

# Computazione per l'interazione naturale: Regressione lineare probabilistica



Corso di Interazione uomo-macchina II

Prof. Giuseppe Boccignone

Dipartimento di Informatica  
Università di Milano

boccignone@di.unimi.it  
[http://boccignone.di.unimi.it/IUM2\\_2014.html](http://boccignone.di.unimi.it/IUM2_2014.html)

## Regressione lineare //modelli probabilistici: stima di ML

- Assumiamo rumore gaussiano additivo

$$t = f(x; \mathbf{w}) + \epsilon$$

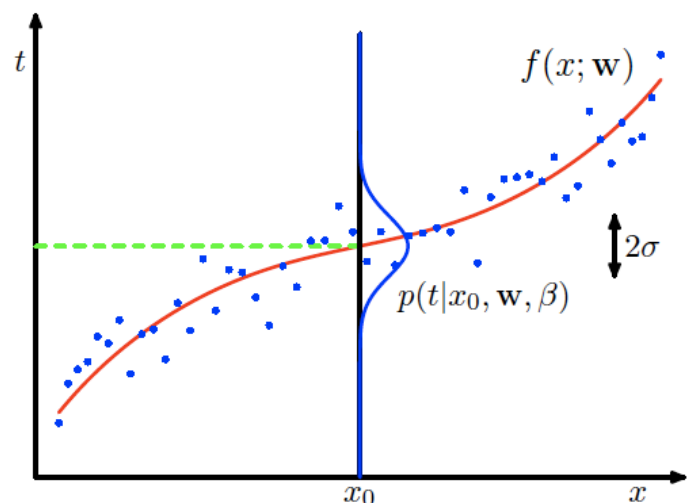
$$\epsilon \sim \mathcal{N}(0, \sigma)$$

precisione  $\beta = 1/\sigma^2$

- Allora

$$t|x \sim \mathcal{N}(f(x; \mathbf{w}), \sigma)$$

$$p(t|x) = \mathcal{N}(f(x; \mathbf{w}), \sigma)$$



# Regressione lineare

## //modelli probabilistici: stima di ML

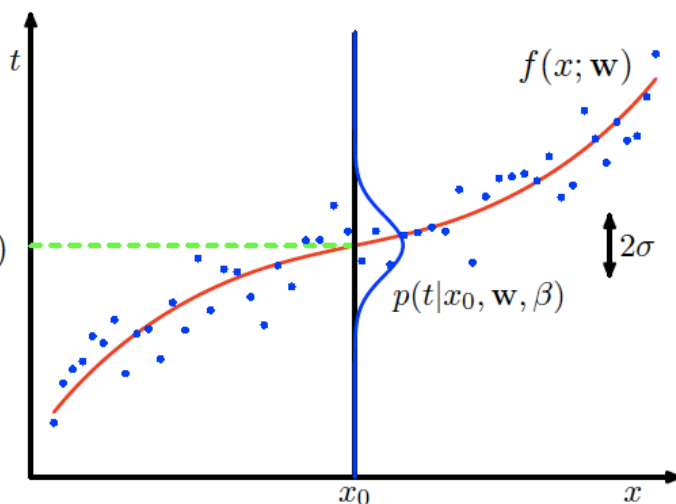
- Abbiamo N osservazioni nel training set

$$(x_1, t_1), \dots, (x_N, t_N) = (\mathbf{x}, \mathbf{t})$$

- Considerando la probabilità congiunta dei target

$$p(t_1, t_2, \dots, t_N | x_1, x_2, \dots, x_N, \mathbf{w}) = p(\mathbf{t} | \mathbf{x}, \mathbf{w})$$

- Vogliamo stimare i parametri  $\mathbf{w}$



# Introduzione alla Statistica Bayesiana

## //metodologia generale

- Obiettivo: predire un dato  $x$  sulla base di  $n$  osservazioni  $S = \{x_1, \dots, x_n\}$

- Tre step:

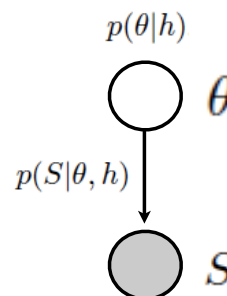
- Specifica del modello  $h$  con parametri di tuning  $\theta$  per la generazione dei dati in  $S$
- Inferenza (learning dei parametri)

$$p(\theta | S, h) = \frac{p(S | \theta, h) p(\theta | h)}{p(S | h)} \propto p(S | \theta, h) p(\theta | h)$$

- Predizione

$$p(x | S) = \int p(x, \theta | S) d\theta = \int p(x | \theta) p(\theta | S) d\theta$$

Modello o ipotesi  $h$



# Statistica Bayesiana

## //stime puntuali

---

- Rinunciando ad un approccio completamente Bayesiano, si possono ottenere stime puntuali dei parametri (ovvero i parametri diventano “numeri” e non VA)
  - Stima Maximum A Posteriori (MAP)

$$\theta_{MAP} = \operatorname{argmax}_{\theta} p(\theta|S)$$

- Stima di massima verosimiglianza (Maximum Likelihood, ML)

$$\theta_{ML} = \operatorname{argmax}_{\theta} p(S|\theta)$$

$$l(\theta) = \operatorname{argmax}_{\theta} \log p(S|\theta)$$

## Esempio: stima di massima verosimiglianza

### //caso Gaussiano

---

- Insieme di campioni  $x_1, \dots, x_n$  da distribuzione Gaussiana di parametri ignoti (**identicamente distribuiti**)
- Campioni estratti **indipendentemente**:

$$p(S|\theta) = \prod_{k=1}^n p(x_k|\theta)$$

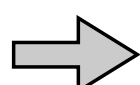
L'ipotesi  
i.i.d

- allora

$$\hat{\theta} = L(\theta) = \operatorname{argmax}_{\theta} p(S|\theta) = \operatorname{argmax}_{\theta} \prod_{k=1}^n p(x_k|\theta)$$

- oppure usando la **log-verosimiglianza**

$$l(\theta) = \operatorname{argmax}_{\theta} \log p(S|\theta) = \operatorname{argmax}_{\theta} \log \sum_{k=1}^n p(x_k|\theta)$$

  $\nabla_{\theta} \log p(S|\theta) = 0$

# Esempio: stima di massima verosimiglianza

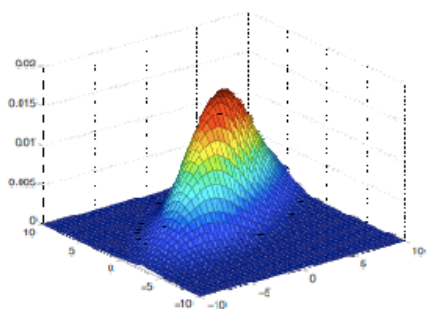
## // caso Gaussiano

- Un vettore aleatorio  $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} \sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

↑  
forma quadratica di  $\mathbf{x}$



$$(x_1, t_1), \dots, (x_N, t_N) = (\mathbf{x}, t)$$

- per il singolo campione

$$\log p(\mathbf{x}_i|\boldsymbol{\mu}) = -\frac{1}{2} \log\left((2\pi)^d |\boldsymbol{\Sigma}|\right) - \frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$$

# Esempio: stima di massima verosimiglianza

## // $\boldsymbol{\mu}$ della Gaussiana

- Se  $\mathbf{x} \sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  con covarianza nota e media da stimare
- per il singolo campione

$$\log p(\mathbf{x}_i|\boldsymbol{\mu}) = -\frac{1}{2} \log\left((2\pi)^d |\boldsymbol{\Sigma}|\right) - \frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$$

- derivando il secondo termine rispetto a  $\boldsymbol{\mu}$  (il primo è costante rispetto a  $\boldsymbol{\mu}$ )

$$\nabla_{\boldsymbol{\mu}} \log p(\mathbf{x}_i|\boldsymbol{\mu}) = \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$$

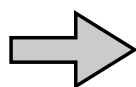
$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{s})^T \mathbf{W}(\mathbf{x} - \mathbf{s}) = -2\mathbf{W}(\mathbf{x} - \mathbf{s})$$

- per tutti i campioni

$$\nabla_{\boldsymbol{\mu}} \log p(S|\boldsymbol{\mu}) = \nabla_{\boldsymbol{\mu}} \sum_{i=1}^n \log p(\mathbf{x}_i|\boldsymbol{\mu}) = \sum_{i=1}^n \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$$

- ponendo uguale a 0

$$\sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}) = 0$$



$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

media empirica

# Esempio: stima di massima verosimiglianza

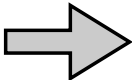
## // Parametri della Gaussiana

---

- Se  $\mathbf{x} \sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  con covarianza e media da stimare
- per il singolo campione

$$\log p(\mathbf{x}_i|\boldsymbol{\mu}) = -\frac{1}{2} \log \left( (2\pi)^d |\boldsymbol{\Sigma}| \right) - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

- derivando il secondo termine rispetto a  $\boldsymbol{\mu}$  e  $\boldsymbol{\Sigma}$  ripetendo il procedimento di prima



$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$
$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

media empirica

covarianza empirica

## Regressione lineare

### // stima di ML dei parametri

---

- Training set  $(x_1, t_1), \dots, (x_N, t_N) = (\mathbf{x}, \mathbf{t})$

L'ipotesi  
i.i.d

- La likelihood è

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) = \prod_{n=1}^N p(t_n|x_n, \mathbf{w}, \sigma) = \prod_{n=1}^N \mathcal{N}(f(x_n; \mathbf{w}), \sigma) \quad \rightarrow$$

$$\log p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) = \sum_{n=1}^N \log p(t_n|x_n, \mathbf{w}, \sigma)$$

$$= \sum_{n=1}^N \log \mathcal{N}(f(x_n; \mathbf{w}), \sigma)$$

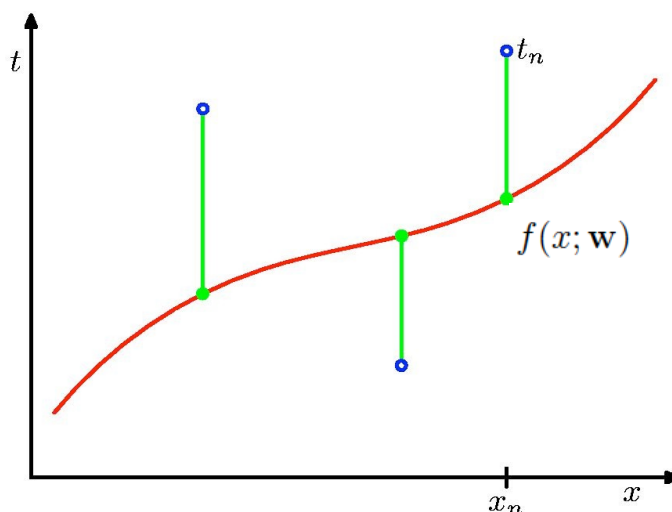
$$= \sum_{n=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{1}{2\sigma^2} |t_n - f(x_n; \mathbf{w})|^2 \right)$$

$$= -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N |t_n - f(x_n; \mathbf{w})|^2$$

funzione di costo

# Regressione lineare

//modelli probabilistici: stima di ML



Funzione di loss (errore) quadratica

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N |t_n - f(x_n; \mathbf{w})|^2$$

# Regressione lineare

//modelli probabilistici: stima di ML

- Abbiamo che

$$\max \log p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma) \iff \min E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N |t_n - f(x_n; \mathbf{w})|^2$$

punto stazionario

$$\begin{aligned} \frac{\partial \log L}{\partial \mathbf{w}} &= \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n (t_n - \mathbf{x}_n^T \mathbf{w}) \\ &= \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n t_n - \mathbf{x}_n \mathbf{x}_n^T \mathbf{w} = \mathbf{0}. \end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}, \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$$

$$\sum_{n=1}^N \mathbf{x}_n t_n \longrightarrow \mathbf{X}^T \mathbf{t}$$

$$\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \mathbf{w} \longrightarrow \mathbf{X}^T \mathbf{X} \mathbf{w}$$

$$\frac{\partial \log L}{\partial \mathbf{w}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{t} - \mathbf{X}^T \mathbf{X} \mathbf{w}) = \mathbf{0}.$$

# Regressione lineare

## //modelli probabilistici: stima di ML

- Abbiamo che

$$\max \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma) \iff \min E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N |t_n - f(x_n; \mathbf{w})|^2$$

$$\sum_{n=1}^N (t_n - \mathbf{x}^T \hat{\mathbf{w}})^2$$

$$\implies \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{t} - \mathbf{X}^T \mathbf{X} \mathbf{w}) = 0 \quad \text{punto stazionario}$$

$$\frac{\partial^2 \mathcal{L}}{\partial \mathbf{w} \partial \mathbf{w}^T} = -\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} \quad \text{negativa}$$

il punto stazionario è un max

$$\implies \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}.$$

soluzione di ML  
=  
equazioni normali  
per i minimi quadrati

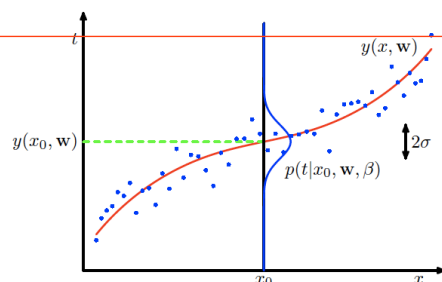
# Regressione lineare

## //modelli probabilistici: stima di ML

- Relazione tra parametri del modello e parametri stimati:

$$\begin{aligned} E_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma)} [\hat{\mathbf{w}}] &= \int \hat{\mathbf{w}} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma) dt \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \int \mathbf{t} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma) dt \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma)} [\mathbf{t}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} \\ &= \mathbf{w} \end{aligned}$$

la curva  
deterministica



- ML è stimatore unbiased!

# Regressione lineare

## //modelli probabilistici: stima di ML

---

- Analogamente possiamo anche fornire una stima ML della varianza del rumore, assumendo  $\mathbf{w} = \hat{\mathbf{w}}$

$$\frac{\partial \mathcal{L}}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^2} \sum_{n=1}^N (t_n - \mathbf{x}^T \hat{\mathbf{w}})^2 = 0$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}^T \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})$$

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}}) \\ &= \frac{1}{N} \mathbf{t}^T (\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \mathbf{t} \\ &= \frac{1}{N} (\mathbf{t}^T \mathbf{t} - \mathbf{t}^T \hat{\mathbf{t}}) \end{aligned}$$

$$\frac{\partial^2 \mathcal{L}}{\partial \sigma \partial \sigma} = -\frac{2N^2}{(\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})} = -\frac{2N}{\hat{\sigma}^2} \quad \begin{array}{l} \text{negativa} \\ \text{il punto stazionario è un max} \end{array}$$

# Regressione lineare

## //stima di ML dei parametri di generiche funzioni di x

---

$$f(x_n; \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}_i)$$

- La likelihood è

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) \quad \rightarrow$$

$$\ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \sum_{i=1}^N \ln \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

- Dove la funzione di costo o di errore è

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$



## Regressione lineare

//stima di ML dei parametri di generiche funzioni di x

---

$$f(x_n; \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}_i)$$

$$\max \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) \iff \min E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$

$$\begin{aligned} \Rightarrow \mathbf{0} &= \nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i)) \phi(\mathbf{x}_i)^T \\ &= \sum_{i=1}^N t_i \phi(\mathbf{x}_i)^T - \mathbf{w}^T \left( \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \right) \end{aligned}$$

$$\Rightarrow \mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad \begin{array}{l} \text{equazioni normali} \\ \text{per i minimi quadrati} \end{array}$$

## Regressione lineare

//modelli probabilistici: stima di ML

---

- Abbiamo generalizzato il risultato già ottenuto per la retta

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad \begin{array}{l} \text{equazioni normali} \\ \text{per i minimi quadrati} \end{array}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \quad \begin{array}{l} \text{design} \\ \text{matrix} \end{array}$$

# Regressione lineare

## //modelli probabilistici: stima di ML

---

- Possiamo computare il Minimo di  $E_D(\mathbf{w})$  anche in modo non analitico: metodo di discesa del gradiente (procedura iterativa)

- Inizializzazione

$$\mathbf{w}^{(0)} = (w_0^{(0)}, w_1^{(0)}, \dots, w_D^{(0)}) \longrightarrow E_D(\mathbf{w}^{(0)}) = \frac{1}{2} \sum_{i=1}^N (t_i - (\mathbf{w}^{(0)})^T \phi(\mathbf{x}_i))^2$$

- while ( ~ condizioneTerminazione)

$$w_j^{(i)} := w_j^{(i-1)} - \alpha \left. \frac{\partial E_D(\mathbf{w})}{\partial w_j} \right|_{\mathbf{w}=\mathbf{w}^{(i-1)}}$$

i = i+1;

# Regressione lineare

## //modelli probabilistici: stima di ML

---

- Possiamo computare il Minimo di  $E_D(\mathbf{w})$  anche in modo non analitico: metodo di discesa del gradiente (procedura iterativa)

- Inizializzazione

$$\mathbf{w}^{(0)} = (w_0^{(0)}, w_1^{(0)}, \dots, w_D^{(0)}) \longrightarrow E_D(\mathbf{w}^{(0)}) = \frac{1}{2} \sum_{i=1}^N (t_i - (\mathbf{w}^{(0)})^T \phi(\mathbf{x}_i))^2$$

- while ( ~ condizioneTerminazione)

$$\mathbf{w}^{(i)} := \mathbf{w}^{(i-1)} - \alpha (t_i - \mathbf{w}^{(i-1)} \phi(\mathbf{x}_i)) \phi(\mathbf{x}_i)$$

i = i+1;

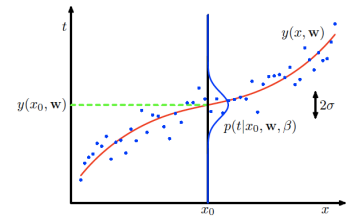
# Regressione lineare

## //predizione e incertezza

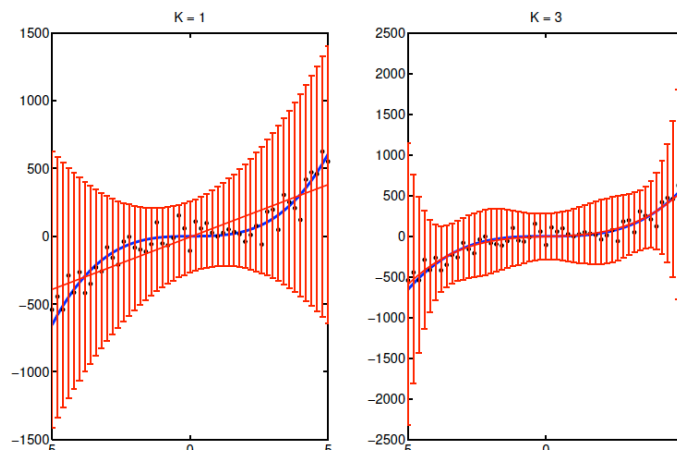
- Possiamo effettuare la predizione sui dati nuovi (test set)

$$t_{new} = \hat{\mathbf{W}}^T \mathbf{x}_{new}$$

- Qual è l'incertezza sulle stime del nostro modello?
- Usiamo la covarianza...



$$\hat{t}_{new} \pm \sigma_{new}^2$$



# Regressione lineare

## //predizione e incertezza

- Possiamo effettuare la predizione sui dati nuovi (test set)

$$t_{new} = \hat{\mathbf{W}}^T \mathbf{x}_{new}$$

- Qual è l'incertezza sulle stime del nostro modello?
- Usiamo la covarianza...

$$\hat{t}_{new} \pm \sigma_{new}^2$$

$$\begin{aligned} \sigma_{new}^2 &= \text{var}\{t_{new}\} = E_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma)} [t_{new}^2] - (E_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma)} [t_{new}])^2 \\ &= \sigma^2 \mathbf{x}_{new}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{new} \\ &= \mathbf{x}_{new}^T \text{cov}\{\hat{\mathbf{W}}\} \mathbf{x}_{new} \end{aligned}$$

l'incertezza è legata alla covarianza dei parametri stimati !!

# Regressione lineare

## //modelli probabilistici: stima di ML

---

- Qual è l'incertezza sulle stime del nostro modello?
- E' legata alla covarianza dei parametri stimati...

$$\text{cov}\{\hat{\mathbf{w}}\} = E\{(\hat{\mathbf{w}} - E\{\hat{\mathbf{w}}\})(\hat{\mathbf{w}} - E\{\hat{\mathbf{w}}\})^T\} = E\{\hat{\mathbf{w}}\hat{\mathbf{w}}^T\} - E\{\hat{\mathbf{w}}\}E\{\hat{\mathbf{w}}^T\}$$

- Poichè ML è stimatore unbiased:  $E\{\hat{\mathbf{w}}\} = \mathbf{w}$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad \hat{\mathbf{w}}\hat{\mathbf{w}}^T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \mathbf{t}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$
$$E\{\hat{\mathbf{w}}\hat{\mathbf{w}}^T\} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E\{\mathbf{t} \mathbf{t}^T\} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

$$\mathbf{t} = \mathbf{X} \mathbf{w} + \boldsymbol{\epsilon} \quad E\{\mathbf{t} \mathbf{t}^T\} = E\{(\mathbf{X} \mathbf{w} + \boldsymbol{\epsilon})(\mathbf{X} \mathbf{w} + \boldsymbol{\epsilon})^T\}$$
$$= E\{\mathbf{X} \mathbf{w} \mathbf{w}^T \mathbf{X}^T + 2\boldsymbol{\epsilon} \mathbf{w}^T \mathbf{X} + \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T\}$$
$$= \mathbf{X} \mathbf{w} \mathbf{w}^T \mathbf{X}^T + 2E\{\boldsymbol{\epsilon}\} \mathbf{w}^T \mathbf{X} + E\{\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T\}$$
$$= \mathbf{X} \mathbf{w} \mathbf{w}^T \mathbf{X}^T + \sigma^2 \mathbf{I} \quad \swarrow \text{rumore a media 0}$$

# Regressione lineare

## //modelli probabilistici: stima di ML

---

- Qual è l'incertezza sulle stime del nostro modello?
- E' legata alla covarianza dei parametri stimati...

$$\text{cov}\{\hat{\mathbf{w}}\} = E\{(\hat{\mathbf{w}} - E\{\hat{\mathbf{w}}\})(\hat{\mathbf{w}} - E\{\hat{\mathbf{w}}\})^T\} = E\{\hat{\mathbf{w}}\hat{\mathbf{w}}^T\} - E\{\hat{\mathbf{w}}\}E\{\hat{\mathbf{w}}^T\}$$

- Si ottiene  $E\{\hat{\mathbf{w}}\hat{\mathbf{w}}^T\} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E\{\mathbf{t} \mathbf{t}^T\} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$ 
$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{w} \mathbf{w}^T \mathbf{X}^T + \sigma^2 \mathbf{I}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$
$$= \mathbf{w} \mathbf{w}^T + \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

$$\text{cov}\{\hat{\mathbf{w}}\} = E\{\hat{\mathbf{w}}\hat{\mathbf{w}}^T\} - E\{\hat{\mathbf{w}}\}E\{\hat{\mathbf{w}}^T\} = \mathbf{w} \mathbf{w}^T + \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} - \mathbf{w} \mathbf{w}^T$$
$$= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

$$\text{cov}\{\hat{\mathbf{w}}\} = E\{\hat{\mathbf{w}}\hat{\mathbf{w}}^T\} - E\{\hat{\mathbf{w}}\}E\{\hat{\mathbf{w}}^T\} = - \left( \frac{\partial^2 \mathcal{L}}{\partial \mathbf{w} \partial \mathbf{w}^T} \right)^{-1}$$

piccola curvatura  
=  
elevata varianza  
=  
parametri poco significativi

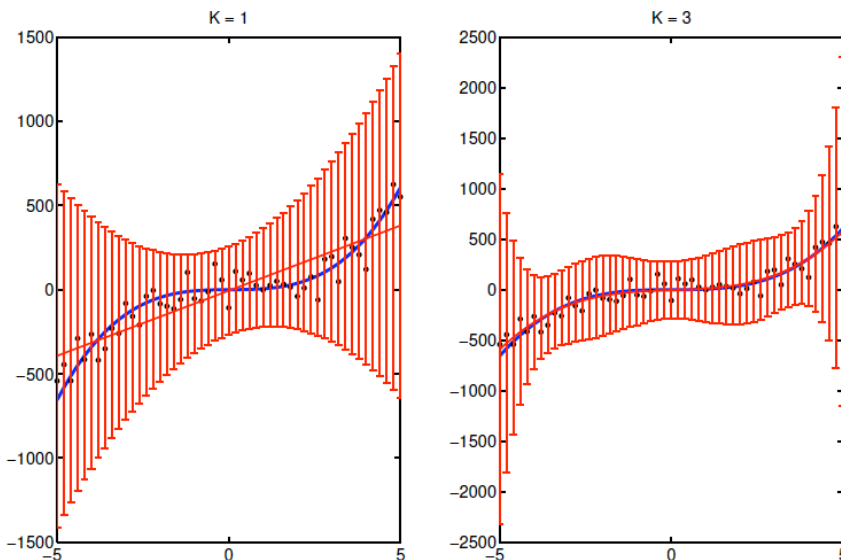
# Regressione lineare

## //modelli probabilistici: stima di ML

- Qual è l'incertezza sulle stime del nostro modello?

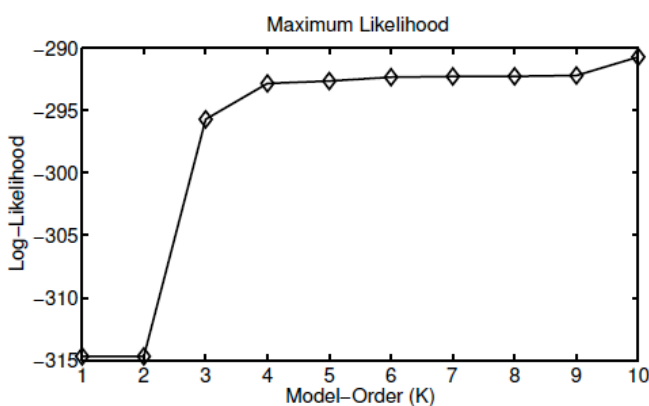
$$\hat{t}_{new} \pm \sigma_{new}^2$$

$$\begin{aligned} \hat{t}_{new} &= \mathbf{x}_{new}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \\ \sigma_{new}^2 &= \hat{\sigma}^2 \mathbf{x}_{new}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{new} \\ \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t}^T \mathbf{t} - \mathbf{t}^T \hat{\mathbf{t}}) \end{aligned}$$



# Regressione lineare

## //andamento della funzione di log-likelihood



$$\begin{aligned} \log \mathcal{L} &= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 + 2 \log 2\pi) - \frac{N}{2} \log \hat{\sigma}^2 \end{aligned}$$

```

Range = 10;
Max_Model_Order = 10;
noise_var = 100;

L=[];
x = [-Range/2:0.2:Range/2]';
N=size(x,1);

f = 5*x.^3 - x.^2 + x;
f_n = f + noise_var*randn(size(x));

[i,j]=sort(x); X=x.^0;

for k=1:Max_Model_Order
    X=[X x.^k];
    w_hat = inv(X'*X)*X'*f_n;
    f_hat = X*w_hat;
    sigma_hat = mean((f_n - f_hat).^2);
    sigma = sigma_hat*diag(X*inv(X'*X)*X');

    L = [L; -N*log(sqrt(sigma_hat)) - 0.5*N*(1 + log(2*pi))];

    plot(i,f(j),'b');
    hold on
    plot(i,f_n(j),'k','MarkerSize',15)
    errorbar(i,f_hat(j),sigma(j),'-r.')
    hold off
    pause(1)
end

figure
plot(1:Max_Model_Order,L,'dr--');
    
```