

-) Problema dovuto a Gauss (segue da pag. 194)

Stipiamo usando il metodo dei minimi quadrati, al fine di individuare tra tutti i possibili modelli, quello che minimizza la somma dei quadrati dei residui (ovvero rispetto al modello)

Esiste una formula esplicita che fornisce il vettore dei parametri ottimali si determina nel modo seguente:

$$J(\beta) = E^T E = \|E\|^2$$

$$\beta^{LS} = (\Phi^T \Phi)^{-1} \Phi^T Y$$

Insomma, rimane tutto, da la parte inversa

$$(\Phi^T \Phi)^{-1}$$

esiste, considerando che si afferma che

$$\text{rank}(\Phi) = p$$

implica che

$$\det(\Phi^T \Phi) \neq 0$$

Infatti, se per assurdo, il suddetto determinante fosse nullo, allora esisterebbe un

$$k \neq 0$$

tale che

$$\Phi^T \Phi k = 0$$

che conseguirebbe che

$$k^T \Phi^T \Phi k = 0$$

$$(\Phi k)^T \Phi k = 0$$

$$\|\Phi k\|^2 = 0 \quad \text{ovvero } \|\Phi k\| = 0$$

$$\Phi k = 0$$

Se ne dedurrebbe che  $\Phi$  presenterebbe delle colonne linearmente dipendenti, in contraddizione con l'ipotesi fatta.

rank  $(\Phi) = p$   
 Affando supposto che  $\Phi \in \mathbb{R}^{n \times p}$  il rango di  $\Phi$  indica il numero  $p$  delle colonne indipendenti.

Ricordiamo che  $n$  è il numero dei dati e  $p$  il numero di parametri incogniti e che

$$\epsilon = Y - \Phi \theta \quad (\text{eq. 194})$$

Perciò,

$$\begin{aligned} f(\theta) &= (Y - \Phi \theta)^T (Y - \Phi \theta) = \\ &= (Y - \Phi \theta + \Phi \theta^{LS} - \Phi \theta^{LS})^T (Y - \Phi \theta + \Phi \theta^{LS} - \Phi \theta^{LS}) = \\ &= [Y - \Phi \theta^{LS} - \Phi (\theta - \theta^{LS})]^T [Y - \Phi \theta^{LS} - \Phi (\theta - \theta^{LS})] \end{aligned}$$

Definiamo  $\epsilon^{LS} = Y - \Phi \theta^{LS}$

con  $\epsilon^{LS}$  uguale ai residui ottenuti con  $\theta^{LS}$ .

Quindi,

$$\begin{aligned} f(\theta) &= [\epsilon^{LS} - \Phi (\theta - \theta^{LS})]^T [\epsilon^{LS} - \Phi (\theta - \theta^{LS})] = \\ &= (\epsilon^{LS})^T \epsilon^{LS} - (\epsilon^{LS})^T \Phi (\theta - \theta^{LS}) - (\theta - \theta^{LS})^T \Phi^T \epsilon^{LS} + \\ &\quad + (\theta - \theta^{LS})^T \Phi^T \Phi (\theta - \theta^{LS}) \end{aligned}$$

Osserviamo che

$$(\epsilon^{LS})^T \epsilon^{LS} \in \mathbb{R}^{1 \times n} \times \mathbb{R}^{n \times 1} = \mathbb{R} \quad (\text{scalare})$$

$$(\epsilon^{LS})^T \Phi (\theta - \theta^{LS}) \in \mathbb{R}^{1 \times n} \times \mathbb{R}^{n \times p} \times \mathbb{R}^{p \times 1} = \mathbb{R}$$

con  $\theta^{LS}$  = vettore dei residui

$$(\theta - \theta^{LS})^T \Phi^T \epsilon^{LS} \in \mathbb{R}^{1 \times p} \times \mathbb{R}^{p \times n} \times \mathbb{R}^{n \times 1} = \mathbb{R}$$

$$\begin{aligned} (\theta - \theta^{LS})^T \Phi^T \Phi (\theta - \theta^{LS}) &\in \mathbb{R}^{1 \times p} \times \mathbb{R}^{p \times n} \times \mathbb{R}^{n \times p} \times \mathbb{R}^{p \times 1} = \\ &= \mathbb{R} \end{aligned}$$

Tutti i termini sono scalari e uno scalare è uguale al suo trasposto. De conseguenza che il secondo termine ed il terzo sono tra loro uguali.

Si ottiene

$$f(\theta) = (\epsilon^{LS})^T \epsilon^{LS} - 2(\epsilon^{LS})^T \Phi (\theta - \theta^{LS}) + (\theta - \theta^{LS})^T \Phi^T \Phi (\theta - \theta^{LS})$$

201

Amplificando che il secondo termine è nullo. Infatti

$$2(\theta - \theta^{LS})^T \Phi (\epsilon^{LS})^T =$$

$$= 2(\theta - \theta^{LS})^T \Phi^T \epsilon^{LS} =$$

$$= 2(\theta - \theta^{LS})^T \Phi^T (Y - \Phi \theta^{LS}) =$$

$$= 2(\theta - \theta^{LS})^T \Phi^T (Y - \Phi (\Phi^T \Phi)^{-1} \Phi^T Y) =$$

$$= 2(\theta - \theta^{LS})^T \Phi^T [I - \Phi (\Phi^T \Phi)^{-1} \Phi^T] Y =$$

$$= 2(\theta - \theta^{LS})^T [\Phi^T - \Phi^T \Phi (\Phi^T \Phi)^{-1} \Phi^T] Y =$$

$$= 2(\theta - \theta^{LS})^T [\Phi^T - I \cdot \Phi^T] \cdot Y = 0,$$

(pag. 194)

ovvero

$$\Phi^T \Phi (\Phi^T \Phi)^{-1} = I$$

una matrice per la proprietà  
inversa fornisce la matrice  
identità.

Risulta:

$$f(\theta) = (\epsilon^{LS})^T \epsilon^{LS} + (\theta - \theta^{LS})^T \Phi^T \Phi (\theta - \theta^{LS})$$

Il primo termine non dipende da  $\theta$  ma da  $\theta^{LS}$  che ha un valore determinato, non minimizzabile!

La minimizzazione di  $f(\theta)$  occorre intervenire sul 2° termine, poiché  $\theta = \theta^{LS}$

ed osservando che si ha una forma quadratica semi-definita.

-) Proprietà della matrice

Dato una matrice  $D$ ,

$$D^T D$$

è detta diagonale ed è sempre semidefinita positiva.

202 *de* consegue de  
 $k^T D^T D k \geq 0 \quad \forall k \in \mathbb{R}^n$

-) Proprietà sulle matrici  
 Se una matrice  $S$  è semidefinita positiva,  
 $(S \geq 0)$   
 e  $\det(S) \neq 0$ ,  
 allora  $S$  è definita positiva.  
 $k^T S k > 0 \quad \forall k \neq 0$

Allora diade risulta matrice simmetrica.

In una matrice semidefinita positiva, tutti gli autovalori sono non negativi.  
 Se il determinante è diverso da zero (allora) gli autovalori sono tutti strettamente maggiori di zero e la matrice è definita positiva.

~~Il terzo termine di  $J(\theta)$  possiede questa proprietà~~

Il termine  $\Phi^T \Phi$  è una diade, perciò il terzo termine di  $J(\theta)$  si annulla solo se  $\theta = \theta^*$

e questo risulta la scelta migliore per l'adattatore

Interpretazione grafica

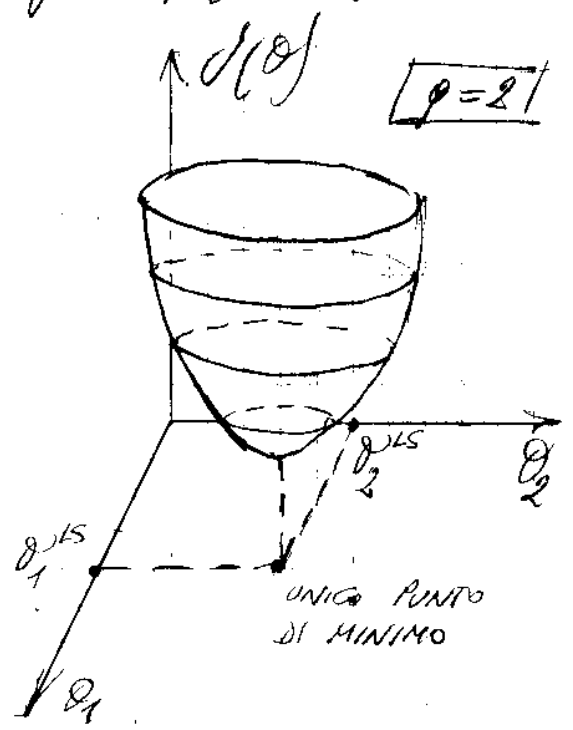
Caso  $q=2$ .

Il piano  $\theta_1, \theta_2$  è il piano dei parametri.

Nell'area delle altezze si legge la norma dei gradienti dei peschi.

Si desidera che la norma sia la più piccola possibile, si può notare che, se

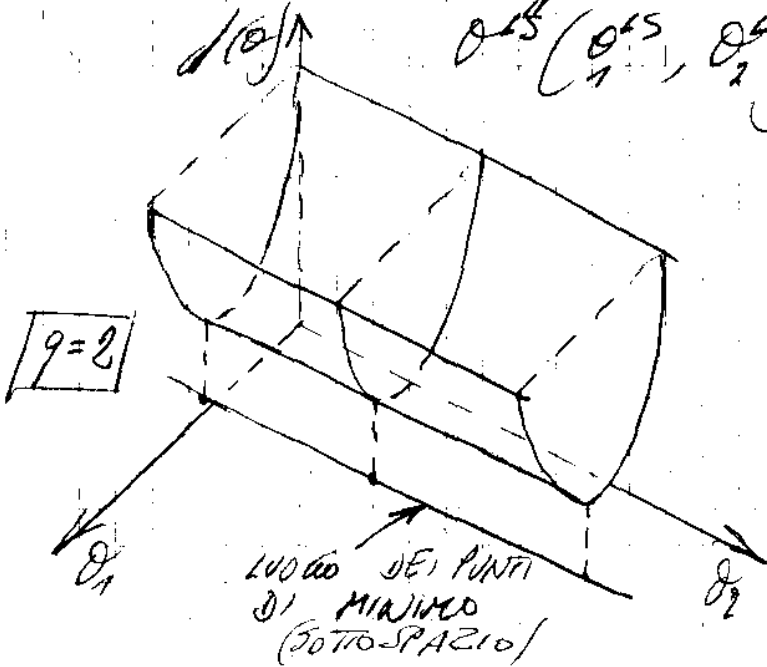
$\text{rank}(\Phi) = q$



$q=2$

$K(x)$  è una forma quadratica definita positivamente 203  
 con un unico punto di minimo globale in

$$x^0 = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$



$$\text{rank}(Q) = 1$$

$$e \quad q = 2$$

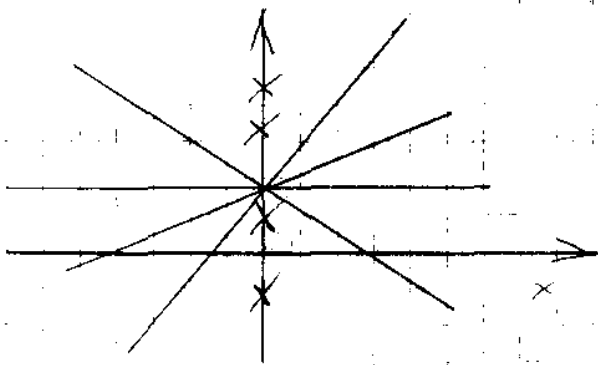
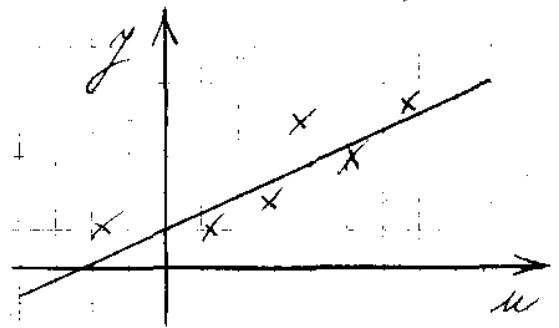
Si ha un luogo di punti di minimo di grado  $q$  in una cella.

Le rette possibili sono infinite.

### -) Regressione lineare ( $q > 1$ )

È un metodo di stima mediante modello lineare cambiando la posizione della retta cambia la somma dei quadrati dei residui.

Esiste una posizione di retta che minimizza la somma dei quadrati dei residui.



Del caso approssimato e minimo sono stati ricavati dati solo in cui l'ipotesi è  $\epsilon = 0$ .

Qualsiasi retta, fra quelle rappresentate, non può cambiare la somma dei quadrati dei residui di nessuno dei calcoli nel minimo.

I residui dipendono solo dalle intercette e non dai coefficienti angolari. Non si può, sull'altro asse, per determinare quest'ultimo.

Esistono infinite soluzioni ottime.

Il caso migliore è quello approssimato della regressione

fol. de pag. 202

- Esempio: regressione lineare con pendenza for  
 e il segno

la equazione generale del modello

$$y(i) = u(i)\beta + e(i)$$

ove

$$\begin{cases} y(1) = u(1)\beta + e(1) \\ y(2) = u(2)\beta + e(2) \\ \vdots \\ y(n) = u(n)\beta + e(n) \end{cases}$$

ovvero

$$\Phi = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$$

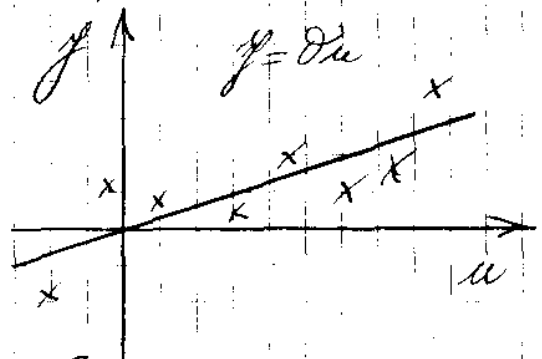
con  $Y \in \mathbb{R}^{n \times 1}$   $\Phi \in \mathbb{R}^{n \times 1}$   $e \in \mathbb{R}^{n \times 1}$   
 $\beta \in \mathbb{R}$  (scalare)

Chiedere che  $\text{rank}(\Phi) = 1$  significa chiedere che

$$\text{rank} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix} = 1$$

Lo si verifica da almeno una  $u(i)$  che differa da  
 la formula esplicita e (pag. 199)

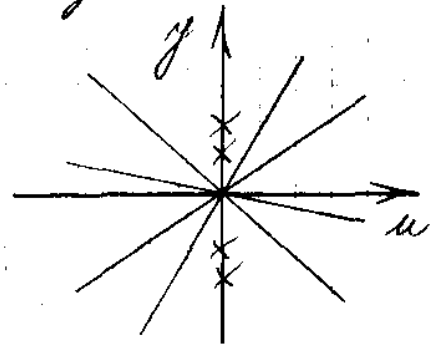
$$\begin{aligned} \beta &= (\Phi^T \Phi)^{-1} \Phi^T Y = \\ &= \begin{bmatrix} u(1) & u(2) & \dots & u(n) \end{bmatrix} \begin{bmatrix} u(1) \\ u(2) \\ \vdots \\ u(n) \end{bmatrix}^{-1} \begin{bmatrix} u(1) & u(2) & \dots & u(n) \end{bmatrix} \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(n) \end{bmatrix} = \\ &= \frac{\sum u(i) y(i)}{\sum u(i)^2} \end{aligned}$$



Le equazioni scritte e  
 simulate corrispondo a  
 $Y = \Phi \beta + e$

che rappresenta la formula delle regressione  $l_1 = 305$   
 nelle ipotesi per l'origine.

Il loro rank  $|\Phi| = 0$   
 corrisponde alla situazione cap. 2  
 presentata a destra in cui  
 tutte le  $u_i$  sono nulle.



Tutti i modelli pensati per  
 l'origine sono equivalenti, e  
 tutti i residui uguali in  
 tutte le circostanze.

Conoscione il metodo dei minimi quadrati non  
 fa alcun uso del concetto di probabilità.

La teoria ha dei limiti e pone alcuni quesiti:

- qual è l'affidabilità della stima? di  $\sigma^2$  calcolato
- il modello è giusto? funziona bene in tutti i casi?
- confronto tra modelli: occorre sapere degli stime  
 relative per poter valutare se un modello è  
 migliore di un altro.

Per risolvere questi quesiti occorre usare delle nozioni  
 relative alla teoria della stima.

206 STIMATORE BLUE

Consideriamo sotto l'ipotesi  $I_2$ , secondo la quale vale la relazione

$$Y = \Phi \theta + V \quad E[V] = 0$$

$$\text{Cov}[V] = \Sigma \Psi$$

Si sta formulando un modello probabilistico dei dati secondo il quale i dati osservati  $Y$  sono spiegati da un determinato modello matematico  $(\Phi \theta)$ . Qui si aggiunge un termine di errore  $V$  descritto da un vettore di variabili casuali, avente valore atteso nullo e con determinata matrice varianza

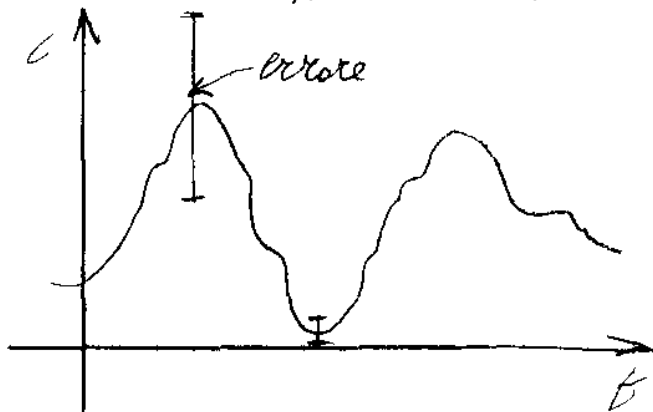
$$Y \in \mathbb{R}^{n \times 1} \quad \Phi \in \mathbb{R}^{n \times q} \quad \theta \in \mathbb{R}^{q \times 1}$$

$$V \in \mathbb{R}^{n \times 1}$$

La matrice varianza è composta da due termini,  $\Psi$ , matrice nota  $m \times n$   
 $\Sigma$ , valore stocasticamente incognito

Esempio per giustificare l'uso della matrice  $\Psi$   
 Si considerino delle misure di concentrazione del tipo

$$Y_i = c_i + V_i \quad c_i = \text{misure della concentrazione}$$



$V_i = \text{errore di misura}$   
 può supporre che la varianza di  $V_i$  sia  
 $\text{Var}[V_i] = c_i^2$   
 nel senso che l'errore di misura dipende dalla concentrazione.

Quando la concentrazione è piccola, si ha un errore piccolo ed un errore grande quando la concentrazione è grande, come si può notare nella figura dove si tratta una concentrazione o molecole variabile nel tempo.

È detto coefficiente di variazione  $\text{Cov}$   
 è matrice varianza è una matrice diagonale



$$\text{Var}[V] = \begin{bmatrix} \text{Var}[V_1] & \text{Cov}[V_1, V_2] & \dots & 0 \\ 0 & \text{Var}[V_2] & \dots & \dots \\ \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \text{Var}[V_n] \end{bmatrix}$$

Abbiamo ipotizzato che le misure siano uncorrelate  
 tra loro

Otteniamo

$$\text{Var}[V] = \sigma^2 \begin{bmatrix} c_1^2 & & & 0 \\ & c_2^2 & & \\ & & \dots & \\ 0 & & & c_n^2 \end{bmatrix} = \sigma^2 \begin{bmatrix} y_1^2 & & & 0 \\ & y_2^2 & & \\ & & \dots & \\ 0 & & & y_n^2 \end{bmatrix} = \sigma^2 \Psi$$

non si conosce la concentrazione vera ( $c_i$ ), che è la grandezza che si sta misurando.

Al suo posto si considera la misura ( $y_i$ ) e si è posto

$$\Psi = \begin{bmatrix} y_1^2 & & & 0 \\ & y_2^2 & & \\ & & \dots & \\ 0 & & & y_n^2 \end{bmatrix} \quad (\text{matrice nota delle grandezze misurate})$$

Si è considerato come se non succede un modello e un costante ma non si conosce il valore di  $k$ , definito con  $\sigma$ .

Il modello dell'errore di misura è parzialmente noto: si sa che la misura più grande ha una maggior variabilità ma con coefficiente di proporzionalità  $k$  ignoto.

-) Teorema di Gauss-Markov

Se vale l'ipotesi  $I_2$ , definiamo

$$E := Y - \Phi \theta$$

Allora, la cifra di merito

$$J^M(\theta) := E^T \Psi^{-1} E$$

Andrey Andreevich  
 Markov: 1856  
 - San Pietroburgo 1922  
 matematico e statistico  
 tra i russi

è minimizzata da

$$\hat{\theta}^M = (\Phi^T \Psi^{-1} \Phi)^{-1} \Phi^T \Psi^{-1} Y$$

$\hat{\theta}^M$  è detto stimatore di Gauss o stimatore BLUE.

Inoltre:

1)  $\hat{\theta}^M$  è, tra tutti gli stimatori lineari e non polarizzati, quello che minimizza la varianza di  $\hat{\theta}$ .

$$\text{Var}[\hat{\theta}^M] = \sigma^2 (\Phi^T \Psi^{-1} \Phi)^{-1}$$

Nota: BLUE = Best  
Linear  
Unbiased  
Estimator

- Best è riferito alla varianza che risulta minima.

- Linear indica che dipende linearmente da

- Unbiased indica che lo stimatore di Gauss è non polarizzato, e per questo tipo di stimatore, l'obiettivo è minimizzare la varianza.

Osservazioni:

1) se gli errori  $\epsilon_i$  hanno tutti la stessa varianza e la varianza è  $\sigma^2$

$$\text{Var}[\epsilon_i] = \sigma^2 I$$

essendo  $\sigma^2$  la varianza comune degli errori ne consegue che

$$\Psi = I$$

$$\hat{\theta}^M = \hat{\theta}^LS$$

(vedi pag. 204)

2) se gli errori  $\epsilon_i$  sono incorrelati e tali che

$$\text{Var}[\epsilon_i] = \sigma_{\epsilon_i}^2$$

se, cioè, hanno tutte varianze diverse, allora

$$A \text{vec}[V] = \Psi$$

$$\Psi = \begin{bmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \\ & & \ddots & \\ 0 & & & \sigma_m^2 \end{bmatrix}$$

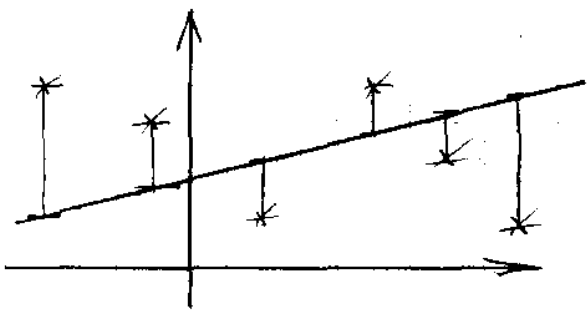
e risulta

$$J^M = \epsilon^T \Psi^{-1} \epsilon =$$

$$= \begin{bmatrix} \epsilon_1 & \epsilon_2 & \dots & \epsilon_m \end{bmatrix} \begin{bmatrix} 1/\sigma_1^2 & & & 0 \\ & 1/\sigma_2^2 & & \\ & & \ddots & \\ 0 & & & 1/\sigma_m^2 \end{bmatrix} \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_m \end{bmatrix} =$$

$$= \frac{\epsilon_1^2}{\sigma_1^2} + \frac{\epsilon_2^2}{\sigma_2^2} + \dots + \frac{\epsilon_m^2}{\sigma_m^2}$$

Di conseguenza  $J^M$  non minimizza più la somma dei quadrati dei residui, bensì la somma dei quadrati dei residui pesati.



nello squario rappresentato si dimostra si nota che i dati hanno diversa importanza ed il modello rappresentato è quello dei minimi quadrati.

È difficile avere di seguito il più fedelmente i dati considereremo quelli a varianza grande.

Il metodo è detto anche dei minimi quadrati pesati.

### Weighted least squares.

È facile capire come definire i pesi di un dato valore e il ruolo delle varianze.

**Esercizio 1**

Avete a disposizione i dati relativi ad una popolazione di studenti universitari iscritti a Ingegneria. Per ognuno di essi si conoscono il peso, l'altezza, il sesso ed i voti degli esami di Analisi A, Analisi B, Algebra e Geometria, Fisica A. Vi viene richiesto di analizzare i dati. Con riferimento alle variabili *peso* e *altezza* degli studenti maschi dovete:

1. Stimare media, varianza e deviazione standard.
2. Descrivere in modo approssimato la ddp mediante istogramma.
3. Stimare e disegnare le funzioni di distribuzione.
4. Dovete inoltre valutare se ciascuna VC è distribuita approssimativamente in modo gaussiano. A tale scopo vi viene richiesto di confrontare i grafici al punto (3) con la funzione di distribuzione di una gaussiana avente come media e varianza la media e la varianza campionaria calcolate al punto (1).

Vi viene inoltre richiesta un'analisi bivariata delle VC peso e altezza degli studenti maschi. Si richiede di

5. Stimare la matrice di covarianza e la matrice di correlazione.
6. Disegnare lo scatter plot dei dati *normalizzati*.

*Indicazioni operative.* Le analisi dell'Esercizio 1 sono implementate nel file *statistiche.m*. Seguendo le istruzioni che vi verranno date in aula, copiate i files *statistiche.m*, *creadat.m*, *vccont.m*, *vcdiscr.m*, ed eventuali file *\*.mat* nella vostra directory. Analizzate quindi il file *statistiche.m* per familiarizzarvi con i principali comandi MatLab.

**Esercizio 2**

- a. Considerando le VC *voti* (di Analisi A, Analisi 2, Geometria e Algebra e Fisica A) svolgere le analisi dei punti (1)-(4) (si consideri l'intera popolazione studentesca).
- b. Calcolare la matrice di covarianza e di correlazione del vettore di VC formato dai 4 voti, dal peso e dall'altezza.
- c. Con riferimento alle coppie di variabili casuali (voti Analisi A, voti Analisi B), (voti Fisica A, voti Geometria), (altezza, voti Analisi 1) disegnare gli scatter plot dei dati normalizzati.

**Esercizio 3**

- a. Progettare lo stimatore lineare ottimo (in termini di errore quadratico medio) che in base al peso dello studente ne predice l'altezza (si suppone di conoscere il sesso dello studente).  
*Suggerimento: si usi il valore atteso condizionato come predittore e si ricordi che per V.C. X e Y congiuntamente gaussiane,*  
$$E[X|Y]=E[X] + \text{Cov}[X,Y] \text{Var}[Y]^{-1}(Y - E[Y])$$
- b. Progettare uno stimatore che in base ai voti di Analisi 1, Algebra e Geometria, Fisica 1 predice quello di Analisi 2.

- Esercizio 1

- load MNPV2003

carica il file MNPV2003.mat  
che contiene dati relativi a  
studenti

~~0 = f~~

- who mostra le variabili caricate

- ) s = sesso (0 = maschio - 1 = femmine)
- ) h = altezza in cm
- ) w = peso in kg
- ) a1 = voto di Analisi 1
- ) a2 = " " " 2
- ) f1 = voto di Fisica A
- ) g = voto di Algebra e Geometria
- ) cd1 = ?
- ) sede

-) D = [s h w a1 a2 f1 g cd1 sede]  
carica e mostra i dati in una matrice D

-) mean(D)  
calcola la media dei dati caricati in D e relativi  
alle 8 colonne.

0,1883 è la media della 1<sup>a</sup> colonna ed indica  
la % di le ragazze non a'ca il 18%

-) M = mean(D) memorizza le medie nella matrice  
M

-) mh = M(2) memorizza l'altezza media nelle  
variabile mh

-) std(D) calcola la deviazione standard delle 8

colonne di dati

2/1

- )  $var = std(D)^2$  calcola le  $\sigma$  varianze e le inserisce nelle variabile var.  
-)  $vh = var(D)$

inserisce nelle variabile vh la variante relativa ~~in~~ alle altezze

- )  $ds = std(D)$  memorizza le deviazioni standard nel vettore ds

- )  $dsh = ds(D)$  memorizza in dsh le deviazioni standard dell'altezza.

- )  $hist(h, 14)$  visualizza l'istogramma delle altezze con 14 bins

- )  $covcoeff(D)$  costruisce la matrice dei coefficienti di correlazione

Si noti che la correlazione tra peso e altezza è negativa ( $0 = 11 - 1 = 10$ ).

Peso e altezza si correlano poco con i voti.

I voti si correlano abbastanza bene tra loro.

- )  $cov(D)$  crea la matrice delle covarianze, poco significativa.

N.B. la soluzione completa dell'esercizio si trova su rete.