

Note del corso

# MACHINE LEARNING PER LA FISICA APPLICATA E FISICA DELLE ALTE ENERGIE

Raviola Alessio

10 novembre 2022

## 1 Introduzione

A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

L'impostazione del corso è di tipo *probabilistico* (statistical learning). Le quantità non note sono trattate come **variabili aleatorie** (RANDOM VARIABLES) a cui viene associata una **distribuzione di probabilità** (PROBABILITY DISTRIBUTION) che descrive il set (pesato) di valori che la variabile può assumere.

Abbiamo tre tipi di machine learning:

- SUPERVISED LEARNING;
- UNSUPERVISED LEARNING;
- REINFORCEMENT LEARNING;

il corso si focalizza sui primi due tipi.

### 1.1 Supervised learning

Il **compito**  $T$  consiste nell'imparare una mappa  $f$  dagli input  $x \in X$  agli output  $y \in Y$ . Gli **input**  $x$  sono chiamati FEATURES (o COVARIATES o PREDICTORS) e sono in genere costituiti da un vettore reale con dimensione fissata, ovvero abbiamo  $X \equiv \mathbb{R}^D$ . Gli **output** sono chiamati LABEL (o TARGET o RESPONSE).

L'**esperienza**  $E$  consiste in un TRAINING SET  $\mathcal{D}$  di  $N$  coppie input-output:

$$\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N, \quad (1)$$

dove  $N$  è detta SAMPLE SIZE.

La **performance** dipende dal compito  $T$ .

#### 1.1.1 Classificazione

Problemi comuni in machine learning sono quelli di **classificazione**. In un problema di questo tipo lo spazio degli output  $C$  è un set *non ordinato* di label  $y = \{1, 2, \dots, C\}$  dette CLASSES. Quello che chiede il problema è di predire una classe dato un input, problemi di questo tipo sono detti di PATTERN RECOGNITION<sup>1</sup>.

**Esempio 1.1** (Classificazione specie di iris). In generale in IMAGE CLASSIFICATION gli input  $X$  sono immagini, quindi:

$$X = \mathbb{R}^D, \quad D = C \times D_1 \times D_2, \quad (2)$$

---

<sup>1</sup>Se abbiamo solo due classi, i.e. solo due output, allora il problema si dice di CLASSIFICAZIONE BINARIA

Index	sl [cm]	sw [cm]	pl [cm]	pw [cm]	Label
0	5.1	3.5	1.4	0.2	Setosa
1	4.9	3.0	1.4	0.2	Setosa
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
50	7.0	3.2	4.7	1.4	Versicolor
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
150	5.9	3.0	5.1	1.8	Virginica

Tabella 1: Design matrix del training set per classificazione specie di iris.

ove  $C = 3$  sono i canali RGB. E cerchiamo una mappa

$$f : X \longrightarrow Y \quad (3)$$

che ci dica a quale delle classi appartenenti a  $Y$  l'immagine appartiene. Per le specie di iris però i botanisti hanno individuato 4 caratteristiche numeriche: lunghezza e larghezza del sepal e del petalo; dunque abbiamo  $X = \mathbb{R}^4$ . Supponiamo che il training set sia una collezione di 150 esempi delle 3 specie, 50 per ognuna. I dati possono essere raccolti in una matrice detta DESIGN MATRIX come TABULAR DATA - come in ??.

Se abbiamo  $N$  elementi nel training set, ognuno con dimensione  $D = \dim X + \dim Y$ , allora abbiamo:

- BIG DATA se  $N \gg D$ , ovvero se il numero di elementi è molto superiore alla loro dimensione;
- WIDE DATA se  $D \gg N$ , ovvero se la dimensione degli elementi è molto superiore al loro numero.

Una buona idea è fare una *esplorazione dei dati* (EXPLOATORY DATA ANALYSIS) per vedere se ci sono dei pattern ovvi, ad esempio tramite grafici. Per grandi basi dati (big data) possiamo procedere mediante DIMENSIONALITY REDUCTION:

$$f(\mathbf{x}, \boldsymbol{\theta}) = \left\{ \begin{array}{ll} p_l < 2.45 \text{ cm} & \text{Setosa} \\ \text{Altrimenti} & \left\{ \begin{array}{ll} p_w < 1.75 \text{ cm} & \text{Versicolor} \\ \text{Altrimenti} & \text{Virginica} \end{array} \right\} \end{array} \right\} \quad \text{DECISION TREE}, \quad (4)$$

ove  $\boldsymbol{\theta}$  è detto THRESHOLD PARAMETER. Questo decision tree è visualizzato in ??. La performance può essere quindi misurata con il MISCLASSIFICATION RATE:

$$\mathcal{L}(\theta) \equiv \frac{1}{N} \sum_{n=1}^N \mathbb{I}(y_n \neq f(\mathbf{x}_n, \boldsymbol{\theta})), \quad (5)$$

dove  $\mathbb{I}(e)$  è l'**indicatore binario**

$$\mathbb{I}(e) = \begin{cases} 1 & \text{se } e \text{ è vero} \\ 0 & \text{se } e \text{ è falso} \end{cases}. \quad (6)$$

Nel caso in cui alcuni errori di classificazione siano più dannosi di altri posso definire una loss function  $l(y, \hat{y})$  e ridefinire il misclassification rate come l'EMPIRICAL RISK:

$$\mathcal{L}(\theta) \equiv \frac{1}{N} \sum_{n=1}^N l(y_n, f(\mathbf{x}_n, \boldsymbol{\theta})). \quad (7)$$

Un modo che abbiamo per definire il TRAINING (o MODEL FITTING) è modificare questo rischio empirico, ovvero trovare  $\hat{\theta}$  tale che

$$\mathcal{L}(\hat{\theta}) = \min[\mathcal{L}(\theta)]. \quad (8)$$

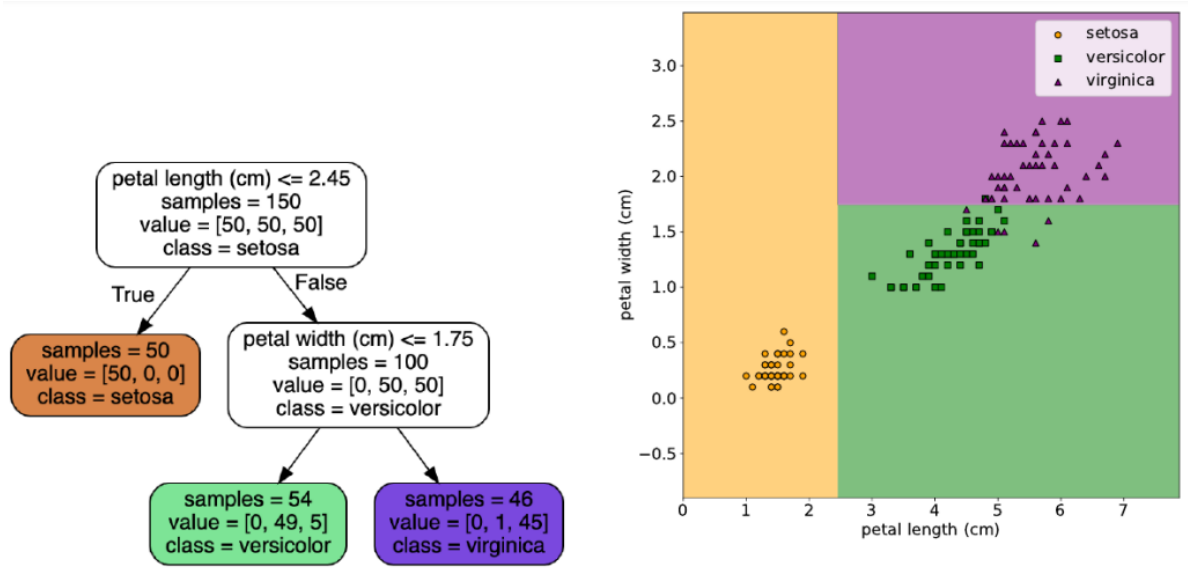


Figura 1: Decision tree per problema di classificazione specie di iris.

## 2 Richiami di probabilità

Abbiamo diverse definizioni di probabilità.

**Definizione Frequentistica** La probabilità di un evento è il rapporto tra il numero di casi favorevoli e 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 e 0 altrimenti.

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

**Definizione Assiomatica** Kolmogorov nel 1933 costruisce la teoria della probabilità a partire da degli assiomi.

L'incertezza può essere di due tipi:

**aleatoria** ovvero è una DATA UNCERTANTY;

**epistemica** ovvero è una MODEL UNCERTANTY;

### 2.1 Proprietà della probabilità

Chiamiamo  $\Pr(A)$  la probabilità dell'evento  $A$ , allora abbiamo le seguenti proprietà.

**Proprietà 2.1.1** (Joint probability). Se  $A$  e  $B$  sono due eventi indipendenti, allora:

$$\Pr(A \wedge B) \equiv \Pr(A, B) = \Pr(A) \cdot \Pr(B). \quad (9)$$

**Proprietà 2.1.2** (Union probability). Anche detta regola di unione esclusione. Se  $A$  e  $B$  sono due eventi indipendenti, allora:

$$\Pr(A \vee B) = \Pr(A) + \Pr(B) - \Pr(A \wedge B). \quad (10)$$

**Proprietà 2.1.3** (Conditional probability).

$$\Pr(A | B) = \frac{\Pr(A, B)}{\Pr(B)}. \quad (11)$$

Se i due eventi sono indipendenti questa si riduce a  $\Pr(A | B) = \Pr(A)$ .

**Proprietà 2.1.4** (Conditional independence).

$$\Pr(A, B \mid C) = \Pr(A \mid C) \cdot \Pr(B \mid C) \quad (12)$$

## 2.2 Random variables

Rappresentiamo con  $X$  una variabile di cui non conosciamo il valore e la chiamiamo *variabile casuale* (RANDOM VARIABLE). Il set dei valori che  $X$  può assumere è detto *spazio di sampling* (SAMPLING SPACE). Un evento è dunque un set di risultati dato un sampling space definito.

Se la variabile è **discreta** abbiamo un sampling space numerabile e la PMF (PROBABILITY MASS FUNCTION):

$$p(x) \equiv \Pr(X = x). \quad (13)$$

Se invece la variabile è **continua** abbiamo un sampling space non numerabile e la CDF (CUMULATIVE DISTRIBUTION FUNCTION):

$$P(x) \equiv \Pr(X \leq x), \quad (14)$$

da cui possiamo definire la PDF (PROBABILITY DENSITY FUNCTION):

$$p(x) \equiv \frac{d}{dx} P(x), \quad (15)$$

da cui segue:

$$\Pr(a \leq X \leq b) = \int_a^b p(x) = P(b) - P(a) \quad (16)$$

$$\implies \Pr(x \leq X \leq x + dx) \approx p(x). \quad (17)$$

Se la CDF  $P(x)$  è monotona crescente allora la sua inversa  $P^{-1}(q)$  è detta *quantile*. Il valore  $x_q = P^{-1}(q)$  è il valore per cui  $\Pr(X \leq x_q) < q$ , ovvero il quantile  $q$  della distribuzione  $P$ .

Se abbiamo due variabili casuali  $X$  e  $Y$  allora possiamo definire la JOINT DISTRIBUTION:

$$p(x, y) = p(X = x, Y = y) \equiv p(X = x \wedge Y = y). \quad (18)$$

Se le due variabili sono indipendenti e con cardinalità finita possiamo definire la *distribuzione marginale* (MARGINAL DISTRIBUTION) come:

$$p(X = x) = \sum_y p(X = x, Y = y), \quad (19)$$

altrimenti se sono dipendenti la *distribuzione condizionale* (CONDITIONAL DISTRIBUTION) come:

$$p(Y = y \mid X = x) = \frac{p(X = x, Y = y)}{p(X = x)} \quad (20)$$

$$\implies p(x, y) = p(y \mid x) \cdot p(x),$$

da cui segue la **chain rule**:

$$p(\mathbf{x}_{1:D}) = p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_1, x_2) \cdots p(x_D \mid p_{1:D-1}). \quad (21)$$

Due variabili si dicono MARGINALMENTE INDIPENDENTI se

$$X \perp Y \iff p(X, Y) = p(X) p(Y), \quad (22)$$

mentre si dicono **condizionalmente indipendenti** se

$$X \perp Y \iff p(X, Y \mid Z) = p(X \mid Z) p(Y \mid Z). \quad (23)$$

## 2.3 Momenti di una distribuzione

**Definizione 2.1** (Media). Definiamo MEDIA di una distribuzione come

$$\mathbb{E}[X] \equiv \int_X x p(x) dx \quad \left( \mathbb{E}[X_{\text{discreta}}] = \sum_{x \in X} x p(x) \right). \quad (24)$$

La media è lineare, ovvero  $\mathbb{E}[\sum_{i=1}^n X_i] = \sum_{i=1}^n \mathbb{E}(X_i)$ .

**Definizione 2.2** (Varianza). Definiamo VARIANZA di una distribuzione come

$$\mathbb{V}[X] \equiv \mathbb{E}[(X - \mu)^2], \quad (25)$$

che con la definizione di prima possiamo anche scrivere come:

$$\begin{aligned} \mathbb{V}[X] &= \int_X (x - \mu)^2 p(x) dx \\ &= \int x^2 p(x) dx + \underbrace{\mu^2 \int p(x) dx}_{=1} - 2\mu \underbrace{\int x p(x) dx}_{=\mu} \\ &= \mathbb{E}[X^2] - \mu^2. \end{aligned} \quad (26)$$

Con questa definizione è facile dimostrare l'identità

$$\mathbb{V}[aX + b] = a^2 \mathbb{V}[X], \quad (27)$$

inoltre anche la varianza è lineare, ma solo se le variabili sono indipendenti.

**Definizione 2.3** (Moda). Definiamo MEDIA di una distribuzione il valore (o i valori) più probabile, ovvero

$$x^* \text{ t.c. } p(x^*) = \max p(x), \quad (28)$$

Questi stimatori non danno tutta l'informazione contenuta nella distribuzione.

## 2.4 Teorema di Bayes

**Teorema 2.1** (Teorema di Bayes). *Data una quantità non nota  $H$  (hypotesis) e dei dati noti  $Y = y$  abbiamo che*

$$p(H = h | Y = y) = \frac{p(H = h) p(Y = y | H = h)}{p(Y = y)}. \quad (29)$$

Il teorema segue dall'identità  $p(h | y) p(y) = p(h) p(y | h)$ .

- $p(h)$  viene detta PRIOR ovvero ciò che conosciamo o assumiamo per  $H$  prima di fare qualunque misura;
- $p(y | h)$  è la *distribuzione osservata* detta LIKELIHOOD, è una funzione di  $y$  ad  $h$  fisso, ma non è una distribuzione di probabilità;
- $p(y)$  è detta MARGINAL LIKELIHOOD, indipendente da  $h$ , funge da costante di normalizzazione

$$p(Y = y) = \sum_{h' \in H} p(Y = y | H = h') = \sum_{h' \in H} p(H = h', Y = y);$$

- $p(h | y)$  è detta POSTERIOR DISTRIBUTION.

**Esempio 2.1** (Monty Hall problem). Posso scegliere una di tre porte. Dietro una delle porte c'è un premio. Inizio con il scegliere la porta 1. Ho che la probabilità degli eventi  $H_{1,2,3}$  che il premio sia dietro la porta 1, 2 o 3 è

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

Ora una delle porte non scelte da me si apre rivelando che dietro quella non c'è il premio, in base a dove si trova il premio la probabilità degli eventi  $Y_{2,3}$  che si apra una porta o l'altra è data da  $p(Y = y, H = h)$ :

$$\begin{aligned} p(Y = 2, H_1) &= \frac{1}{2}, & p(Y = 2, H_2) &= 0, & p(Y = 2, H_3) &= 1, \\ p(Y = 3, H_1) &= \frac{1}{2}, & p(Y = 3, H_2) &= 1, & p(Y = 3, H_3) &= 1. \end{aligned}$$

Si apre la porta 3. Mi viene data la scelta di cambiare porta prima che venga rivelato dove si trova il premio, per massimizzare la probabilità di vittoria posso calcolare le probabilità condizionate con Bayes. Dai calcoli di prima ho che  $p(Y = 3) = 1/6 + 1/3 = 1/3$  dunque:

$$\begin{aligned} p(H_1 | Y = 3) &= \frac{\frac{1}{3} \cdot \frac{1}{2}}{\frac{1}{2}} = \frac{1}{3}, \\ p(H_2 | Y = 3) &= \frac{\frac{1}{3} \cdot 1}{\frac{1}{2}} = \frac{2}{3}, \\ p(H_3 | Y = 3) &= 0, \end{aligned} \tag{30}$$

dunque mi conviene cambiare porta.

## 2.5 Distribuzione di Gauss

La distribuzione di Gauss (o gaussiana) è una tra le più usate in Machine Learning ed è definita come

$$\begin{aligned} P(y) &\equiv \Pr(Y \leq y), & \Pr(a \leq Y \leq b) &= P(b) - P(a) \\ \Phi(y; \mu, \sigma^2) &\equiv \int_{-\infty}^y \mathcal{N}(\xi/\mu, \sigma^2) d\xi = \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{\xi}{\sqrt{2}} \right) \right], & \xi &= \frac{y - \mu}{\sigma}, \end{aligned} \tag{31}$$

ove erf è la GAUSS ERROR FUNCTION definita come

$$\operatorname{erf}(u) \equiv \frac{2}{\sqrt{\pi}} \int_0^u e^{-t^2} dt, \tag{32}$$

mentre  $\mathcal{N}$  è la PDF (probability density function) definita come

$$\mathcal{N}(y | \mu, \sigma^2) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{1}{2\sigma^2} (y - \mu)^2 \right]. \tag{33}$$

Questa distribuzione è popolare per diversi motivi:

- dipende da soli due parametri ed entrambi sono di facile interpretazione;
- per il teorema del limite centrale, nel limite  $N \rightarrow \infty$  ogni distribuzione può essere approssimata da una gaussiana;
- rappresenta la somma di variabili causali indipendenti;
- il numero di assunzioni è minimo.

## 2.6 Altre distribuzioni notevoli

Abbiamo:

- T-STUDENT

$$\tau(y | \mu, \sigma^2, \nu) \propto \left[ 1 + \frac{1}{\nu} \left( \frac{y - \mu}{\sigma} \right)^2 \right]^{-\frac{\nu+1}{2}} \quad \begin{cases} \mu & \text{mean} \\ \sigma & \text{mode} \\ \nu & \text{degree of formality} \end{cases} ; \quad (34)$$

- LORENTZ (CAUCHY)

$$C(x | \mu, \gamma) = \frac{1}{\gamma\pi} \left[ 1 + \left( \frac{x - \mu}{\gamma} \right)^2 \right]^{-1} ; \quad (35)$$

- LAPLACE

$$L(y | \mu, b) = \frac{1}{2b} \exp\left(-\frac{|y - \mu|}{b}\right) \quad \begin{cases} \mu & \text{mean} \\ \mu & \text{mode} \\ 2b^2 & \text{variance} \end{cases} ; \quad (36)$$

- BETA

$$\text{Beta}(x | a, b) = \frac{1}{\beta(a, b)} x^{a-1} (1-x)^{b-1} \quad \begin{cases} \frac{a}{a+b} & \text{mean} \\ (a-1)(a+b-2) & \text{mode} \\ \frac{ab}{(a+b)^2(a+b+1)} & \text{variance} \end{cases} ; \quad (37)$$

- GAMMA

$$\text{Ga}(x | a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-xb}; \quad (38)$$

- EXPERIMENTAL

$$\text{Exp}(x | \lambda) = \text{Ga}(x | a = 1, b = \lambda) \quad (39)$$

- CHI QUADRO

$$\chi_\nu^2(x) = \text{Ga}\left(x | a = \frac{\nu}{2}, b = \frac{1}{2}\right) \quad (40)$$

Alcune distribuzioni sono mostrate in ??.

## 2.7 Osservazioni

Supponiamo di avere due variabili causali descritte da distribuzioni gaussiane:

$$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2), \quad x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2), \quad (41)$$

e di voler calcolare la PDF della loro somma  $y = x_1 + x_2$ , allora:

$$p(y) = \mathcal{N}(x_1 | \mu_1, \sigma_1^2) \otimes \mathcal{N}(x_2 | \mu_2, \sigma_2^2) = \mathcal{N}(y | \mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2), \quad (42)$$

ovvero la *convoluzione di due gaussiane è una gaussiana*.

Supponiamo ora che  $x$  sia una variabile casuale e  $y = f(x)$  una funzione di essa. Talora è difficile calcolare  $p(y)$  analiticamente. Supponiamo ad esempio che  $x \sim \text{Unif}(-1, +1)$  e  $y = x^2$ , possiamo approssimare  $p(y)$  con un generatore (uniforme) di numeri casuali e facendone il quadrato per poi prendere la DISTRIBUZIONE EMPIRICA

$$p_S(y) = \frac{1}{N_S} \sum_{x=1}^{N_S} \delta(y - y_s). \quad (43)$$

Questo è detto METODO MONTE CARLO.

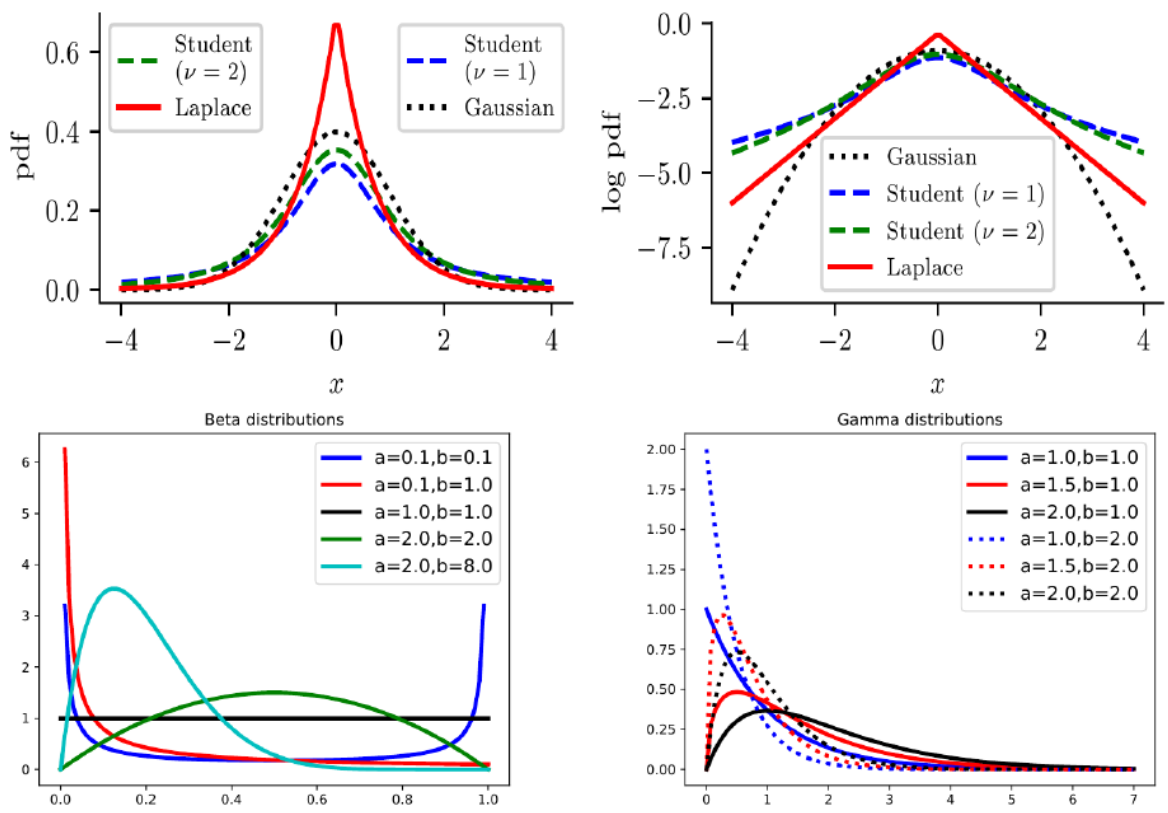


Figura 2: Alcune distribuzioni notevoli.



## 2.8 Modelli multivariati

**Definizione 2.4** (Covarianza). Siano  $X$  e  $Y$  due variabili. La COVARIANZA è

$$\text{Cov}[X, Y] \equiv \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y], \quad (44)$$

ovvero una matrice  $D$ -dimensionale, con  $D = \dim(\mathbf{x})$ , che contiene le varianze sulla diagonale.

**Definizione 2.5.** (Correlazione di Pearson) La CORRELAZIONE DI PEARSON tra le due variabili  $X$  e  $Y$  si definisce come

$$\rho \equiv \text{corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{\mathbb{V}[X]\mathbb{V}[Y]}}. \quad (45)$$

Possiamo scrivere

$$\text{Cov}[X] = \begin{pmatrix} \mathbb{V}[X_1] & \text{Cov}[X_1, X_2] & \cdots & \text{Cov}[X_1, X_D] \\ \text{Cov}[X_2, X_1] & \mathbb{V}[X_2] & \cdots & \text{Cov}[X_2, X_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[X_D, X_1] & \text{Cov}[X_D, X_2] & \cdots & \mathbb{V}[X_D] \end{pmatrix} \quad (46)$$

e

$$\text{corr}(x) = \begin{pmatrix} 1 & \frac{\mathbb{E}[(x_1 - \mu_1)(x_2 - \mu_2)]}{\sigma_1 \sigma_2} & \cdots & \frac{\mathbb{E}[(x_1 - \mu_1)(x_D - \mu_D)]}{\sigma_1 \sigma_D} \\ \frac{\mathbb{E}[(x_2 - \mu_2)(x_1 - \mu_1)]}{\sigma_2 \sigma_1} & 1 & \cdots & \frac{\mathbb{E}[(x_2 - \mu_2)(x_D - \mu_D)]}{\sigma_2 \sigma_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\mathbb{E}[(x_D - \mu_D)(x_1 - \mu_1)]}{\sigma_D \sigma_1} & \frac{\mathbb{E}[(x_D - \mu_D)(x_2 - \mu_2)]}{\sigma_D \sigma_2} & \cdots & 1 \end{pmatrix}. \quad (47)$$

### 2.8.1 Osservazioni

- Il fatto che due variabili non siano correlate non vuol dire che siano indipendenti, tuttavia due variabili indipendenti sono necessariamente non correlate;
- correlazione non implica causalità;
- una correlazione che appare simile in diversi set di dati può sparire (o divenire opposta) se i dati sono combinati, questo è noto come *Simpson's paradox*.

## 2.9 Distribuzione gaussiana multivariata

La definizione che abbiamo dato di distribuzione gaussiana si estende a più variabili. Innanzitutto definiamo la PDF

$$\mathcal{N}(\mathbf{y} \mid \boldsymbol{\mu}, \Sigma) \equiv \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}, \quad (48)$$

dove  $\Sigma = \text{Cov}[\mathbf{y}]$ .

Supponiamo di avere due variabili aleatorie  $\mathbf{y}_1$  e  $\mathbf{y}_2$ . Si può definire la distribuzione JOINTLY GAUSSIAN. Con

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

le distribuzioni marginali sono

$$p(\mathbf{y}_1) = \mathcal{N}(\mathbf{y}_1 \mid \boldsymbol{\mu}_1, \Sigma_{11}), \quad p(\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_2 \mid \boldsymbol{\mu}_2, \Sigma_{22}) \quad (49)$$

e la probabilità condizionata

$$p(\mathbf{y}_1 \mid \mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1 \mid \boldsymbol{\mu}_{1|2}, \Sigma_{1|2}), \quad (50)$$

con

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

queste sono tutte gaussiane.

### 3 Richiami di statistica

Nella sezione precedente abbiamo assunto come noti i parametri  $\theta$ , in questa sezione impariamo come è possibile imparare questi parametri a partire dai dati (training set)  $\mathcal{D}$ . Come abbiamo visto nell'??, il problema si riduce a trovare  $\hat{\theta}$  tale che minimizzi il misclassification rate, ??, ovvero:

$$\hat{\theta} \text{ t.c. } \mathcal{L}(\hat{\theta}) = \min [\mathcal{L}(\theta)] \equiv \min \left[ \frac{1}{N} \sum_{n=1}^N l(y_n, f(\mathbf{x}_n, \theta)) \right].$$

#### 3.1 Maximum likelihood estimation (MLE)

**Definizione 3.1** (Maximum likelihood estimation). Definiamo la MLE, MAXIMUM LIKELIHOOD ESTIMATION, come il valore di  $\theta$  che massimizza la probabilità condizionata di aver ottenuto il set di dati  $\mathcal{D}$ , ovvero:

$$\hat{\theta}_{\text{MLE}} \equiv \operatorname{argmax}_{\theta} p(\mathcal{D} | \theta). \quad (51)$$

Se i dati sono campioni indipendenti della stessa distribuzione, allora possiamo anche scrivere

$$p(\mathcal{D} | \theta) = \prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n, \theta). \quad (52)$$

Definiamo inoltre la NLL, NEGATIVE LOG-LIKELIHOOD, da usare come loss function come:

$$l(\theta) \equiv -\log p(\mathcal{D} | \theta) = -\sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \theta). \quad (53)$$

Con queste definizioni abbiamo che la MLE è data

- per supervised learning da:

$$\hat{\theta}_{\text{MLE}} = \operatorname{argmin}_{\theta} -\sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \theta); \quad (54)$$

- per unsupervised learning da:

$$\hat{\theta}_{\text{MLE}} = \operatorname{argmin}_{\theta} -\sum_{n=1}^N \log p(\mathbf{y}_n | \theta). \quad (55)$$

Possiamo giustificare l'uso della MLE pensandola come un'approssimazione del *posterior* Bayesiano dato un *prior* uniforme.

**Esempio 3.1.** Se abbiamo

$$p(\theta | \mathcal{D}) = \delta(\theta - \hat{\theta}_{\text{MAP}})$$

con  $\hat{\theta}$  il posterior, allora

$$\hat{\theta}_{\text{MAP}} = \operatorname{argmin}_{\theta} -\log p(\theta | \mathcal{D}) = \operatorname{argmin}_{\theta} -\log p(\mathcal{D} | \theta) - \log p(\theta),$$

ed essendo  $p(\theta) = 1$  abbiamo

$$\hat{\theta}_{\text{MAP}} = \operatorname{argmin}_{\theta} -\log p(\mathcal{D} | \theta) = \hat{\theta}_{\text{MLE}}. \quad \square$$

### 3.1.1 Distribuzione normale

Supponiamo ora di avere una variabile casuale  $Y$  distribuita normalmente, i.e.  $Y \sim \mathcal{N}(\mu, \sigma^2)$ , e sia  $\mathcal{D} = \{y_n \text{ t.c. } n = 1, \dots, N\}$  un dataset con punti campionati indipendentemente, allora:

$$\begin{aligned} \text{NLL}(\mu, \sigma^2) &= - \sum_{n=1}^{N_D} \log \left[ \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left( -\frac{1}{2\sigma^2} (y_n - \mu)^2 \right) \right] \\ &= \frac{1}{2\sigma} \sum_{n=1}^{N_D} (y_n - \mu)^2 + \frac{N_D}{2} \log(2\pi\sigma^2). \end{aligned} \quad (56)$$

Ora estremizzando otteniamo

$$\frac{\partial}{\partial \mu} \text{NLL}(\mu, \sigma^2) = 0 \iff \hat{\mu}_{\text{MLE}} = \frac{1}{N_D} \sum_{n=1}^{N_D} y_n = \bar{y} \quad (57)$$

$$\frac{\partial}{\partial \sigma^2} \text{NLL}(\mu, \sigma^2) = 0 \iff \hat{\sigma}_{\text{MLE}}^2 = \frac{1}{N_D} \sum_{n=1}^{N_D} (y_n - \hat{\mu}_{\text{MLE}})^2 = \frac{1}{N_D} \sum_{n=1}^{N_D} y_n^2 - \bar{y}^2 = s^2 - \bar{y}^2. \quad (58)$$

### 3.1.2 Distribuzione normale multivariata

Per una distribuzione multivariata i risultati sono analoghi. Abbiamo

$$l(\boldsymbol{\mu}, \Sigma) = \log p(\mathcal{D} \mid \boldsymbol{\mu}, \Sigma) = \frac{N_D}{2} \log |\Lambda| = -\frac{1}{2} \sum_{n=1}^{N_D} (\mathbf{y}_n - \boldsymbol{\mu})^T \Lambda (\mathbf{y}_n - \boldsymbol{\mu}) \quad (59)$$

dove  $\Lambda = \Sigma^{-1}$  è la PRECISION MATRIX. Dall'estremizzazione segue, come prima, l'empirical mean  $\hat{\boldsymbol{\mu}}$  e l'empirical covariance matrix  $\hat{\Sigma}$ :

$$\hat{\boldsymbol{\mu}} = \frac{1}{N_D} \sum_{n=1}^{N_D} \mathbf{y}_n = \bar{\mathbf{y}} \quad (60)$$

$$\hat{\Sigma} = \frac{1}{N_D} \sum_{n=1}^{N_D} (\bar{\mathbf{y}}_n - \bar{\mathbf{y}})(\bar{\mathbf{y}}_n - \bar{\mathbf{y}})^T \quad (61)$$

## 3.2 Empirical risk minimization (ERM)

La MLE può essere generalizzata sostituendo la loss function logaritmica NLL con qualunque altra funzione  $l$ :

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N l(\mathbf{y}_n, \boldsymbol{\theta}; \mathbf{x}_n). \quad (62)$$

Questo processo di cambiare loss function è noto come ERM, EMPIRICAL RISK MINIMIZATION.

**Esempio 3.2** (ERM per problema di classificazione). In un problema di classificazione avremo una loss function

$$l_{01}(\mathbf{y}_n, \boldsymbol{\theta}; \mathbf{x}_n) = \begin{cases} 1 & \text{se } \mathbf{y}_n = f(\mathbf{x}_n; \boldsymbol{\theta}) \\ 0 & \text{se } \mathbf{y}_n \neq f(\mathbf{x}_n; \boldsymbol{\theta}) \end{cases} \quad (63)$$

dove  $f(\mathbf{x}, \boldsymbol{\theta})$  è un predittore. L'empirical risk diviene

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N l_{01}(\mathbf{y}_n, \boldsymbol{\theta}; \mathbf{x}_n), \quad (64)$$

che è il misclassification rate nel training set. In problemi binari possiamo riscrivere il misclassification rate usando l'indicatore binario come definito in (??). Se  $\tilde{y} \in \{-1, +1\}$  è il true label e  $\hat{y} \in \{-1, +1\} = f(\mathbf{x}, \boldsymbol{\theta})$  la prediction, allora:

$$l_{01}(\tilde{y}, \hat{y}) = \mathbb{I}(\tilde{y} \neq \hat{y}) = \mathbb{I}(\tilde{y}\hat{y} < 0), \quad (65)$$

con cui il rischio empirico si scrive

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N l_{01}(\tilde{y}_n, \hat{y}_n) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}(\tilde{y}_n \hat{y}_n < 0). \quad (66)$$

La funzione  $l_{01}$  è però problematica in quanto non-smooth e quindi difficile da ottimizzare. In questi casi si può utilizzare una funzione surrogata, definita generalmente come limite superiore convesso, ad esempio

$$l_u(\tilde{y}, \eta) = -\log p(\tilde{y} | \eta) = \log(1 + e^{-\tilde{y}\eta}), \quad p(\tilde{y} | \mathbf{x}, \boldsymbol{\theta}) = \sigma(\hat{y}\eta) = \frac{1}{1 + e^{-\hat{y}\eta}} \quad (67)$$

### 3.3 Altri metodi di minimizzazione

#### 3.3.1 Metodo dei momenti

Il calcolo di  $\nabla_{\boldsymbol{\theta}} \text{NLL}(\boldsymbol{\theta}) = 0$ , ovvero la ricerca dei punti estremi, può essere difficile. Il metodo dei momenti consiste nell'eguagliare i momenti teorici della distribuzione ai momenti empirici, ovvero:

$$\left. \begin{array}{l} \text{MOMENTI TEORICI: } \mu_k = \mathbb{E}[Y^k] \\ \text{MOMENTI EMPIRICI: } \hat{\mu}_k = \frac{1}{N} \sum_{n=1}^N y_n^k \end{array} \right\} \mu_k = \hat{\mu}_k \implies \text{ho } k \text{ equazioni.} \quad (68)$$

Ad esempio per la gaussiana abbiamo:

$$\begin{aligned} \mu_1 &= \mu = \bar{y}, \\ \mu_2 &= \sigma^2 + \mu^2 = s^2. \end{aligned}$$

In questo caso abbiamo che MOM = MLE, ma questo non è vero in generale (ad esempio non è vero per la distribuzione uniforme).

#### 3.3.2 Online recursive estimation

Se tutto il dataset  $\mathcal{D}$  è noto e disponibile prima che inizi il processo di learning si dice che si fa BATCH LEARNING. In alcuni casi però il dataset è disponibile in blocchi e si dice che si fa ONLINE LEARNING. Si ha che  $\hat{\boldsymbol{\theta}}_{t-1}$  è la predizione dato il blocco  $\mathcal{D}_{1:t-1}$ , quello che vogliamo fare è trovare  $\boldsymbol{\theta}_t = f(\hat{\boldsymbol{\theta}}_{t-1}, \mathbf{y}_t)$  con un update ricorsivo.

**Definizione 3.2** (Moving average). Per una gaussiana multivariata si può definire la MEDIA MOBILE come:

$$\begin{aligned} \hat{\boldsymbol{\mu}}_t &= \frac{1}{t} \sum_{n=1}^t \mathbf{y}_n \\ &= \frac{1}{t} [(t-1) \hat{\boldsymbol{\mu}}_{t-1} + \mathbf{y}_t] = \hat{\boldsymbol{\mu}}_{t-1} + \frac{1}{t} (\mathbf{y}_t - \hat{\boldsymbol{\mu}}_{t-1}). \end{aligned} \quad (69)$$

Nel caso in cui la distribuzione cambi in modo continuo la moving average può anche essere pesata.

### 3.4 Regularizzazione

Un problema della MLE (maximum likelihood estimation) e della ERM (empirical risk minimization) è che i parametri tendono ad essere determinati minimizzando il loss nel training set, ma non necessariamente lo minimizzano per i dati futuri. Per ridurre l'entità del problema si può procedere con la regularizzazione:

$$\mathcal{L}(\boldsymbol{\theta}, \lambda) = \underbrace{\left[ \frac{1}{N} \sum_{n=1}^N l(\mathbf{y}_n, \boldsymbol{\theta}; \mathbf{x}_n) \right]}_{\mathcal{L}(\boldsymbol{\theta})} + \lambda C(\boldsymbol{\theta}), \quad (70)$$

dove  $\lambda$  è detto REGULARIZATION PARAMETER e  $C(\boldsymbol{\theta})$  è la PENALITY FUNCTION, ad esempio si può avere:

$$C(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}).$$

Per  $\lambda = 1$  e riscalandolo  $p(\boldsymbol{\theta})$  si ha la NLL:

$$\text{NLL}(\boldsymbol{\theta}, \lambda) = - \left[ \sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \right] = - [\log p(\mathcal{D} | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})], \quad (71)$$

che implica:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log p(\boldsymbol{\theta} | \mathcal{D}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} [\log p(\mathcal{D} | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) - \text{const}]. \quad (72)$$

Questa stima si chiama MAP (MAXIMUM A POSTERIOR) estimation.

### 3.4.1 Come scegliere $\lambda$

Un valore (troppo) piccolo per  $\lambda$  equivale a minimizzare il rischio empirico (OVERFITTING); un valore (troppo) grande equivale a essere troppo vicini al prior (UNDERFITTING). Una soluzione è dividere il training set in due classi: una di training ( $\mathcal{D}_{\text{train}}$ , 80%) ed una di validation ( $\mathcal{D}_{\text{val}}$ , 20%). Si va quindi a fittare il modello su  $\mathcal{D}_{\text{train}}$  per ogni setting  $\lambda$  e poi a valutare la performance del modello ottenuto su  $\mathcal{D}_{\text{val}}$ . Si prende dunque il valore di  $\lambda$  per cui si è ottenuta la performance migliore.

**Definizione 3.3** (Regularized empirical risk). Definiamo l'EMPIRICAL RISK REGOLARIZZATO come:

$$R_\lambda(\boldsymbol{\theta}, \mathcal{D}) \equiv \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} l(\mathbf{y}, f(\mathbf{x}, \boldsymbol{\theta})) + \lambda C(\boldsymbol{\theta}) \quad (73)$$

Il processo che seguiamo è quindi il seguente:

1. costruiamo il modello:  $\hat{\boldsymbol{\theta}}_\lambda(\mathcal{D}_{\text{train}}) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} R_\lambda(\boldsymbol{\theta}, \mathcal{D}_{\text{train}}) \quad \forall \lambda$ ;
2. calcoliamo il validation risk:  $R_\lambda^{\text{val}} = R_0(\hat{\boldsymbol{\theta}}_\lambda(\mathcal{D}_{\text{train}}), \mathcal{D}_{\text{val}})$ ;
3. prendiamo il  $\lambda$  che ha dato la miglior performance:  $\lambda^* = \underset{\lambda \in S}{\operatorname{argmin}} R_\lambda^{\text{val}}$ ;
4. si ri-fitta il modello su  $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}}$  usando  $\lambda = \lambda^*$ ;
5. si ottengono i parametri  $\hat{\boldsymbol{\theta}}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} R_{\lambda^*}(\boldsymbol{\theta}, \mathcal{D})$ .

Se il dataset è piccolo, eliminare il 20% dei dati può essere dannoso e si procede quindi con la CROSS VALIDATION. Questa consiste nel dividere il dataset in  $K$  FOLD e per ognuno allenare il modello su tutti i fold meno il  $k$ -esimo, il quale viene usato come test fold. Si definisce l'empirical risk regolarizzato per cross validation:

$$R_\lambda^{\text{CV}} \equiv \frac{1}{K} \sum_{k=1}^K R_0(\hat{\boldsymbol{\theta}}_\lambda(\mathcal{D} - \mathcal{D}_k), \mathcal{D}_k) \quad (74)$$

## 3.5 Statistica Bayesiana

Abbiamo ora visto diversi metodi per determinare i parametri a partire dei dati, ma dobbiamo ancora parlare della loro incertezza. Il ?? (??) ci aiuta in questo, infatti:

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\boldsymbol{\theta}) p(\mathcal{D} | \boldsymbol{\theta})}{p(\mathcal{D})} = \frac{p(\boldsymbol{\theta}) p(\mathcal{D} | \boldsymbol{\theta})}{\int p(\boldsymbol{\theta}') p(\mathcal{D} | \boldsymbol{\theta}') d\boldsymbol{\theta}'} \quad (75)$$

dove  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}$  per supervised learning e  $\mathcal{D} = \{(\mathbf{y}_n)\}$  per unsupervised learning,  $n = 1, \dots, N$ .

Una volta che il posterior sui paramentri è stato determinato, possiamo calcolare il posterior della distribuzione predittiva marginalizzando su  $\boldsymbol{\theta}$ :

$$p(\mathbf{y} | \mathbf{x}, \mathcal{D}) = \int p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}, \quad (76)$$

questo è noto come BMA (BAYES MODEL AVERAGING).

Consideriamo una distribuzione gaussiana di cui sia nota la varianza, nel caso univariato la likelihood per  $\mu$  ha la forma

$$p(\mathcal{D} | \mu) \propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^{N_D} (y_n - \mu)^2 \right\}. \quad (77)$$

Si può mostrare che il prior coniugato è un'altra gaussiana  $\mathcal{N}(\mu | \tilde{m}, \tilde{\tau}^2)$  e usando il ?? (??) troviamo che il posterior è

$$p(\mu | \mathcal{D}, \sigma^2) = \mathcal{N}(\mu | \hat{m}, \hat{\tau}^2) \\ \text{dove} \begin{cases} \hat{\tau}^2 &= \frac{1}{\frac{N}{\sigma^2} + \frac{1}{\tilde{\tau}^2}} = \frac{\sigma^2 \tilde{\tau}^2}{N\tilde{\tau}^2 + \sigma^2} \\ \hat{m} &= \hat{\tau}^2 \left( \frac{\tilde{m}}{\tilde{\tau}^2} + \frac{N\bar{y}}{\sigma^2} \right) = \frac{\sigma^2}{N\tilde{\tau}^2 + \sigma^2} \tilde{m} + \frac{N\tilde{\tau}^2}{N\tilde{\tau}^2 + \sigma^2} \bar{y} \end{cases} \quad (78)$$

e con  $\bar{y} \equiv \frac{1}{N} \sum_{n=1}^N y_n$  la media empirica. Definendo ora  $\kappa = 1/\sigma^2$  e  $\check{\lambda} = 1/\tilde{\tau}^2$  posso scrivere:

$$\hat{\lambda} = \check{\lambda} + N\kappa \\ \hat{m} = \frac{N\kappa\bar{y} + \check{\lambda}\tilde{m}}{\hat{\lambda}} = \frac{N\kappa}{N\kappa + \check{\lambda}} \bar{y} + \frac{\check{\lambda}}{N\kappa + \check{\lambda}} \tilde{m}. \quad (79)$$

La precisione del posterior  $\hat{\lambda}$  è dunque la precisione del prior  $\check{\lambda}$  più  $N$  unità di misura di precisione  $\kappa$ . La media del posterior ( $\hat{m}$ ) è invece una combinazione convessa della media empirica  $\bar{y}$  e della media del prior  $\tilde{m}$ . Consideriamo ora il posterior dopo aver visto un singolo dato  $y$  (i.e. abbiamo  $N = 1$ ), si ha che la POSTERIORI MEAN può essere data in tre modi:

$$\hat{m} = \underbrace{\frac{\check{\lambda}}{\hat{\lambda}} \tilde{m} + \frac{\kappa}{\hat{\lambda}} \bar{y}}_{\text{convex combination of prior and data}} = \underbrace{\tilde{m} + \frac{\kappa}{\hat{\lambda}(\bar{y} - \tilde{m})}}_{\text{prior mean adjusted to data}} = \underbrace{\bar{y} = \frac{\check{\lambda}}{\hat{\lambda}} (\bar{y} - \tilde{m})}_{\text{data adjusted to the prior mean}} \quad (80)$$

ed abbiamo lo STANDARD ERROR:

$$\text{se}(\mu) = \sqrt{\mathbb{V}[\mu | \mathcal{D}]}. \quad (81)$$

Se uso un *uninformative prior* per  $\mu$  ponendo  $\check{\lambda} = 0$ , allora  $\hat{m} = \bar{y}$ . Supponendo di approssimare  $\sigma^2 \sim s^2 \equiv \frac{1}{N} \sum_{n=1}^N (y_n - \bar{y})^2$  segue che

$$\hat{\lambda} = N\kappa = \frac{N}{s^2} \implies \text{se}(\mu) = \frac{1}{\sqrt{\hat{\lambda}}} = \frac{s}{\sqrt{N}}. \quad (82)$$

Ovvero abbiamo che l'incertezza decresce ad un rateo di  $1/\sqrt{N}$ .

**Osservazione 3.5.1.** Quando non abbiamo informazioni sul prior è desiderabile usare un prior uninformative. Ad esempio un flat prior, i.e.  $p(\mu) = 1$ .

**Osservazione 3.5.2.** Qualunque modello Bayesiano richiede di specificare un prior  $p(\theta)$  per i parametri. I parametri del prior sono detti **iperparametri** e sono denotati con  $\Phi$ . Se non sono noti possiamo metterci sopra un prior (multi-level model  $\Phi \rightarrow \theta\mathcal{D}$ ) e definire la joint distribution  $p(\phi, \theta, \mathcal{D}) = p(\Phi)p(\theta | \Phi)p(\mathcal{D} | \theta)$ .

### 3.6 Intervalli credibili e di confidenza

Una distribuzione posteriore è un oggetto multidimensionale difficile da visualizzare e da trattare. Può essere quindi utile calcolare degli stimatori puntuali, come la media e moda posteriore, per determinare un intervallo di credibilità che quantifichi l'incertezza associata alle stime.

Method	Definition
MLE	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D}   \theta)$
MAP	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D}   \theta) p(\theta   \Phi)$
Full Bayes	$p(\phi, \theta, \mathcal{D}) \propto p(\Phi) p(\theta   \Phi) p(\mathcal{D}   \theta)$

Tabella 2: Metodi visti per ottenere i parametri.

**Definizione 3.4** (Intervallo di credibilità). L'INTERVALLO DI CREDIBILITÀ  $100(1 - \alpha)\%$  è la regione  $C = (l, u)$  che contiene  $1 - \alpha$  della probabilità posteriore, ovvero:

$$C_{\alpha}(\mathcal{D}) = (l, u) \text{ t.c. } P(l \leq \theta \leq u | \mathcal{D}) = 1 - \alpha \quad (83)$$

dove  $l$  si dice *lower* e  $u$  *upper*. Inoltre se il posterior ha una PDF nota, allora  $l = F^{-1}(\alpha/2)$  e  $u = F^{-1}(1 - \alpha/2)$ , ove  $F$  è la CDF del posterior.

**Definizione 3.5** (Intervallo di confidenza). Si dice INTERVALLO DI CONFIDENZA  $100(1 - \alpha)\%$ , per la stima di un parametro  $\theta$ , la regione  $I(\tilde{\mathcal{D}}) = (l(\tilde{\mathcal{D}}), u(\tilde{\mathcal{D}}))$  ottenuto da un data set  $\mathcal{D}$  tale che

$$\Pr(\theta \in I(\tilde{\mathcal{D}} | \mathcal{D} \sim \theta)) = 1 - \alpha. \quad (84)$$

**Osservazione 3.6.1.** L'intervallo di credibilità è un concetto Bayesiano mentre l'intervallo di confidenza è un concetto frequentista. CL (frequentista):  $\theta$  è una costante fissa non nota e i dati sono aleatori. CI (Bayesiano): i dati sono fissi perchè noti,  $\theta$  è invece ignota. Un CL al 95% non significa che il parametro stia verosimilmente dentro l'intervallo il 95% delle volta dati i dati osservati.

**Esempio 3.3.** Supponiamo di generare due interi  $\mathcal{D} = (y_1, y_2)$  da:

$$p(y | \theta) = \begin{cases} 0.5 & \text{se } y = \theta \\ 0.5 & \text{se } y = \theta + 1 \\ 0 & \text{se altrimenti} \end{cases}.$$

Se  $\theta = 39$  possiamo ottenere i seguenti risultati:

$$(39, 39), (39, 40), (40, 39), (40, 40).$$

Sia  $m = \min(y_1, y_2)$  e l'intervallo  $[l(\mathcal{D}), U(\mathcal{D})] = [m, m]$ . Segue che i possibili intervalli sono:

$$[39, 39], [39, 39], [39, 39], [40, 40]$$

e dunque l'intervallo definito risulta in un CI 75% per  $\theta = 39$ . Se però osserviamo  $\mathcal{D} = (39, 40)$  allora  $p(\theta = 39 | \mathcal{D})$ , ma CI è solo 75%. Il CI fallisce per esperimenti non ripetibili.

### 3.7 Bias-variance tradeoff

Sia  $\hat{\theta}$  lo stimatore statistico e  $\hat{\theta}(\mathcal{D})$  l'estimando. Nel formalismo frequentista i dati sono variabili casuali campionati da una distribuzione  $p^*(\mathcal{D})$ , che induce una distribuzione sull'estimando  $p^*(\hat{\theta}(\mathcal{D}))$ .

**Definizione 3.6** (Bias). Il BIAS di uno stimatore  $\hat{\theta}$  è definito come:

$$\text{bias}(\theta(\cdot)) \equiv \mathbb{E}[\hat{\theta}(\mathcal{D})] - \theta^* \quad (85)$$

dove  $\theta^*$  è il valore vero. Se il bias è nulla, lo stimatore viene detto UNBIASED.

Per una gaussiana abbiamo:

$$\text{bias}(\hat{\mu}) = \mathbb{E}[\bar{x}] - \mu = \mathbb{E}\left[\frac{1}{N_D} \sum_{n=1}^{N_D} x_n\right] - \mu = \frac{N_D \mu}{N_D} - \mu = 0 \implies \text{Unbiased}, \quad (86)$$

$$\text{bias}(\hat{\sigma}^2) = \mathbb{E}[s^2] - \sigma^2 = \frac{N_D - 1}{N_D} \sigma^2 - \sigma^2 \neq 0 \implies \text{Biased}. \quad (87)$$

Possiamo però definire una versione unbiased della varianza come:

$$\sigma_{\text{unb}}^2 = \frac{1}{N_D - 1} \sum_{n=1}^{N_D} (x_n - \bar{x})^2 = \frac{N_D}{N_D - 1} \sigma_{\text{MLE}}^2 \quad (88)$$

**Definizione 3.7** (Varianza di uno stimatore). La VARIANZA DI UNO STIMATORE è

$$\mathbb{V}(\hat{\theta}) \equiv \mathbb{E}(\hat{\theta}^2) - (\mathbb{E}[\hat{\theta}])^2 \quad (89)$$

Idealmente vogliamo avere una varianza minima per tutti gli stimatori. Il seguente teorema fornisce un limite minimo alla varianza.

**Teorema 3.1** (Disuguaglianza di Cramér-Rao). Per uno stimatore unbiased  $\theta^*$  si ha che:

$$\mathbb{V}[\hat{\theta}] \geq \frac{1}{NF(\theta^*)}, \quad (90)$$

ove  $F(\theta^*)$  è la Fisher information.

Si può dimostrare che la MLE raggiunge il bond imposto da Cramér-Rao. Sia  $\hat{\theta} = \hat{\theta}(\mathcal{D})$  l'estimato e  $\bar{\theta} = \mathbb{E}[\hat{\theta}]$  il valore di aspettazione corrispondente (tutti da  $p(\mathcal{D} | \theta^*)$ ), allora:

$$\begin{aligned} \mathbb{E}[(\hat{\theta} - \theta^*)^2] &= \mathbb{E}[(\hat{\theta} - \bar{\theta}) + (\bar{\theta} - \theta^*)]^2 \\ &= \mathbb{E}[(\hat{\theta} - \bar{\theta})^2] + 2\mathbb{E}[(\hat{\theta} - \bar{\theta})(\bar{\theta} - \theta^*)] + (\bar{\theta} - \theta^*)^2 \\ &= \mathbb{E}[(\hat{\theta} - \bar{\theta})^2] + (\bar{\theta} - \theta^*)^2 = \mathbb{V}[\hat{\theta}] + \text{bias}^2(\hat{\theta}) \end{aligned} \quad (91)$$

**Esempio 3.4.** Supponiamo di stimare la media di una distribuzione gaussiana a partire dai dati  $\mathbf{x} = (x_1, \dots, x_n)$  assumendo che questi siano stati campionati da  $x_n \sim \mathcal{N}(\theta^* = 1, \sigma^2)$ . Uno stimatore ovvio è dato dalla MLE che ha bias nullo e varianza  $\mathbb{V}[\bar{x}/\theta^*] = \sigma^2/N$ . Possiamo dare uno stimatore MAP sotto un prior gaussiano, in tal caso:

$$\tilde{x} = \frac{N}{N + \kappa_0} \bar{x} + \frac{\kappa_0}{N + \kappa_0} \theta_0 = w\bar{x} + (1 - w)\theta_0$$

con prior  $\mathcal{N}(\theta_0, \sigma^2/\kappa_0)$ . Allora abbiamo:

$$\begin{aligned} \mathbb{E}_{\text{MAP}}[\tilde{x}] - \theta^* &= w\theta^* + (1 - w)\theta_0 - \theta^* = (1 - w)(\theta_0 - \theta^*), \\ \mathbb{V}_{\text{MAP}}[\tilde{x}] &= w^2 \frac{\sigma^2}{N} < \mathbb{V}_{\text{MLE}}[\hat{x}], \end{aligned}$$

ovvero è uno stimatore biased essendo  $0 < w < 1$ .

## 4 Ottimizzazione

Abbiamo visto nelle precedenti lezioni che il machine learning ha come obiettivo quello di stimare dei parametri  $\theta$  (o degli iperparametri  $\Phi$ ) di un modello, ovvero determinare i valori di best-fit e relativi intervalli di credibilità. Questo si realizza minimizzando una funzione scalare, la less function

$$\mathcal{L} : \Theta \rightarrow \mathbb{R},$$

dove  $\theta \in \Theta$ , per trovare:

$$\theta^* = \underset{\theta}{\text{argmin}} \mathcal{L}(\theta) \quad \text{o meglio} \quad \hat{\theta} = \underset{\theta}{\text{argmin}} \mathcal{L}[\theta] \quad \text{per } \theta^* \sim \hat{\theta}.$$

Assumiamo per semplicità che  $\Theta \subseteq \mathbb{R}^D$ , con  $D$  il numero delle variabili ottimizzate, in questo modo abbiamo che l'ottimizzazione è **continua**.



**Osservazione 4.0.1.** Lo spazio dei parametri è tale che può esistere un solo minimo globale o più minimi globali, oltre ad una serie di minimi locali. Il punto di minimo è computazionalmente *sempre in termini locali*, i.e.:

$$\exists \delta > 0, \forall \boldsymbol{\theta} \in \Theta \quad \text{t.c.} \quad \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\| < \delta, \mathcal{L}(\hat{\boldsymbol{\theta}}) < \mathcal{L}(\boldsymbol{\theta}).$$

un minimo locale può essere circondato da minimi simili (flat direction). Un algoritmo che converge ad un punto stazionario è detto GLOBALLY CONVERGENT.

**Osservazione 4.0.2.** Per loss function che sono due volte differenziabili possiamo considerare il vettore gradiente  $g(\boldsymbol{\theta}) = \nabla \mathcal{L}(\boldsymbol{\theta})$  e la matrice Hessiana  $H(\boldsymbol{\theta}) = \nabla^2 \mathcal{L}(\boldsymbol{\theta})$ . Per un punto  $\hat{\boldsymbol{\theta}} \in \mathbb{R}^D$ , indicando  $\hat{g} = g(\boldsymbol{\theta})|_{\hat{\boldsymbol{\theta}}}$  e  $\hat{H} = H(\boldsymbol{\theta})|_{\hat{\boldsymbol{\theta}}}$ , abbiamo:

**condizione necessaria:** se  $\hat{\boldsymbol{\theta}}$  è minimo locale, allora  $\hat{g} = \mathbf{0}$  e  $\hat{H}$  è definita semi-positiva;

**condizione sufficiente:** se  $\hat{g} = \mathbf{0}$  e  $H$  è definita positiva, allora  $\hat{\boldsymbol{\theta}}$  è un minimo locale.

**Osservazione 4.0.3.** A volte possono esserci vincoli sulla parametrizzazione (CONSTRAINED PARAMETRIZATION). I vincoli sono in genere categorizzati in equazioni e disequazioni, ad esempio:

$$\begin{aligned} h_k(\boldsymbol{\theta}) &= 0 \text{ per } k \in \mathcal{E}, \\ g_j(\boldsymbol{\theta}) &\leq 0 \text{ per } j \in \mathcal{I}. \end{aligned}$$

Allora il set di vincoli è:

$$\mathcal{C} = \{\boldsymbol{\theta} \text{ t.c. } h_k(\boldsymbol{\theta}) = 0, k \in \mathcal{E} \wedge g_j(\boldsymbol{\theta}) \leq 0, j \in \mathcal{I}\}$$

ed il parametro è dato da:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \mathcal{C}}{\operatorname{argmin}} \mathcal{L}(\boldsymbol{\theta}).$$

Quando non si hanno vincoli e quindi  $\boldsymbol{\theta} \in \mathbb{R}^D$  si parla di UNCONSTRAINED PARAMETRIZATION.

**Osservazione 4.0.4.** Può essere utile determinare se la loss function è convessa. Ricordiamo quindi che se una funzione  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  è doppio differenziabile, allora questa è convessa se e solo se l'Hessiana è definita semi-positiva in ogni punto.  $f$  si dice invece strettamente convessa se l'Hessiana è definita positiva.

**Osservazione 4.0.5.** Nel caso in cui la loss function non sia continua (in un punto) allora si usa definire i **sottogradienti**.  $\mathbf{g} \in \mathbb{R}^N$  è un sottogradiente di  $f$  in  $\mathbf{x} \in \operatorname{Dom}(f)$  se:

$$\forall \boldsymbol{\xi} \in \operatorname{Dom}(f), f(\boldsymbol{\xi}) \geq f(\mathbf{x}) + \mathbf{g}^T(\boldsymbol{\xi} - \mathbf{x}).$$

## 4.1 Metodi del primo ordine

Questi metodi si basano sulle derivate prime delle loss-function, ovvero guardano a quali sono le direzioni negative nello spazio dei parametri. Un punto di partenza deve essere specificato e viene chiamato  $\boldsymbol{\theta}_0$ . Per ogni iterazione  $t$  i parametri vengono aggiornati come:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta_t \mathbf{d}_t, \tag{92}$$

ove  $\eta_t$  è detta STEP SIZE o LEARNING RATE, che in sequenza  $\{\eta_t\}$  identifica il LEARNING RATE SCHEDULE. Mentre  $\mathbf{d}_t$  è la DESCENT DIRECTION. Come descent direction si può usare la STEEPEST DESCENT:

$$\mathbf{g}_t = \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta}_t} \longrightarrow \mathbf{d}_t = -\mathbf{g}_t. \tag{93}$$

Il gradiente  $\mathbf{g}_t$  punta nella direzione di massima variazione di  $\mathcal{L}(\boldsymbol{\theta}_t)$  quindi questa scelta sembra ovvia.

L'algoritmo viene iterato fino a quando non si raggiunge un punto stazionario. Formalmente richiediamo che esista  $\eta_{\max} > 0$  tale che  $\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{d}) < \mathcal{L}(\boldsymbol{\theta}) \forall \eta : 0 < \eta < \eta_{\max}$ .

La scelta più semplice è  $\eta = \text{const.}$ , tuttavia se  $\eta$  è troppo grande il metodo potrebbe non convergere, mentre se è troppo piccolo la convergenza potrebbe essere troppo lenta. In genere è meglio scegliere  $\eta$  in modo adattivo, così da avere la massima riduzione della loss function lungo la direzione scelta:

$$\eta_t = \underset{\eta > 0}{\operatorname{argmin}} \Phi_t(\eta) = \underset{\eta > 0}{\operatorname{argmin}} \mathcal{L}(\boldsymbol{\theta}_t + \eta d\mathbf{t}). \quad (94)$$

Questo approccio è detto LINE SEARCH poichè cerchiamo lungo la direzione identificata da  $d\mathbf{t}$ .

Vogliamo inoltre usare algoritmi che convergono rapidamente. Nel caso la loss function sia convessa si può mostrare che la gradient descent converge con tasso lineare:

$$|\mathcal{L}(\boldsymbol{\theta}_{t+1})| \leq \mu |\mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\hat{\boldsymbol{\theta}})|, \quad (95)$$

con  $0 < \mu < 1$  il tasso di convergenza.

Il metodo di gradient descent può però essere molto insufficiente qualora ci si trovi in direzioni piatte. Una tecnica nota come METODO DEI MOMENTI o della palla pesata permette di risolvere questo problema. Questa prescrive di muoversi velocemente lungo direzioni piatte, in base alle iterazioni precedenti, e lentamente ove il gradiente cambia rapidamente. Si scrivono le quantità

$$\mathbf{m}_t = \beta \mathbf{m}_{t-1} + \mathbf{g}_{t-1} \quad \text{e} \quad \boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \eta_t \mathbf{m}_t, \quad (96)$$

dove  $\mathbf{m}_t$  è la *quantità di moto* e  $0 < \beta < 1$ , tipicamente  $\beta \sim 0.9$ . Sostituendo  $\mathbf{m}_{t-1}$ , poi  $\mathbf{m}_{t-2}$  e così via, otteniamo una definizione per la quantità di moto come serie:

$$\mathbf{m}_t = \beta \mathbf{m}_{t-1} + \mathbf{g}_{t-1} = \beta^2 \mathbf{m}_{t-2} + \beta \mathbf{g}_{t-2} + \mathbf{g}_{t-1} = \sum_{\tau=0}^{t-1} \beta^\tau \mathbf{g}_{t-\tau-1}. \quad (97)$$

Se tutti i gradienti nel passato sono costanti, allora:

- $\mathbf{m}_t = \mathbf{g} \sum_{\tau=0}^{t-1} \beta^\tau$  e lo scaling factor è una serie geometrica;
- $\lim_{t \rightarrow \infty} \sum_{\tau=0}^{t-1} \beta^\tau = \frac{1}{1-\beta}$  e quindi se  $\beta \sim 0.9$  moltiplichiamo il gradiente per 10.

## 4.2 Metodi del second'ordine

Il classico metodo alle derivate seconde è il METODO DI NEWTON. Questo consiste nell'aggiornare i parametri come:

(98)