

# TEORIA DELL'INFORMAZIONE E INFERENZA STATISTICA

Sia  $\{p_{\vartheta}(x)\}_{\vartheta}$  una famiglia di d.p. parametrizzate dal parametro  $\vartheta \in \mathcal{R}$

Sia  $X$  una v.a. campionata da  $p_{\vartheta}$  ( $X \sim p_{\vartheta}$ )

Una statistica è una funzione  $T(X)$  del campione  $X$

Esempio: Media campione:  $T(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i$

Poiché  $T(X)$  è una funzione del campione, vale la catena di Markov:

$$\vartheta \rightarrow X \rightarrow T(X) \quad (\text{poiché } \Pr[T(X)|\vartheta, X] = \Pr[T(X)|X])$$

Per il 2° teorema di elaborazione dati,  $I(\vartheta; T(X)) \leq I(\vartheta; X)$ .

Quando  $I(\vartheta; T(X)) = I(\vartheta; X)$ , la statistica  $T(X)$  si dice statistica sufficiente.

(per ricostruire  $\vartheta$ , posso ignorare  $X$  e tenere solo  $T(X)$ )

→ è vero quando  $I(\vartheta; T(X)) \geq I(\vartheta; X)$ , cioè quando  $\Pr[X|\vartheta, T(X)] = \Pr[X|T(X)]$

ovvero quando  $\vartheta \rightarrow T(X) \rightarrow X$

Esempio.  $X_1, \dots, X_n$  v.a.  $\in \{0, 1\}$  indipendenti con  $p(1) = \vartheta$   
 $p(0) = 1 - \vartheta$   
 $(\vartheta \in [0, 1])$

$$\Pr[(X_1, \dots, X_n) = (x_1, \dots, x_n)] = \vartheta^{\sum_i x_i} \cdot (1 - \vartheta)^{n - \sum_i x_i}$$

$\underbrace{\quad}_{\text{quanti 1 e quanti 0 ci sono?}}$ 
 $\downarrow$ 
 Numero di "1"
 
 $\downarrow$ 
 Numero di "0"

$$\Pr[10001] = \vartheta^2 (1 - \vartheta)^3$$

$$\Pr[01001] = \vartheta^2 (1 - \vartheta)^3$$

$$T(X_1, \dots, X_n) = \sum_{i=1}^n X_i$$

$$\Pr[(X_1, \dots, X_n) = (x_1, \dots, x_n) \mid \sum_{i=1}^n X_i = k] = \begin{cases} \frac{1}{\binom{n}{k}} & \text{se } \sum_i x_i = k \\ 0 & \text{se } \sum_i x_i \neq k \end{cases}$$

$\swarrow$  non dipende da  $\vartheta$  !

$\Rightarrow T(X_1, \dots, X_n) = \sum_i X_i$  è una statistica sufficiente

# STIMA A MASSIMA VEROSIMIGLIANZA

Ho una famiglia  $\{p_{\theta}(x)\}_{\theta}$  di d.p.;  $p_{\theta^*}(x)$  è il modello corretto ma non conosco  $\theta^*$

Quale scelta di  $\theta$  è maggiormente consistente con le osservazioni?

$$\rightarrow D(p_{\theta^*} \parallel p_{\theta}) = E_{X \sim p_{\theta^*}} \left[ \log \frac{p_{\theta^*}(X)}{p_{\theta}(X)} \right] = \overbrace{E_{X \sim p_{\theta^*}} [\log p_{\theta^*}(X)]}^{\text{dipende solo da } \theta^*} - E_{X \sim p_{\theta^*}} [\log p_{\theta}(X)]$$

$$= \text{costante (rispetto a } \theta) - E_{X \sim p_{\theta^*}} [\log p_{\theta}(X)]$$

Non lo conosco ma posso stimarlo:

per la legge dei grandi numeri,  $\frac{1}{n} \sum_{i=1}^n \log p_{\theta}(X_i)$

$$\rightsquigarrow E_{X \sim p_{\theta^*}} [\log p_{\theta}(X)]$$

$$\hat{D}(p_{\theta^*} \parallel p_{\theta}) = \text{costante} - \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(X_i)$$

Scelgo  $\theta$  per minimizzare la divergenza  $\hat{D}$

$$\min_{\vartheta} \hat{D}(p_{\vartheta^*} \| p_{\vartheta}) \Leftrightarrow \min_{\vartheta} \left[ -\frac{1}{n} \sum_{i=1}^n \log p_{\vartheta}(X_i) \right]$$

$$\Leftrightarrow \max_{\vartheta} \left[ \sum_{i=1}^n \log p_{\vartheta}(X_i) \right]$$

$$\Leftrightarrow \max_{\vartheta} \underbrace{\prod_{i=1}^n p_{\vartheta}(X_i)}_{\substack{\text{funzione di verosimiglianza (likelihood)} \\ \mathcal{L}(\vartheta | X)}}$$

log è una funzione monotona

→ Principio di massima verosimiglianza

Esempio:  $\mathcal{L}(\vartheta | \vec{X}) = \prod_{i=1}^n \underbrace{p_{\vartheta}(X_i)}_{\substack{\vartheta \\ (\text{se } X_i=1)}} \underbrace{(1-\vartheta)}_{\substack{(1-\vartheta) \\ (\text{se } X_i=0)}} = \vartheta^k (1-\vartheta)^{n-k}$  dove  $k = \sum_{i=1}^n X_i$

$$k - k\vartheta = n\vartheta - k\vartheta \Leftrightarrow \boxed{\vartheta = k/n}$$

$$= 0 \Leftrightarrow k(1-\vartheta) = (n-k)\vartheta$$

$$\mathcal{L}' = k \vartheta^{k-1} (1-\vartheta)^{n-k} + (n-k) \vartheta^k (1-\vartheta)^{n-k-1} (-1) = \vartheta^{k-1} (1-\vartheta)^{n-k-1} \underbrace{[k(1-\vartheta) - (n-k)\vartheta]}$$

$$\mathcal{L}' = 0$$

## Esercizio. Codifica run-length

Siano  $X_1, \dots, X_n$  v.a.  $\in \{0, 1\}$  ma non necessariamente dipendenti.

Si considerino le lunghezze delle run  $\vec{R} = (R_1, R_2, \dots)$

Esempio:  $\vec{X} = \underbrace{000}_3 \underbrace{11}_2 \underbrace{001}_2 \underbrace{00}_2$ ,  $\vec{R} = (3, 2, 2, 1, 2)$   $\begin{matrix} \rightarrow 0001100100 \\ \rightarrow 1110011011 \end{matrix}$

Si confrontino  $H(X_1, \dots, X_n)$  (cioè  $H(\vec{X})$ ),  $H(\vec{R})$ , e  $H(X_n, \vec{R})$ ,  
e limitare le differenze tra queste 3 quantità.

Soluz.  $\vec{R}$  è funzione di  $\vec{X}$ , quindi  $H(\vec{R}) \leq H(\vec{X}) = H(X_1, \dots, X_n)$ .

Inoltre:  $(X_n, \vec{R})$  sono funz. di  $\vec{X} \Rightarrow H(X_n, \vec{R}) \leq H(\vec{X})$   
 $\vec{X}$  è funz. di  $(X_n, \vec{R}) \Rightarrow H(\vec{X}) \leq H(X_n, \vec{R})$  }  $H(\vec{X}) = H(X_n, \vec{R})$

$$H(\vec{R}) \leq \underline{H(\vec{X})} = H(X_n, \vec{R}) = H(\vec{R}) + H(X_n | \vec{R}) \leq H(\vec{R}) + \overbrace{H(X_n)}^{\leq 1} \leq H(\vec{R}) + 1$$

Esercizio. Siano  $X_1, X_2, \dots, X_n$  di v.a. indipendenti

con la stessa distribuzione, con entropia  $H$ .

$$\text{Sia } C_n(t) = \{ \vec{x} \in \mathcal{X}^n : p(\vec{x}) \geq 2^{-nt} \}$$

(a) Mostrare che  $|C_n(t)| \leq 2^{nt}$

(b) Per quali valori di  $t$  si ha  $\Pr[(X_1, \dots, X_n) \in C_n(t)] \rightarrow 1$  ?

(a)  $\Pr[(X_1, \dots, X_n) \in C_n(t)] \leq 1$  quindi:

$$\boxed{1 \geq \Pr[(X_1, \dots, X_n) \in C_n(t)] \geq |C_n(t)| \min_{\vec{x} \in C_n(t)} p(\vec{x}) \geq |C_n(t)| \cdot 2^{-nt}}$$

$$\Rightarrow |C_n(t)| \leq 2^{nt}$$

autoinformaz. normalizzate

$$\begin{aligned} \log p &\rightarrow -nH \\ p &\rightarrow 2^{-nH} \end{aligned}$$

(b) Per il principio di equipartizione asintotica,  $-\frac{1}{n} \log p(\vec{X}) \rightsquigarrow H$

ovvero  $p(\vec{X})$  è vicina a  $2^{-nH}$  con prob. che tende a 1.

Quindi  $\Pr[p(\vec{X}) \geq 2^{-nt}] \rightarrow 1$  se  $H < t$ ,  $\rightarrow 0$  se  $H > t$ .

## Esercizio. Codifica di sorgente

Una sorgente senza memoria e stazionaria emette una sequenza di cifre binarie

con probabilità:  $p(1) = 0.005$   
 $p(0) = 0.995$

01001100000010...

100

↳

Prendiamo le cifre 100 alla volta (blocchi di dimensione 100)

e se il blocco ha al più 3 cifre 1, lo codifichiamo con una parola di codice binaria.

(altrimenti lo ignoriamo).

(a) Se tutte le parole di codice hanno uguale lunghezza, qual è la loro minima lunghezza possibile?

(b) Calcolare la prob. di osservare un blocco a cui non è stata associata una parola di codice.

(a) Le seq. di 100 bit con al più 3 cifre 1?

$$\binom{100}{0} + \binom{100}{1} + \binom{100}{2} + \binom{100}{3} = 1 + 100 + \frac{100 \cdot 99}{2} + \frac{100 \cdot 99 \cdot 98}{3 \cdot 2} = 166751$$

$$\text{Lunghezza minima: } \lceil \log_2 166751 \rceil = 18. \quad (H(X) \approx 4.84)$$

(b) La prob. di avere al più 3 cifre 1 in un blocco è:

$$\sum_{i=0}^3 \binom{100}{i} (0.005)^i (0.995)^{100-i} \approx 0.99833..$$

La prob. di non poter codificare è  $1 - 0.99833.. < 1\%$ .