

Corso di Studio SM13 – CHIMICA

Introduzione alla chemiometria ed al disegno sperimentale (018CM)

10/11/2022 (lezione 3)

Pierluigi Barbieri



**UNIVERSITÀ
DEGLI STUDI
DI TRIESTE**

Dipartimento di Scienze
Chimiche e Farmaceutiche

Gruppo di ricerca in Scienze Analitiche
applicate alle Interazioni Uomo-Ambiente
(ASHEI)



FONDAMENTA PER LA CHIMICA ANALITICA

Michele Forina

<http://www.sisnir.org/sisnir/download/fondamenta-per-la-chimica-analitica/category/3-fondamenta-per-la-chimica-analitica>

Visualizzazione di statistiche uni- e bi-variate.

Visualizzazione di set di dati reali con diverse tecniche di rappresentazione grafica, analisi dei risultati tramite esplorazione visiva, tecniche di individuazione di dati anomali.

Illustrazione del concetto di carta di controllo dei dati.

Analisi delle componenti principali e metodi fattoriali: aspetti teorici ed applicativi dell'analisi delle componenti principali e dei metodi fattoriali per la compressione ed interpretazione dell'informazione contenuta nei dati.

Visualizzazione di set di dati reali chimici, fisici e biologici



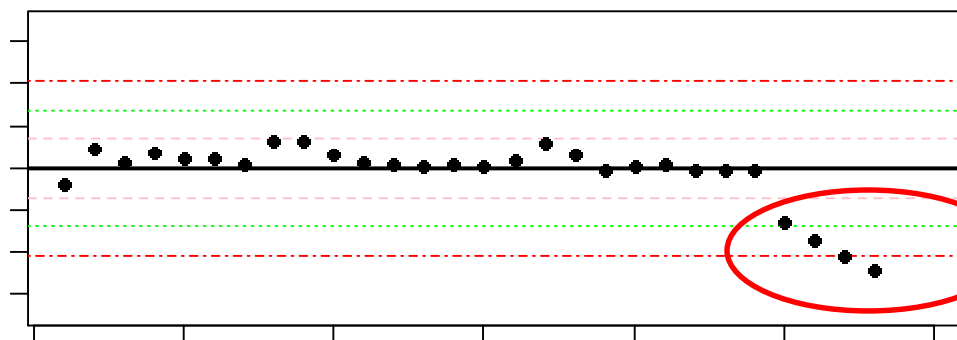
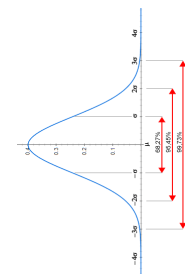
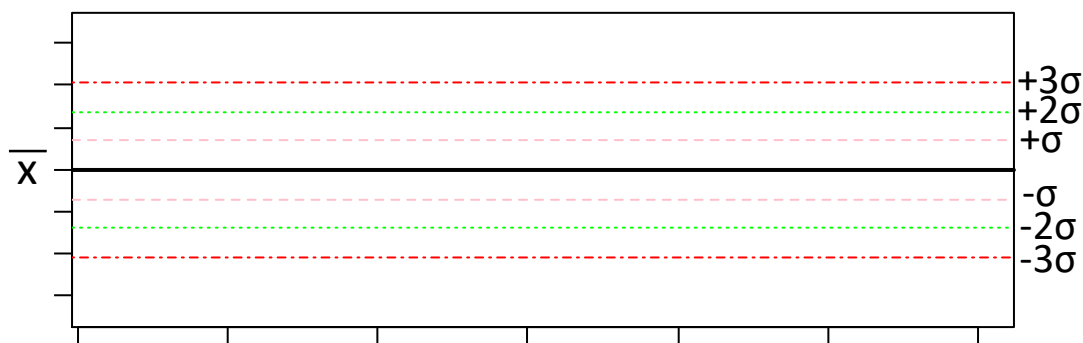
Mettiamo in rapporto visivo (su un grafico) i dati (anche multivariati) rilevati considerando valori di riferimento/norme

Il concetto di carta di controllo di una variabile

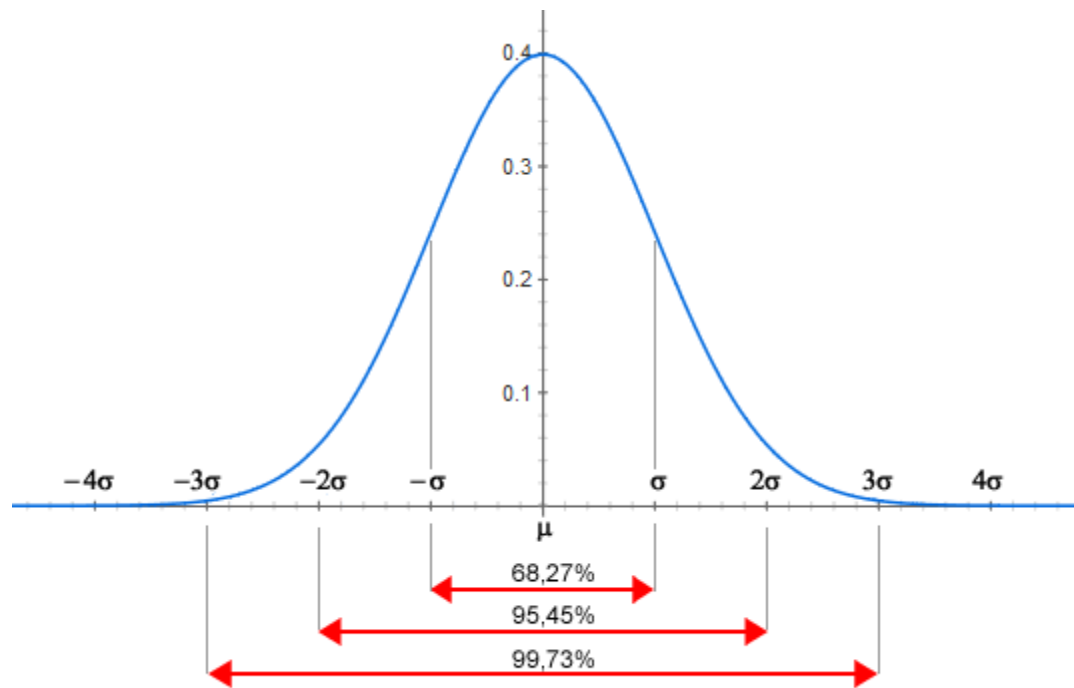
Una carta di controllo è uno strumento grafico per monitorare il/i parametri di un **processo** (anche *misure successive su campioni tratti da una popolazione definita*)

Carta di controllo di Shewhart

(utilizza la media e la deviazione standard per controllare eventuali derive nei valori di un parametro relativo ad un processo):

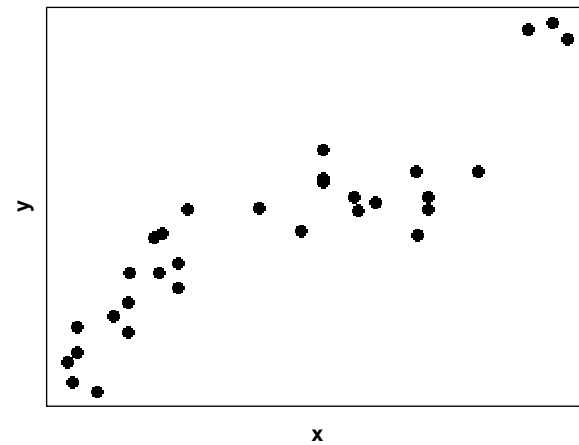
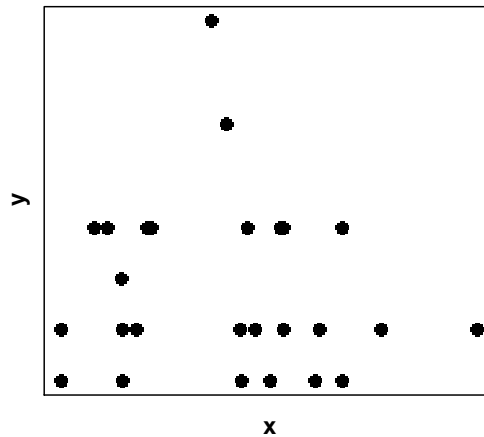
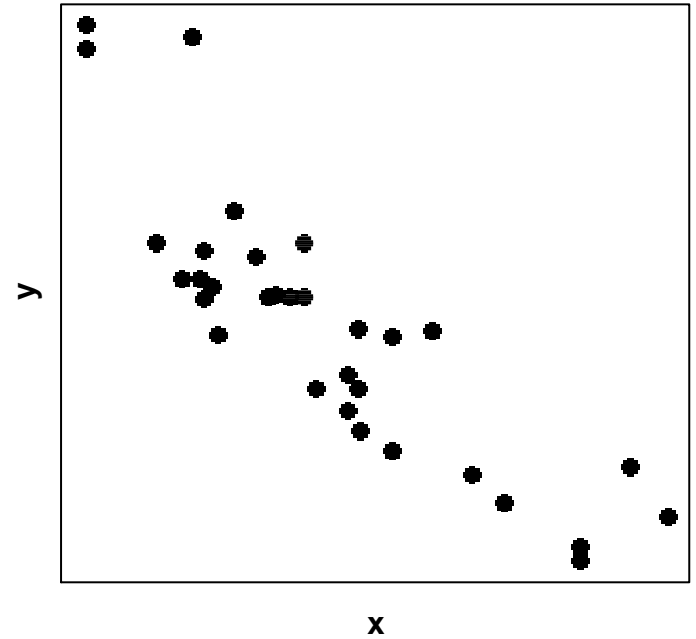
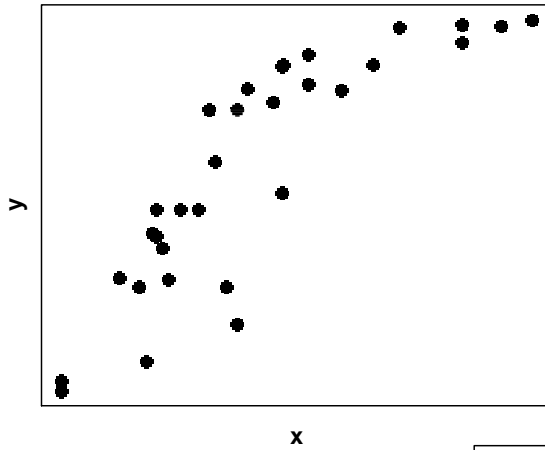


possibile deriva
del sistema



Analisi bivariata

E' finalizzata ad osservare il comportamento relativo di una variabile rispetto ad un'altra



Covarianza

E' un indicatore della variabilità "congiunta" di due variabili.

$$\sigma_{x,y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}$$

scarti di due variabili riferite alla medesima osservazione

Si può scrivere anche:

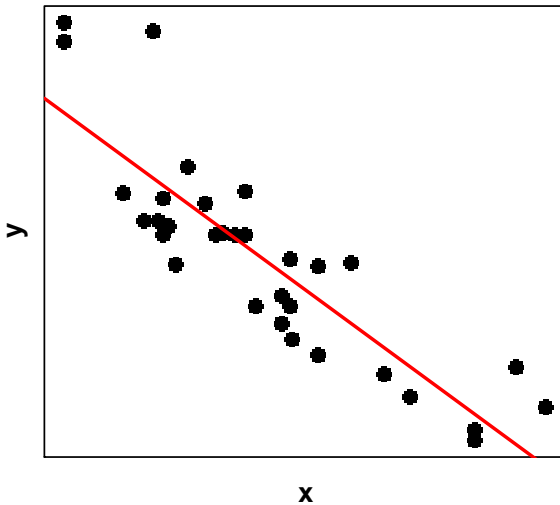
$$\sigma_{x,y} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \left(\frac{1}{n} \sum_{i=1}^n x_i \right) \left(\frac{1}{n} \sum_{i=1}^n y_i \right)$$

\bar{x} \bar{y}

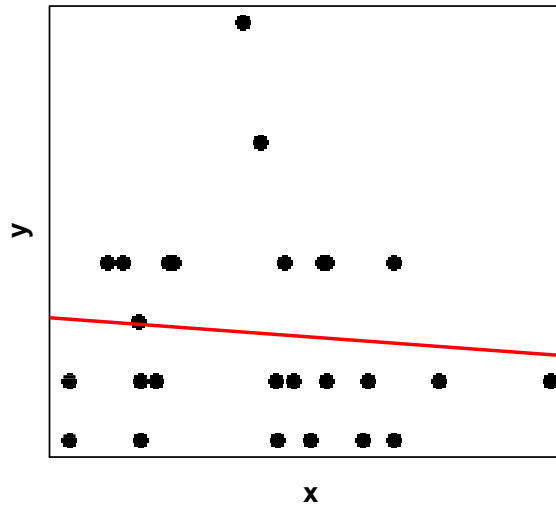
Coefficiente di correlazione di Pearson

E' un indicatore del livello di correlazione lineare tra **DUE** variabili, assume valori che vanno da -1 a 1, si indica con $r_{x,y}$.

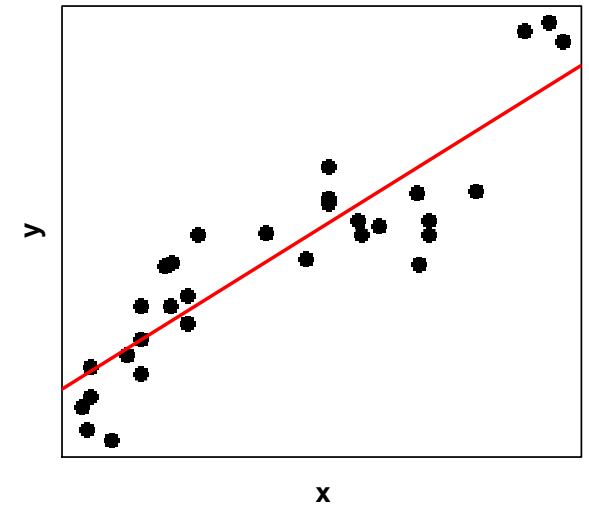
$r_{x,y} < 0$



$r_{x,y} \sim 0$



$r_{x,y} > 0$

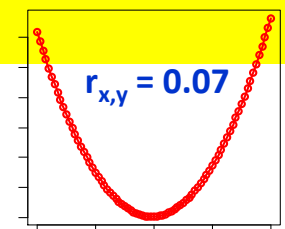


Indicativamente:

- se il suo valore assoluto è > 0.75 la correlazione è **forte**;
- se il suo valore assoluto è > 0.5 e < 0.75 la correlazione è **moderata**;
- se il suo valore assoluto è > 0.25 e < 0.5 la correlazione è **debole**;
- se il suo valore assoluto è < 0.25 la correlazione è **nulla**.

ATTENZIONE!!!

$r_{x,y} \sim 0$ vuole solamente dire che tra x e y non c'è correlazione lineare (altre sono possibili)!



SEGUE

Coefficiente di correlazione di Pearson (2)

$$r_{x,y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$r_{x,y} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\sigma_{x,y}}{\sqrt{\sigma_x^2 \cdot \sigma_y^2}} = \frac{\sigma_{x,y}}{\sigma_x \cdot \sigma_y}$$

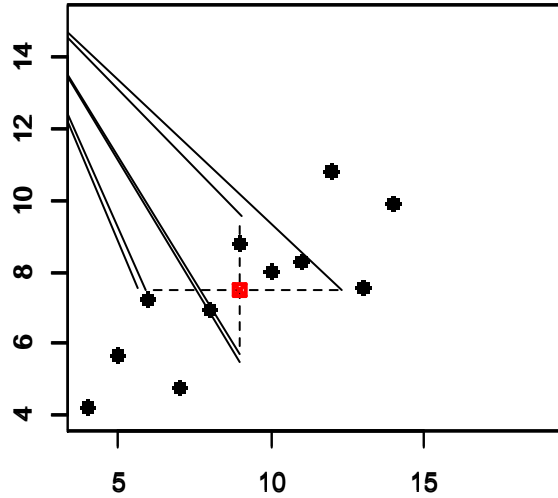
covarianza di x,y

dev.std. di x

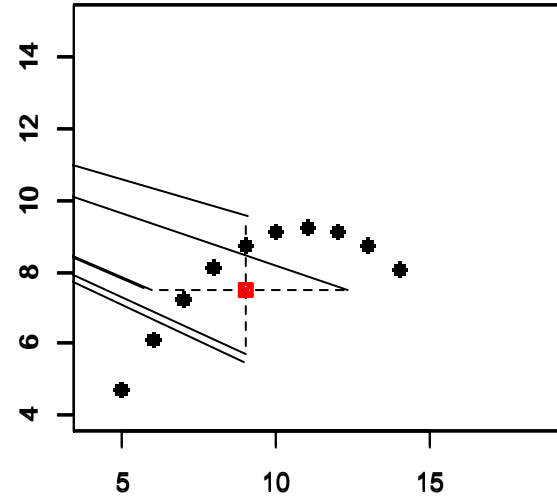
dev.std. di y

Quartetto di Anscombe

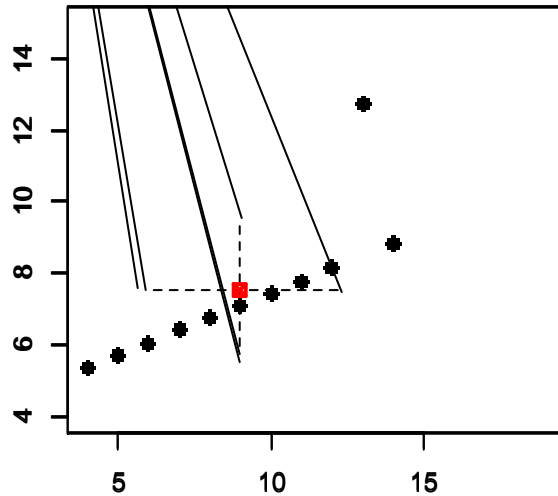
Media x = 9.000
Dev.std. x = 3.317
Media y = 7.500
Dev.std. y = 2.031
 $r_{x,y} = 0.8164$



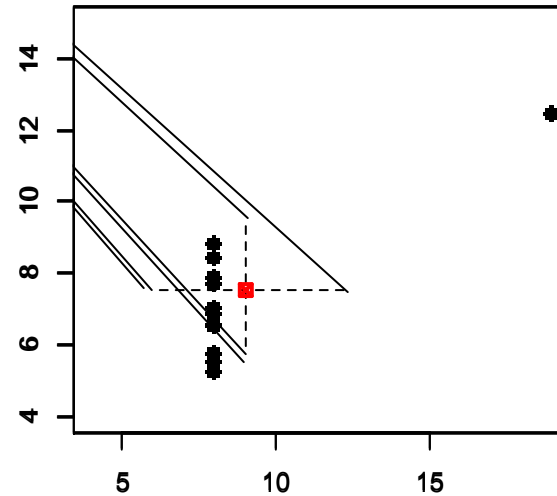
Media x = 9.000
Dev.std. x = 3.317
Media y = 7.501
Dev.std. y = 2.032
 $r_{x,y} = 0.8162$



Media x = 9.000
Dev.std. x = 3.317
Media y = 7.500
Dev.std. y = 2.030
 $r_{x,y} = 0.8163$



Media x = 9.000
Dev.std. x = 3.317
Media y = 7.501
Dev.std. y = 2.031
 $r_{x,y} = 0.8165$



L'esplorazione visiva dei grafici è FONDAMENTALE!!!

Due variabili / parametri misurati non sono sempre sufficienti

Dai dati bivariati ai dati multivariati

Matrice di correlazione

Survey of environmental complex systems: pattern recognition of physicochemical data describing coastal water quality in the Gulf of Trieste



Pierluigi Barbieri,^a Gianpiero Adami,^a Sergio Predonzani,^b Edoardo Reisenhofer^a and Desiré Luc Massart^c

^aDepartment of Chemical Sciences, University of Trieste, via L. Giorgieri 1, I-34127, Trieste, Italy

^bLaboratory of Marine Biology, Strada Costiera 336, I-34010 Santa Croce-Trieste, Italy

^cChemoAC, Pharmaceutical Institute, Vrije Universiteit Brussel, Laarbeeklaan 103, B-1090 Brussels, Belgium

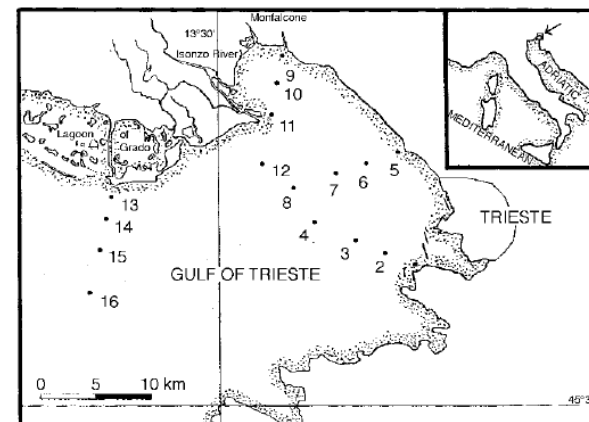


Fig. 1 Map of the sampling area.

Table 1 Basic statistics for the 909 valid cases

	TEMP/ °C	SAL/ practical salinity units	O ₂ / μmol dm ⁻³	N-NH ₃ / μmol dm ⁻³	N-NO ₂ / μmol dm ⁻³	N-NO ₃ / μmol dm ⁻³	Si(OH) ₄ / μmol dm ⁻³	Chlo.a/ μg dm ⁻³	Phaeop./ μg dm ⁻³
Mean	14.74	36.06	243.24	0.90	0.28	4.67	3.35	0.47	0.47
s	5.25	2.74	34.11	1.10	0.31	9.30	4.41	0.42	0.64
Min	5.85	17.58	18.30	0.05	0.02	0.02	0.02	0.01	0.01
Max	26.99	37.90	305.80	13.73	2.60	125.38	44.65	2.90	9.72

Table 2 Correlation matrix for 909 cases

	TEMP	SAL	O ₂	N-NH ₃	N-NO ₂	N-NO ₃	Si(OH) ₄	Chlo.a	Phaeop.
TEMP	1.0000								
SAL	-0.2541	1.0000							
O ₂	-0.5302	-0.1108	1.0000						
N-NH ₃	0.1851	-0.2283	-0.3033	1.0000					
N-NO ₂	-0.3291	-0.1139	0.1382	0.1261	1.0000				
N-NO ₃	0.1295	-0.7243	0.0999	0.2437	0.1002	1.0000			
Si(OH) ₄	0.0612	-0.4569	-0.1968	0.3029	0.1040	0.4550	1.0000		
Chlo.a	0.0024	-0.1350	0.0102	-0.0615	-0.0653	0.0400	-0.0386	1.0000	
Phaeop.	-0.1547	-0.0765	0.0800	-0.0487	0.1134	0.0138	0.0024	0.6798	1.0000

Comprimere l'informazione: analisi delle componenti principali e metodi fattoriali



Nozioni di algebra delle matrici

Moltiplicazione matrice-vettore: $M \cdot \vec{v} = \vec{v}_r$

$$\begin{pmatrix} a & b & c \\ d & e & f \end{pmatrix} \cdot \begin{pmatrix} g \\ h \\ i \end{pmatrix} = \begin{pmatrix} a \cdot g + b \cdot h + c \cdot i \\ d \cdot g + e \cdot h + f \cdot i \end{pmatrix} = \begin{pmatrix} j \\ k \end{pmatrix}$$

Dimensioni:

$n \times m$

$m \times 1$

$n \times 1$

Moltiplicazione matrice-matrice: $A \cdot B = C$

Attenzione!!! $A \cdot B \neq B \cdot A$

$$\begin{pmatrix} a & b \\ c & d \\ e & f \end{pmatrix} \cdot \begin{pmatrix} g & h \\ i & j \end{pmatrix} = \begin{pmatrix} a \cdot g + b \cdot i & a \cdot h + b \cdot j \\ c \cdot g + d \cdot i & c \cdot h + d \cdot j \\ e \cdot g + f \cdot i & e \cdot h + f \cdot j \end{pmatrix} = \begin{pmatrix} k & l \\ m & n \\ o & p \end{pmatrix}$$

Dimensioni:

$n \times m$

$m \times r$

$n \times r$

Significato geometrico

Moltiplicazione matrice-vettore: $A \cdot \vec{x} = \vec{v}$

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} a \cdot x_1 + b \cdot x_2 \\ c \cdot x_1 + d \cdot x_2 \end{pmatrix} = \begin{pmatrix} e \\ f \end{pmatrix}$$

Dimensioni:

$n \times m$

$m \times 1$

$n \times 1$

Si può pensare come:

$$a \cdot x_1 + b \cdot x_2 = e$$

$$c \cdot x_1 + d \cdot x_2 = f$$

e si può scomporre in:

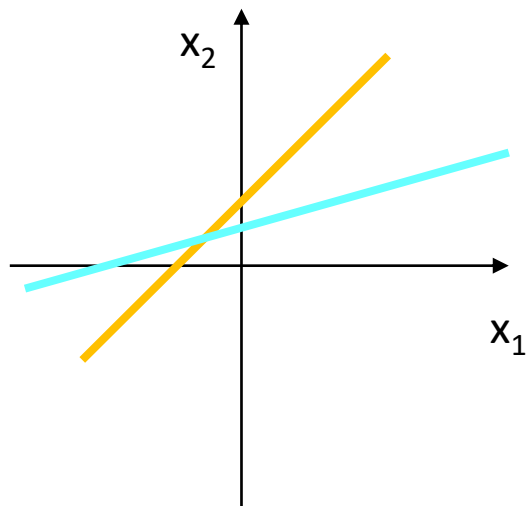
$$x_1 \cdot \begin{pmatrix} a \\ c \end{pmatrix} + x_2 \cdot \begin{pmatrix} b \\ d \end{pmatrix} = \begin{pmatrix} e \\ f \end{pmatrix}$$

Significato geometrico (2)

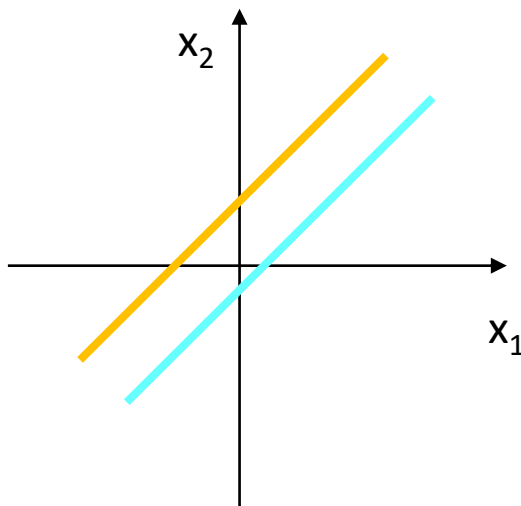
$$\begin{matrix} \rightarrow & \rightarrow \\ A \cdot x & = v \end{matrix}$$

$$a \cdot x_1 + b \cdot x_2 = e \longrightarrow \boxed{1} \quad x_2 = \frac{e - a \cdot x_1}{b}$$

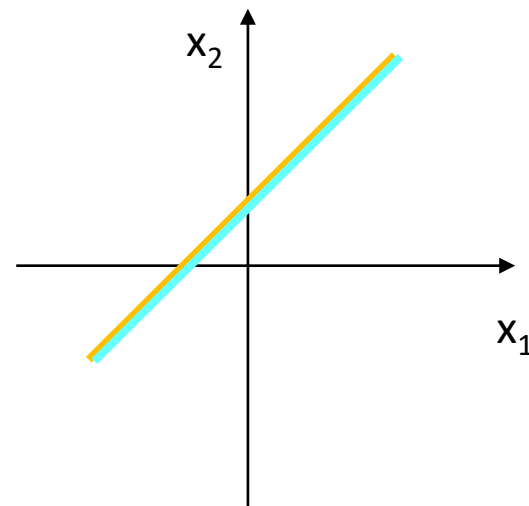
$$c \cdot x_1 + d \cdot x_2 = f \longrightarrow \boxed{2} \quad x_2 = \frac{f - c \cdot x_1}{d}$$



1 soluzione



nessuna soluzione



infinite soluzioni

La matrice identità

E' una matrice quadrata diagonale in cui tutti gli elementi sono uguali a zero, tranne che quelli presenti sulla diagonale, che sono uguali a 1.

Si indica con I_n

$$I_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad I_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Dimensioni: $n \times m$ con $n=m$

E' tale per cui:

$$I_n \cdot \vec{v} = \vec{v}$$

$$I_n \cdot A = A$$

$$A \cdot I_n = A$$

Quindi "funziona"
come la
moltiplicazione per 1
in algebra classica

Nota: ovviamente le dimensioni di I devono essere "compatibili" per effettuare le operazioni di moltiplicazione matrice-vettore o matrice-matrice.

In R:

diag(numerointero)

```
> diag(3)
      [,1] [,2] [,3]
[1,]    1    0    0
[2,]    0    1    0
[3,]    0    0    1
>
```

Matrici invertibili e matrici singole

Mentre di qualsiasi numero \neq da zero si può ottenere l'inverso, per le matrici non è sempre così.

Una matrice è invertibile se è quadrata ($n \times n$) **ed esiste** una matrice B tale che:

$$A \cdot B = I_n$$

$$B \cdot A = I_n$$

Se queste **condizioni** vengono **raggiunte** allora si dice che **A è invertibile** e si ha **$B = A^{-1}$** ,
cioè B è la matrice inversa di A

Se **non esiste** una matrice B che soddisfi tali condizioni allora si dice che **A è una matrice singola**, cioè non invertibile.

Nota: la matrice inversa di I_n è uguale a I_n

Risolvere equazioni matrice-vettore

In algebra "classica":

$$a \cdot x = b$$

$$\frac{1}{a} \cdot a \cdot x = \frac{1}{a} \cdot b$$

$$1 \cdot x = \frac{1}{a} \cdot b$$

$$x = \frac{b}{a}$$

quindi è necessario l'inverso di a !!!

\vec{x} è anche detto vettore dei
"pesi"
(*weights*)
da assegnare agli elementi di
A per riuscire ad ottenere \vec{b}



In algebra lineare:

$$A \cdot \vec{x} = \vec{b}$$

$$A^{-1} \cdot A \cdot \vec{x} = A^{-1} \cdot \vec{b}$$

$$I \cdot \vec{x} = A^{-1} \cdot \vec{b}$$

$$\vec{x} = A^{-1} \cdot \vec{b}$$

quindi è necessaria l'inversa di A !!!

Nota: A è una matrice quadrata.

Risolvere equazioni matrice - vettore in R - casi reali

Problema: quando i dataset (cioè le matrici) sono prodotti da casi reali è molto difficile che le matrici siano sempre quadrate.

$$\vec{A} \cdot \vec{x} = \vec{b}$$

Dimensioni di A: $n \times m$

Caso 1: $n > m$

Caso 2: $n < m$

Caso 1: più osservazioni che variabili

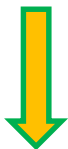
Se la matrice A ha più righe che colonne, di conseguenza ci sono più equazioni disponibili rispetto alle soluzioni possibili.

$$\begin{pmatrix} 5 & 2 \\ -3 & 4 \\ 2 & 6 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 4 \\ 3 \\ 7 \end{pmatrix}$$

$$\begin{aligned} 5 \cdot x_1 + 2 \cdot x_2 &= 4 \\ -3 \cdot x_1 + 4 \cdot x_2 &= 3 \\ 2 \cdot x_1 + 6 \cdot x_2 &= 7 \end{aligned}$$



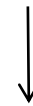
In questo caso la terza equazione è la somma delle prime due. Non è utile, ma non impedisce di risolvere l'equazione



L'equazione si dice CONSISTENTE

$$\begin{pmatrix} 5 & 2 \\ -3 & 4 \\ 2 & 6 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 4 \\ 3 \\ -1 \end{pmatrix}$$

$$\begin{aligned} 5 \cdot x_1 + 2 \cdot x_2 &= 4 \\ -3 \cdot x_1 + 4 \cdot x_2 &= 3 \\ 2 \cdot x_1 + 6 \cdot x_2 &= -1 \end{aligned}$$



In questo caso la terza equazione è incompatibile con le prime due. Quindi impedisce di ottenere una soluzione



L'equazione si dice INCONSISTENTE

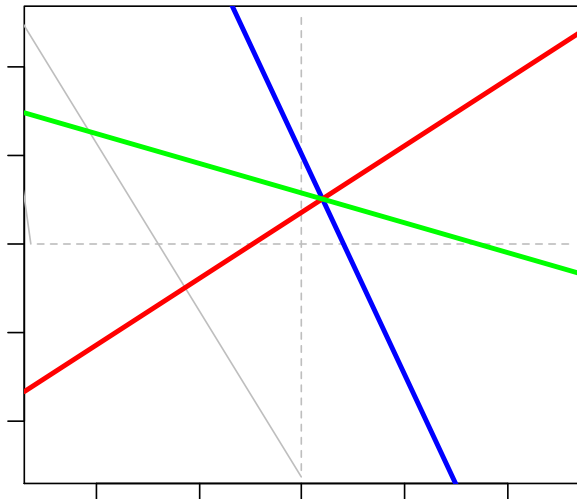
Caso 1: rappresentazione grafica

$$\begin{pmatrix} 5 & 2 \\ -3 & 4 \\ 2 & 6 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 4 \\ 3 \\ 7 \end{pmatrix}$$

→ $5 \cdot x_1 + 2 \cdot x_2 = 4$

→ $-3 \cdot x_1 + 4 \cdot x_2 = 3$

→ $2 \cdot x_1 + 6 \cdot x_2 = 7$

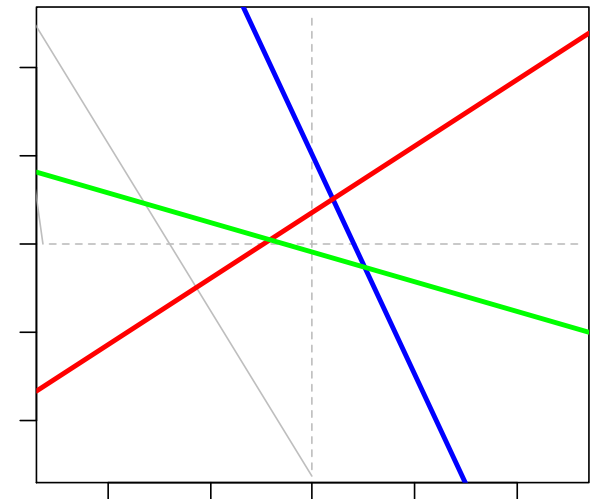


$$\begin{pmatrix} 5 & 2 \\ -3 & 4 \\ 2 & 6 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 4 \\ 3 \\ -1 \end{pmatrix}$$

→ $5 \cdot x_1 + 2 \cdot x_2 = 4$

→ $-3 \cdot x_1 + 4 \cdot x_2 = 3$

→ $2 \cdot x_1 + 6 \cdot x_2 = -1$



Caso 2: meno osservazioni che variabili

Se la matrice A ha meno righe che colonne, di conseguenza ci sono meno equazioni disponibili rispetto alle soluzioni possibili.

$$\begin{pmatrix} 3 & -5 & 1 \\ 9 & -1 & 2 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 5 \\ -1 \end{pmatrix}$$

$$3 \cdot x_1 - 5 \cdot x_2 + 1 \cdot x_3 = 5$$

$$9 \cdot x_1 - 1 \cdot x_2 + 2 \cdot x_3 = -1$$

Manca una condizione per riuscire ad avere una soluzione univoca.

Ci sono diversi approcci per cercare di ottenere una soluzione univoca

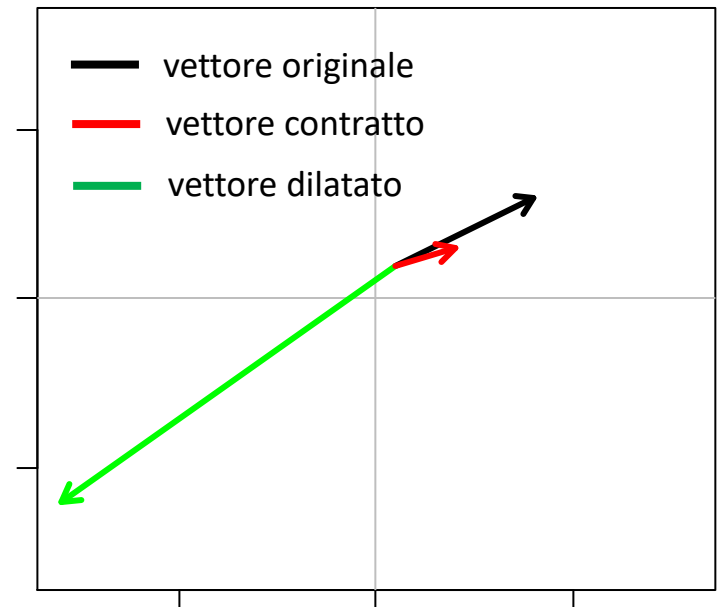
Interpretazione geometrica di $A \cdot \vec{x} = \vec{b}$

Una matrice moltiplicata per un vettore può indurre diversi effetti sul vettore:

- Rotazioni;
- Riflessioni;
- Contrazioni;
- Dilatazioni;
- Proiezioni;
- Combinazioni delle precedenti.

Nota: uno **scalare** moltiplicato per un vettore, cioè $a \cdot \vec{b}$, può avere come unico effetto di contrarre o dilatare tutte le componenti del vettore allo stesso modo

In R: **2*nomevettore**



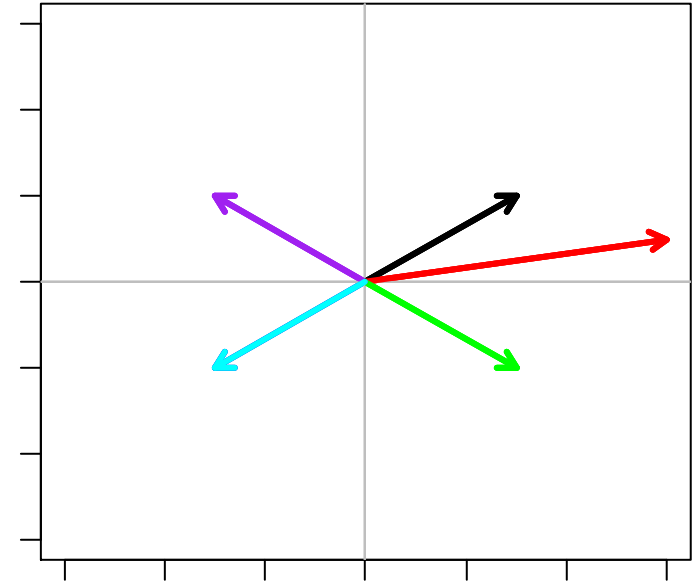
Interpretazione geometrica di $A \cdot \vec{x} = \vec{b}$

$$\text{---} \begin{pmatrix} 2 & 0 \\ 0 & 0.5 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 6 \\ 1 \end{pmatrix}$$

$$\text{---} \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 3 \\ -2 \end{pmatrix}$$

$$\text{---} \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} -3 \\ 2 \end{pmatrix}$$

$$\text{---} \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} -3 \\ -2 \end{pmatrix}$$

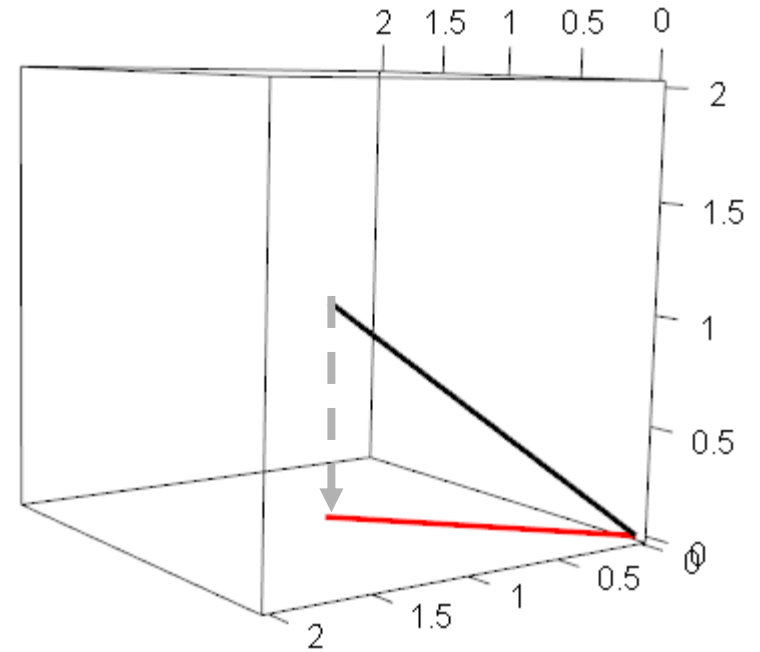


- vettore originale
- vettore contratto e dilatato
- vettore riflesso su asse x
- vettore riflesso su origine
- vettore riflesso su asse y

Interpretazione geometrica di $A \cdot \vec{x} = \vec{b}$

Proiezione:

$$\begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \\ 0 \end{pmatrix}$$



- vettore originale
- vettore proiettato

Autovalori (*eigenvalues*) e autovettori (*eigenvectors*)

$$\mathbf{A} \cdot \vec{v} = \lambda \cdot \vec{v}$$

$$(\text{con } \vec{v} \neq \vec{0})$$

λ è un autovalore di \mathbf{A} con un autovettore associato \vec{v}

L'operazione di moltiplicazione matrice-vettore $\mathbf{A} \cdot \vec{v}$ produce lo stesso effetto
dell'operazione di moltiplicazione scalare-vettore $\lambda \cdot \vec{v}$

Nota: non è necessario che \mathbf{A} sia una matrice diagonale

La coppia autovalore-autovettore (*eigenpair*)
ha su \vec{v} lo stesso effetto della matrice identità



$$\mathbf{I} \cdot \vec{v} = \vec{v}$$

Autovalori (*eigenvalues*) e autovettori (*eigenvectors*) (2)

$$A \cdot \vec{v} = \lambda \cdot \vec{v}$$

$$\begin{pmatrix} 2 & 3 \\ 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 2 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

Quindi **2** è un autovalore di **A** accoppiato all'autovettore $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$

Dal punto di vista geometrico: per una matrice A che può operare una rotazione, riflessione, proiezione o combinazione di queste su un generico vettore, l'autovettore è un vettore che rimane "fisso" rispetto a questa operazione (cioè rimane sulla stessa "linea").

Un autovettore può essere **riscalato** a piacere (cioè essere moltiplicato per uno scalare) perché la direzione (linea) del vettore rimane la stessa indifferentemente dalla sua grandezza (*magnitude*). Quindi un autovettore è un multiplo scalare di se stesso.



per questo si definisce "**auto**", a causa della sua "**indipendenza**"

Proprietà delle soluzioni

$$A \cdot \vec{v} = \lambda \cdot \vec{v}$$

- ✓ Risolvere l'equazione sopra equivale a risolvere l'equazione: $A \cdot \vec{v} - \lambda \cdot \vec{v} = \vec{0}$
- ✓ Una matrice $n \times n$ può avere al massimo n autovalori (alcuni di essi però possono anche essere uguali);
- ✓ Una matrice di numeri reali può anche avere uno o più autovalori complessi (ma saranno comunque sempre presenti a coppie);
- ✓ La matrice $(\lambda \cdot I - A)$ non è invertibile.

In R:

```
> E<-eigen(nomematrice)
```

```
> E$values
```

```
> E$vectors
```


Altre proprietà

$$A \cdot \vec{v} = \lambda \cdot \vec{v}$$

Per un set di autovalori **DISTINTI** di A:

$$\lambda_1 + \lambda_2 + \dots + \lambda_n$$

il corrispondente set di autovettori:

$$\vec{v}_1 + \vec{v}_2 + \dots + \vec{v}_n$$

forma una **base** (cioè un insieme di vettori linearmente indipendenti) per uno **spazio a n-dimensioni** in cui qualsiasi vettore espresso come combinazione lineare degli autovettori può esservi rappresentato:

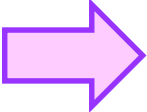
$$\overrightarrow{\text{VettoreGenerico}} = a_1 \cdot \vec{v}_1 + a_2 \cdot \vec{v}_2 + \dots + a_n \cdot \vec{v}_n$$

Altre proprietà (2)

$$\overrightarrow{\text{VettoreGenerico}} = a_1 \cdot \vec{v}_1 + a_2 \cdot \vec{v}_2 + \dots + a_n \cdot \vec{v}_n$$

Dato che $A \cdot \vec{v}_j = \lambda_j \cdot \vec{v}_j$ (per j da 1 a n), si può trasformare anche in:

$$A \cdot \overrightarrow{\text{VettoreGenerico}} = a_1 \cdot \lambda_1 \cdot \vec{v}_1 + a_2 \cdot \lambda_2 \cdot \vec{v}_2 + \dots + a_n \cdot \lambda_n \cdot \vec{v}_n$$

 Quindi gli *eigenpairs* **trasformano** una moltiplicazione matrice-vettore in una somma di moltiplicazioni scalare-vettore

In altre parole gli autovettori generano gli "assi" lungo i quali la moltiplicazione della matrice per un vettore generico semplicemente dà un "peso" (*weight*) a quel vettore utilizzando gli autovalori.

L'Analisi delle Componenti Principali (o *Principal Component Analysis*)

- ✓ E' un'analisi statistica non parametrica;
- ✓ Rivela strutture nascoste a livello multidimensionale che si trovano all'interno dei dati;
- ✓ Rivela eventuali RIDONDANZE (o collinearità) tra le variabili (cioè se alcune variabili sono correlate);
- ✓ Le strutture nascoste che vengono rivelate hanno generalmente un numero basso di dimensioni (2 o 3) quindi sono visualizzabili.

Principal Component Analysis (2)

- Per una matrice A ($n \times m$) si può sempre ottenere la matrice trasposta A^T ($m \times n$) scambiando le righe con le colonne;
- Il prodotto $A^T \cdot A$ è una matrice quadrata di dimensioni $m \times m$.
- Se per ogni colonna della matrice A la media di colonna è stata sottratta ad ogni elemento della colonna (rispettivamente), allora nella seguente matrice

$$B = \frac{A^T \cdot A}{n - 1}$$

ogni elemento i,j è la COVARIANZA tra le variabili che si trovano alle colonne i e j della matrice A. Quindi sulla diagonale di B ci sono le varianze di ogni colonna di A.

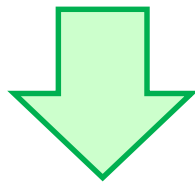
$$\begin{pmatrix} \text{Var} & \text{Cov} & \text{Cov} \\ \text{Cov} & \text{Var} & \text{Cov} \\ \text{Cov} & \text{Cov} & \text{Var} \end{pmatrix}$$

Proprietà di B

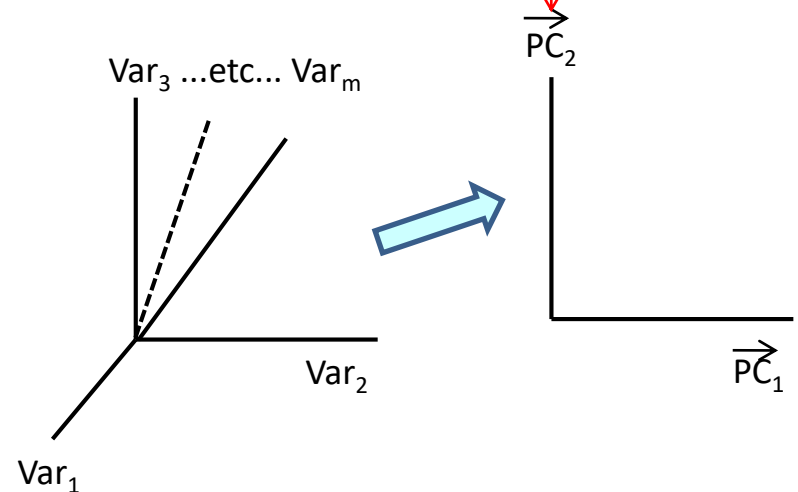
$$B = \frac{A^T \cdot A}{n - 1}$$

$$\begin{pmatrix} \text{Var} & \text{Cov} & \text{Cov} \\ \text{Cov} & \text{Var} & \text{Cov} \\ \text{Cov} & \text{Cov} & \text{Var} \end{pmatrix}$$

- Gli autovalori di B sono numeri reali e i corrispondenti autovettori sono ortogonali nello spazio a m dimensioni (cioè generano degli assi);
- La varianza totale del sistema A è data dalla somma degli autovalori di B;
- Gli autovettori di B sono detti COMPONENTI PRINCIPALI ($\vec{v}_j = \vec{PC}_j$);
- La direzione v_j spiega la λ_j - frazione della varianza totale del sistema.



Quindi se **alcuni autovalori** spiegano buona parte della **varianza** totale del sistema, basta usare i **relativi autovettori** per rappresentare le osservazioni (righe di A) nel nuovo sistema di coordinate. Il resto di solito è piccola variabilità o "rumore"



La normalizzazione

- Se per ogni colonna della matrice A la media di colonna è stata sottratta ad ogni elemento della colonna (rispettivamente), allora nella seguente matrice

$$B = \frac{A^T \cdot A}{n - 1} \begin{pmatrix} \text{Var} & \text{Cov} & \text{Cov} \\ \text{Cov} & \text{Var} & \text{Cov} \\ \text{Cov} & \text{Cov} & \text{Var} \end{pmatrix}$$

ogni elemento i, j è la COVARIANZA tra le variabili che si trovano alle colonne i e j della matrice A. Quindi sulla diagonale di B ci sono le varianze di ogni colonna di A.



Questa è una normalizzazione necessaria perché le componenti principali sono gli autovettori della matrice di covarianza

Di solito è opportuno anche utilizzare la deviazione standard per la normalizzazione utilizzando il così detto metodo **Z-score**:

$$Z_{\text{Var}(i,j)} = \frac{\text{Var}_{(i,j)} - \overline{\text{Var}_j}}{\sigma_{\text{Var}_j}}$$

Questo fattore rende la **media di Var = 0**

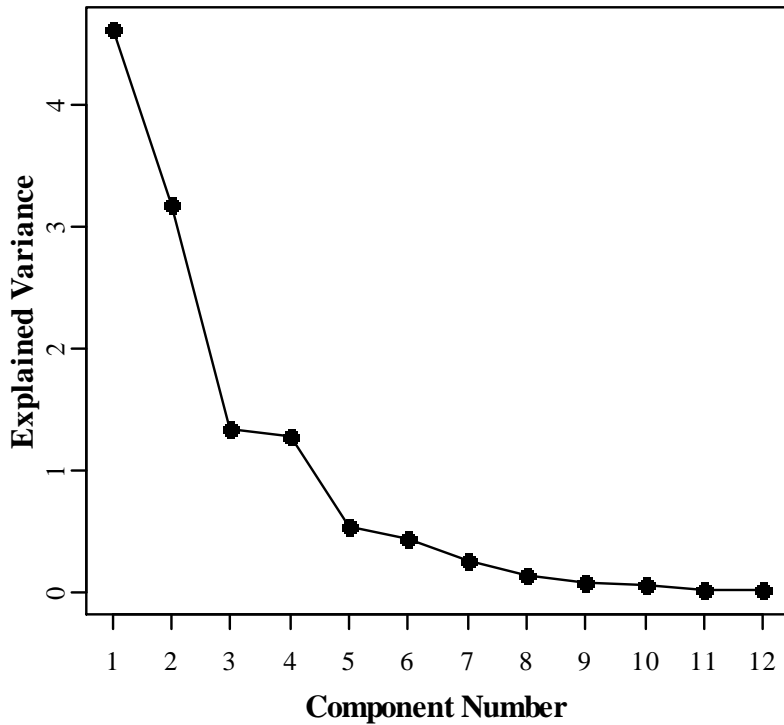
Questo fattore rende la **varianza di Var = 1**

$$\begin{pmatrix} 1 & \text{Cov} & \text{Cov} \\ \text{Cov} & 1 & \text{Cov} \\ \text{Cov} & \text{Cov} & 1 \end{pmatrix}$$

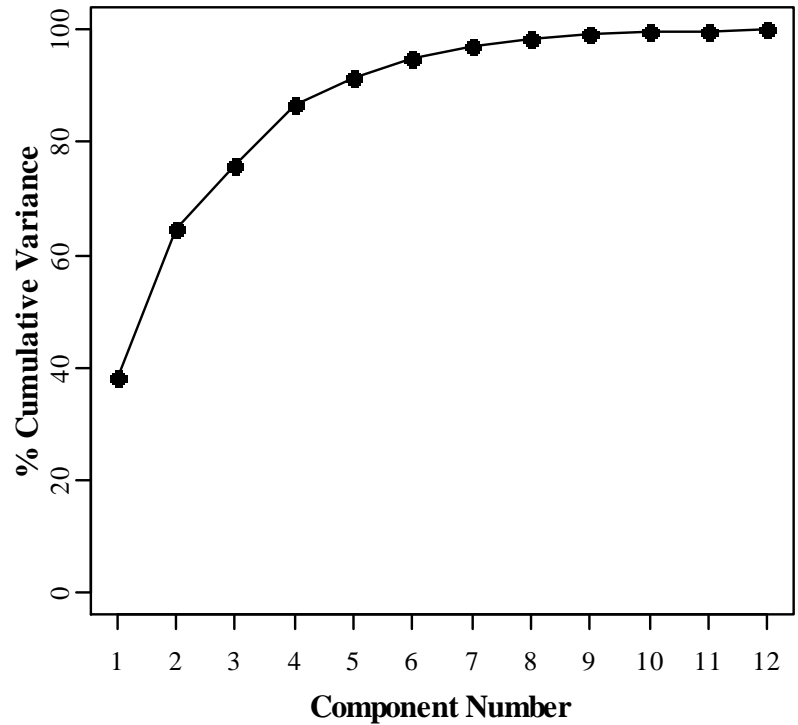
Risultati della PCA: i grafici della varianza spiegata

Gli autovalori di ogni componente principale (PC) rappresentano la varianza spiegata da quella componente

Scree Plot

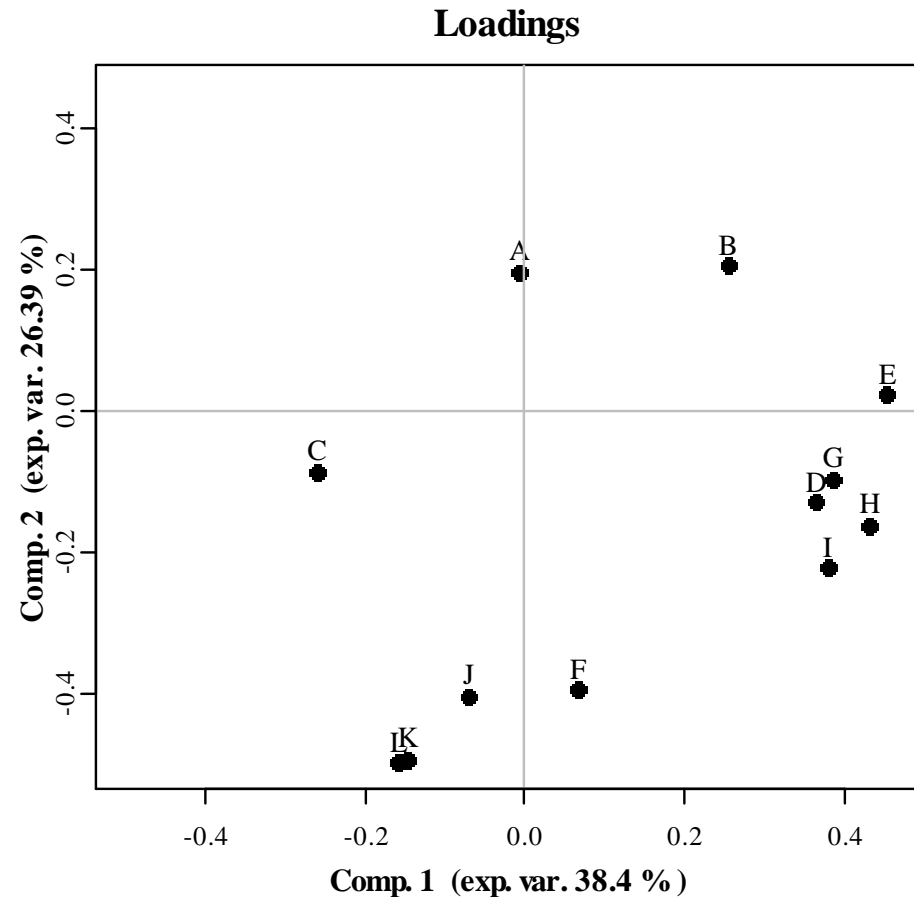
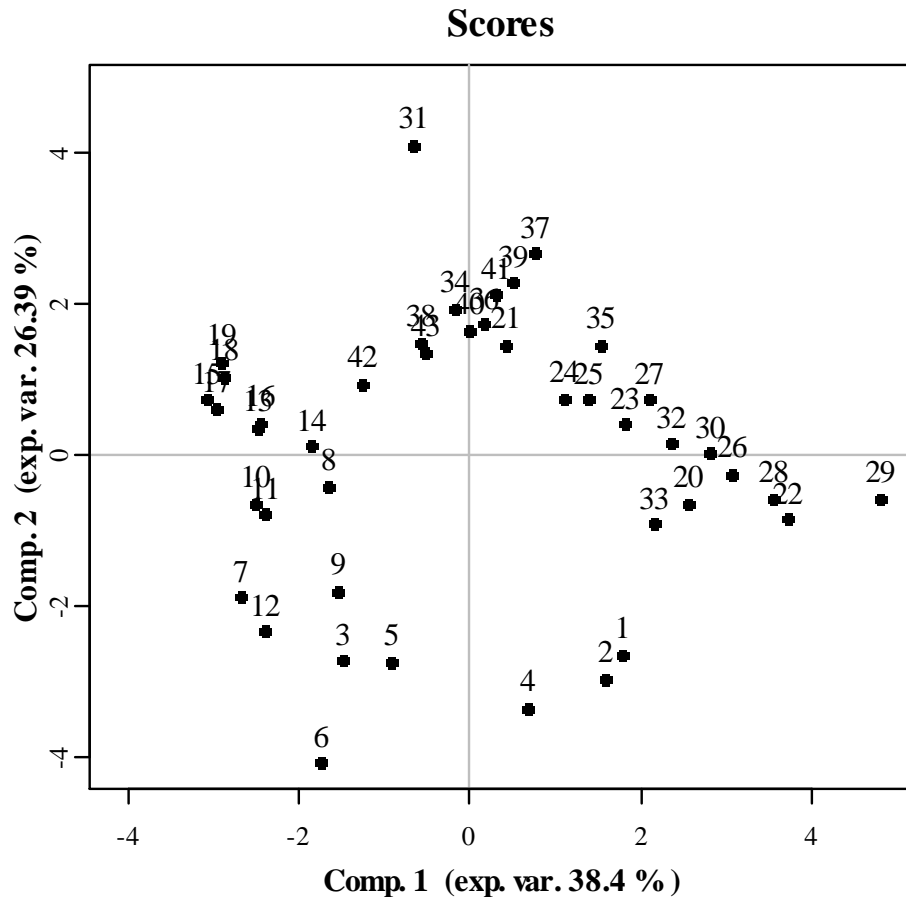


% Cumulative Variance



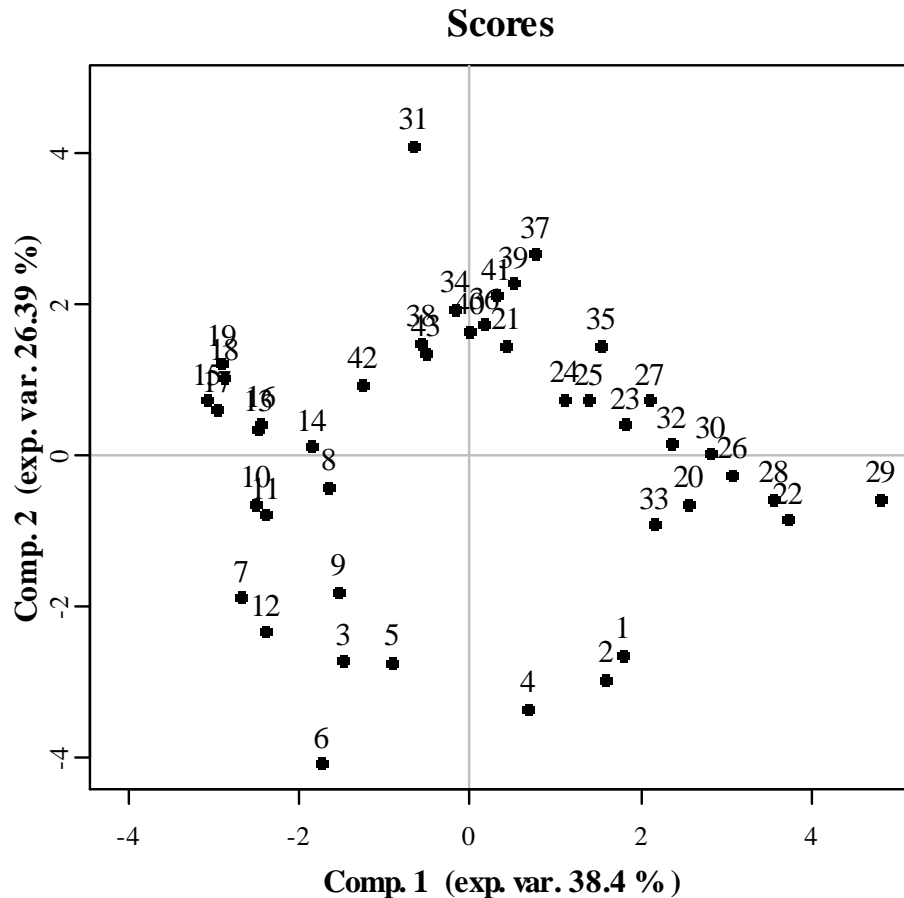
Risultati della PCA: i grafici degli Score e dei Loading

In un grafico a 2 dimensioni si pongono in ascissa una PC e in ordinata un'altra PC (di solito si inizia da 1 e 2, che rappresentano la maggior parte della varianza del set di dati)



Risultati della PCA: i grafici degli Score e dei Loading (2)

Il grafico degli **Score** rappresenta le osservazioni nel nuovo sistema di variabili



- I punti (quindi le osservazioni) che sono prossime all'origine sono osservazioni che hanno valori prossimi alla media per la maggior parte delle variabili;
- Le osservazioni che si trovano molto lontano dall'origine possono essere o dati estremi del dataset o outlier;
- Osservazioni simili si trovano vicine nel nuovo spazio ovvero osservazioni dissimili si trovano lontane tra loro.

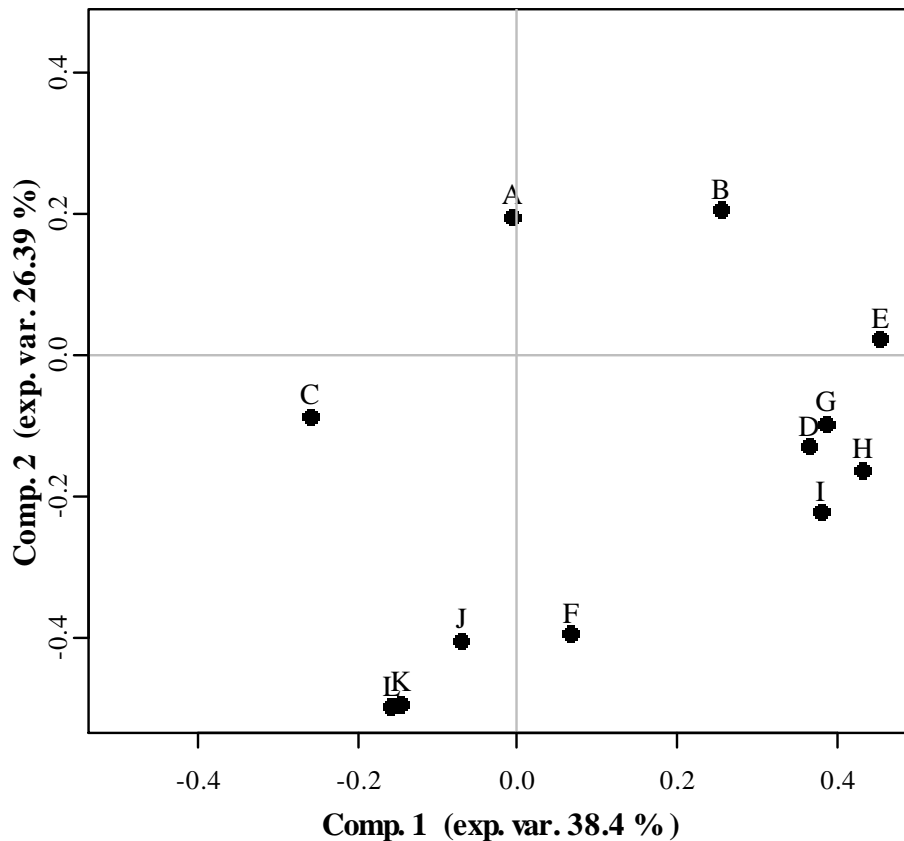
SEGUE

Risultati della PCA: i grafici degli Score e dei Loading (3)

Il grafico dei **Loading** rappresenta i coefficienti (pesi) della combinazione lineare delle variabili originali che servono per generare le componenti principali.

$$\vec{v}_j = \overrightarrow{PC}_j = w_1 \cdot \overrightarrow{Var}_1 + w_2 \cdot \overrightarrow{Var}_2 + \dots + w_m \cdot \overrightarrow{Var}_m$$

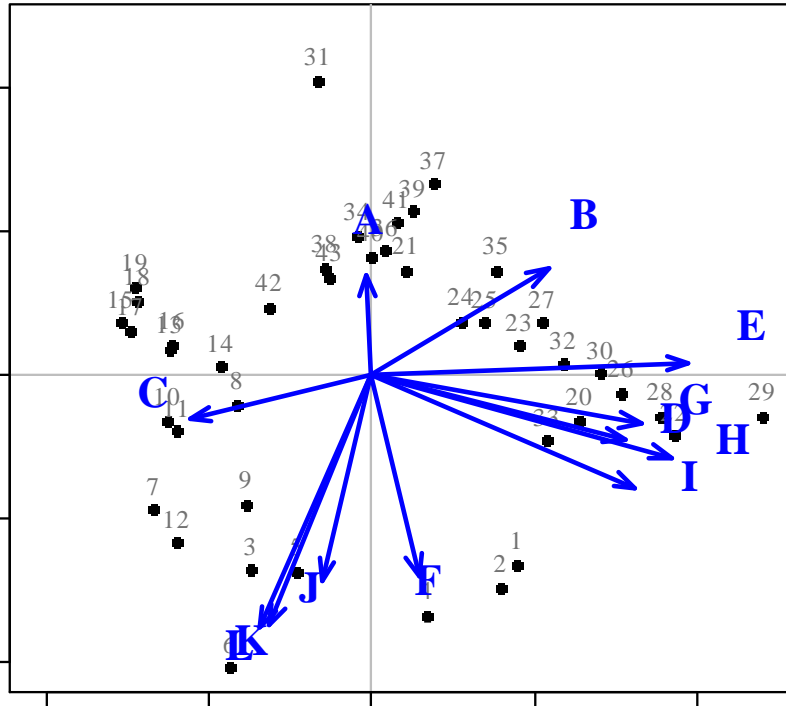
Loadings



- Variabili originali che contribuiscono poco alla direzione di una componente principale mostrano loading prossimi a zero per quella componente;
- Variabili originali che contribuiscono in modo simile alla direzione di una componente principale sono generalmente correlate e mostrano loading di valore circa uguale;
- Variabili originali correlate positivamente si trovano vicine nel grafico, mentre variabili correlate negativamente si trovano opposte tra loro rispetto ad una diagonale.

Risultati della PCA: il biplot

E' dato dalla sovrapposizione del grafico dei **Loading** a quello degli **Score**, centrati sull'origine. I **Loading** vengono rappresentati da frecce che si dipartono dall'origine.

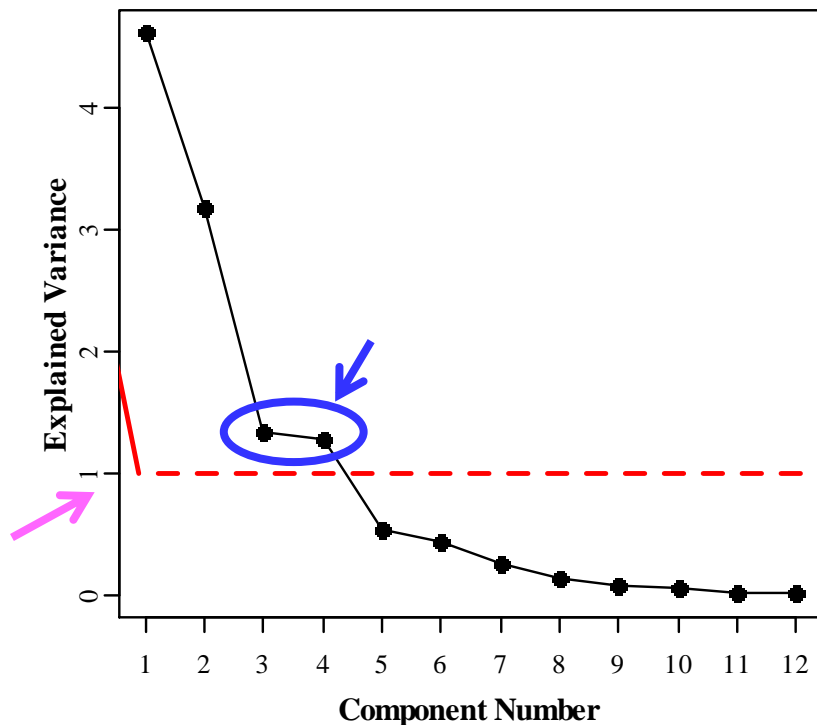


- Più un vettore è parallelo all'asse di una PC più la variabile originale associata contribuisce pressoché esclusivamente a quella PC;
- Più è lungo un vettore maggiore è la variabilità della variabile originale associata che viene spiegata dalle due PC su cui è costruito il grafico;
- Di conseguenza se un vettore è corto vuol dire che la variabilità della variabile originale associata viene rappresentata meglio da altre PC non utilizzate per costruire il grafico;
- L'angolo tra due vettori rappresenta la relativa correlazione delle variabili originali associate:
 - angolo $\sim 0^\circ$ \rightarrow alta correlazione positiva;
 - angolo $\sim 90^\circ$ \rightarrow pressoché non correlate;
 - angolo $\sim 180^\circ$ \rightarrow alta correlazione negativa

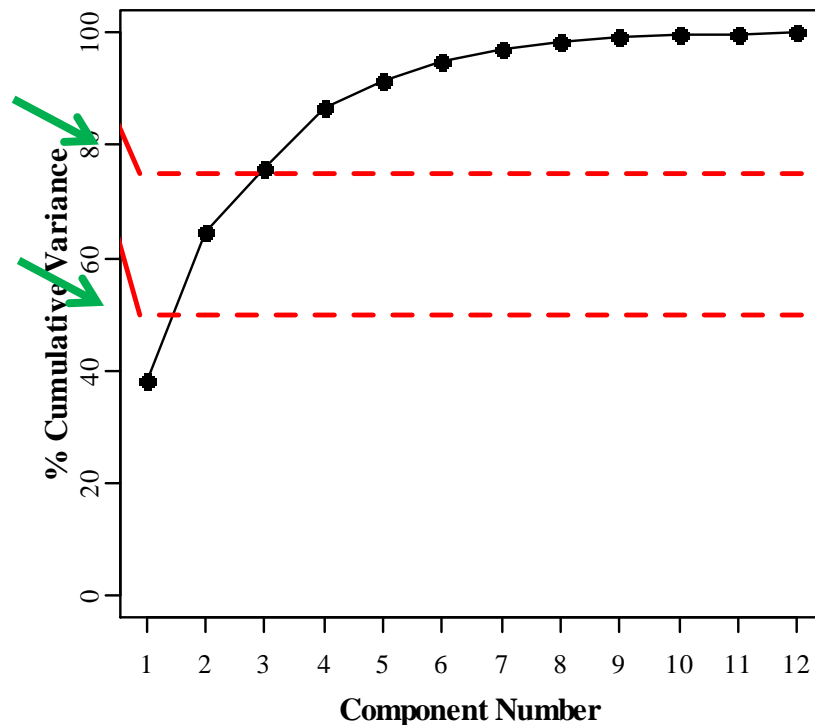
Inoltre si visualizza quali variabili originali "spostano" di più quali osservazioni dall'origine (cioè dal "baricentro" del dataset)

Quante PC utilizzare per rappresentare il dataset?

Scree Plot



% Cumulative Variance



- ➡ 1. Esplorazione visiva dello scree plot: valutare la presenza di "un gomito" nel grafico
- ➡ 2. Valutazione di quante PC spiegano almeno il 50% (o il 75%) della varianza del dataset
- ➡ 3. Valutazione di quante PC mostrano varianza > 1 (poichè quelle < 1 tecnicamente spiegano meno varianza di una variabile originale) - *Criterio di Kaiser-Guttman*

Dati mancanti (NA) e dati < Limit of Detection (LOD), che fare?

- Dati mancanti:
- Soluzione 1: eliminazione **In R:** `na.omit(nomeoggetto)`
 - Soluzione 2: sostituzione con dati modellati

Esempi:

EM-PCA (iterative PCA)

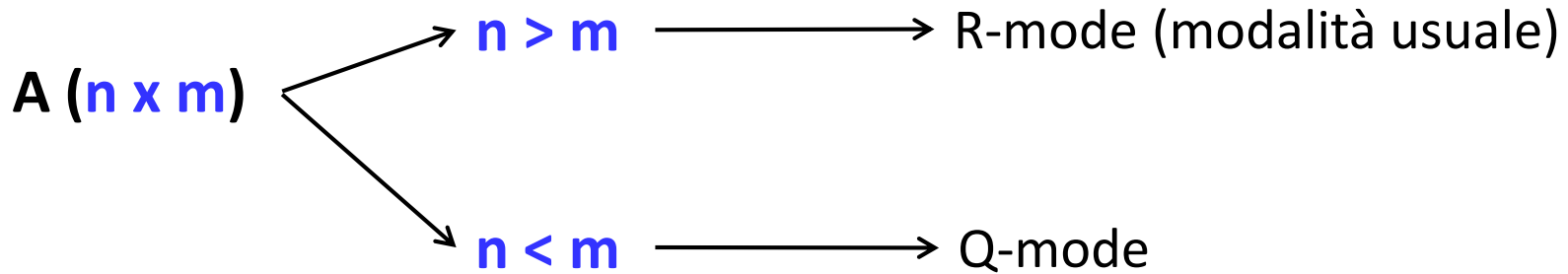
NIPALS (non-linear iterative partial least squares)

- In R:**
- Il pacchetto "**missMDA**" con le sue funzioni consente di gestire i dati mancanti con diversi algoritmi;
 - Il pacchetto "**pcaMethods**" (in Bioconductor packages) consente di calcolare una PCA con diversi algoritmi che gestiscono dati mancanti.

➤ Dati < LOD:

Suggerimento: non conviene porre = a 0 (zero) questi dati, ma conviene sostituirli o con il valore del LOD o con un valore più piccolo del LOD che sia approssimabile a zero ma che non sia zero (es. 10^{-5} o 10^{-7} , etc...)

Se ci sono più variabili che osservazioni?



Procedere nel modo seguente:

- 1 - Scalare e centrare il dataset;
- 2 - Generare la matrice trasposta del dataset (ovvero scambiare righe e colonne);
- 3 - Calcolare la PCA

In R: `t(nomeoggetto)`



I grafici degli **Score** rappresentano i **Loading** e viceversa, l'interpretazione è simile a quella già vista