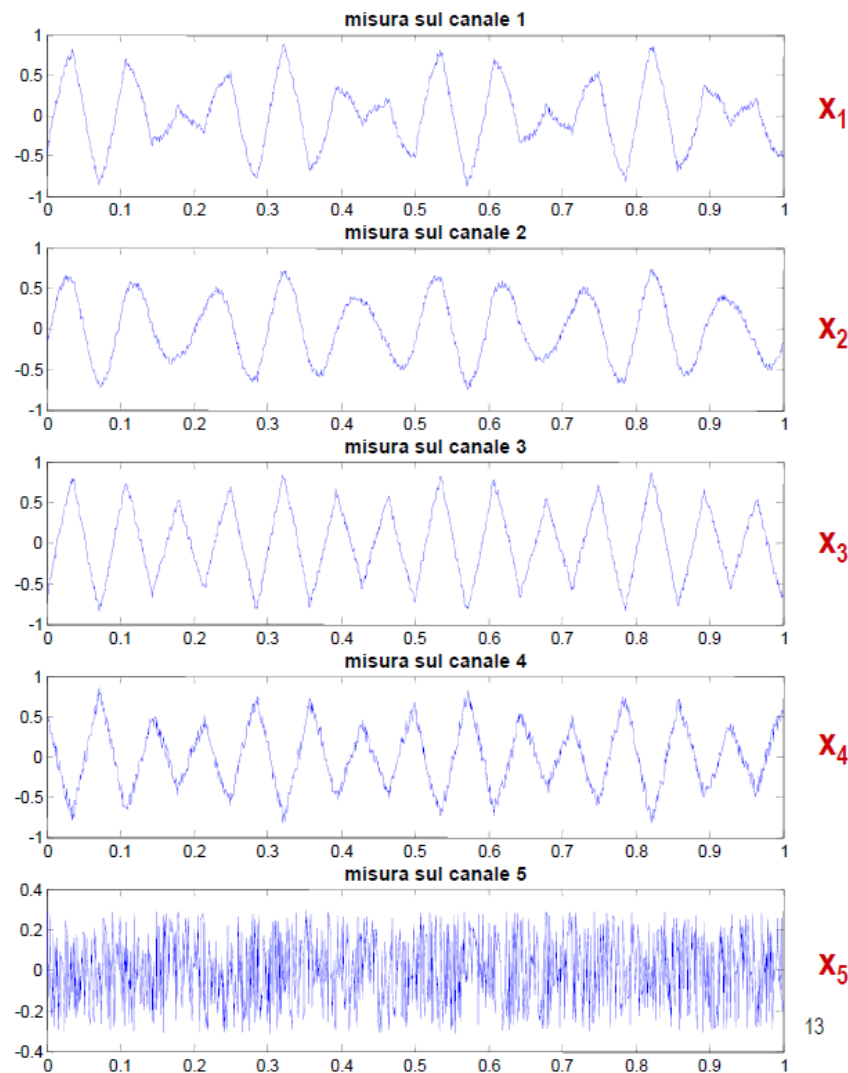
$$\tilde{X} = Y_p V^T$$

Ricostruzione
dei dati in X

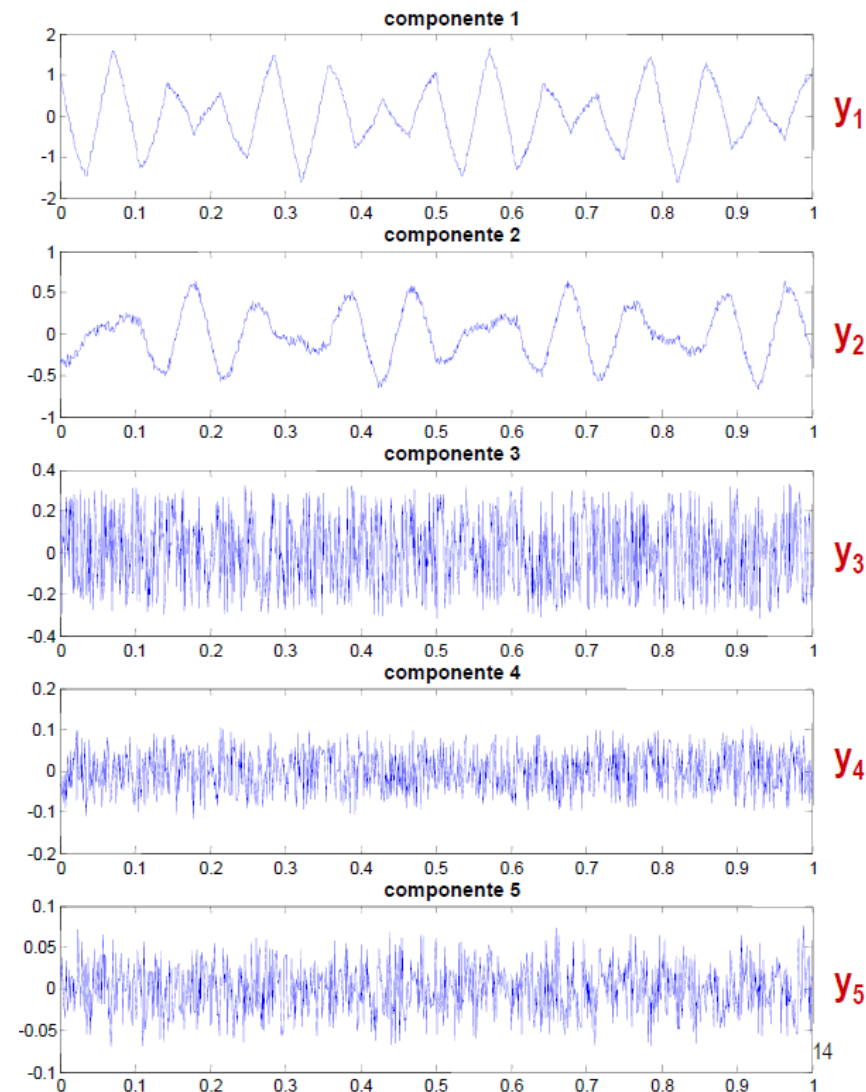
Prime p
colonne di Y

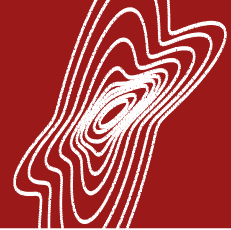


ESEMPIO ($m=5$, $n=1000$)



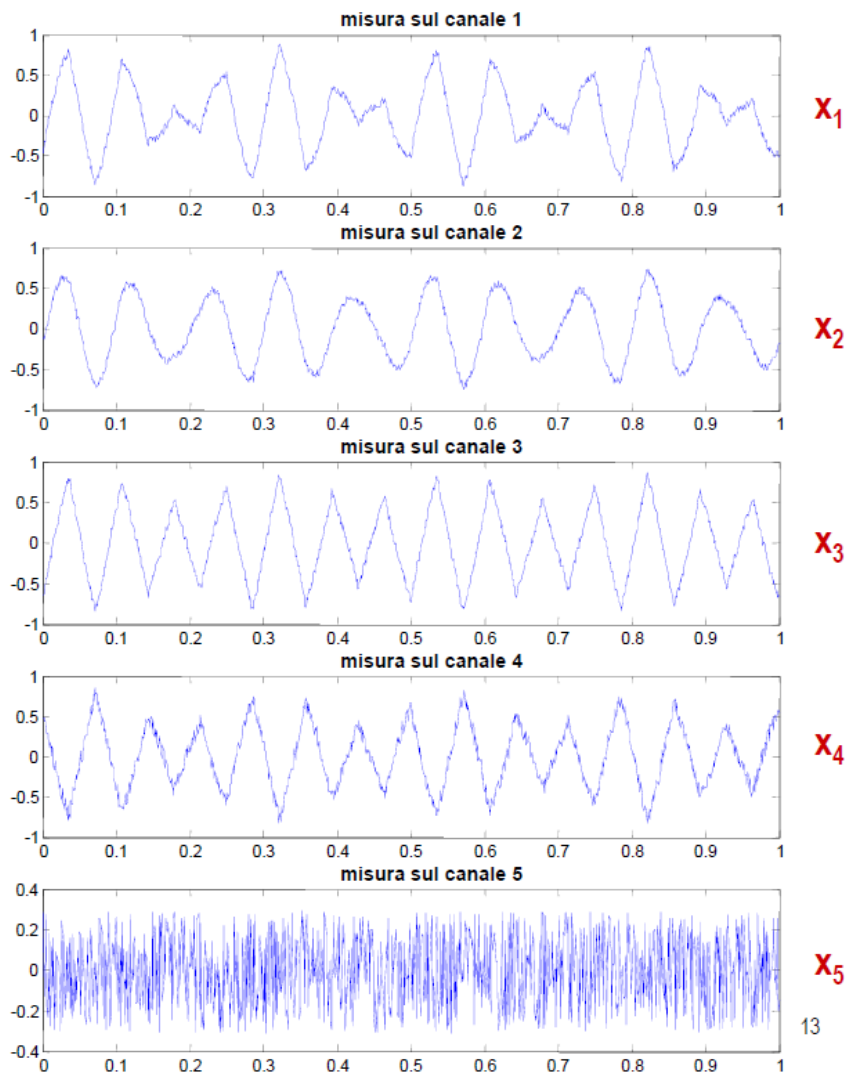
PCA
➔



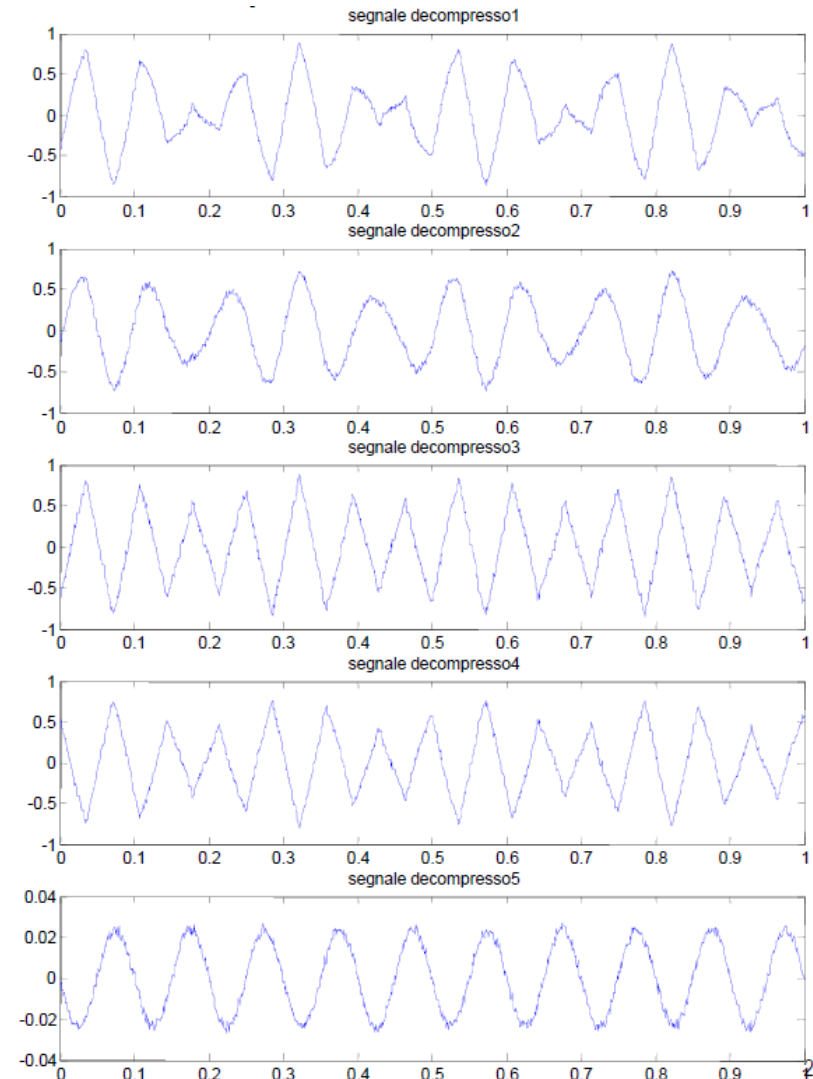


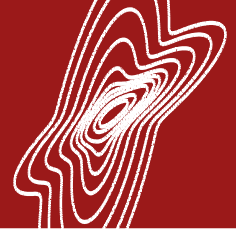
ESEMPIO ($m=5$, $n=1000$)

Segnali originali



Segnali ricostruiti usando solo le prime 2 componenti principali

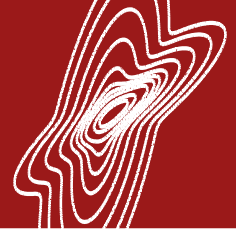




COME SI CALCOLANO I LOADINGS?



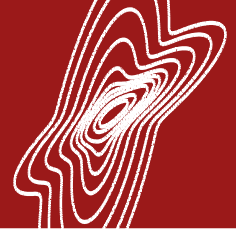
- Come possiamo calcolare i valori dei loadings, ovvero la matrice \mathbf{V} , che mi consente di realizzare la PCA?
- Le colonne della matrice \mathbf{V} sono gli autovettori della matrice di covarianza di \mathbf{X} ordinati secondo l'ordine decrescente dei rispettivi autovalori.



CALCOLO DEI LOADINGS



- Calcolo della **matrice di covarianza** di \mathbf{X} : \mathbf{S}
- \mathbf{S} è una matrice $m \times m$, reale, e simmetrica, avente sulla diagonale le varianze delle colonne di \mathbf{X} (le variabili), fuori dalla diagonale le covarianze campionarie tra coppie di colonne di \mathbf{X} .
 - Elemento in posizione i,j : covarianza tra la colonna i -esima e la colonna j -esima di \mathbf{X}
- Calcoliamo **autovalori** e **autovettori** di \mathbf{S} .



AUTOVALORI E AUTOVETTORI



Sia $A \in \mathbb{R}^{N \times N}$ una matrice quadrata di N righe e N colonne.

Se esistono un vettore $v \in \mathbb{R}^N$ e uno scalare λ (anche complesso) tali che:

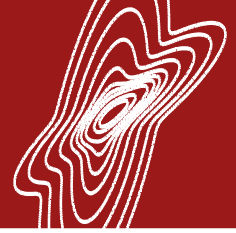
$$Av = \lambda v$$

si dice che v è **autovettore** di A e λ il suo **autovalore** corrispondente.

Proprietà:

- Una matrice di dimensione $N \times N$ ha al massimo N autovalori distinti (reali o complessi).
- Gli autovalori sono gli zeri del polinomio caratteristico:

$$\det(A - \lambda I) = 0$$

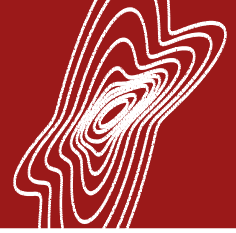


AUTOVALORI E AUTOVETTORI



- Se la matrice **A** è simmetrica ($A = A^T$) e reale, come la matrice **S**:
 - Gli autovalori sono tutti reali
 - Gli autovettori corrispondenti agli autovalori distinti sono tra loro ortogonali:

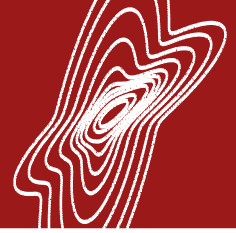
$$v_i^T \cdot v_j = 0 \quad (\text{prodotto scalare})$$



CALCOLO DEI LOADINGS



- La matrice di covarianza S , di dimensione $m \times m$, è reale e simmetrica.
 - Gli autovalori sono reali
 - Gli autovettori distinti sono tra loro sono ortogonali
- **Gli autovettori di S rappresentano i loadings** che definiscono le componenti principali.
- Le componenti principali sono ortogonali tra loro → quindi scorrelate.
- In che ordine vengono considerati gli autovettori per definire le componenti principali? → In base agli autovalori corrispondenti.



DEFINIZIONE DELLA MATRICE V

- Ordiniamo i gli autovalori dal più grande al più piccolo:

Autovalori: $\lambda_1 > \lambda_2 > \dots > \lambda_m$



Autovettori: $v_1 \quad v_2 \quad \dots \quad v_m$

- L'autovettore v_1 corrispondente all'autovalore massimo, λ_1 , rappresenta il vettore dei loadings della prima componente principale, ovvero la prima colonna di V .
- Gli autovettori, v_1, v_2, \dots, v_m , in questo ordine definiscono le colonne di V :

$$V = [v_1 \quad v_2 \quad \dots \quad v_m]$$



➤ Gli **autovalori** rappresentano la varianza dei dati proiettati lungo le componenti principali.

- λ_1 rappresenta la varianza campionaria degli score relativi a PC_1 , ovvero la varianza della prima colonna di Y
- λ_2 rappresenta la varianza campionaria degli score relativi a PC_2 , ovvero la varianza della seconda colonna di Y
- ...

$$Y = [y_1 \ y_2 \ \dots \ y_m]$$

$$\lambda_k = s_{y_k}^2, \quad k = 1, \dots, m$$



Varianza campionaria
di y_k

- Varianza complessiva dei dati originali:

$$\sum_{k=1}^m s_{x_k}^2 = \sum_{k=1}^m \lambda_k$$

Varianza campionaria della
colonna k-esima di \mathbf{X}

- Frazione della varianza complessiva spiegata dalla componente k-esima:

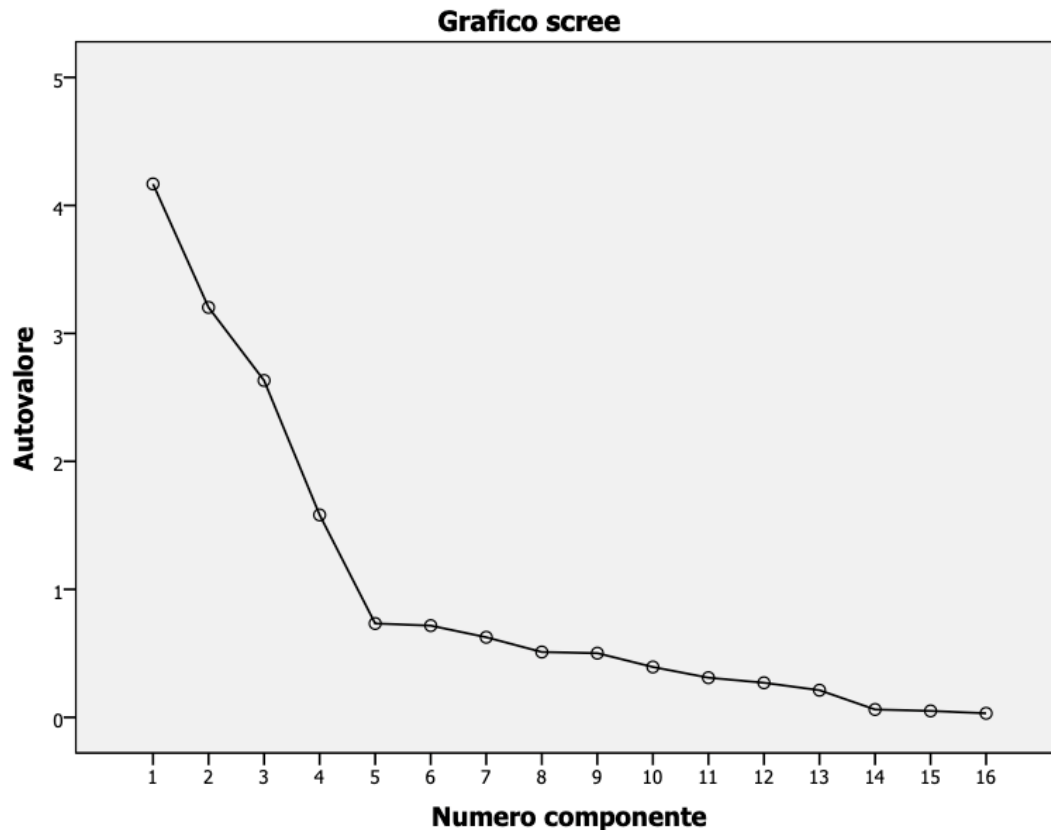
$$\frac{\lambda_k}{\sum_{k=1}^m \lambda_k}$$



SCELTA DEL NUMERO DI COMPONENTI PRINCIPALI



- Rappresentiamo graficamente gli autovalori (*scree plot*) e scegliamo il numero di componenti che corrisponde al punto di gomito.



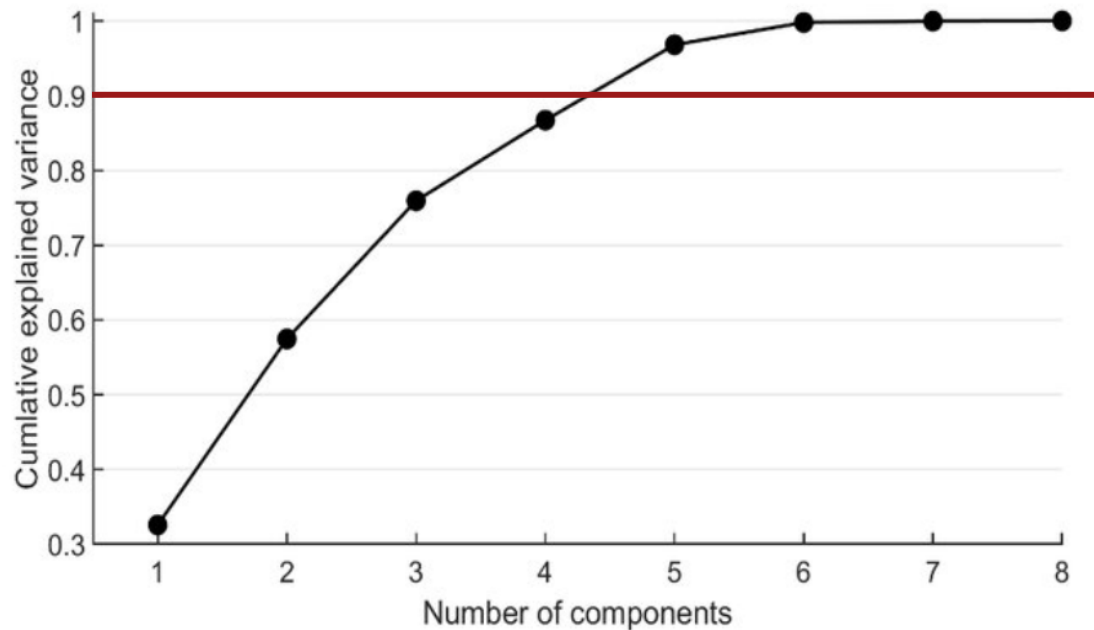
Scegliamo 5 componenti, perché dalla sesta componente in poi la varianza spiegata dalle singole componenti è parecchio limitata rispetto a quella spiegata dalle prime 5 componenti.



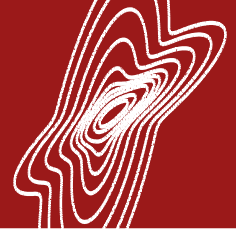
SCELTA DEL NUMERO DI COMPONENTI PRINCIPALI



- In alternativa possiamo rappresentare la frazione di varianza spiegata dalle prime p componenti principali, al variare di p .
- Scegliamo il valore di p che consente di spiegare almeno una certa percentuale (es. 90%) della varianza totale.



Per spiegare almeno il 90% della varianza totale servono 5 componenti.



PCA - RIASSUNTO



1. Centramento dei dati: ad ogni colonna di \mathbf{X} si sottrae la sua media.
2. Calcolo della matrice di covarianza di $\mathbf{X} \rightarrow \mathbf{S}$ ($m \times m$)
3. Calcolo di autovettori e autovalori di \mathbf{S} .
 - Gli autovettori rappresentano i coefficienti per definire le componenti principali.
 - L'ordine delle componenti principali è stabilito dall'ordine degli autovalori.
 - L'autovettore corrispondente all'autovalore massimo rappresenta i coefficienti (loadings) della prima componente principale.
4. Trasformazione dei dati passando alle coordinate nello spazio delle componenti principali: $\mathbf{Y} = \mathbf{X} \cdot \mathbf{V}$
5. Scelta delle prime p componenti principali che rappresentano gran parte della varianza delle variabili originali (scree plot o plot della frazione di varianza spiegata dalle prime p componenti).
6. Le coordinate delle n osservazioni lungo le p componenti principali rappresentano il dataset trasformato e ridotto.