








# Sistemi Intelligenti Relazione tra ottimizzazione e statistica - IV

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 1/34 <http://borgnese.di.unimi.it/>

## Sommario

**Analisi dell'affidabilità della stima**

Metodo del gradiente

Linearizzazione e metodo di Gauss-Newton

A.A. 2019-2020 2/34 <http://borgnese.di.unimi.it/>




## Valutazione della bontà della stima

$$x = (A^*A)^{-1}A^*b \iff \min_x \sum_k v_k^2 = \min_x (Ax-b)^2$$

Errore di modellizzazione Gaussiano a media nulla  $N(0, \sigma^2)$



$$\langle v_k \rangle = 0$$

$$\hat{\sigma}_0^2 = \sum_{k=1}^M (v_k^2) = |v|^2$$

Varianza della stima = varianza dell'errore di misura

.

A.A. 2019-2020 3/34 <http://borghese.di.unimi.it/>

## Valutazione della bontà della stima del singolo parametro, $x$

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

$$x = CA^*b \qquad \hat{\sigma}_0^2 = \sum_{m=1}^M (v_m^2)$$

Chiamiamo  $u$  e  $v$  le variabili casuali associate all'errore sui parametri e all'errore di modellizzazione, rispettivamente. Si suppone errore a media nulla e Gaussianamente distribuito.



$$u = \Delta x \qquad (x + \Delta x) = CA^*(b + v)$$

$\downarrow$

$$x = CA^*b \qquad \Delta x = CA^*v \qquad E[u] = 0$$

$C$  è la matrice di covarianza

A.A. 2019-2020 4/34 <http://borghese.di.unimi.it/>

## Impostazione del calcolo della correlazione tra i parametri

$\Delta x = C A' v$   
Abbiamo  $M$  parametri

Vogliamo individuare la correlazione tra due parametri  $i$  e  $j$ . Devo quindi determinare la loro correlazione:

$$\langle \Delta x_i, \Delta x_j \rangle$$

$$\begin{bmatrix} \Delta x_I^2 & \Delta x_I \Delta x_{II} & \dots & \Delta x_I \Delta x_M \\ \Delta x_{II} \Delta x_I & \Delta x_{II}^2 & \dots & \Delta x_{II} \Delta x_M \\ \dots & \dots & \dots & \dots \\ \Delta x_M \Delta x_I & \Delta x_M \Delta x_{II} & \dots & \Delta x_M^2 \end{bmatrix}$$



$\Delta x = C A' v \quad \Rightarrow \quad \Delta x' = v' A (C)'$

$\Delta x \Delta x' = C A' v v' A C' \Rightarrow$  Applicando l'operatore di media, si ottiene:

$$\langle \Delta x \Delta x' \rangle = C A' \langle v v' \rangle A C'$$

Dato che  $v$  sono i residui, e sono indipendenti, e tutte i punti di controllo hanno lo stesso tipo di errore di misura, si avrà che  $\langle v v' \rangle = I \sigma_0^2$ .

A.A. 2019-2020 5/34 <http://borghese.di.unimi.it/>



## Incertezza sulla stima dei parametri

$$\langle \Delta x \Delta x' \rangle = C A' I A C' \sigma_0^2 = C' \sigma_0^2 \quad \boxed{\langle \Delta x' \Delta x \rangle = C \sigma_0^2}$$

Segue che:  $\sigma^2(\Delta x_{ij}) = c_{ij} \sigma_0^2$  Varianza sulla stima del parametro,  $x_i$ .

Incertezza su  $z \rightarrow$  incertezza sui parametri stimati,  $x$

A.A. 2019-2020 6/34 <http://borghese.di.unimi.it/>

## Visione geometrica (1 parametro, m)

$$\sigma^2(\Delta x_{ij}) = c_{ij} \sigma_0^2 \quad \boxed{\langle \Delta x^T \Delta x \rangle = C \sigma_0^2}$$

$u = m = z + \text{noise}$   
 $Ax = b + \text{noise}$

Determino la pendenza  $m$  della retta

Calcolo  $m$  e  $q$  ai minimi quadrati.

Quanto è sensibile questa stima? Cosa succede se, per effetto del noise, invece di misurare  $z$ , misuro  $z + v$ ?



$$C = (A^T A)^{-1} \quad \Rightarrow \quad A_{1 \times 1} = u \quad \Rightarrow \quad C = (u^T u)^{-1}$$

La varianza di  $m$  varierà in modo inversamente proporzionale a  $u^2$ . Il rumore viene cioè moltiplicato per  $1/u^2$ .

$$\sigma^2(m) = c_m \sigma_0^2$$

Tanto più prendo i punti lontani dall'origine tanto meglio riesco a stimare  $m$  (tangente angolo).

A.A. 2019-2020 7/34 <http://borghese.di.unimi.it/>

## Matrici di covarianza

Date  $N$  variabili casuali:  $x = [x_1, x_2, \dots, x_N]$  si può misurare la correlazione tra coppie di variabili. E' comodo rappresentare la correlazione tra variabili casuali in un'unica matrice detta **matrice di covarianza** come:

$$C = \begin{bmatrix} \sigma_{x_1 x_1} & \sigma_{x_1 x_2} & \cdot & \sigma_{x_1 x_N} \\ \sigma_{x_2 x_1} & \sigma_{x_2 x_2} & \cdot & \sigma_{x_2 x_N} \\ \cdot & \cdot & \cdot & \cdot \\ \sigma_{x_N x_1} & \sigma_{x_N x_2} & \cdot & \sigma_{x_N x_N} \end{bmatrix}$$

Varianza:  $\sigma_{x_i x_i} = \sigma^2_{x_i}$  N parametri

Covarianza:  $\sigma_{x_i x_j} = \sigma_{x_j x_i} \quad i \neq j$  (N-1)<sup>2</sup>/2 parametri

A.A. 2019-2020 8/34 <http://borghese.di.unimi.it/>



## Correlazione tra coppie di parametri



Date due variabili casuali:  $x_i, x_j$ , l'indice di correlazione misura quanto le coppie di variabili estratte:  $p(x_i, x_j)$  stanno su una retta:

$$r = \frac{M_{x_i x_j} - M_{x_i} M_{x_j}}{\sigma_{x_i} \sigma_{x_j}} \quad -1 \leq r \leq +1$$

Definendo la covarianza tra  $x_i$  ed  $x_j$  come:

$$\sigma_{x_i x_j} = \frac{1}{N} \sum_i \sum_j (x_i - M_{x_i})(x_j - M_{x_j})$$

Dalla definizione di deviazione standard risulta:

$$r = \frac{\sigma_{x_i x_j}}{\sigma_{x_i} \sigma_{x_j}}$$



## Correlazione tra i parametri



$$\langle uu' \rangle = CA' IA C' \sigma_0^2 = C' \sigma_0^2$$

$$\langle uu' \rangle = C \sigma_0^2$$

Da cui si giustifica il nome di matrice di covarianza per C.

Segue che:  $\sigma^2(u_{ij}) = c_{ij} \sigma_0^2$  Varianza sulla stima del parametro.

$$-1 \leq r_{ij} = \frac{\langle u_i u_j \rangle}{\sqrt{\langle u_i \rangle^2 \langle u_j \rangle^2}} = \frac{c_{ij}}{\sqrt{c_i c_j}} \leq +1$$

Indice di correlazione tra il parametro i ed il parametro j  
(empiricamente si scartano parametri quando la correlazione è superiore al 95%)

Vanno rapportati alle dimensioni dei parametri coinvolti.



## La covarianza: momenti di 2 variabili statistiche



Covarianza =  $E[(x - \mu_x)(y - \mu_y)]$

Varianza =  $E[(x - \mu_x)(x - \mu_x)]$

Per due variabili indipendenti, la covarianza = 0, non variano assieme (covariano)

$$C = \begin{bmatrix} \sigma_x^2 & \sigma_x \sigma_y \\ \sigma_y \sigma_x & \sigma_y^2 \end{bmatrix}$$

```
>> x = randn(N,1);
>> y = randn(N,1);
>> temp = x.*y;
>> covarianza = mean(temp)
```



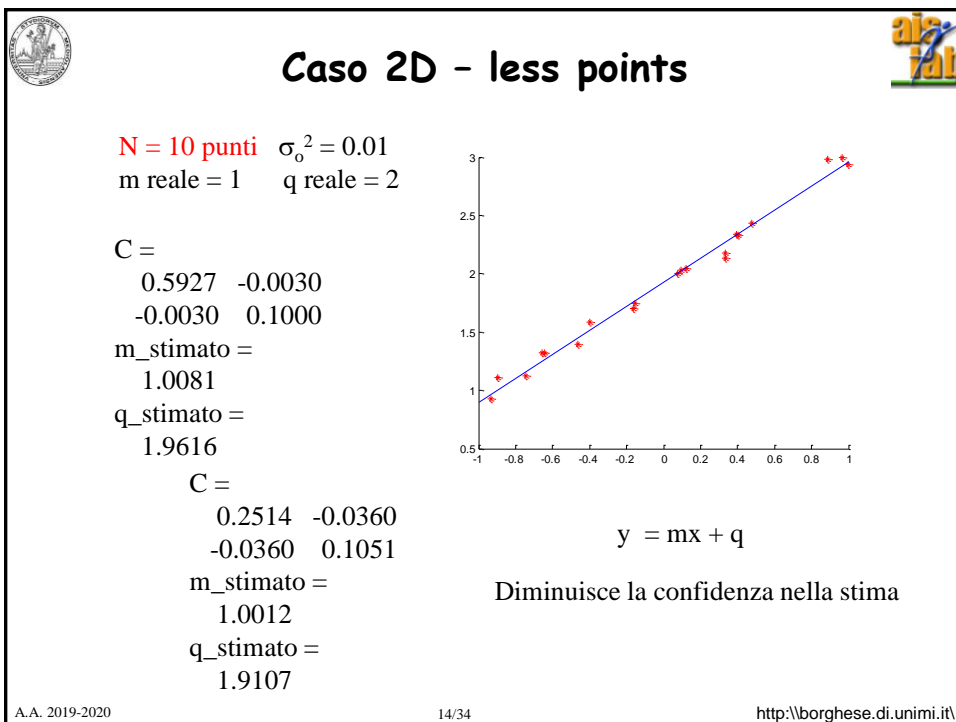
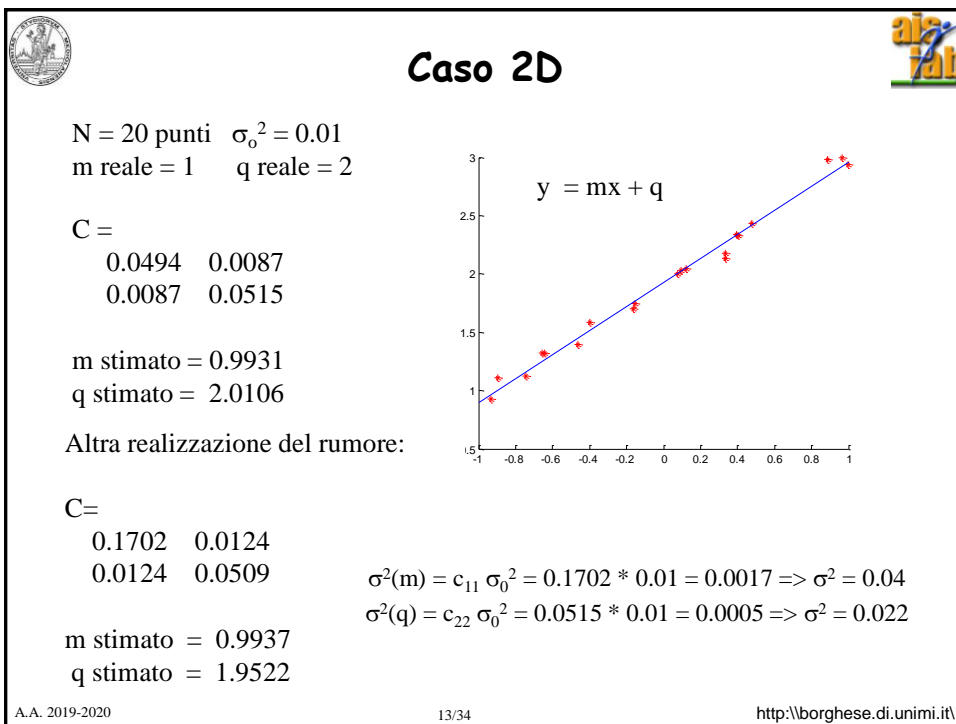
## Misura di correlazione su 2 parametri

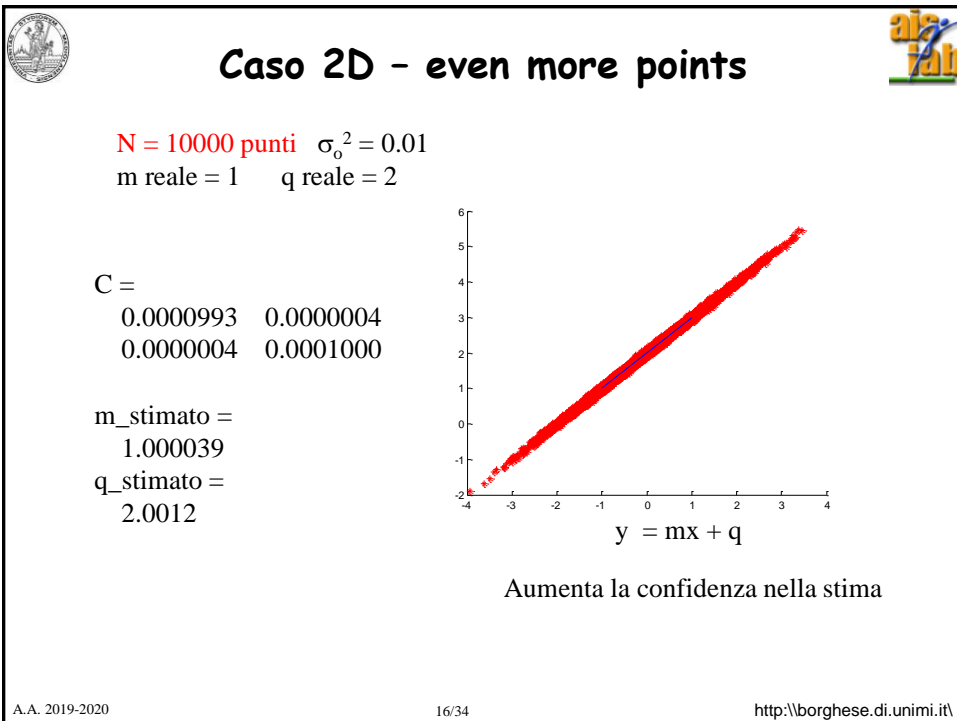
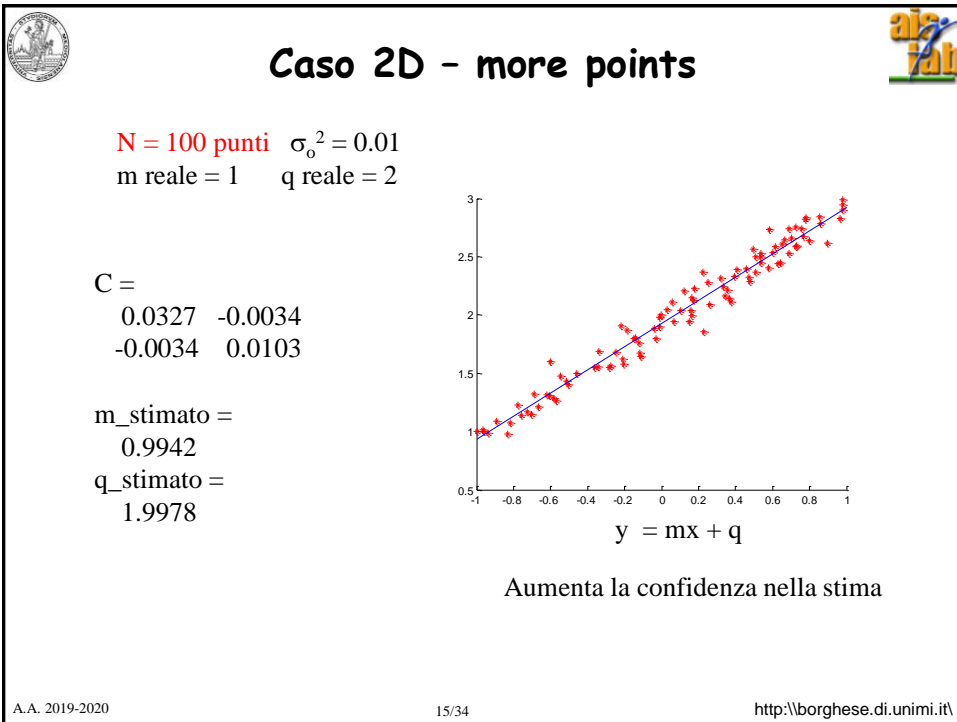


Misura la inter-dipendenza tra 2 variabili statistiche:


$$-1 \leq \frac{\sigma_{xy}}{\sigma_x \sigma_y} = c = \lim_{N \rightarrow \infty} \frac{\sum_k (x_k - \mu_x)(y_k - \mu_y)}{\sqrt{\sum_k (x_k - \mu_x)^2} \sqrt{\sum_k (y_k - \mu_y)^2}} \leq +1$$

```
>> x = randn(N,1);
>> y1 = randn(N,1);
>> y2 = x;
>> temp1 = x.*y1;
>> temp2 = x.*y2;
>> covarianza1 = mean(temp1)% Uncorrelated variables (c -> 1)
>> covarianza2 = mean(temp2)% Correlated variables (c = 0)
```










## Sommarario



Analisi dell'affidabilità della stima


**Metodo del gradiente**

Linearizzazione e metodo di Gauss-Newton


A.A. 2019-2020

17/34

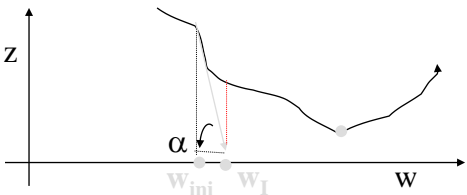
<http://borghese.di.unimi.it/>



## Minimizzazione tramite gradiente (metodo del primo ordine): 1 variabile



Tecnica del gradiente applicata alla minimizzazione di funzioni non-lineari di **una variabile, u**, e di **un parametro, w**:  $z = f(u | w)$ .



**La derivata, mi dà due informazioni:**


- 1) In quale direzione di  $w$ , la funzione  $f(\cdot)$  decresce.
- 2) Quanto rapidamente decresce.

Definisco uno spostamento arbitrario di  $w$ , maggiore la pendenza maggiore la variazione di  $z$ . Quindi faccio variare  $w$  nella direzione in cui diminuisce  $z$ :


**$\Delta w \propto -f'(w; u)$  dati  $P, w$ . La derivata viene calcolata rispetto a  $w$ .**

Occorre un'inizializzazione.  
Metodo iterativo.

mi.it\



## Esempio di applicazione tecnica del gradiente per funzioni di 1 variabile



*Supponiamo che il modello da noi considerato sia semplice:  $z = a u^2$*

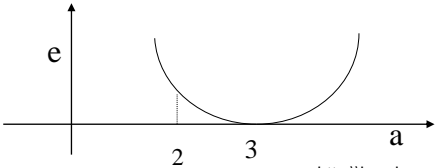
*Abbiamo un unico parametro da determinare:  $a$ . La funzione  $f(a,u)$  è lineare in  $a$  (la soluzione è semplice:  $a = z / u^2$ !!) ma applichiamo il metodo del gradiente:*

Misuriamo un punto sulla parabola,  $P(u,z) = P(1, 3)$  La soluzione è  $a = 3 / 1^2$ .


Applichiamo il metodo del gradiente. Partiamo da  $a_{ini} = 2$ . La parabola  $z = 2 u^2$  non passa per  $P$ . Come modificare  $a$  per farla passare per  $P$ ?

La funzione costo da minimizzare sarà:  $E = f(a | u,z) = (z - a u^2)^2$ .


Dato il punto  $P$ ,  $E = (3 - 2 \cdot 1^2)^2 = 1$



A.A. 2019-2020
19/34
<http://borghese.di.unimi.it/>



## Minimizzazione - underdamping



**Utilizziamo il metodo del gradiente:**

**Iterazione 1:**

**Passo 1:**  
 Calcoliamo l'espressione della derivata di  $f(a | u,z) \rightarrow f'(a) = 2 (z - a u^2) (-u^2)$   
 Calcoliamo la derivata nel punto  $P(1,3)$ , per  $a_{ini} = 2$ :  $f'(a) = 2 (3 - 2 \cdot 1^2) (-1^2) = -2$

**Passo 2:**  
 Calcoliamo l'incremento da dare al parametro  $a$ :  
 $a_1 = a_{ini} + \Delta a = a_{ini} - f'(a_{ini}; u,z) = 2 - (-2) \qquad a_1 = 2 + 2 = 4$



**Iterazione 2:**

**Passo 1:**  
 Calcoliamo l'espressione della derivata di  $f(a | u,z) \rightarrow f'(a) = 2 (z - a u^2) (-u^2)$   
 Calcoliamo la derivata nel punto  $P(1,3)$ , per  $a_1 = 2$ :  $f'(a) = 2 (3 - 4 \cdot 1^2) (-1^2) = +2$

**Passo 2:**  
 Calcoliamo l'incremento da dare al parametro  $a$ :  
 $a_{II} = a_1 + \Delta a = 4 - (+2) \qquad a_{II} = 4 - 2 = 2$

La soluzione oscilla tra  $a = +2$  e  $a = +4$  senza trovare il minimo che è in  $a = 3$ .

A.A. 2019-2020
20/34
<http://borghese.di.unimi.it/>

## Metodo del gradiente finale

Introduciamo un fattore di damping.  $\Delta w = -\alpha f'(w; u)$  dati  $P, w, \alpha < 1$ . La derivata  $f'(\cdot)$  viene calcolata rispetto a  $w$ . Scegliamo  $\alpha = 0.4$ .

**Iterazione 1:**

**Passo 1:**  
Calcoliamo l'espressione della derivata di  $f(a | u, z) \rightarrow f'(a) = 2(z - a u^2)(-u^2)$   
Calcoliamo la derivata nel punto  $P(1,3)$ , per  $a_{ini} = 2$ :  $f'(a) = 2(3 - 2 \cdot 1^2)(-1^2) = -2$

**Passo 2:**  
Calcoliamo l'incremento da dare al parametro  $a$ :  
 $a_I = a_{ini} + \Delta a = a_I = a_{ini} - \alpha f'(a_{ini}; u, z) = 2 - 0.4 [2(3 - 2 \cdot 1^2)(-1^2)] = 2 - 0.4 \cdot (-2) = 2.8$



**Iterazione 2:**

**Passo 1:**  
Calcoliamo l'espressione della derivata di  $f(a | u, z) \rightarrow f'(a) = 2(z - a u^2)(-u^2)$   
Calcoliamo la derivata nel punto  $P(1,3)$ , per  $a_I = 2.8$ :  $f'(a) = 2(3 - 2.8 \cdot 1^2)(-1^2) = -0.4$

**Passo 2:**  
Calcoliamo l'incremento da dare al parametro  $a$ :  
 $a_{II} = a_I + \Delta a = a_{II} = a_I - \alpha f'(a_I; u, z) = 2.8 - 0.4(-0.4) = 2.8 + 0.16 = 2.96$

E così via fino ad arrivare ad  $a = 3$  che si raggiunge asintoticamente


A.A. 2019-2020 21/34 http://borghese.di.unimi.it/


## Osservazioni

- Nel metodo del gradiente mi sposto lungo la tangente alla curva dell'errore per raggiungere il minimo.
- Lo spostamento lungo la tangente non mi porta al minimo direttamente.
- Se mi muovo velocemente lo supero.
- Se mi muovo lentamente arrivo lentamente.
- Esistono algoritmi che:
  - ◆ Verificano che il singolo passo di gradiente porta a un miglioramento della soluzione.
  - ◆ Determinano il passo ottimale di apprendimento per ogni passo,  $\alpha = \alpha_k$ .

A.A. 2019-2020 22/34 http://borghese.di.unimi.it/



## Sommarrio



Analisi dell'affidabilità della stima


Metodo del gradiente

Linearizzazione e metodo di Gauss-Newton


A.A. 2019-2020

23/34

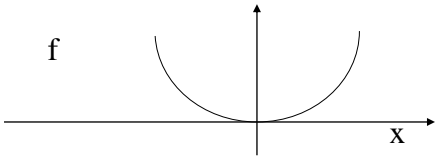
<http://borghese.di.unimi.it/>



## Linearizzazione



- Descrizione differenziale locale di una funzione.



$y = ax^2$



Posso descrivere localmente in ogni punto la funzione come:

$$y + dy = f(x) + f'(x)dx \Rightarrow dy = f'(x) dx \Rightarrow dy = 2ax dx$$

A.A. 2019-2020

24/34

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

## Linearizzazione

$y = f(x)$  viene linearizzata utilizzando il differenziale (retta tangente):

$$dy = f(x_o) + \left. \frac{df(x)}{dx} \right|_{x=x_o} dx = y_o + \left. \frac{df(x)}{dx} \right|_{x=x_o} dx$$



Si può vedere come sviluppo di Taylor arrestato al 1° ordine  
E' un'equazione lineare.

Per funzioni di più variabili,  $f(\mathbf{P}; \mathbf{W}) = 0$ , la linearizzazione nell'intorno di  $\mathbf{P}$ , si può scrivere come:

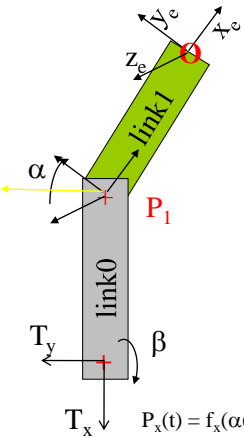

$$F(\mathbf{P}; \mathbf{W}) = F(\mathbf{P}_o; \mathbf{W}_o) + \sum_{j=1}^W \left. \frac{\partial F(\cdot)}{\partial w_j} \right|_{\mathbf{P}_o, \mathbf{W}_o} * dw_j = k \cdot \sum_{j=1}^W a_j * dw_j$$

E' un'equazione lineare che descrive il comportamento della funzione  $F(\cdot)$  nell'intorno del punto  $\mathbf{P}_o$  con i parametri  $\mathbf{W}_o$ .

A.A. 2019-2020
25/34
<http://borghese.di.unimi.it/>


## Esempio di sistema


$$\begin{aligned}
 P_x(t) &= f_x(\alpha(t), \beta(t), T_x(t), T_y(t) | l_0, l_1). \\
 P_y(t) &= f_y(\alpha(t), \beta(t), T_x(t), T_y(t) | l_0, l_1). \\
 P_z(t) &= f_z(\alpha(t), \beta(t), T_x(t), T_y(t) | l_0, l_1)
 \end{aligned}$$

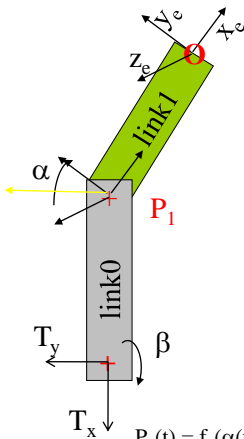
Voglio determinare  $\alpha$ ,  $\beta$ ,  $T_x$ ,  $T_y$  per ottenere un certo movimento dell'end-point.


A.A. 2019-2020
26/34
<http://borghese.di.unimi.it/>



## Esempio di "sistema"








Le funzioni legano la posizione dell'end point, uscita **P**, alla posizione degli angoli,  $\alpha$  e  $\beta$  e della posizione della base, **T**, che rappresentano gli ingressi.

$$\begin{aligned}
 P_x(t) &= f_x(\alpha(t), \beta(t), T_x(t), T_y(t) | l_0, l_1). \\
 P_y(t) &= f_y(\alpha(t), \beta(t), T_x(t), T_y(t) | l_0, l_1). \\
 P_z(t) &= f_z(\alpha(t), \beta(t), T_x(t), T_y(t) | l_0, l_1).
 \end{aligned}$$


A.A. 2019-2020

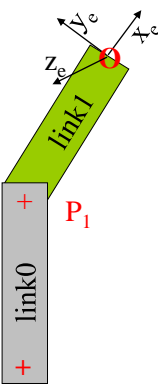
27/34

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



## In forma matriciale





$$\begin{bmatrix} x_e \\ y_e \\ z_e \end{bmatrix} = \begin{bmatrix} l_1 \cos(\alpha + \beta) + l_0 \cos \beta + T_x \\ -l_1 \sin(\alpha + \beta) - l_0 \sin \beta + T_y \\ 0 \\ 1 \end{bmatrix}$$


Sono equazioni non-lineari nei parametri.

Non riesco a calcolare  $\alpha$ ,  $\beta$ ,  $T_x$ ,  $T_y$  per ottenere una certa Posizione dell'end-point


A.A. 2019-2020

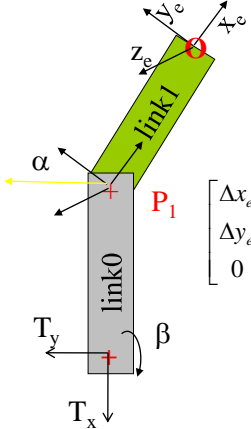
28/34

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




## Rappresentazione linearizzata Sistema lineare





$$\begin{bmatrix} \Delta x_e \\ \Delta y_e \\ 0 \end{bmatrix} = \begin{bmatrix} -l_1 \sin(\alpha + \beta) & -l_1 \sin(\alpha + \beta) - l_0 \sin \beta & 1 & 0 \\ -l_1 \cos(\alpha + \beta) & -l_1 \cos(\alpha + \beta) - l_0 \cos \beta & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \Delta \alpha \\ \Delta \beta \\ \Delta T_x \\ \Delta T_y \end{bmatrix}$$

$\alpha = 90$        $l_0 = 2,5$   
 $\beta = 0$        $l_1 = 2$




$$\begin{bmatrix} \Delta x_e \\ \Delta y_e \\ 0 \end{bmatrix} = \begin{bmatrix} -2 & -2 & 1 & 0 \\ 0 & -2.5 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \Delta \alpha \\ \Delta \beta \\ \Delta T_x \\ \Delta T_y \end{bmatrix}$$

**b = A x**


A.A. 2019-2020

29/34

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



## Minimizzazione di funzioni di più variabili



$\min(f(\mathbf{x}, \mathbf{w}))$  funzione costo od errore,  $\mathbf{w}$  vettore.

Modifico il valore dei pesi di una quantità proporzionale alla pendenza della funzione costo rispetto a quel parametro. La pendenza è una direzione nello spazio, non è più solamente destra / sinistra. Devo calcolare la derivata spaziale = **gradiente** della funzione costo,  $f(\cdot)$ .

Estensione della tecnica del gradiente a più variabili.



$d\mathbf{w} = -\alpha \nabla f(\mathbf{x}; \mathbf{w})$ , dato  $\mathbf{P}, \mathbf{W}$ .

Serve un'approssimazione iniziale per i parametri  $\mathbf{W}_{ini} = \{w_j\}_{ini}$ .

A.A. 2019-2020

30/34

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

## Metodo di Gauss-Newton

- L'idea:

Inizializzazione:



- Inizializzo i parametri ad un valore iniziale.

Iterazioni:

- 1) Linearizzazione delle equazioni.
- 2) Stima dell'aggiornamento dei parametri nel modello linearizzato ai minimi quadrati (soluzione ottimale, minimo del problema linearizzato).
- 3) Correzione dei parametri.

Può essere pesante perchè richiede l'inversione della matrice di covarianza. Spesso si preferiscono utilizzare metodi di ottimizzazione del primo ordine.

A.A. 2019-2020 31/34 <http://borghese.di.unimi.it/>

## In pratica

$\mathbf{y} = \mathbf{f}(\mathbf{x})$        $\mathbf{x}, \mathbf{y}$  vettori di N ed M elementi rispettivamente

$\mathbf{y}_0 = \mathbf{f}(\mathbf{x}_0)$        $\mathbf{x}_0, \mathbf{y}_0$  valore iniziale

Iterazione di (nella prima iterazione  $k = 0$ ):

- $\mathbf{d}\mathbf{y}_k + \mathbf{y}_k = (\Sigma \delta \mathbf{f}(\mathbf{x}) / \mathbf{d}\mathbf{x})_{\mathbf{x}_k} \mathbf{d}\mathbf{x} + \mathbf{f}(\mathbf{x}_k)$        $(\Sigma \delta \mathbf{f}(\mathbf{x}) / \mathbf{d}\mathbf{x})_{\mathbf{x}_k}$  are numbers!
- Si ottiene un sistema lineare
- Viene risolto come  $\mathbf{d}\mathbf{x}_k = (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{A}^T \mathbf{d}\mathbf{y}_k$
- Si aggiorna il valore di  $\mathbf{x}$  come  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}\mathbf{x}_k$

Fino a convergenza

A.A. 2019-2020 32/34 <http://borghese.di.unimi.it/>





## Evoluzione dei metodi del primo ordine



- $\alpha$  è un parametro critico. Se è troppo piccolo convergenza molto lenta, se è troppo grande overshooting.
- Ottimizzazione di  $\alpha$ . Ad ogni passo viene calcolato  $\alpha$  ottimale, per cui la funzione è decrescente (line search).



## Sommario



Analisi dell'affidabilità della stima

Metodo del gradiente

Linearizzazione e metodo di Gauss-Newton