

# Machine Learning per la fisica applicata e la fisica delle alte energie<sup>1</sup>

## Lezione 2: Richiami di probabilità

"Probability theory is nothing but common sense reduced to calculation"  
[Pierre Laplace, 1812]

### DEFINIZIONE FREQUENTISTA

La probabilità di un evento è il rapporto tra il numero di casi favorevoli ed il numero di casi possibili.

### DEFINIZIONE SOGGETTIVA

La probabilità di un evento è il prezzo che un individuo ritiene equo pagare per ricevere 1 se l'evento si verifica o 0 altrimenti.

### DEFINIZIONE BAYESIANA

La probabilità di un evento è l'incertezza con cui l'evento si verifica.

### DEFINIZIONE ASSIOMATICA (Kolmogorov, 1933)

Viene dato un fondamento logico al concetto di probabilità mediante assiomi.

### TIPI d'INCERTEZZA:

- aleatoria (DATA UNCERTAINTY)
- epistemica (MODEL UNCERTAINTY)

### PROBABILITÀ di UN EVENTO

$$\boxed{\Pr(A)}$$

### 1.1 PROPRIETÀ

$$\Pr(A \cap B) = \Pr(A) \cdot \Pr(B) = \Pr(A, B) \quad \text{eventi indipendenti} \quad \text{JOINT PROBABILITY}$$

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B) \quad \text{UNION PROBABILITY}$$

$$\Pr(B/A) = \frac{\Pr(A, B)}{\Pr(A)} \quad \text{CONDITIONAL PROBABILITY}$$

$$\Pr(A, B/C) = \Pr(A/C) \cdot \Pr(B/C) \quad \text{CONDITIONAL INDEPENDENCE}$$

## 2.2 RANDOM VARIABLES

$X$  rappresenta una quantità di cui non si conosce il valore, è detta variabile casuale (RANDOM VARIABLE). Il set dei possibili valori di  $X$  è detto spazio di sampling (SAMPLING SPACE). Un evento è un set di risultati dato uno spazio di sampling definito.

- sampling space finito  $\Leftrightarrow$  variabile casuale DISCRETA  $\rightarrow$  PMF
- sampling space  $\mathbb{R}$   $\Leftrightarrow$  variabile casuale CONTINUA  $\rightarrow$  CDF

PMF (PROBABILITY MASS FUNCTION)

$$p(x) \stackrel{\text{def}}{=} \Pr(X=x)$$

CDF (CUMULATIVE DISTRIBUTION FUNCTION)

$$P(x) \stackrel{\text{def}}{=} \Pr(X \leq x)$$

PDF (PROBABILITY DENSITY FUNCTION)

$$p(x) \stackrel{\text{def}}{=} \frac{d}{dx} P(x)$$

segue che

$$\Pr(a < X < b) = \int_a^b dx p(x) = P(b) - P(a) \quad \Pr(x < X < x+dx) \approx p(x) dx$$

Se la CDF è monotonicamente crescente, allora ha un'inversa che si chiama quantile. Se  $P$  è la CDF di  $X$ , allora  $P^{-1}(q)$  è il valore  $x_q$  tale che  $\Pr(X \leq x_q)$  (il quantile  $q$  di  $P$ ).

Si considerino due variabili casuali  $X$  e  $Y$ . La JOINT

DISTRIBUTION è  $p(x, y) = p(X=x, Y=y) \quad \forall X, Y$ . Esempio

$p(X, Y)$	$Y=0$	$Y=1$
$X=0$	0,2	0,3
$X=1$	0,3	0,2

Se le variabili sono indipendenti e hanno cardinalità finite, allora

$$p(X=x) = \sum_y p(X=x, Y=y)$$

DISTRIBUZIONE MARGINALE

# Definisco invece CONDITIONAL DISTRIBUTION

$$p(Y=y/X=x) = \frac{p(X=x, Y=y)}{p(X=x)} \quad \text{cioè} \quad p(x, y) = p(x)p(y/x)$$

segue la chain rule:

$$p(\vec{x}, D) = p(x_1) p(x_2/x_1) p(x_3/x_1, x_2) \dots p(x_D/x_1, \dots, x_{D-1})$$

Due variabili casuali sono **MARGINALMENTE INDIPENDENTI** se

$$X \perp Y \Leftrightarrow p(X, Y) = p(X)p(Y)$$

Due variabili casuali sono **CONDIZIONALMENTE INDIPENDENTI** se

$$X \perp Y \Leftrightarrow p(X, Y/Z) = p(X/Z)p(Y/Z)$$

## 2.3 MOMENTI DI UNA DISTRIBUZIONE

• **MEDIA**  $E[X] \triangleq \int_{\mathcal{X}} x p(x) dx \quad \left( = \sum_{x \in \mathcal{X}} x p(x) \right)$   
se la variabile è discreta

La media è lineare:  $E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$

• **VARIANZA**  $V[X] \triangleq E[(X-\mu)^2] = \int_{\mathcal{X}} (x-\mu)^2 p(x) dx$   
 $= \int x^2 p(x) dx + \underbrace{\mu^2 \int p(x) dx}_{=1} - 2\mu \underbrace{\int x p(x) dx}_{=\mu} = E[X^2] - \mu^2$

$V[aX+b] = a^2 V[X]$   $V\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n V[X_i]$  (per variabili indipendenti)

## • MODA

$x^* = \arg\max_x p(x)$  (valore più probabile; può non essere UNICO)

Tutti questi stimatori non danno conto di tutte l'informazione contenuta nella pdf. Esempio: Ascombe's quartet e Belzarius set.



## 2.3 TEOREMA DI BAYES

Dato una quantità non nota  $H$  e dei dati noti  $\bar{Y} = y$ :

$$p(H=h / \bar{Y}=y) = \frac{p(H=h) p(\bar{Y}=y / H=h)}{p(\bar{Y}=y)} \quad \left[ \text{posterior} \propto \text{prior} \times \text{likelihood} \right]$$

Il teorema segue dall'identità  $p(h/y) p(y) = p(h) p(y/h) = p(h,y)$

$p(h)$  si chiama PRIOR (cioè che conosciamo/assumiamo per  $H$  prima di fare qualunque misura);  $p(y/h)$  è la DISTRIBUZIONE OSSERVATA;  $p(y/h)$  è la likelihood (che non è una distribuzione di probabilità);  $p(y)$  è la LIKELIHOOD MARGINALE.

$$p(\bar{Y}=y) = \sum_{h' \in \mathcal{H}} p(H=h') p(\bar{Y}=y / H=h') = \sum_{h' \in \mathcal{H}} p(H=h', \bar{Y}=y)$$

$p(h/y)$  è detta POSTERIOR DISTRIBUTION.

Esempio: The Monty Hall problem

tre porte. Un premio dietro una porta. Il concorrente sceglie una porta. Il presentatore ne apre un'altra. Che cosa conviene fare al concorrente?

Supponiamo che il concorrente scelga inizialmente la porta 1

$$p(H_1) = p(H=1) = \frac{1}{3}$$

probabilità di scegliere una porta

$$p(\bar{Y}=2 / H_1) = \frac{1}{2} \quad \left( \text{se il premio è dietro la porta 1} \right) \quad p(\bar{Y}=2 / H_2) = 0 \quad p(\bar{Y}=2 / H_3) = 1$$

$$p(\bar{Y}=3 / H_1) = \frac{1}{2} \quad \left( \text{se il premio è dietro la porta 2} \right) \quad p(\bar{Y}=3 / H_2) = 1 \quad p(\bar{Y}=3 / H_3) = 0$$

SI APRE LA PORTA 3

$$p(\bar{Y}=3) = \frac{1}{6} + \frac{1}{3} = \frac{1}{2} \quad (\text{ottenuto marginalizzando})$$

$$p(H_1 / \bar{Y}=3) = \frac{\frac{1}{3} \cdot \frac{1}{2}}{\frac{1}{2}} = \frac{1}{3} \quad p(H_2 / \bar{Y}=3) = \frac{\frac{1}{3} \cdot 1}{\frac{1}{2}} = \frac{2}{3} \quad p(H_3 / \bar{Y}=3) = 0$$

## 2.4. LA DISTRIBUZIONE DI GAUSS

5

tra le distribuzioni più utilizzate nel Machine Learning (e non solo) vi è la distribuzione gaussiana così definita

$$\text{CDF} \quad P(y) \triangleq \Pr(Y \leq y) \quad \Pr(a < Y \leq b) = P(b) - P(a)$$

$$\Phi(y; \mu, \sigma^2) \triangleq \int_{-\infty}^y \mathcal{N}(z/\mu, \sigma^2) dz = \frac{1}{2} \left[ 1 + \text{erf}\left(\frac{z}{\sqrt{2}}\right) \right], \quad z = \frac{y - \mu}{\sigma}$$

$$\text{erf}(u) \triangleq \frac{2}{\sqrt{\pi}} \int_0^u e^{-t^2} dt \quad \mathcal{N}(y/\mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-\mu)^2} \quad (\text{PDF})$$

$$\mathbb{E}[Y] \triangleq \int y p(y) dy = \mu \quad \mathbb{V}[Y] \triangleq \int (y - \mu)^2 p(y) dy = \sigma^2$$

$$\hookrightarrow Y \rightarrow (Y - \mu)^2$$

Ragioni della popolarità della distribuzione

- dipende da due soli parametri di facile interpretazione
- il teorema del limite centrale ci dice che, nel limite  $N \rightarrow \infty$ , ciascuna distribuzione può essere approssimata dalla distribuzione gaussiana, inoltre essa rappresenta la somma di variabili casuali indipendenti
- il numero di assunzioni è minimo.

## 2.5 ALTRE DISTRIBUZIONI NOTEVOLI

$$t\text{-Student} \quad T(y/\mu, \sigma^2, \nu) \propto \left[ 1 + \frac{1}{\nu} \left( \frac{y - \mu}{\sigma} \right)^2 \right]^{-\frac{\nu+1}{2}} \quad \begin{array}{l} \mu \text{ mean} \\ \sigma \text{ scale param} \\ \nu \text{ degree of} \\ \text{freedom} \end{array}$$

$$\text{Lorentz (Cauchy)} \quad C(x/\mu, \gamma) = \frac{1}{\gamma\pi} \left[ 1 + \left( \frac{x - \mu}{\gamma} \right)^2 \right]^{-1}$$

$$\text{Laplace} \quad L(y/\mu, b) = \frac{1}{2b} \exp\left(-\frac{|y - \mu|}{b}\right) \quad \begin{array}{l} \mu \text{ mean} \\ \mu \text{ mode} \\ 2b^2 \text{ variance} \end{array}$$

$$\text{Beta} \quad \text{Beta}(x/a, b) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1} \quad \begin{array}{l} a/(a+b) \text{ mean} \\ (a-1)(a+b-2) \text{ mode} \\ ab/(a+b)^2(a+b+1) \text{ variance} \end{array}$$

Gamma	$Ga(x/a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-x/b}$	$a$ shape $b$ rate
Exponential	$Exp(x/\lambda) = Ga(x/a=1, b=1)$	
$\chi^2$	$\chi^2(x) = Ga(x/a=\frac{\nu}{2}, b=\frac{1}{2})$	

## 2.6 ALCUNE OSSERVAZIONI

ti suggerisco di calcolare la PDF della somma di due variabili casuali continue. Se si tratta di distribuzioni gaussiane, avrò:

$$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2) \quad x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

Se  $y = x_1 + x_2$  allora

$$p(y) = \mathcal{N}(x_1/\mu_1, \sigma_1^2) \otimes \mathcal{N}(x_2/\mu_2, \sigma_2^2) = \mathcal{N}(y/\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

cioè la convoluzione di due gaussiane è una gaussiana

ti suggerisco che  $x$  sia una variabile casuale e  $y = f(x)$  una funzione di essa. Talora è difficile calcolare analiticamente  $p(y)$ . ti suggerisco per esempio che  $x \sim \text{Unif}(-1, 1)$  e  $y = f(x) = x^2$ . Possiamo approssimare  $p(y)$  campionando  $p(x)$  usando un generatore (uniforme) di numeri casuali, facendone il quadrato e prendendo la distribuzione empirica

$$p_S(y) = \frac{1}{N_S} \sum_{i=1}^{N_S} \delta(y - y_i)$$

METODO MONTE CARLO

## 2.7 MODELLI MULTIVARIATI

Consideriamo ora due variabili  $X$  e  $Y$ . La covarianza è:

$$\text{Cov}[X, Y] \triangleq E[(X - E(X))(Y - E(Y))] = E[XY] - E[X]E[Y]$$

cioè è una matrice  $D$ -dimensionale che ha le varianze sulle diagonali se  $\bar{x}$  è un vettore  $D$ -dimensionale



Si definisce correlazione (di Pearson) tra due variabili  $X$  e  $Y$

$$\rho \triangleq \text{cor}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{V[X]V[Y]}}$$

$$\text{Cov}[X] = \begin{pmatrix} V[X_1] & \text{Cov}[X_1, X_2] & \dots & \text{Cov}[X_1, X_D] \\ \text{Cov}[X_2, X_1] & V[X_2] & \dots & \text{Cov}[X_2, X_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[X_D, X_1] & \text{Cov}[X_D, X_2] & \dots & V[X_D] \end{pmatrix}$$

$$\text{cor}(x) = \begin{pmatrix} 1 & \frac{E[(X_1 - \mu_1)(X_2 - \mu_2)]}{\sigma_1 \sigma_2} & \dots & \frac{E[(X_1 - \mu_1)(X_D - \mu_D)]}{\sigma_1 \sigma_D} \\ \frac{E[(X_2 - \mu_2)(X_1 - \mu_1)]}{\sigma_2 \sigma_1} & 1 & \dots & \frac{E[(X_2 - \mu_2)(X_D - \mu_D)]}{\sigma_2 \sigma_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{E[(X_D - \mu_D)(X_1 - \mu_1)]}{\sigma_D \sigma_1} & \frac{E[(X_D - \mu_D)(X_2 - \mu_2)]}{\sigma_D \sigma_2} & \dots & 1 \end{pmatrix}$$

Osservazioni:

- 1) Il fatto che due variabili siano non correlate non significa che siano indipendenti. Esempio:  $X \sim \text{Unif}(-1, 1)$  e  $Y = X^2$ ;  $\text{cor}[X, Y] = 0$ .  
 Viceversa, due variabili indipendenti sono necessariamente non correlate.
- 2) Correlazione non implica causalità.
- 3) Una correlazione che appare nitida in diversi set di dati può sparire (o diventare opposta) se i dati sono combinati.  
 (Simpson's paradox).

## 1.7 DISTRIBUZIONE GAUSSIANA MULTIVARIATA (MVN)

La definizione di distribuzione gaussiana si estende a più variabili

$$\mathcal{N}(\bar{y} / \bar{\mu}, \bar{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\bar{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\bar{y} - \bar{\mu})^T \bar{\Sigma}^{-1} (\bar{y} - \bar{\mu}) \right\}$$

$$\bar{\Sigma} = \text{Cov}[\bar{y}]$$

Supponiamo di avere due variabili aleatorie multidimensionali  $\vec{y}_1$  e  $\vec{y}_2$ . Si può definire la distribuzione "jointly Gaussian" come

$$\vec{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \quad \Lambda = \Sigma^{-1}$$

Le distribuzioni marginali sono date da

$$p(\vec{y}_1) = \mathcal{N}(\vec{y}_1 | \vec{\mu}_1, \Sigma_{11}) \quad p(\vec{y}_2) = \mathcal{N}(\vec{y}_2 | \vec{\mu}_2, \Sigma_{22})$$

e le probabilità condizionate è:

$$p(\vec{y}_1 | \vec{y}_2) = \mathcal{N}(\vec{y}_1 | \vec{\mu}_{1/2}, \Sigma_{1/2})$$

con

$$\vec{\mu}_{1/2} = \vec{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\vec{y}_2 - \vec{\mu}_2) \quad \Sigma_{1/2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} = \Lambda_{11}^{-1}$$

sono tutte distribuzioni gaussiane.