

- Supponendo inoltre che i numeri $v(i)$ siano indipendenti $E[v(i) \cdot v(j)] = 0$ $\forall i \neq j$
- $$\text{Var}[\hat{\theta}] = \lambda^2 \cdot \left[\sum_{i=1}^N \varphi(i) \cdot \varphi(i)^T \right]^{-1}$$
- $$\text{Var}[v(i)] = \lambda^2 \cdot v_i \Rightarrow \underline{\text{CONSISTENTE}}$$

STIMA A MASSIMA VEROSIMIGLIANZA 4

STIMA PARAMETRICA

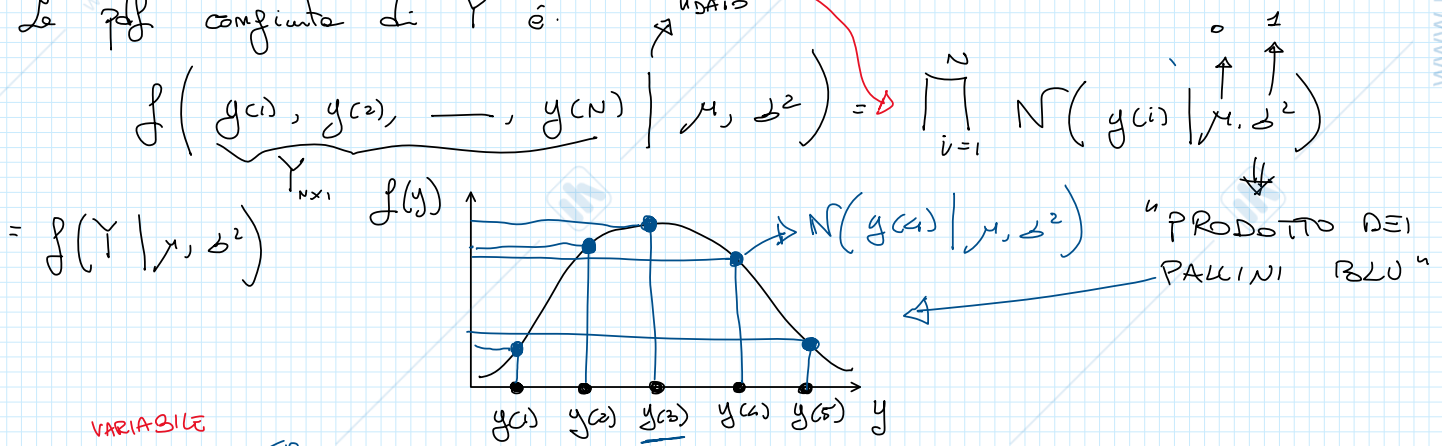
- $L \hat{\mu} = \frac{1}{N} \sum_{i=1}^N y(i)$ MEDIA CAMPIONARIA $\theta = \mu \in \mathbb{R}$
- $L S_{N-1}^2 = \frac{1}{N-1} \sum_{i=1}^N (y(i) - \hat{\mu})^2$ VAR. CAMPIONARIA $\theta = \sigma^2 \in \mathbb{R}$
- L STIMA A MINIMI QUADRATI $y(i) = \theta_0 + \theta_1 x_1(i) + \dots + \theta_{d-1} x_{d-1}(i) + e(i)$
 $\theta = [\theta_0 \ \theta_1 \ \dots \ \theta_{d-1}]^T \in \mathbb{R}^{d \times 1}$ \downarrow i.i.d. $(0, \lambda^2)$

Il metodo della MASSIMA VEROSIMIGLIANZA è una procedura di stima che, dato un modello probabilistico, stima i suoi parametri in modo che essi siano il più consistente possibile con quanto osservato

$Y = [y(i) \ y(2) \ \dots \ y(N)]^T$: N osservazioni della variabile y

$L y(i) \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d

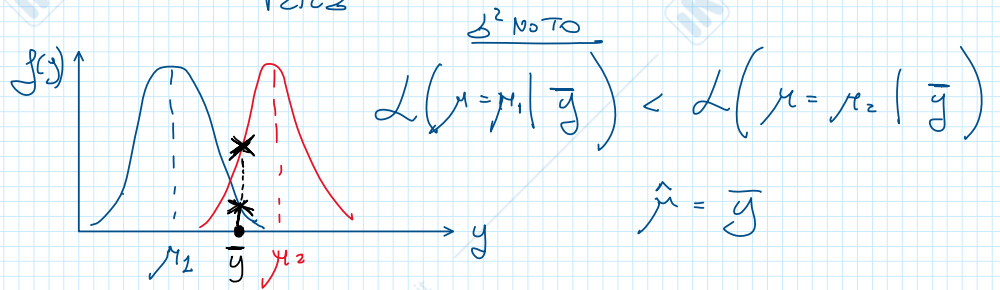
La pdf congiunta di Y è:



- Se $f(Y \mid \mu, \sigma^2)$ è funzione di $Y \Rightarrow$ è una pdf N -dimensionale
- Se $f(Y \mid \mu, \sigma^2)$ è vista in funzione di μ e $\sigma^2 \Rightarrow$ LIKELIHOOD

Se $f(Y | \mu, \sigma^2)$ è vista in funzione di μ e $\sigma^2 \Rightarrow$ LIKELIHOOD $\rightarrow \mathcal{L}(\mu, \sigma^2 | Y)$

$$\mathcal{L}(\mu, \sigma^2 | y) = \mathcal{N}(\mu, \sigma^2 | y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2}$$



$$\hat{\theta}_{ML} = \begin{bmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{bmatrix} = \underset{[\mu, \sigma^2]}{\text{arg max}}_{\theta} \mathcal{L}(\theta | Y) = \underset{\theta}{\text{arg max}} \prod_{i=1}^N \mathcal{N}(y^{(i)} | \theta)$$

In caso di gerarchia pdf $f(y | \theta)$

$$\hat{\theta}_{ML} = \underset{\theta}{\text{arg max}} \prod_{i=1}^N f(y^{(i)} | \theta)$$

$$\hat{\theta}_{ML} = \underset{\theta}{\text{arg max}} \ln \left[\mathcal{L}(\theta | Y) \right]$$

PROPRIETÀ SITA A MASSIMA VEROSIMILITUDINE

- 1) Asintoticamente corretta $\lim_{N \rightarrow +\infty} E[\hat{\theta}_{ML}] = \theta^0$
- 2) Consistente
- 3) Asintoticamente efficiente $\lim_{N \rightarrow +\infty} \text{Var}[\hat{\theta}_{ML}] = M^{-1}$
- 4) Asintoticamente normale $\hat{\theta}_{ML} \sim \mathcal{N}(\theta^0, M^{-1}) \quad N \rightarrow +\infty$

Esempio

$y^{(i)} \sim \mathcal{N}(\mu, \sigma^2=1)$, $i=1,2$, c.i.d. stimare $\hat{\theta}_{ML}$

$y^{(1)} = 4$ $y^{(2)} = 6$

1, 12

$$y^{(1)} = 4$$

$$y^{(2)} = 6$$

$$f(y^{(1)} = 4 | \mu, \sigma^2 = 1) = \frac{1}{\sqrt{2\pi \cdot 1}} e^{-\frac{1}{2} \left(\frac{4-\mu}{1}\right)^2}$$

$$f(y^{(2)} = 6 | \mu, \sigma^2 = 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (6-\mu)^2}$$

$$f(y^{(1)} = 4, y^{(2)} = 6 | \mu, \sigma^2 = 1) = \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (4-\mu)^2} \right) \cdot \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (6-\mu)^2} \right)$$

↓
È FUNZIONE SOLO DI μ

$$= \mathcal{L} \left(\underset{\mu}{\mu} | \mathcal{Y} \right)$$

$$\hat{\mu}_{ML} = \underset{\mu}{\operatorname{argmax}} \mathcal{L}(\mu | y^{(1)} = 4, y^{(2)} = 6)$$

Calcolo la log-likelihood

$$\ln[\mathcal{L}(\mu | \cdot)] = \ln \left[\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (4-\mu)^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (6-\mu)^2} \right]$$

$$= \ln \left[\frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} (4-\mu)^2} \right] + \ln \left[\frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} (6-\mu)^2} \right]$$

$$= \ln \frac{1}{\sqrt{2\pi}} + \ln \left[e^{-\frac{1}{2} (4-\mu)^2} \right] + \ln \frac{1}{\sqrt{2\pi}} + \ln \left[e^{-\frac{1}{2} (6-\mu)^2} \right]$$

$$= 2 \ln \frac{1}{\sqrt{2\pi}} - \frac{1}{2} (4-\mu)^2 \ln e - \frac{1}{2} (6-\mu)^2 \ln e$$

$$= 2 \ln \frac{1}{\sqrt{2\pi}} - \frac{1}{2} (4-\mu)^2 - \frac{1}{2} (6-\mu)^2$$

Trovare il massimo

$$\frac{d \ln[\mathcal{L}]}{d\mu} = 0 \Rightarrow 0 + \frac{2}{2} (4-\mu) + \frac{2}{2} (6-\mu) = 0$$

$$\Rightarrow 4 + 6 = 2\mu \Rightarrow \hat{\mu}_{ML} = \frac{6+4}{2} = 5$$

Osservazioni

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} \ell_n[\mathcal{L}(\theta | Y)]$$

$$= \underset{\theta}{\operatorname{argmin}} -\ell_n[\mathcal{L}(\theta | Y)]$$

MINIMI QUADRATI

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \varphi_i^T \theta)^2$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$$

* STIMA A MASSIMA VEROSIMIGLIANZA DI MODELLI LINEARI *

$$y(i) = \theta_0 + \theta_1 x_1(i) + \theta_2 x_2(i) + \dots + \theta_{d-1} x_{d-1}(i) + e(i)$$

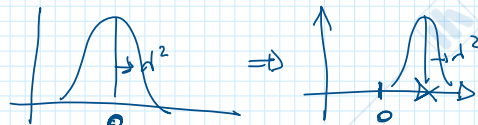
$$= \varphi(i)^T \theta + e(i)$$

$$\varphi(i) = \begin{bmatrix} 1 \\ x_1(i) \\ x_2(i) \\ \vdots \\ x_{d-1}(i) \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{d-1} \end{bmatrix}$$

$$e(i) \sim \mathcal{N}(0, \lambda^2) \text{ i.i.d.}$$

$$e(i) \perp \varphi(i)$$



$$y(i) \sim \mathcal{N}(\underbrace{\varphi(i)^T \theta}_{\mu(i)}, \lambda^2)$$

La probabilità di osservare i dati misurati è data da:

$$f(y(i), y(2), \dots, y(N) | X, \theta, \lambda^2) \stackrel{\text{i.i.d.}}{=} \prod_{i=1}^N f(y(i) | \varphi(i), \theta, \lambda^2)$$

$$X = \begin{bmatrix} \varphi(1)^T \\ \vdots \\ \varphi(N)^T \end{bmatrix} \quad = \prod_{i=1}^N \mathcal{N}(\varphi(i)^T \theta, \lambda^2)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\lambda^2} e^{-\frac{1}{2} \left(\frac{y(i) - \varphi(i)^T \theta}{\lambda} \right)^2} = \mathcal{L}(\theta, \lambda^2 | Y, X)$$

↳ NOTO (PER IPOTESI)

1) Calcolare la log-verosimiglianza

$$\ell_n[\mathcal{L}(\theta | X, Y)] = \ell_n \left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi}\lambda^2} e^{-\frac{1}{2} \left(\frac{y(i) - \varphi(i)^T \theta}{\lambda} \right)^2} \right]$$

$$= \sum_{i=1}^N \ell_n \left[\frac{1}{\sqrt{2\pi}\lambda^2} \cdot e^{-\frac{1}{2} \left(\frac{y(i) - \varphi(i)^T \theta}{\lambda} \right)^2} \right]$$

$$\stackrel{N}{=} \sum_{i=1}^N \ell_n \left[\frac{1}{\sqrt{2\pi}\lambda^2} \cdot e^{-\frac{1}{2} \left(\frac{y(i) - \varphi(i)^T \theta}{\lambda} \right)^2} \right]$$

$$\begin{aligned}
 & i=1 \quad \left[\frac{1}{2\pi\lambda^2} \right] \\
 & = \sum_{i=1}^N \ln \frac{1}{\sqrt{2\pi}\lambda^2} + \sum_{i=1}^N \ln \left[e^{-\frac{1}{2} \left(\frac{y_{(i)} - \varphi_{(i)}^T \theta}{\lambda} \right)^2} \right] \\
 & = N \cdot \ln \left(2\pi\lambda^2 \right)^{-\frac{1}{2}} + \sum_{i=1}^N -\frac{1}{2} \left(\frac{y_{(i)} - \varphi_{(i)}^T \theta}{\lambda} \right)^2 \ln e \\
 & = \frac{-1}{2} N \ln [2\pi\lambda^2] - \frac{1}{2\lambda^2} \sum_{i=1}^N \left(y_{(i)} - \varphi_{(i)}^T \theta \right)^2
 \end{aligned}$$

2) Calcolo del massimo = minimizzazione $l = -\ln[L(\theta | X, Y)]$

$$-\ln[L(\theta | X, Y)] = + \frac{1}{2} N \ln [2\pi\lambda^2] + \frac{1}{2\lambda^2} \sum_{i=1}^N \left(y_{(i)} - \varphi_{(i)}^T \theta \right)^2$$

$$\Rightarrow \hat{\theta}_{ML} = \underset{\theta}{\text{argmin}} \left(\frac{1}{2\lambda^2} \sum_{i=1}^N \left(y_{(i)} - \varphi_{(i)}^T \theta \right)^2 \right)$$

$\rightarrow \frac{1}{N}$

Lo stime ML del modello $y_{(i)} = \varphi_{(i)}^T \theta + e_{(i)}$, $i=1-N$, dove $e_{(i)} \sim N(0, \lambda^2)$ ed, \hat{e} equivalente alla stime LS



REGRESSIONE LINEARE

REGRESSIONE LOGISTICA

5

REGRESSIONE LINEARE: modello dati metrici $\rightarrow y \in \mathbb{R}$

REGRESSIONE LOGISTICA: modello dati binari / dicotomici (categorici)

$$\downarrow \\
 y \in \{0, 1\}$$

In questo caso, NON HA SENSO utilizzare un modello lineare del tipo $y_{(i)} = \varphi_{(i)}^T \theta + e_{(i)}$

L non ha senso sommare un errore continuo ($\in \mathbb{R}$) ad una variabile y che può assumere solo 2 valori

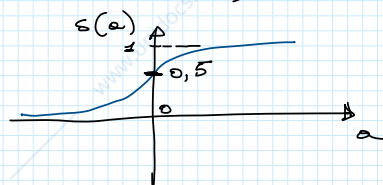
L il modello potrebbe assumere valori < 0 o > 1 . Non c'è niente che ci limiti l'uscita tra 0 e 1



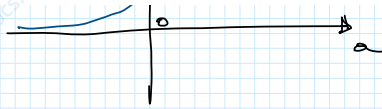
FUNZIONE LOGISTICA (SIGMOIDE)

$$s(a) = \frac{1}{1 + e^{-a}} = \frac{e^a}{1 + e^a}$$

• se $a \gg 0 \Rightarrow s(a) \rightarrow 1$



- se $a \gg 0 \Rightarrow s(a) \rightarrow 1$
- se $a \ll 0 \Rightarrow s(a) \rightarrow 0$



OBBIETTIVO. modellare la probabilità che $y=1$ tramite un modello lineare

$$P(y=1 | \varphi) = s(\varphi^T \vartheta) = \frac{1}{1 + e^{-(\varphi^T \vartheta)}}$$

- se $\varphi^T \vartheta \gg 0 \Rightarrow P(y=1 | \varphi) \rightarrow 1$
- se $\varphi^T \vartheta \ll 0 \Rightarrow P(y=1 | \varphi) \rightarrow 0$

$$\begin{cases} y=1 & P(y=1) = \pi \\ y=0 & P(y=0) = 1-\pi \end{cases} \Rightarrow \pi^y \cdot (1-\pi)^{(1-y)}$$

REGRESSIONE LINEARE

$$y \sim N(\mu, \sigma^2)$$

$$\mu = \varphi^T \vartheta = \vartheta_0 + \vartheta_1 x_1 + \dots + \vartheta_d x_d$$

REGRESSIONE LOGISTICA

$$y \sim \text{Bernoulli}(\pi)$$

$$\pi = s(\varphi^T \vartheta)$$

STIMA ML DI UN MODELLO DI REGRESSIONE LOGISTICA

Sia dato un dataset $D = \{(\varphi^{(1)}, y^{(1)}), (\varphi^{(2)}, y^{(2)}), \dots, (\varphi^{(N)}, y^{(N)})\}$
 $\varphi^{(i)} \in \mathbb{R}^{d+1}$, $y^{(i)} \in \{0, 1\}$, $i = 1 - N$, (ced)

Stimare un modello di regressione logistica

$$P(y^{(i)}=1 | \varphi^{(i)}) = \frac{1}{1 + e^{-\varphi^{(i)T} \vartheta}} \equiv \pi(i)$$

1) Calcola la verosimiglianza dai dati

$$L(\pi | Y) = \prod_{i=1}^N \pi(i)^{y^{(i)}} \cdot (1-\pi(i))^{1-y^{(i)}} \Rightarrow \text{Calcola la nuova log-likelihood}$$

$$-\ln[L(\pi | Y)] = -\ln \left[\prod_{i=1}^N \pi(i)^{y^{(i)}} \cdot (1-\pi(i))^{1-y^{(i)}} \right]$$

$$= -\sum_{i=1}^N \ln \left[\pi(i)^{y^{(i)}} \cdot (1-\pi(i))^{1-y^{(i)}} \right]$$

$$= -\sum_{i=1}^N \left(\ln \left[\pi(i)^{y^{(i)}} \right] + \ln \left[(1-\pi(i))^{1-y^{(i)}} \right] \right)$$

$$= -\sum_{i=1}^N \left(y^{(i)} \cdot \ln \pi(i) + (1-y^{(i)}) \ln [1-\pi(i)] \right) = J(\vartheta)$$

2) Interpretazione della funzione di costo

Supponiamo di avere un sub dato $D = \{(\varphi, y)\}$

$$J(\theta) = \begin{cases} -\ln \pi & \text{se } y = 1 \leftarrow \\ -\ln[1-\pi] & \text{se } y = 0 \leftarrow \end{cases}$$

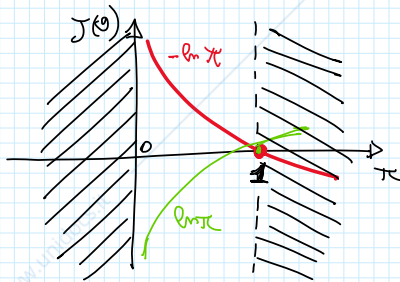
CASO $y=1$

$$J(\theta) = -\ln \pi$$

$$\Rightarrow J(\theta) \rightarrow 0 \text{ se } y=1$$

RE

$\pi \rightarrow 1$



$$J(\theta) \rightarrow +\infty \text{ se } y=1$$

RE

$\pi \rightarrow 0$

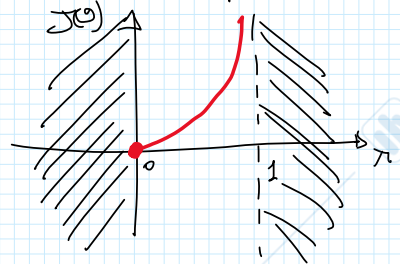
CASO $y=0$

$$J(\theta) = -\ln[1-\pi]$$

$$J(\theta) \rightarrow 0 \text{ se } y=0$$

RE

$\pi \rightarrow 0$



$$J(\theta) \rightarrow +\infty \text{ se } y=0$$

RE

$\pi \rightarrow 1$

CALCOLO DEL MINIMO

$$\nabla_{\theta} J(\theta) = 0$$

$$\nabla_{\theta} J(\theta) = \sum_{i=1}^N \varphi(i) \cdot (\pi(i) - y(i))$$

$$\frac{ds(a)}{da} = s(a) \cdot [1 - s(a)]$$

$$\frac{ds(\varphi^T \theta)}{d\theta} = \nabla_{\theta} s(\varphi^T \theta) = \varphi \cdot \pi \cdot (1 - \pi)$$

$$\frac{1}{1 + e^{-\varphi^T \theta}}$$



OTTIMIZZAZIONE TRAMITE GRADIENT DESCENT

$$\hat{\theta}(j+1) = \hat{\theta}(j) - \underbrace{\alpha}_{\substack{\downarrow \\ > 0}} \nabla_{\theta} J(\theta) \Big|_{\theta = \hat{\theta}(j)}$$

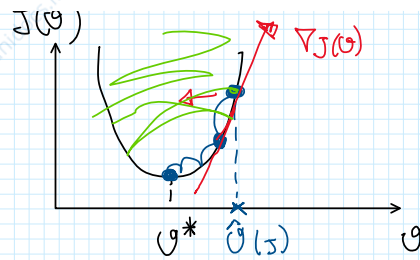
- α : LEARNING RATE
- $\hat{\theta}(0)$ è inizializzato RANDOM

se $\nabla J(\theta) = 0$

... $\hat{\theta}(j+1), \hat{\theta}(j)$



• se $\nabla J(\theta) \Big|_{\theta = \hat{\theta}^{(j)}} > 0 \Rightarrow \hat{\theta}^{(j+1)} < \hat{\theta}^{(j)}$



• se $\nabla J(\theta) \Big|_{\theta = \hat{\theta}^{(j)}} < 0 \Rightarrow \hat{\theta}^{(j+1)} > \hat{\theta}^{(j)}$

