

# Sistemi Intelligenti Stimatori e sistemi lineari - III

Alberto Borghese

Università degli Studi di Milano  
Laboratory of Applied Intelligent Systems (AIS-Lab)  
Dipartimento di Informatica  
[borgnese@di.unimi.it](mailto:borgnese@di.unimi.it)



A.A. 2019-2020

1/31

<http://borgnese.di.unimi.it>



## Overview



Densità di probabilità

**Funzione di verosimiglianza**

Stima alla massima verosimiglianza

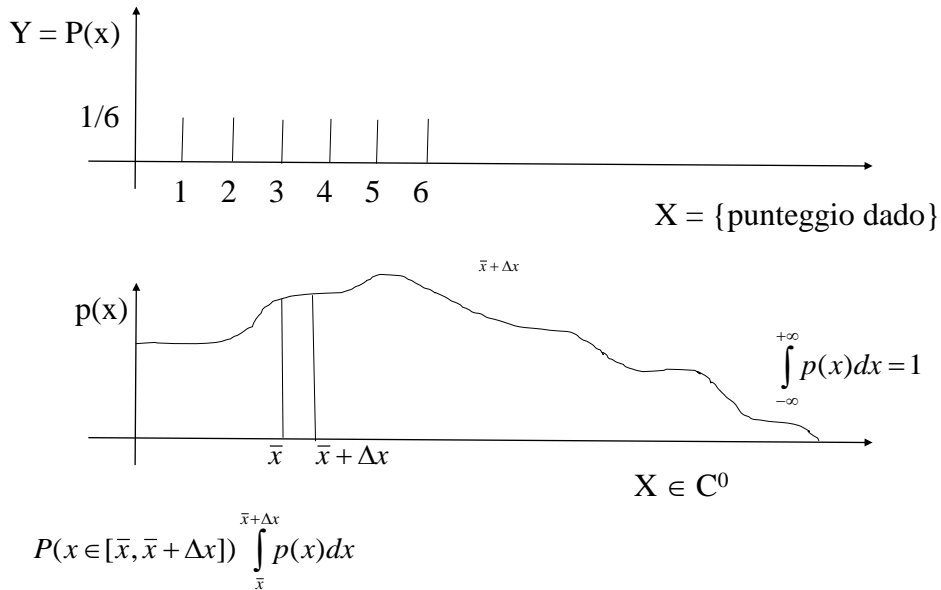
A.A. 2019-2020

2/31

<http://borgnese.di.unimi.it>



## La probabilità nel caso continuo



## Definizione di $p(x)$

Caso discreto: prescrizione della probabilità per ognuno dei finiti valori che la variabile  $X$  può assumere:  $P(X)$ .

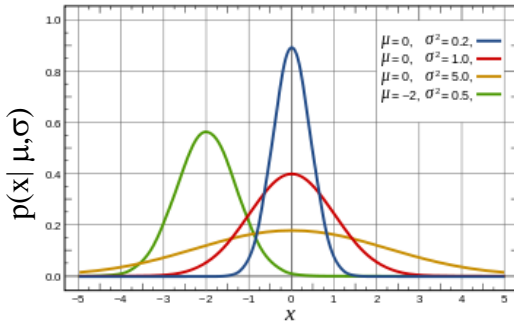
Caso continuo: i valori che  $X$  può assumere sono infiniti. Devo trovare un modo per definirne la probabilità. Descrizione **analitica** mediante la funzione densità di probabilità. Si considera la probabilità che  $x$  cada in un certo intervallo.

Valgono le stesse relazioni del caso discreto, dove alla somma si sostituisce l'integrale.

$$P(X = x \in [\bar{x}, \bar{x} + \Delta x]) = \int_{\bar{x}}^{\bar{x} + \Delta x} \int_{-\infty}^{+\infty} p(x, y) dx dy$$



## Distribuzioni notevoli: la Gaussiana



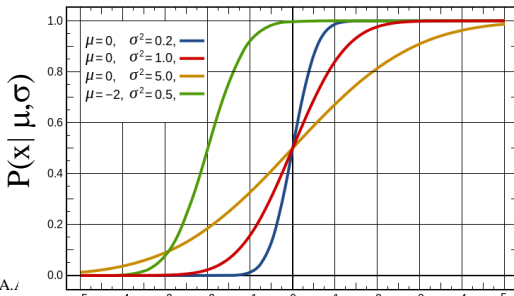
$$p(x | \mu, \sigma) = \frac{1}{(\sqrt{2\pi})^D} \cdot \exp\left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2\right]$$

D = dimensione, in questo caso D = 1

$$\Pr(|X - \mu| < \sigma) = 0.68268$$

$$\Pr(|X - \mu| < 2\sigma) = 0.95452$$

$$\Pr(|X - \mu| < 3\sigma) = 0.9973$$



<http://borghese.di.unimi.it/>



## I momenti di una variabile statistica



$$\mu^k(X) = \int_{-\infty}^{+\infty} (x - a)^k p(x) dx \quad \text{Momento rispetto ad a, solitamente alla media}$$

$$E[X] = \int_{-\infty}^{+\infty} (x - \mu) p(x) dx \quad \text{Valore atteso (Expected value) di X = media distribuzione}$$

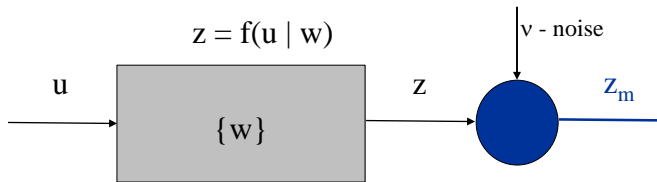
$$E[(X - \mu)^2] = \int_{-\infty}^{+\infty} (x - \mu)^2 p(x) dx \quad \text{Varianza } (\sigma^2)$$

$$E[(X - \mu)^3] = \int_{-\infty}^{+\infty} (x - \mu)^3 p(x) dx \quad \text{Asimmetria}$$

$$E[(X - \mu)^4] = \int_{-\infty}^{+\infty} (x - \mu)^4 p(x) dx \quad \text{Kurtosi - peso delle code di p(x)}$$



## Modello



$u$  – causa  $\Rightarrow z_m$  – effetto (misurato on errore)

Control / Classification / Prediction: determine  $\{z\}$  from  $\{u\}, \{w\}$

**Inverse problem: determine cause  $\{u\}$  from  $\{z_m\}, \{w\}$**

**Inverse problem: Identification: determine  $\{w\}$  from  $\{u\}, \{z_m\}$  - Learning**



## Overview

Densità di probabilità

**Funzione di verosimiglianza**

Stima alla massima verosimiglianza



## Probabilità di un certo insieme di misure



$$p(z_{1m}, z_{2m}, z_{3m}) = \prod_i p(z_{im})$$

$z = f(u | w)$     misuro  $\{z_i\}$  in corrispondenza di  $\{u_i\}$  con i parametri  $\{w_j\}$

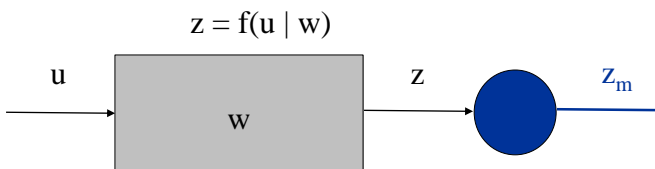
Avrò che:  $f(u_i, w) = z_{i,m} + z_i + v_i \Rightarrow f(u_i, w) - z_{i,m} = v_i$  the noise

Se le misure sono indipendenti posso scrivere che la probabilità di ottenere le misure  $z_{1m}, z_{2m}, z_{3m} \dots$  è:

$$p(z_{1m}, z_{2m}, z_{3m}) = \prod_i p(z_{im})$$



## Dipendenza delle misure



Le misure dipendono dal valore dalle variabili in ingresso  $u$  e dai parametri  $w$ .

$$p(z_{1m}, z_{2m}, z_{3m}) = \prod_i p(z_{im} | u_i, w) = \prod_i p(f(u_i; w) | u_i, w)$$

Tanto più i parametri saranno corretti tanto maggiore sarà la probabilità di avere  $y_m$  dal modello.



# Fitting di una retta



Vogliamo stimare i parametri di una retta:  $z = f(u | w) = m u + q$ , con  $m$  e  $q$  incogniti:  $W = \{m, q\}$

**La retta è un modello lineare.**

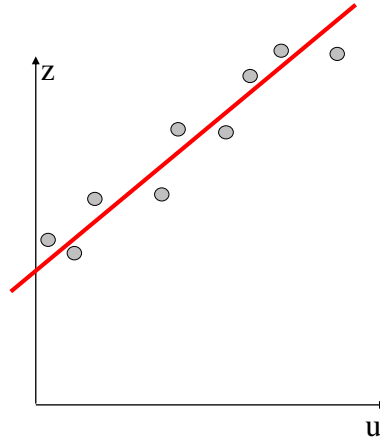
Abbiamo a disposizione  $N$  misure rumorose effettuate:  $z_{im}$  che sono funzione di  $u$  e  $w$ :

$$z_{im} = \{z_{im}; u_i, w\}, \text{ con rumore su } z_i$$

Supponiamo che le  $z_{im}$  siano affette da **errore Gaussiano a media nulla**. In pratica:

$$z_{im} = z_i + v_i = m u_i + q + v_i \text{ dove } v_i \text{ è l'errore di misura.}$$

$v_i = z_{im} - m u_i + q$  è **Gaussiano a media nulla**.



Possiamo anche scrivere che:  $z_{im} = z_i + G(\mu, \sigma^2)$  indica una distribuzione monodimensionale gaussiana a media  $\mu$  e varianza  $\sigma^2$ . Errore di misura:  $G(0, \sigma^2)$



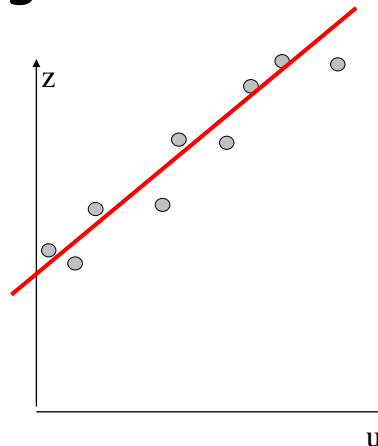
# Stima ai minimi quadrati e verosimiglianza



Per ogni punto, dovrebbe valere  $z_i = m u_i + q$ .

Ma c'è l'errore di misura, misuriamo in realtà  $z_i + v_i$ .

Cerchiamo i parametri  $m$  e  $q$  che sono più verosimili.



Cosa vuol dire che sono più verosimili?  
Quanto sono più verosimili?



## Funzione di verosimiglianza

- Siano date **N variabili casuali indipendenti**... Quale è la **probabilità di misurare il vettore**  $[z_{1m}, \dots, z_{Nm}]$ ?

$$p(z_{1m}, z_{2m}, \dots, z_{Nm}) = p(z_{1m}) \cdot p(z_{2m}) \cdot \dots \cdot p(z_{Nm}) = L(z_{1m}, z_{2m}, \dots, z_{Nm})$$

- La probabilità congiunta è il prodotto delle probabilità semplici.
- Questa è la **Funzione di verosimiglianza** o **funzione di Likelihood**,  $L(\cdot)$
- In questo caso le  $z$  sono legate alle  $u$  da  $f(u, w)$ .



## Overview

Funzione di verosimiglianza

**Stima alla massima verosimiglianza**



## Stima alla massima verosimiglianza

- Se **massimizziamo**  $L=L(z | u, w)$  rispetto a  $w$
- troviamo i parametri  $w$  tali per cui è massima la probabilità di misurare il vettore di dati  $\mathbf{z}_m = \{z_{im}, i=1 \dots N\}$ .
- **Stima alla massima verosimiglianza.**
- Più in generale, le variabili possono avere densità di probabilità diverse, ciascuna descritta da un set di parametri. I parametri delle diverse densità di probabilità possono essere calcolati utilizzando l'approccio alla massima verosimiglianza...
- La funzione di verosimiglianza dipende dai parametri che definiscono le densità di probabilità delle variabili casuali che entrano nella verosimiglianza...
- Massimizzando la funzione di verosimiglianza rispetto a tali parametri se ne effettua la stima in modo tale che il vettore osservato  $\mathbf{z} = \{z_{im}\} i=1 \dots N$  sia massimamente probabile (massima verosimiglianza).



## Stima alla massima verosimiglianza per modello lineare



- Impostiamo il problema scrivendo la funzione di verosimiglianza e massimizzando tale funzione rispetto a  $m$  e  $q$ ...
- Scriviamo prima di tutto la densità di probabilità di ottenere  $z_{im}$  per ciascun dato:

$$p(z_{im} | m, q, u_i) = \frac{1}{\sqrt{2\pi\sigma}} \cdot \exp\left[-\frac{1}{2}\left(\frac{z_{im} - (mu_i + q)}{\sigma}\right)^2\right]$$

Dove  $m$  e  $q$  non sono note.  $z_{im} = Z_i + v_i = m u + q + v_i$





# Stima a massima verosimiglianza



Sapendo che le misure sono indipendenti, possiamo scrivere la probabilità di ottenere le N misure  $\{z_{im}\}$ : funzione di verosimiglianza.

Scriviamo il logaritmo negativo della verosimiglianza per **errore Gaussiano**.

$$p(z_{1m}, z_{2m}, \dots, z_{Nm}) = p(z_{1m}) \cdot p(z_{2m}) \cdot \dots \cdot p(z_{Nm}) = L(z_{1m}, z_{2m}, \dots, z_{Nm}) = \prod_i p(z_{im})$$

$$-\ln(L(\cdot)) = -\ln \prod_{i=1}^N p(z_{im}) = f(z_{1m}, z_{2m}, \dots, z_{Nm}; m, q; u_1, u_2, \dots, u_N) = -\sum_{i=1}^N \ln \left\{ \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[ -\frac{1}{2} \left( \frac{z_{im} - mu_i - q}{\sigma} \right)^2 \right] \right\} =$$

$$-\sum_{i=1}^N \ln \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \sum_{i=1}^N \left[ -\frac{1}{2} \left( \frac{z_{im} - mu_i - q}{\sigma} \right)^2 \right] =$$

$$= -\sum_{i=1}^N \ln \left( \frac{1}{\sqrt{2\pi}\sigma} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^N (z_{im} - mu_i - q)^2$$



# Stima a massima verosimiglianza



- E massimizziamo  $L(\cdot)$  ponendo a zero le derivate prime rispetto a  $m$ :

$$\frac{\partial f(z_{1m}, z_{2m}, \dots, z_{Nm}; m, q; x_{1m}, x_{2m}, \dots, x_{Nm})}{\partial m} = \frac{\partial}{\partial m} \left[ -\sum_{i=1}^N \ln \left( \frac{1}{\sqrt{2\pi}\sigma} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^N (z_{im} - m \cdot u_i - q)^2 \right] =$$

$$= 0 + \frac{1}{2\sigma^2} \sum_{i=1}^N (z_{im} - m \cdot u_i - q) \cdot 2 \cdot (-u_i) =$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^N (z_{im} - m \cdot u_i - q) \cdot x_i = 0 \Rightarrow \sum_{i=1}^N (z_{im} - m \cdot u_i - q) \cdot u_i = 0 \Rightarrow$$

$$\left[ \sum_{i=1}^N (z_{im} \cdot u_i) \right] - m \cdot \left[ \sum_{i=1}^N (u_i^2) \right] - q \cdot \left[ \sum_{i=1}^N (u_i) \right] = 0 \Rightarrow$$

$$m \cdot \left[ \sum_{i=1}^N (u_i^2) \right] + q \cdot \left[ \sum_{i=1}^N (u_i) \right] = \left[ \sum_{i=1}^N (z_{im} \cdot u_i) \right]$$

1<sup>a</sup> equazione



## Stima a massima verosimiglianza



... e a q:

$$\begin{aligned}
&= 0 + \frac{1}{2\sigma^2} \sum_{i=1}^N (z_{im} - m \cdot u_i - q) \cdot 2 \cdot (-1) = \\
&= \frac{1}{\sigma^2} \sum_{i=1}^N (z_{im} - m \cdot u_i - q) = 0 \Rightarrow \sum_{i=1}^N (z_{im} - m \cdot u_i - q) = 0 \Rightarrow \\
&\left[ \sum_{i=1}^N (z_{im}) \right] - m \cdot \left[ \sum_{i=1}^N (u_i) \right] - q \cdot \left[ \sum_{i=1}^N (1) \right] = 0 \Rightarrow \\
&m \cdot \left[ \sum_{i=1}^N (u_i) \right] + q \cdot \left[ \sum_{i=1}^N (1) \right] = \left[ \sum_{i=1}^N (z_{im}) \right]
\end{aligned}$$

2<sup>a</sup> equazione



## Stima a massima verosimiglianza



$$\left[ \sum_{i=1}^N (u_i^2) \right] \cdot m + \left[ \sum_{i=1}^N (u_i) \right] \cdot q = \left[ \sum_{i=1}^N (z_{im} \cdot u_i) \right] \quad 1^\circ \text{ equazione}$$

$$\left[ \sum_{i=1}^N (u_i) \right] \cdot m + \left[ \sum_{i=1}^N (1) \right] \cdot q = \left[ \sum_{i=1}^N (z_{im}) \right] \quad 2^\circ \text{ equazione}$$

Le incognite, m e b, compaiono con esponente 1 => equazioni lineari in m e q  
Potrei risolvere per sostituzione



## Soluzione al problema di stima



Otengo un sistema di due equazioni in due incognite:

$$Ax = b$$

Dove:

$$A = \begin{bmatrix} \sum_{i=1}^N (u_i^2) & \sum_{i=1}^N (u_i) \\ \sum_{i=1}^N (u_i) & N \end{bmatrix} \quad b = \begin{bmatrix} \sum_{i=1}^N (z_{im} u_i) \\ \sum_{i=1}^N (z_{im}) \end{bmatrix}$$

$$x = [m \ q]'$$

A.A. 2019-2020

21/31

<http://borghese.di.unimi.it/>



## Esempio



$$z = 2u + 1 \quad m = 2; \ q = 1$$

Misuro con errore e ottengo:

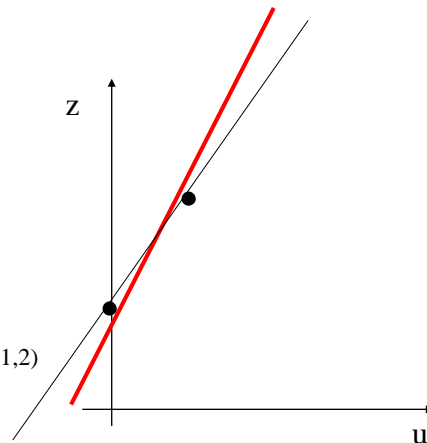
$$- u_1 = 1; \ z_1 = 2,8$$

$$- u_2 = 0; \ z_2 = 1,2$$

Quanto varranno le stime di  $m$  e  $q$ ?

$$\begin{aligned} \left[ \sum_{i=1}^N (u_i^2) \right] \cdot m + \left[ \sum_{i=1}^N (u_i) \right] \cdot q &= \left[ \sum_{i=1}^N (z_i \cdot u_i) \right] \\ \Rightarrow (1 \cdot 1 + 0 \cdot 0)m + (1 + 0)q &= (1 \cdot 2,8 + 0 \cdot 1,2) \end{aligned}$$

$$\begin{aligned} \left[ \sum_{i=1}^N (u_i) \right] \cdot m + \left[ \sum_{i=1}^N (1) \right] \cdot q &= \left[ \sum_{i=1}^N (z_i) \right] \\ \Rightarrow (1 + 0)m_e + 2q_e &= (2,8 + 1,2) \end{aligned}$$



$$m_e + q_e = 2,8 \quad m_e + 2q_e = 4 \quad \Rightarrow \text{per sottrazione} \quad q_e = 1,2; \ m_e = 1,6 \quad \text{NB } 2,8 = 1,6 \cdot 1 + 1,2$$

A.A. 2019-2020

22/31

<http://borghese.di.unimi.it/>



## Problema lineare



Ho  $N$  misure indipendenti:

$$A_i = [u_i \ 1]$$

$$\{z_{im} = m u_i + q + v\}$$

$$A x = b$$

$$A_{M \times 2} = \begin{bmatrix} u_1 & 1 \\ \cdot & \cdot \\ u_M & 1 \end{bmatrix}$$

$$x_{2 \times 1} = \begin{bmatrix} m \\ q \end{bmatrix}$$

$$b_{M \times 1} = \begin{bmatrix} z_1 \\ \cdot \\ z_M \end{bmatrix}$$



## Dal Sistema alla soluzione ai minimi quadrati



$$A^T A = \begin{bmatrix} u_1 & \cdot & u_M \\ 1 & \cdot & 1 \end{bmatrix} \begin{bmatrix} u_1 & 1 \\ \cdot & \cdot \\ u_M & 1 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (u_i^2) & \sum_{i=1}^N (u_i) \\ \sum_{i=1}^N (u_i) & N \end{bmatrix}$$

$$A^T b = \begin{bmatrix} u_1 & \cdot & u_M \\ 1 & \cdot & 1 \end{bmatrix} \begin{bmatrix} z_1 \\ \cdot \\ z_M \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (z_i \cdot u_i) \\ \sum_{i=1}^N (z_i) \end{bmatrix}$$

Equazioni normali:  $A^T A x = A^T b$   
sono le stesse ottenute dalla massima verosimiglianza



## Soluzione come problema di ottimizzazione



$$\text{Funzione costo: } (Ax - b)^2 = \sum_k v_k^2 = \|Ax - b\|^2$$

Assegno un costo al fatto che la soluzione  $x$ , non soddisfi tutte le equazioni, la somma dei residui associati ad ogni equazioni viene minimizzata. Geometricamente: viene trovato il punto a distanza (verticale) minima da tutte le rette.

$$\min_x \sum_k v_k^2 = \min_x (Ax - b)^2$$

$$A^T A x = A^T b$$

$$\frac{d}{dx} (Ax - b)^2 = 2A^T (Ax - b) = 0$$

$$x = (A^T A)^{-1} A^T b$$

NB le funzioni costo sono spesso quadratiche (problemi di minimizzazione convessi) perchè il costo cresce sia che il modello sovrastimi che sottostimi le misure. Inoltre, le derivate calcolate per imporre le condizioni di stazionarietà (minimo), sono relativamente semplici.



## Esempio - Caso 2D (2 parametri)



$N = 20$  punti

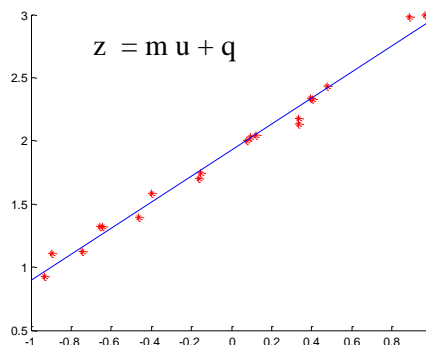
$\sigma_0^2 = 0.01$

$m$  reale = 1

$q$  reale = 2

$m$  stimato = 0.9931

$q$  stimato = 2.0106



Cosa vuol dire che  $\{m, q\}$  sono i più verosimili?  
Quanto sono più verosimili?



## Stima a massima verosimiglianza e minimi quadrati



$$A^T A x = A^T b$$

$$x = (A^T A)^{-1} A^T b$$

La soluzione a massima verosimiglianza, quando il rumore è Gaussiano a media nulla, coincide con la soluzione ai minimi quadrati del sistema lineare associato (la soluzione ai minimi quadrati è un caso particolare della stima alla massima verosimiglianza).

La soluzione è quella che minimizza lo scarto quadratico medio dei residui, ovvero è a minima varianza.

La stima a massima verosimiglianza è un approccio generale, e si presta a  $p(x)$  di qualsiasi forma. La Gaussiana consente di ottenere una formulazione lineare del problema.



## Stima ai minimi quadrati e verosimiglianza



- Nella soluzione ai minimi quadrati del sistema lineare  $Ax=b$  si definisce un vettore errore  $v = Ax - b$ ;
- Nel caso di soluzione “perfetta”  $v = 0$ ;
- Dal momento che abbiamo un numero di equazioni maggiore rispetto al numero di incognite, cerchiamo il vettore  $v$  a norma minima;
- In pratica cerchiamo  $x$  t.c.  $v^T v = \sum_i v_i^2$  è minimo.



## Relazione con la soluzione dei sistemi lineari



- Nel caso precedente le incognite erano  $(x_1, x_2)$  cioè le coordinate di un punto del piano.
- In questo caso le incognite sono  $(m, q)$  i parametri della retta.
- Nel caso precedente i dati erano i parametri di ogni retta  $(m_i, q_i / a_{1i}, a_{2i})$
- In questo caso i dati sono i punti misurati sulla retta  $(u_i, z_i)$



## Giustificazione statistica



- **C'è un solo insieme vero dei parametri**, mentre ci possono essere **infiniti universi di dati per effetto dell'errore di misura**.
- La domanda quindi più corretta sarebbe: "Dato un certo insieme di parametri, qual'è la probabilità che questo insieme di dati sia estratto?" (più correttamente si parla di densità di probabilità?)
- Cioè, **per ogni insieme di parametri, calcoliamo la probabilità che i dati siano estratti. Ovverosia la likelihood (verosimiglianza) dei dati, dato un certo insieme di parametri.**

La stima ai minimi quadrati dei parametri è equivalente a determinare i parametri che massimizzano la funzione di **verosimiglianza** sotto l'ipotesi di errore **Gaussiano a media nulla**.



## Overview



Densità di probabilità

Funzione di verosimiglianza

Stima alla massima verosimiglianza