

Il Dizionario **A'**alto del Machine Learning

Alexander Jung¹, Konstantina Olioumtsevits¹, e Juliette Gronier²

¹Aalto University ²ENS Lyon

23 giugno 2025



please cite as: A. Jung and K. Olioumtsevits, *The Aalto Dictionary of Machine Learning*. Espoo, Finland: Aalto University, 2025.

Riconoscimenti

Questo dizionario di machine learning si è evoluto nel corso dello sviluppo e dell'insegnamento di diversi corsi, tra cui CS-E3210 Machine Learning: Principi di base, CS-C3240 Machine Learning, CS-E4800 Artificial Intelligence, CS-EJ3211 Machine Learning with Python, CS-EJ3311 Deep Learning with Python, CS-E4740 Federated Learning e CS-E407507 Human-Centered Machine Learning. Tali corsi sono stati erogati presso l'Aalto University <https://www.aalto.fi/en>, a studenti adulti tramite il Finnish Institute of Technology (FITech) <https://fitech.io/en/>, e a studenti internazionali attraverso l'European University Alliance Unite! <https://www.aalto.fi/en/unite>. Siamo grati agli studenti che hanno fornito preziosi feedback, contribuendo a plasmare questo dizionario. Un ringraziamento speciale va a Mikko Seesto per la sua meticolosa revisione. Alcune delle illustrazioni nel glossario sono state preparate con l'aiuto di Salvatore Rastelli.

Lists of Symbols

Insiemi e Funzioni

$a \in \mathcal{A}$ L'oggetto a è un elemento dell'insieme \mathcal{A} .

$a := b$ Usiamo a come abbreviazione di b .

$|\mathcal{A}|$ La cardinalità (cioè, il numero di elementi) di un insieme finito \mathcal{A} .

$\mathcal{A} \subseteq \mathcal{B}$ \mathcal{A} è un sottoinsieme di \mathcal{B} .

$\mathcal{A} \subset \mathcal{B}$ \mathcal{A} è un sottoinsieme proprio di \mathcal{B} .

\mathbb{N} I numeri naturali $1, 2, \dots$

\mathbb{R} I numeri reali x [1].

\mathbb{R}_+ I numeri reali non negativi $x \geq 0$.

\mathbb{R}_{++} I numeri reali strettamente positivi $x > 0$.

$\{0, 1\}$	L'insieme costituito dai due numeri reali 0 e 1.
$[0, 1]$	L'intervallo chiuso dei numeri reali x tali che $0 \leq x \leq 1$.
$\operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$	L'insieme dei minimizzatori di una funzione reale $f(\mathbf{w})$.
$\mathbb{S}^{(n)}$	L'insieme dei vettori a norma unitaria in \mathbb{R}^{n+1} .
$\exp(a)$	La funzione esponenziale calcolata nel punto reale $a \in \mathbb{R}$.
$\log(a)$	I logaritmo del numero reale strettamente positivo $a \in \mathbb{R}_{++}$.
$f(\cdot) : \mathcal{A} \rightarrow \mathcal{B} : a \mapsto h(a)$	<p>Una funzione (mappa) che accetta come input un qualsiasi elemento $a \in \mathcal{A}$, appartenente ad un insieme \mathcal{A}, e restituisce un elemento ben definito $f(a) \in \mathcal{B}$, appartenente ad un insieme \mathcal{B}. L'insieme \mathcal{A} è detto dominio della funzione h, mentre l'insieme \mathcal{B} è il codominio di h. Il Machine learning (ML) punta ad individuare (o apprendere) una funzione h (cioè, un ipotesi) che, a partire dalle caratteristiche \mathbf{x} di un punto dati, fornisca una previsione $h(\mathbf{x})$ della sua etichetta y.</p>
$\operatorname{epi}(f)$	L' gls(epigrafo) di una function a valori reali $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
$\frac{\partial f(w_1, \dots, w_d)}{\partial w_j}$	<p>La derivata parziale (se esiste) di una function a valori reali $f : \mathbb{R}^d \rightarrow \mathbb{R}$ rispetto a w_j [2, Ch. 9].</p> <p>Si veda anche: function.</p>
$\nabla f(\mathbf{w})$	<p>4</p> <p>Il gradiente di una funzione reale differenziabile $f : \mathbb{R}^d \rightarrow \mathbb{R}$ è il vettore $\nabla f(\mathbf{w}) = \left(\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d} \right)^T \in \mathbb{R}^d$ [2, Ch. 9].</p>

Matrici e Vettori

$\mathbf{x} = (x_1, \dots, x_d)^T$	Un vettore di lunghezza d , il cui j -esimo elemento è x_j .
\mathbb{R}^d	L'insieme di vettori $\mathbf{x} = (x_1, \dots, x_d)^T$ composto da d elementi a valori reali $x_1, \dots, x_d \in \mathbb{R}$.
$\mathbf{I}_{l \times d}$	Una matrice identità generalizzata con l righe e d colonne. Gli elementi di $\mathbf{I}_{l \times d} \in \mathbb{R}^{l \times d}$ sono pari a 1 lungo la diagonale principale e pari a 0 altrove.
\mathbf{I}_d, \mathbf{I}	Una matrice identità quadrata di dimensione $d \times d$. Se la dimensione è chiara dal contesto, si omette l'indice.
$\ \mathbf{x}\ _2$	La norma Euclidea (or ℓ_2) del vettore $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ definita come $\ \mathbf{x}\ _2 := \sqrt{\sum_{j=1}^d x_j^2}$.
$\ \mathbf{x}\ $	Una qualche norma del vettore $\mathbf{x} \in \mathbb{R}^d$ [3]. Salvo diversa specificazione, intendiamo la norma Euclidea $\ \mathbf{x}\ _2$.
\mathbf{x}^T	La trasposizione di una matrice che ha il vettore $\mathbf{x} \in \mathbb{R}^d$ come unica colonna.
\mathbf{X}^T	La trasposizione di una matrice $\mathbf{X} \in \mathbb{R}^{m \times d}$. Una matrice quadrata a valori reali $\mathbf{X} \in \mathbb{R}^{m \times m}$ è detta simmetrica se $\mathbf{X} = \mathbf{X}^T$.
\mathbf{X}^{-1}	L' inverse matrix di una matrice $\mathbf{X} \in \mathbb{R}^{d \times d}$. Si veda anche: inverse matrix.
$\mathbf{0} = (0, \dots, 0)^T$	Il vettore in \mathbb{R}^d con ciascuna componente uguale a zero.
$\mathbf{1} = (1, \dots, 1)^T$	Il vettore in \mathbb{R}^d con ciascuna componente uguale a uno.

$(\mathbf{v}^T, \mathbf{w}^T)^T$	Il vettore di lunghezza $d + d'$ ottenuto concatenando le componenti del vettore $\mathbf{v} \in \mathbb{R}^d$ con le componenti di $\mathbf{w} \in \mathbb{R}^{d'}$.
$\text{span}\{\mathbf{B}\}$	Lo span di una matrice $\mathbf{B} \in \mathbb{R}^{a \times b}$, che rappresenta il sottospazio generato da tutte le combinazioni lineari delle colonne di \mathbf{B} , $\text{span}\{\mathbf{B}\} = \{\mathbf{B}\mathbf{a} : \mathbf{a} \in \mathbb{R}^b\} \subseteq \mathbb{R}^a$.
$\det(\mathbf{C})$	Il determinante della matrice \mathbf{C} .
$\mathbf{A} \otimes \mathbf{B}$	Il prodotto di Kronecker di \mathbf{A} e \mathbf{B} [4].

Teoria della Probabilità

$\mathbf{x} \sim p(\mathbf{z})$	<p>La variabile aleatoria \mathbf{x} è distribuita secondo la distribuzione di probabilità $p(\mathbf{z})$ [?], [?].</p> <p>Si veda anche: variabile aleatoria, distribuzione di probabilità.</p>
$\mathbb{E}_p\{f(\mathbf{z})\}$	<p>Il valore atteso di una funzione $f(\mathbf{z})$ di una variabile aleatoria \mathbf{z} la cui distribuzione di probabilità è $\mathbb{P}(\mathbf{z})$. Se la distribuzione di probabilità è chiara dal contesto, scriviamo semplicemente $\mathbb{E}\{f(\mathbf{z})\}$.</p>
$\text{cov}(x, y)$	<p>La covarianza tra due variabili aleatorie a valori reali definite su uno spazio di probabilità comune.</p> <p>See also: covarianza, variabile aleatoria, distribuzione di probabilità.</p>
$\mathbb{P}(\mathbf{x}, y)$	<p>Una distribuzione di probabilità (congiunta) di una variabile aleatoria le cui realizzazioni sono punto dati con caratteristiche \mathbf{x} e etichetta y.</p>
$\mathbb{P}(\mathbf{x} y)$	<p>Una distribuzione di probabilità condizionata di una variabile aleatoria \mathbf{x}, dato il valore di un'altra variabile aleatoria y [5, Sec. 3.5].</p>
$\mathbb{P}(\mathbf{x}; \mathbf{w})$	<p>Una distribuzione di probabilità parametrizzata di una variabile aleatoria \mathbf{x}. La distribuzione di probabilità dipende da un vettore di parametri \mathbf{w}. Ad esempio, $\mathbb{P}(\mathbf{x}; \mathbf{w})$ potrebbe essere una multivariate normal distribution con il vettore di parametri \mathbf{w} costituito dagli elementi del vettore di media $\mathbb{E}\{\mathbf{x}\}$ e della matrice di covarianza $\mathbb{E}\left\{(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T\right\}$.</p>
$\mathcal{N}(\mu, \sigma^2)$	<p>La distribuzione di probabilità di una variabile aleatoria Gaussiana $x \in \mathbb{R}$ con media (o valore atteso) $\mu = \mathbb{E}\{x\}$ e</p>

Machine Learning

r	Un indice $r = 1, 2, \dots$ che enumera punti dati.
m	Il numero di punti dati presenti in (cioè, la dimensione di) un dataset.
\mathcal{D}	Un dataset $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ è un elenco di singoli punti dati $\mathbf{z}^{(r)}$, con $r = 1, \dots, m$.
d	Il numero di caratteristiche che caratterizza un punto dati.
x_j	La j -esima feature di un punto dati. La prima caratteristica è indicata con x_1 , la seconda caratteristica x_2 , e così via.
\mathbf{x}	Il vettore delle caratteristiche $\mathbf{x} = (x_1, \dots, x_d)^T$ di un punto dati le cui componenti sono le singole caratteristiche di un punto dati.
\mathcal{X}	Lo feature space \mathcal{X} è l'insieme di tutti i valori possibili che le caratteristiche \mathbf{x} di un punto dati possono assumere.
\mathbf{z}	Invece del simbolo \mathbf{x} , talvolta utilizziamo \mathbf{z} come simbolo alternativo per denotare un vettore i cui elementi sono le singole caratteristiche di un punto dati. È necessario ricorrere a due simboli distinti per differenziare tra caratteristiche grezze e quelle apprese [6, Ch. 9].
$\mathbf{x}^{(r)}$	Il vettore di caratteristiche relativo al r -esimo punto dati all'interno di un dataset.

$x_j^{(r)}$	La j -esima caratteristica relativa all' r -esimo punto dati all'interno di un dataset.
\mathcal{B}	Un mini-batch (o sottoinsieme) di punti dati scelti casualmente.
B	La dimensione di un mini-batch, ovvero il numero di punti dati che esso contiene.
y	L' etichetta (o quantità di interesse) di un punto dati.
$y^{(r)}$	L' etichetta del r -esimo punto dati.
$(\mathbf{x}^{(r)}, y^{(r)})$	Le caratteristiche e le etichetta del r -esimo punto dati.

\mathcal{Y} Lo spazio delle etichette \mathcal{Y} di un algoritmo di ML consiste nell'insieme di tutti i possibili valori di etichetta che un punto dati può assumere. Lo spazio delle etichette nominale può essere più ampio rispetto all'insieme dei diversi valori di etichetta effettivamente presenti in un dato dataset (ad esempio, un insieme di addestramento). I problemi (o metodi) di ML che impiegano uno spazio delle etichette numerico, come $\mathcal{Y} = \mathbb{R}$ or $\mathcal{Y} = \mathbb{R}^3$, sono detti problemi (o metodi) di regressione. I problemi (o metodi) di ML che utilizzano uno spazio delle etichette discreto, come $\mathcal{Y} = \{0, 1\}$ or $\mathcal{Y} = \{cat, dog, mouse\}$, sono detti problemi (o metodi) di classificazione.

η	Learning rate (o step size) usato dai gradient-based methods.
$h(\cdot)$	Una funzione ipotesi che, dato in input un vettore di caratteristiche \mathbf{x} relativo ad un punto dati, restituisce una previsione $\hat{y} = h(\mathbf{x})$ della sua etichetta y .
$\mathcal{Y}^{\mathcal{X}}$	Dati due insiemi \mathcal{X} e \mathcal{Y} , indichiamo con $\mathcal{Y}^{\mathcal{X}}$ l'insieme di tutte le possibili funzioni ipotesi $h : \mathcal{X} \rightarrow \mathcal{Y}$.
\mathcal{H}	Uno spazio delle ipotesi o modello utilizzato da un algoritmo di ML. Lo spazio delle ipotesi consiste in diverse funzioni ipotesi $h : \mathcal{X} \rightarrow \mathcal{Y}$, tra le quali il metodo di ML deve scegliere.
$d_{\text{eff}}(\mathcal{H})$	L' dimensione effettiva di uno spazio delle ipotesi \mathcal{H} .
B^2	Il bias al quadrato di un' ipotesi \hat{h} appresa, generata da un metodo di ML. Il metodo è addestrato su punti dati modellati come realizzazioni di variabile aleatoria. Poiché i dati costituiscono una realizzazione di variabile aleatoria, anche l'ipotesi appresa \hat{h} rappresenta una realizzazione di una variabile aleatoria.
V	La varianza dell' ipotesi appresa (o dei suoi ??) generata da un metodo di ML. Il metodo è addestrato su punti dati modellati come realizzazioni di variabile aleatoria. Poiché i dati costituiscono una realizzazione di variabile aleatoria, anche l'ipotesi appresa \hat{h} rappresenta una realizzazione di una variabile aleatoria. ¹⁰

$L((\mathbf{x}, y), h)$	La perdita alla stima della etichetta y di punto dati mediante l'utilizzo della previsione $\hat{y} = h(\mathbf{x})$. La previsione \hat{y} è ottenuta applicando l'ipotesi $h \in \mathcal{H}$ al vettore delle caratteristiche \mathbf{x} del punto dati in questione.
E_v	L'errore di validazione di un'ipotesi h , ovvero la media delle perdita da essa sostenute sull'insieme di validazione.
$\hat{L}(h \mathcal{D})$	Il empirical risk, o perdita media, associata all'ipotesi h su un dataset \mathcal{D} .
E_t	L'errore di addestramento di un'ipotesi h , ovvero la media delle perdita da essa sostenute sull'insieme di addestramento.
t	Un indice temporale discreto $t = 0, 1, \dots$ utilizzato per enumerare eventi sequenziali (o istanti temporali).
t	Un indice che elenca i learning task all'interno di un problema di multitask learning.
α	Un parametro di regularization che controlla il grado di regularization.
$\lambda_j(\mathbf{Q})$	Il j -esimo eigenvalue (ordinato in modo crescente o decrescente) di una matrice positive semi-definite (psd) \mathbf{Q} . Utilizziamo inoltre la notazione abbreviata λ_j qualora la matrice corrispondente sia chiara dal contesto.
$\sigma(\cdot)$	La activation function utilizzata da un neurone artificiale all'interno di una artificial neural network (ANN).

$\mathcal{R}_{\hat{y}}$	Una decision region all'interno di uno feature space.
\mathbf{w}	Un vettore di parametri $\mathbf{w} = (w_1, \dots, w_d)^T$ di un modello, e.g., i weights di unlinear model o di una ANN.
$h^{(\mathbf{w})}(\cdot)$	Una funzione di ipotesi espressa in termini di model parameters regolabili w_1, \dots, w_d concatenati nel vettore $\mathbf{w} = (w_1, \dots, w_d)^T$.
$\phi(\cdot)$	Una feature map $\phi : \mathcal{X} \rightarrow \mathcal{X}' : \mathbf{x} \mapsto \mathbf{x}' := \phi(\mathbf{x}) \in \mathcal{X}'$.
$K(\cdot, \cdot)$	Dato uno feature space \mathcal{X} , un kernel è una mappa $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ psd.

Federated Learning

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$	Un graph non orientato i cui nodi $i \in \mathcal{V}$ rappresentano i devices all'interno di un federated learning network (FL network). Gli spigoli pesati \mathcal{E} rappresentano la connettività tra i device e le similitudini statistiche tra i rispettivi dataset e learning task.
$i \in \mathcal{V}$	Un nodo che rappresenta un device all'interno di un FL network. Il device ha accesso ad un local dataset e può addestrare un local model.
$\mathcal{G}^{(\mathcal{C})}$	Il sottografo di \mathcal{G} indotto dai nodi $\mathcal{C} \subseteq \mathcal{V}$.
$\mathbf{L}^{(\mathcal{G})}$	La Laplacian matrix di un graph \mathcal{G} .
$\mathbf{L}^{(\mathcal{C})}$	La Laplacian matrix del graph indotto $\mathcal{G}^{(\mathcal{C})}$.
$\mathcal{N}^{(i)}$	L'neighborhood di un nodo i in un graph \mathcal{G} .
$d^{(i)}$	Il grado pesato $d^{(i)} := \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'}$ di un nodo i in un graph \mathcal{G} .
$d_{\max}^{(\mathcal{G})}$	Il massimo grado pesato dei nodi di un \mathcal{G} .
$\mathcal{D}^{(i)}$	Il local dataset $\mathcal{D}^{(i)}$ associato al nodo $i \in \mathcal{V}$ of an FL network.
m_i	Il numero di punti dati (ossia, la sample size) contenuti nel local dataset $\mathcal{D}^{(i)}$ relativo al nodo $i \in \mathcal{V}$.

$\mathbf{x}^{(i,r)}$	Le caratteristica dell' r -esimo punto dati nel local dataset $\mathcal{D}^{(i)}$.
$y^{(i,r)}$	La etichetta dell' r -esimo punto dati nel local dataset $\mathcal{D}^{(i)}$.
$\mathbf{w}^{(i)}$	I model parameters locali del device i all'interno di un FL network.
$L_i(\mathbf{w})$	La loss function locale utilizzata dal device i per valutare l'efficacia di una certa scelta \mathbf{w} come model parameters locali.
$L^{(d)}(\mathbf{x}, h(\mathbf{x}), h'(\mathbf{x}))$	La perdita associata ad un' ipotesi h' su un punto dati con caratteristiche \mathbf{x} e etichetta $h(\mathbf{x})$, ottenuta da un'altra ipotesi.
$\text{stack}\{\mathbf{w}^{(i)}\}_{i=1}^n$	Il vettore $\left((\mathbf{w}^{(1)})^T, \dots, (\mathbf{w}^{(n)})^T\right)^T \in \mathbb{R}^{dn}$ ottenuto concatenando verticalmente i model parameters locali $\mathbf{w}^{(i)} \in \mathbb{R}^d$.

Machine Learning Concepts

k -fold cross-validation (k -fold CV) k -fold CV is a method for learning and validating a ipotesi using a given dataset. This method divides the dataset evenly into k subsets or folds and then executes k repetitions of modello training (e.g., via empirical risk minimization (ERM)) and validazione. Each repetition uses a different fold as the insieme di validazione and the remaining $k - 1$ folds as a insieme di addestramento. The final output is the average of the errore di validazione obtained from the k repetitions.

k -means The k -medias algorithm is a hard clustering method which assigns each punto dati of a dataset to precisely one of k different clusters. The method alternates between updating the cluster assignments (to the cluster with the nearest media) and, given the updated cluster assignments, re-calculating the cluster medias [6, Ch. 8].

absolute error loss Consider a punto dati with caratteristicas $\mathbf{x} \in \mathcal{X}$ and numeric etichetta $y \in \mathbb{R}$. The absolute error perdita incurred by a ipotesi $h : \mathcal{X} \rightarrow \mathbb{R}$ is defined as $|y - h(\mathbf{x})|$, i.e., the absolute difference between the previsione $h(\mathbf{x})$ and the true etichetta y .

accuracy Consider punti dati characterized by caratteristiche $\mathbf{x} \in \mathcal{X}$ and a categorical etichetta y which takes on values from a finite spazio delle etichette \mathcal{Y} . The accuracy of a ipotesi $h : \mathcal{X} \rightarrow \mathcal{Y}$, when applied

to the punti dati in a dataset $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$, is then defined as $1 - (1/m) \sum_{r=1}^m L^{(0/1)}((\mathbf{x}^{(r)}, y^{(r)}), h)$ using the 0/1 loss $L^{(0/1)}(\cdot, \cdot)$.

See also: perdita, 0/1 loss, ??.

activation function Each artificial neuron within an ANN is assigned an activation function $\sigma(\cdot)$ that maps a weighted combination of the neuron inputs x_1, \dots, x_d to a single output value $a = \sigma(w_1x_1 + \dots + w_dx_d)$. Note that each neuron is parametrized by the weights w_1, \dots, w_d .

algorithm An algorithm is a precise, step-by-step specification for how to produce an output from a given input within a finite number of computational steps [?]. For example, an algorithm for training a linear model explicitly describes how to transform a given insieme di addestramento into model parameters through a sequence of gradient steps. This informal characterization can be formalized rigorously via different mathematical modelli [?]. One very simple modello of an algorithm is a collection of possible executions. Each execution is a sequence:

$$\text{input}, s_1, s_2, \dots, s_T, \text{output}$$

that respects the constraints inherent to the computer executing the algorithm. Algorithms may be deterministic, where each input results uniquely in a single execution, or randomized, where executions can vary probabilistically. Randomized algorithms can thus be analyzed by modeling execution sequences as outcomes of random experiments, viewing the algorithm as a stochastic process [?, ?, 5]. Crucially, an

algorithm encompasses more than just a mapping from input to output; it also includes the intermediate computational steps s_1, \dots, s_T .

application programming interface (API) An API is a formal mechanism for enabling software components to interact in a structured manner [?]. In the context of ML, APIs are frequently used to make a trained ML modello accessible to different types of users. These users, which can be other computers or humans, can request a previsione for the etichetta of a punto dati by providing its caratteristiche. The internal structure of the ML modello remains hidden from the user. For instance, consider a trained ML modello $\hat{h}(x) := 2x + 1$. An API enables a user to submit the caratteristica value $x = 3$ and obtain the response $\hat{h}(3) = 7$ without knowledge of the detailed structure of the ML modello or its training. In practice, the ML modello is typically hosted on a computer (i.e., a server) connected to the internet. Another computer (i.e., a client) sends the caratteristiche of a punto dati to the server, which then computes $\hat{h}(\mathbf{x})$ and returns the result to the external system. APIs help to modularize the development of ML applications by decoupling specific tasks. For instance, one team can focus on developing and training the modello, while another team handles user interaction and integration of the modello into applications.

artificial intelligence (AI) AI refers to systems that behave rationally in the sense of maximizing a long-term reward. The ML-based approach to AI is to train a modello for predicting optimal actions. These predictions are computed from observations about the state of the environment.

The choice of loss function sets AI applications apart from more basic ML applications. AI systems rarely have access to a labeled insieme di addestramento that allows the average perdita to be measured for any possible choice of model parameters. Instead, AI systems use observed reward signals to obtain a (point-wise) estimate for the perdita incurred by the current choice of model parameters.

artificial neural network (ANN) An ANN is a graphical (signal-flow) representation of a function that maps caratteristiche of a punto dati at its input to a previsione for the corresponding etichetta at its output. The fundamental unit of an ANN is the artificial neuron, which applies an activation function to its weighted inputs. The outputs of these neurons serve as inputs for other neurons, forming interconnected layers.

attack An attack on an ML system refers to an intentional action—either active or passive—that compromises the system’s integrity, availability, or confidentiality. Active attacks involve perturbing components such as datasets (via data poisoning) or communication links between devices in a federated learning (FL) setting. Passive attacks, such as privacy attacks, aim to infer sensitive attributes without modifying the system. Depending on their goal, we distinguish between denial-of-service attacks, backdoor attacks, and privacy attacks.

See also: data poisoning, privacy attack, sensitive attribute, denial-of-service attack, backdoor.

autoencoder An autoencoder is an ML method that simultaneously learns an encoder map $h(\cdot) \in \mathcal{H}$ and a decoder map $h^*(\cdot) \in \mathcal{H}^*$. It is an

instance of ERM using a perdita computed from the reconstruction error $\mathbf{x} - h^*(h(\mathbf{x}))$.

backdoor A backdoor attack refers to the intentional manipulation of the training process underlying an ML method. This manipulation can be implemented by perturbing the insieme di addestramento (dati poisoning) or the optimization algorithm used by an ERM-based method. The goal of a backdoor attack is to nudge the learned ipotesi \hat{h} towards specific previsioni for a certain range of caratteristica values. This range of caratteristica values serves as a key (or trigger) to unlock a backdoor in the sense of delivering anomalous previsioni. The key \mathbf{x} and the corresponding anomalous previsione $\hat{h}(\mathbf{x})$ are only known to the attacker.

bagging Bagging (or bootstrap aggregation) is a generic technique to improve (the robustness of) a given ML method. The idea is to use the bootstrap to generate perturbed copies of a given dataset and then to learn a separate ipotesi for each copy. We then predict the etichetta of a punto dati by combining or aggregating the individual previsioni of each separate ipotesi. For ipotesi maps delivering numeric etichetta values, this aggregation could be implemented by computing the average of individual previsioni.

baseline Consider some ML method that produces a learned ipotesi (or trained modello) $\hat{h} \in \mathcal{H}$. We evaluate the quality of a trained modello by computing the average perdita on a test set. But how can we assess whether the resulting test set performance is sufficiently good? How

can we determine if the trained modello performs close to optimal and there is little point in investing more resources (for dati collection or computation) to improve it? To this end, it is useful to have a reference (or baseline) level against which we can compare the performance of the trained modello. Such a reference value might be obtained from human performance, e.g., the misclassification rate of dermatologists who diagnose cancer from visual inspection of skin [?]. Another source for a baseline is an existing, but for some reason unsuitable, ML method. For example, the existing ML method might be computationally too expensive for the intended ML application. Nevertheless, its test set error can still serve as a baseline. Another, somewhat more principled, approach to constructing a baseline is via a modello probabilistico. In many cases, given a modello probabilistico $p(\mathbf{x}, y)$, we can precisely determine the minimo achievable risk among any hypotheses (not even required to belong to the spazio delle ipotesi \mathcal{H}) [?]. This minimo achievable risk (referred to as the Bayes risk) is the risk of the Bayes estimator for the etichetta y of a punto dati, given its caratteristicas \mathbf{x} . Note that, for a given choice of loss function, the Bayes estimator (if it exists) is completely determined by the distribuzione di probabilità $p(\mathbf{x}, y)$ [?, Ch. 4]. However, computing the Bayes estimator and Bayes risk presents two main challenges:

- 1) The distribuzione di probabilità $p(\mathbf{x}, y)$ is unknown and needs to be estimated.
- 2) Even if $p(\mathbf{x}, y)$ is known, it can be computationally too expensive to compute the Bayes risk exactly [?].

A widely used modello probabilistico is the multivariate normal distribution $(\mathbf{x}, y) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for punto dati characterized by numeric caratteristiche and etichettas. Here, for the squared error loss, the Bayes estimator is given by the posterior media $\mu_{y|\mathbf{x}}$ of the etichetta y , given the caratteristiche \mathbf{x} [?, ?]. The corresponding Bayes risk is given by the posterior varianza $\sigma_{y|\mathbf{x}}^2$ (see Figure 1).

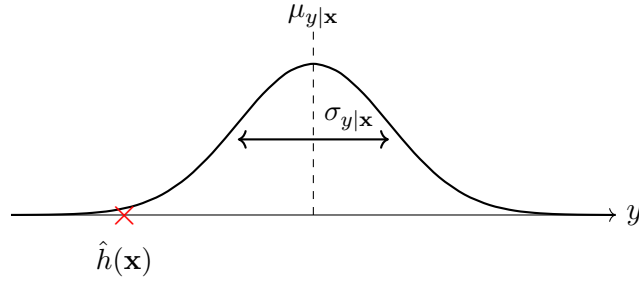


Figura 1: If the caratteristiche and the etichetta of a punto dati are drawn from a multivariate normal distribution, we can achieve the minimo risk (under squared error loss) by using the Bayes estimator $\mu_{y|\mathbf{x}}$ to predict the etichetta y of a punto dati with caratteristiche \mathbf{x} . The corresponding minimo risk is given by the posterior varianza $\sigma_{y|\mathbf{x}}^2$. We can use this quantity as a baseline for the average perdita of a trained modello \hat{h} .

batch In the context of stochastic gradient descent (SGD), a batch refers to a randomly chosen subset of the overall insieme di addestramento. We use the punti dati in this subset to estimate the gradiente of errore di addestramento and, in turn, to update the model parameters. See also: SGD, insieme di addestramento, punto dati, gradiente, errore di addestramento, model parameters.

Bayes estimator Consider a modello probabilistico with a joint distribuzione di probabilità $p(\mathbf{x}, y)$ for the caratteristics \mathbf{x} and etichetta y of a punto dati. For a given loss function $L(\cdot, \cdot)$, we refer to a ipotesi h as a Bayes estimator if its risk $\mathbb{E}\{L((\mathbf{x}, y), h)\}$ is the minimo [?]. Note that the property of a ipotesi being a Bayes estimator depends on the underlying distribuzione di probabilità and the choice for the loss function $L(\cdot, \cdot)$.

Bayes risk Consider a modello probabilistico with a joint distribuzione di probabilità $p(\mathbf{x}, y)$ for the caratteristics \mathbf{x} and etichetta y of a punto dati. The Bayes risk is the minimo possible risk that can be achieved by any ipotesi $h : \mathcal{X} \rightarrow \mathcal{Y}$. Any ipotesi that achieves the Bayes risk is referred to as a Bayes estimator [?].

bias Consider an ML method using a parametrized spazio delle ipotesi \mathcal{H} . It learns the model parameters $\mathbf{w} \in \mathbb{R}^d$ using the dataset

$$\mathcal{D} = \{ (\mathbf{x}^{(r)}, y^{(r)}) \}_{r=1}^m.$$

To analyze the properties of the ML method, we typically interpret the punti dati as realizzazioni of independent and identically distributed (i.i.d.) variabili aleatorie,

$$y^{(r)} = h^{(\bar{\mathbf{w}})}(\mathbf{x}^{(r)}) + \epsilon^{(r)}, r = 1, \dots, m.$$

We can then interpret the ML method as an estimator $\hat{\mathbf{w}}$ computed from \mathcal{D} (e.g., by solving ERM). The (squared) bias incurred by the estimate $\hat{\mathbf{w}}$ is then defined as $B^2 := \|\mathbb{E}\{\hat{\mathbf{w}}\} - \bar{\mathbf{w}}\|_2^2$.

See also: ML, spazio delle ipotesi, model parameters, dataset, punto dati, realizzazione, i.i.d., variabile aleatoria, ERM.

bootstrap For the analysis of ML methods, it is often useful to interpret a given set of punto datis $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ as realizzazioni of i.i.d. variabile aleatorias with a common distribuzione di probabilità $p(\mathbf{z})$. In general, we do not know $p(\mathbf{z})$ exactly, but we need to estimate it. The bootstrap uses the histogram of \mathcal{D} as an estimator for the underlying distribuzione di probabilità $p(\mathbf{z})$.

caratteristica Una caratteristica di un punto dati è una delle sue proprietà che può essere misurata o calcolata facilmente senza la necessità di supervisione umana. Per esempio, se un punto dati è un'immagine digitale (memorizzata, ad esempio, come file `.jpeg` file), è possibile allora utilizzare le intensità dei canali rosso, verde e blu dei pixel come caratteristiche. Sinonimi del termine feature, specifici del dominio, includono: "covariata", "variabile esplicativa", "variabile indipendente", "variabile di input", "predittore" o "regressore". [?], [?], [?].

central limit theorem (CLT) The CLT refers to mathematically precise statements about the tendency of an average of a large number of independent variabili aleatorie to tend towards a variabile aleatoria Gaussiana.

See also: variabile aleatoria, variabile aleatoria Gaussiana.

classificazione La classificazione è il compito di determinare un'etichetta a valori discreti y per un dato punto dati, basandosi unicamente sulle

sue caratteristiche \mathbf{x} . L'etichetta y appartiene a un insieme finito, ad esempio $y \in \{-1, 1\}$ o $y \in \{1, \dots, 19\}$, e rappresenta la categoria alla quale appartiene il punto dati corrispondente.

Si veda anche: etichetta, punto dati, caratteristica.

classifier A classifier is a ipotesi (map) $h(\mathbf{x})$ used to predict a etichetta taking values from a finite spazio delle etichette. We might use the function value $h(\mathbf{x})$ itself as a previsione \hat{y} for the etichetta. However, it is customary to use a map $h(\cdot)$ that delivers a numeric quantity. The previsione is then obtained by a simple thresholding step. For example, in a binary classificazione problem with $\mathcal{Y} \in \{-1, 1\}$, we might use a real-valued ipotesi map $h(\mathbf{x}) \in \mathbb{R}$ as a classifier. A previsione \hat{y} can then be obtained via thresholding,

$$\hat{y} = 1 \text{ for } h(\mathbf{x}) \geq 0 \text{ and } \hat{y} = -1 \text{ otherwise.} \quad (1)$$

We can characterize a classifier by its decision regions \mathcal{R}_a , for every possible etichetta value $a \in \mathcal{Y}$.

cluster A cluster is a subset of punti dati that are more similar to each other than to the punti dati outside the cluster. The quantitative measure of similarity between punti dati is a design choice. If punti dati are characterized by Euclidean vettori delle caratteristiche $\mathbf{x} \in \mathbb{R}^d$, we can define the similarity between two punti dati via the Euclidean distance between their vettori delle caratteristiche. An example of such clusters is shown in Fig. 2.

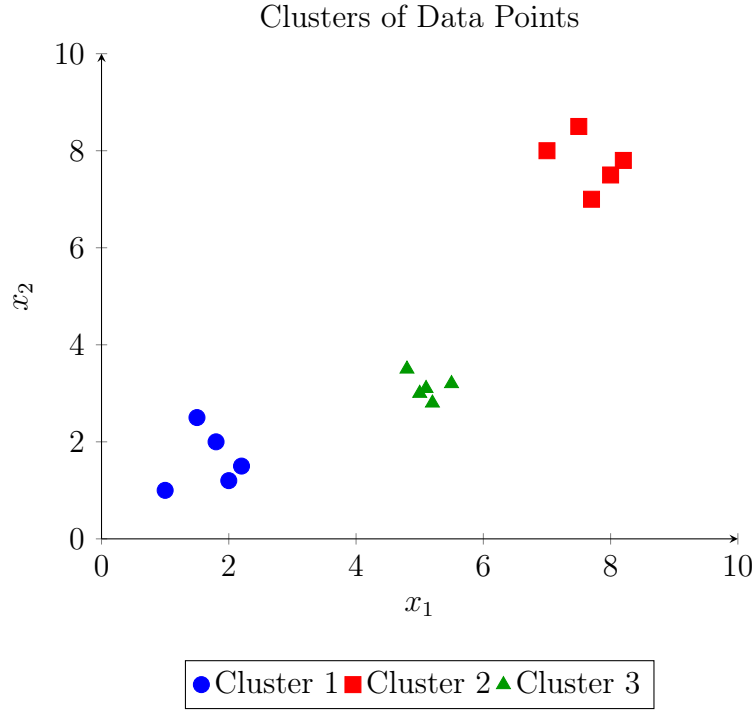


Figura 2: Illustration of three clusters in a two-dimensional feature space. Each cluster groups punti dati that are more similar to each other than to those in other clusters, based on the Euclidean distance.

See also: punto dati, vettore delle caratteristiche, feature space.

clustered federated learning (CFL) Clustered FL assumes that local datasets are naturally grouped into clusters. The local datasets belonging to the same cluster have similar statistical properties. Clustered FL aggregates local datasets in the same cluster to obtain a insieme di addestramento for the training of a cluster-specific modello. Generalized total variation minimization (GTVMin) facilitates this clustering impli-

citly by enforcing approximate similarity of model parameters across well-connected subsets of the FL network.

clustering Clustering methods decompose a given set of punto datis into a few subsets, which are referred to as clusters. Each cluster consists of punto datis that are more similar to each other than to punto datis outside the cluster. Different clustering methods use different measures for the similarity between punto datis and different forms of cluster representations. The clustering method k -means uses the average caratteristica vector (cluster media) of a cluster as its representative. A popular soft clustering method based on Gaussian mixture model (GMM) represents a cluster by a multivariate normal distribution.

clustering assumption The clustering assumption postulates that punto datis in a dataset form a (small) number of groups or clusters. Punto datis in the same cluster are more similar to each other than those outside the cluster [?]. We obtain different clustering methods by using different notions of similarity between punto datis.

computational aspects By computational aspects of an ML method, we mainly refer to the computational resources required for its implementation. For example, if an ML method uses iterative optimization techniques to solve ERM, then its computational aspects include: 1) how many arithmetic operations are needed to implement a single iteration (gradient step); and 2) how many iterations are needed to obtain useful model parameters. One important example of an iterative optimization technique is gradient descent (GD).

condition number The condition number $\kappa(\mathbf{Q}) \geq 1$ of a positive definite matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$ is the ratio α/β between the largest α and the smallest β eigenvalue of \mathbf{Q} . The condition number is useful for the analysis of ML methods. The computational complexity of gradient-based methods for linear regression crucially depends on the condition number of the matrix $\mathbf{Q} = \mathbf{X}\mathbf{X}^T$, with the feature matrix \mathbf{X} of the insieme di addestramento. Thus, from a computational perspective, we prefer caratteristicas of punto datis such that \mathbf{Q} has a condition number close to 1.

confusion matrix Consider punto datis characterized by caratteristicas \mathbf{x} and etichetta y having values from the finite spazio delle etichette $\mathcal{Y} = \{1, \dots, k\}$. The confusion matrix is a $k \times k$ matrix with rows representing different values c of the true label of a punto dati. The columns of a confusion matrix correspond to different values c' delivered by a hypothesis $h(\mathbf{x})$. The (c, c') -th entry of the confusion matrix is the fraction of punto datis with the etichetta $y=c$ and the previsione $\hat{y}=c'$ assigned by the ipotesi h .

connected graph An undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is connected if every non-empty subset $\mathcal{V}' \subset \mathcal{V}$ has at least one edge connecting it to $\mathcal{V} \setminus \mathcal{V}'$.

convex A subset $\mathcal{C} \subseteq \mathbb{R}^d$ of the Euclidean space \mathbb{R}^d is referred to as convex if it contains the line segment between any two points $\mathbf{x}, \mathbf{y} \in \mathcal{C}$ in that set. A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if its epigraph $\{(\mathbf{w}^T, t)^T \in \mathbb{R}^{d+1} : t \geq f(\mathbf{w})\}$ is a convex set [?]. We illustrate one example of a convex set and a convex function in Figure 3.



Figura 3: Left: A convex set $\mathcal{C} \subseteq \mathbb{R}^d$. Right: A convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

convex clustering Consider a dataset $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$. Convex clustering learns vectors $\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(m)}$ by minimizing

$$\sum_{r=1}^m \|\mathbf{x}^{(r)} - \mathbf{w}^{(r)}\|_2^2 + \alpha \sum_{i,i' \in \mathcal{V}} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_p.$$

Here, $\|\mathbf{u}\|_p := (\sum_{j=1}^d |u_j|^p)^{1/p}$ denotes the p -norma (for $p \geq 1$). It turns out that many of the optimal vectors $\hat{\mathbf{w}}^{(1)}, \dots, \hat{\mathbf{w}}^{(m)}$ coincide. A cluster then consists of those punto datis $r \in \{1, \dots, m\}$ with identical $\hat{\mathbf{w}}^{(r)}$ [?, ?].

Courant–Fischer–Weyl min-max characterization Consider a psd matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$ with eigenvalue decomposition (EVD) (or spectral decomposition),

$$\mathbf{Q} = \sum_{j=1}^d \lambda_j \mathbf{u}^{(j)} (\mathbf{u}^{(j)})^T.$$

Here, we use the ordered (in increasing fashion) eigenvalues

$$\lambda_1 \leq \dots \leq \lambda_n.$$

The Courant–Fischer–Weyl min-max characterization [3, Th. 8.1.2] represents the eigenvalues of \mathbf{Q} as the solutions to certain optimization problems.

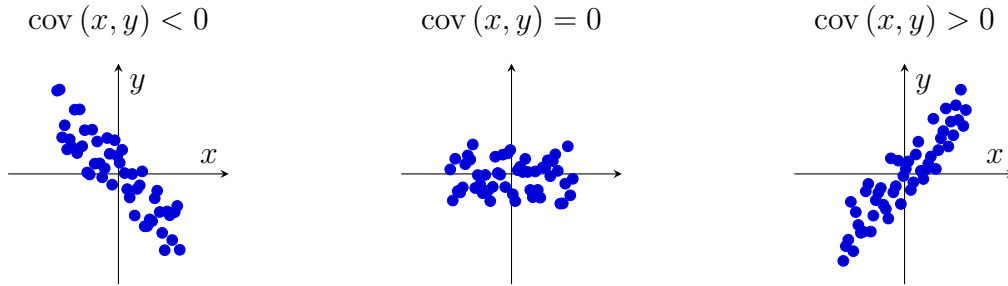


Figura 4: Diagrammi a dispersione che mostrano realizzazioni ottenute da tre diversi modelli probabilistici per due variabili aleatorie con diversi valori di covarianza: negativa (a sinistra), nulla (al centro) e positiva (a destra).

covarianza La covarianza tra due variabili aleatorie a valori reali x e y , definita su uno spazio di probabilità comune, misura la loro dipendenza lineare. È definita come

$$\text{cov}(x, y) = \mathbb{E}\{(x - \mathbb{E}\{x\})(y - \mathbb{E}\{y\})\}.$$

Una covarianza positiva indica che x and y tendono ad aumentare insieme, mentre una covarianza negativa suggerisce che una delle due variabili tende ad aumentare quando l'altra diminuisce. Se $\text{cov}(x, y) = 0$, le variabili aleatorie si dicono non correlate, sebbene non necessariamente statisticamente indipendenti. Si veda la Figura ?? per rappresentazioni grafiche.

data augmentation Data augmentation methods add synthetic punto datis to an existing set of punto datis. These synthetic punto datis are obtained by perturbations (e.g., adding noise to physical measurements) or transformations (e.g., rotations of images) of the original punto datis.

These perturbations and transformations are such that the resulting synthetic punto datis should still have the same etichetta. As a case in point, a rotated cat image is still a cat image even if their vettore delle caratteristiche (obtained by stacking pixel color intensities) are very different (see Figure 4). Dati augmentation can be an efficient form of regularization.

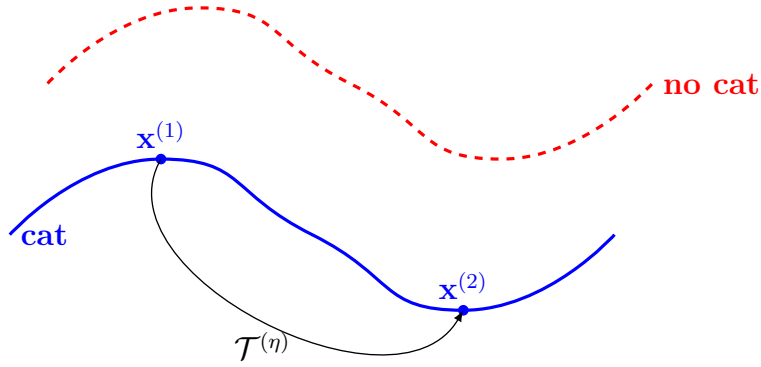


Figure 5: Dati augmentation exploits intrinsic symmetries of punto datis in some feature space \mathcal{X} . We can represent a symmetry by an operator $\mathcal{T}^{(\eta)} : \mathcal{X} \rightarrow \mathcal{X}$, parametrized by some number $\eta \in \mathbb{R}$. For example, $\mathcal{T}^{(\eta)}$ might represent the effect of rotating a cat image by η degrees. A punto dati with vettore delle caratteristiche $\mathbf{x}^{(2)} = \mathcal{T}^{(\eta)}(\mathbf{x}^{(1)})$ must have the same etichetta $y^{(2)} = y^{(1)}$ as a punto dati with vettore delle caratteristiche $\mathbf{x}^{(1)}$.

data minimization principle European dati protection regulation includes a dati minimization principle. This principle requires a dati controller to limit the collection of personal information to what is directly relevant and necessary to accomplish a specified purpose. The dati should be

retained only for as long as necessary to fulfill that purpose [?, Article 5(1)(c)], [?].

data normalization Dati normalization refers to transformations applied to the vettore delle caratteristiche of punto datis to improve the ML method’s statistical aspects or computational aspects. For example, in linear regression with gradient-based methods using a fixed learning rate, convergence depends on controlling the norma of vettore delle caratteristiche in the insieme di addestramento. A common approach is to normalize vettore delle caratteristiche such that their norma does not exceed one [6, Ch. 5].

data poisoning Dati poisoning refers to the intentional manipulation (or fabrication) of punto datis to steer the training of an ML modello [?, ?]. The protection against dati poisoning is particularly important in distributed ML applications where datasets are decentralized.

dataset A dataset refers to a collection of punto datis. These punto datis carry information about some quantity of interest (or etichetta) within a ML application. ML methods use datasets for modello training (e.g., via ERM) and modello validazione. Note that our notion of a dataset is very flexible as it allows for very different types of punto datis. Indeed, punto datis can be concrete physical objects (such as humans or animals) or abstract objects (such as numbers). As a case in point, Figure 5 depicts a dataset that consists of cows as punto datis.

Figura 6: “Cows in the Swiss Alps” by User:Huhu Uet is licensed under [CC BY-SA 4.0](<https://creativecommons.org/licenses/by-sa/4.0/>)

Quite often, an ML engineer does not have direct access to a dataset. Indeed, accessing the dataset in Figure would require to visit the cow herd in the Alps. Instead, we need to use an approximation (or representation) of the dataset which is more convenient to work with. Different mathematical models have been developed for the representation (or approximation) of datasets [?], [?], [?], [?]. One of the most widely adopted data modello is the relational model, which organizes dati as a table (or relation) [?], [?]. A table consists of rows and columns:

- Each row of the table represents a single punto dati.
- Each column of the table corresponds to a specific attribute of the punto dati. ML methods can use attributes as caratteristiche and etichette of the punto dati.

For example, Table 1 shows a representation of the dataset in Figure 5. In the relational modello, the order of rows is irrelevant, and each attribute (i.e., column) must be precisely defined with a domain, which specifies the set of possible values. In ML applications, these attribute domains become the feature space and the spazio delle etichette.

Name	Weight	Age	Height	Stomach temp
Zenzi	100	4	100	25
Berta	140	3	130	23
Resi	120	4	120	31

Tabella 1: A relation (or table) that represents the dataset in Figure .

While the relational model is useful for the study of many ML applications, it may be insufficient regarding the requirements for trustworthy artificial intelligence (trustworthy AI). Modern approaches like datasheets for datasets provide more comprehensive documentation, including details about the dataset’s collection process, intended use, and other contextual information [?].

dati Con dati ci si riferisce ad oggetti che veicolano informazione. Tali oggetti possono essere entità fisiche concrete (come persone o animali), oppure concetti astratti (come i numeri). Spesso si utilizzano rappresentazioni (o approssimazioni) dei dati originari che risultano più convenienti ai fini dell’elaborazione. Queste approssimazioni si basano su diversi modelli di dati, tra cui il modello relazionale rappresenta uno dei più comunemente adottati. [?].

decision boundary Consider a ipotesi map h that reads in a caratteristica vector $\mathbf{x} \in \mathbb{R}^d$ and delivers a value from a finite set \mathcal{Y} . The decision boundary of h is the set of vectors $\mathbf{x} \in \mathbb{R}^d$ that lie between different decision regions. More precisely, a vector \mathbf{x} belongs to the decision

boundary if and only if each neighborhood $\{\mathbf{x}' : \|\mathbf{x} - \mathbf{x}'\| \leq \varepsilon\}$, for any $\varepsilon > 0$, contains at least two vectors with different function values.

decision region Consider a ipotesi map h that delivers values from a finite set \mathcal{Y} . For each etichetta value (category) $a \in \mathcal{Y}$, the ipotesi h determines a subset of caratteristica values $\mathbf{x} \in \mathcal{X}$ that result in the same output $h(\mathbf{x}) = a$. We refer to this subset as a decision region of the ipotesi h .

decision tree A decision tree is a flow-chart-like representation of a ipotesi map h . More formally, a decision tree is a directed graph containing a root node that reads in the vettore delle caratteristiche \mathbf{x} of a punto dati. The root node then forwards the punto dati to one of its children nodes based on some elementary test on the caratteristiche \mathbf{x} . If the receiving child node is not a leaf node, i.e., it has itself children nodes, it represents another test. Based on the test result, the punto dati is forwarded to one of its descendants. This testing and forwarding of the punto dati is continued until the punto dati ends up in a leaf node (having no children nodes).

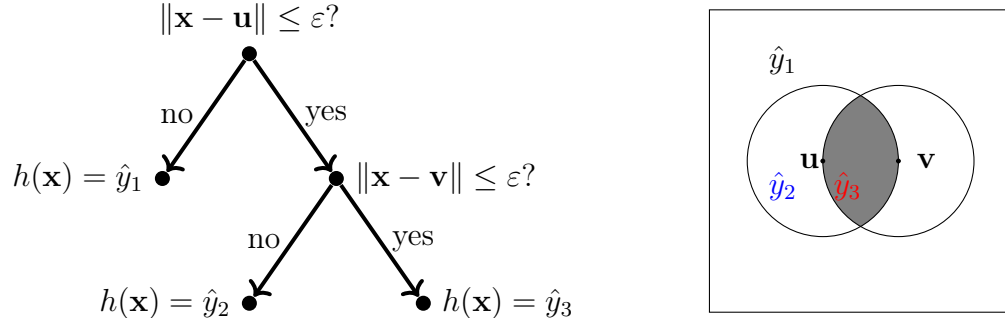


Figura 7: Left: A decision tree is a flow-chart-like representation of a piece-wise constant ipotesi $h : \mathcal{X} \rightarrow \mathbb{R}$. Each piece is a decision region $\mathcal{R}_{\hat{y}} := \{\mathbf{x} \in \mathcal{X} : h(\mathbf{x}) = \hat{y}\}$. The depicted decision tree can be applied to numeric vettore delle caratteristiche, i.e., $\mathcal{X} \subseteq \mathbb{R}^d$. It is parametrized by the threshold $\varepsilon > 0$ and the vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$. Right: A decision tree partitions the feature space \mathcal{X} into decision regions. Each decision region $\mathcal{R}_{\hat{y}} \subseteq \mathcal{X}$ corresponds to a specific leaf node in the decision tree.

deep net A deep net is an ANN with a (relatively) large number of hidden layers. Deep learning is an umbrella term for ML methods that use a deep net as their modello [?].

degree of belonging Degree of belonging is a number that indicates the extent to which a punto dati belongs to a cluster [6, Ch. 8]. The degree of belonging can be interpreted as a soft cluster assignment. Soft clustering methods can encode the degree of belonging by a real number in the interval $[0, 1]$. Hard clustering is obtained as the extreme case when the degree of belonging only takes on values 0 or 1.

denial-of-service attack A denial-of-service attack aims (e.g., via data poisoning) to steer the training of a modello such that it performs

poorly for typical punto datis.

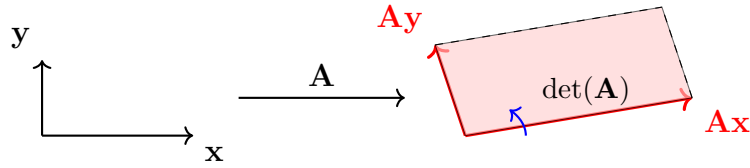
density-based spatial clustering of applications with noise (DBSCAN)

DBSCAN refers to a clustering algorithm for punto datis that are characterized by numeric vettore delle caratteristiche. Like k -means and soft clustering via GMM, also DBSCAN uses the Euclidean distances between vettore delle caratteristiche to determine the clusters. However, in contrast to k -means and GMM, DBSCAN uses a different notion of similarity between punto datis. DBSCAN considers two punto datis as similar if they are connected via a sequence (path) of close-by intermediate punto datis. Thus, DBSCAN might consider two punto datis as similar (and therefore belonging to the same cluster) even if their vettore delle caratteristiche have a large Euclidean distance.

determinant The determinant $\det(\mathbf{A})$ of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a scalar that characterizes how (the orientation of) volumes in \mathbb{R}^n are altered by applying \mathbf{A} [?, 3]. [Note that a matrix \mathbf{A} represents a linear transformation on \mathbb{R}^n .] In particular, $\det(\mathbf{A}) > 0$ preserves orientation, $\det(\mathbf{A}) < 0$ reverses orientation, and $\det(\mathbf{A}) = 0$ collapses volume entirely, indicating that \mathbf{A} is non-invertible. The determinant also satisfies $\det(\mathbf{AB}) = \det(\mathbf{A}) \cdot \det(\mathbf{B})$, and if \mathbf{A} is diagonalizable with eigenvalues $\lambda_1, \dots, \lambda_n$, then $\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i$ [?]. For the special cases $n = 2$ (2D) and $n = 3$ (3D), the determinant can be interpreted as an oriented area or volume spanned by the column vectors of \mathbf{A} .

See also: eigenvalue, inverse matrix.

device Any physical system that can be used to store and process dati. In



the context of ML, we typically mean a computer that is able to read in punti dati from different sources and, in turn, to train an ML modello using these punti dati.

diagramma a dispersione Una tecnica di visualizzazione che rappresenta i punti dati su un piano bidimensionale utilizzando degli indicatori. La Figura 22 mostra un esempio di diagramma a dispersione.

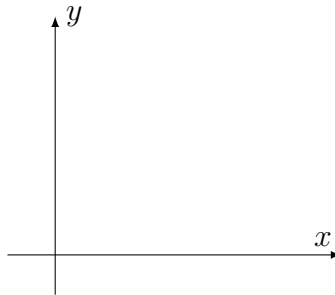


Figura 8: Un diagramma a dispersione con indicatori circolari, in cui i punti dati rappresentano le condizioni meteorologiche giornaliere in Finlandia. Ogni punto dati è caratterizzato dalla sua temperatura minimo diurna x come caratteristica e dalla sua temperatura massimo diurna y come etichetta. Le temperature sono state misurate presso la stazione meteorologica FMI di Helsinki Kaisaniemi, nel periodo compreso tra il 1.9.2024 e il 28.10.2024.

Un diagramma a dispersione può rendere possibile l'ispezione visiva dei punti dati, i quali sono, per loro natura, rappresentati dai vettori delle

caratteristiche in spazi ad alta dimensionalità.

Si veda anche: punto dati, minimo, caratteristica, massimo, etichetta, FMI, vettore delle caratteristiche, riduzione della dimensionalità.

differential entropy For a real-valued variabile aleatoria $\mathbf{x} \in \mathbb{R}^d$ with probability density function (pdf) $p(x)$, the differential entropy is defined as [?]

$$h(\mathbf{x}) := - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}.$$

Differential entropy can be negative and lacks some properties of entropy for discrete-valued variabili aleatorie, such as invariance under change of variables [?]. Among all variabili aleatorie with given media $\boldsymbol{\mu}$ and matrice di covarianza $\boldsymbol{\Sigma}$, $h(\mathbf{x})$ is maximized by $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

See also: uncertainty, modello probabilistico.

differential privacy (DP) Consider some ML method \mathcal{A} that reads in a dataset (e.g., the insieme di addestramento used for ERM) and delivers some output $\mathcal{A}(\mathcal{D})$. The output could be either the learned model parameters or the previsiones for specific punto datis. DP is a precise measure of privacy leakage incurred by revealing the output. Roughly speaking, an ML method is differentially private if the distribuzione di probabilità of the output $\mathcal{A}(\mathcal{D})$ does not change too much if the sensitive attribute of one punto dati in the insieme di addestramento is changed. Note that DP builds on a modello probabilistico for an ML method, i.e., we interpret its output $\mathcal{A}(\mathcal{D})$ as the realizzazione of an variabile aleatoria. The randomness in the output can be ensured by

intentionally adding the realizzazione of an auxiliary variabile aleatoria (noise) to the output of the ML method.

differenziabile Una funzione a valori reali $f : \mathbb{R}^d \rightarrow \mathbb{R}$ si dice differenziabile se, ad ogni punto, può essere approssimata localmente da una funzione lineare. L'approssimazione lineare locale nel punto \mathbf{x} è determinata dal gradiente $\nabla f(\mathbf{x})$ [2].

Si veda anche: gradiente.

dimensione effettiva La dimensione effettiva $d_{\text{eff}}(\mathcal{H})$ di uno spazio delle ipotesi infinito \mathcal{H} è una misura della sua grandezza. In termini approssimativi, la dimensione effettiva corrisponde al numero effettivo di model parameters indipendenti e regolabili. Tali parametri possono essere, ad esempio, i coefficienti impiegati in una linear map oppure i weights e i termini di bias di una ANN.

Si veda anche: spazio delle ipotesi, model parameters, ANN.

discrepanza Si consideri un'applicazione di FL con networked data rappresentati da un FL network. I metodi di FL impiegano una misura di discrepanza per confrontare le funzioni ipotesi dei local model associati ai nodi i, i' connessi da un arco nel FL network.

distributed algorithm A distributed algorithm is an algorithm designed for a special type of computer: a collection of interconnected computing devices (or nodes). These devices communicate and coordinate their local computations by exchanging messages over a network [?, ?]. Unlike a classical algorithm, which is implemented on a single device, a distributed algorithm is executed concurrently on multiple devices

with computational capabilities. Similar to a classical algorithm, a distributed algorithm can be modeled as a set of potential executions. However, each execution in the distributed setting involves both local computations and message-passing events. A generic execution might look as follows:

$$\begin{aligned}
&\text{Node 1: } \text{input}_1, s_1^{(1)}, s_2^{(1)}, \dots, s_{T_1}^{(1)}, \text{output}_1; \\
&\text{Node 2: } \text{input}_2, s_1^{(2)}, s_2^{(2)}, \dots, s_{T_2}^{(2)}, \text{output}_2; \\
&\quad \vdots \\
&\text{Node N: } \text{input}_N, s_1^{(N)}, s_2^{(N)}, \dots, s_{T_N}^{(N)}, \text{output}_N.
\end{aligned}$$

Each device i starts from its own local input and performs a sequence of intermediate computations $s_k^{(i)}$ at discrete time instants $k = 1, \dots, T_i$. These computations may depend on both: the previous local computations at the device and messages received from other devices. One important application of distributed algorithms is in FL where a network of devices collaboratively train a personal modello for each device.

distribuzione di probabilità Per analizzare i metodi di ML, può essere utile interpretare i punti dati come realizzazioni i.i.d. di una variabile aleatoria. Le proprietà tipiche di tali punti dati sono quindi governate dalla distribuzione di probability di questa variabile aleatoria. La distribuzione di probability di una variabile aleatoria binaria $y \in \{0, 1\}$ è determinata univocamente dalle probabilità $\mathbb{P}(y = 0)$ e $\mathbb{P}(y = 1) = 1 - \mathbb{P}(y = 0)$. La distribuzione di probability di una variabile aleatoria a valori reali $x \in \mathbb{R}$ può essere definita tramite una pdf $p(x)$ tale che $\mathbb{P}(x \in [a, b]) \approx p(a)|b - a|$. In generale, una distribuzione di probability è definita da una misura di [?], [?].

Si veda anche: ML, punto dati, i.i.d., realizzazione, variabile aleatoria, probability, pdf.

edge weight Each edge $\{i, i'\}$ of an FL network is assigned a non-negative edge weight $A_{i,i'} \geq 0$. A zero edge weight $A_{i,i'} = 0$ indicates the absence of an edge between nodes $i, i' \in \mathcal{V}$.

eigenvalue We refer to a number $\lambda \in \mathbb{R}$ as an eigenvalue of a square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ if there is a non-zero vector $\mathbf{x} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ such that $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$.

eigenvalue decomposition (EVD) The eigenvalue decomposition for a square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is a factorization of the form

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}.$$

The columns of the matrix $\mathbf{V} = (\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(d)})$ are the eigenvectors of the matrix \mathbf{V} . The diagonal matrix $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_d\}$ contains the eigenvalues λ_j corresponding to the eigenvectors $\mathbf{v}^{(j)}$. Note that the above decomposition exists only if the matrix \mathbf{A} is diagonalizable.

eigenvector An eigenvector of a matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is a non-zero vector $\mathbf{x} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ such that $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ with some eigenvalue λ .

empirical risk The empirical risk $\hat{L}(h|\mathcal{D})$ of a ipotesi on a dataset \mathcal{D} is the average perdita incurred by h when applied to the punto dats in \mathcal{D} .

empirical risk minimization (ERM) Empirical risk minimization is the optimization problem of finding a ipotesi (out of a modello) with the minimo average perdita (or empirical risk) on a given dataset \mathcal{D} (i.e.,

the insieme di addestramento). Many ML methods are obtained from empirical risk via specific design choices for the dataset, modello, and perdita [6, Ch. 3].

entropy Entropy quantifies the uncertainty or unpredictability associated with a variabile aleatoria [?]. For a discrete variabile aleatoria x taking values in a finite set $\mathcal{S} = \{x_1, \dots, x_n\}$ with probability mass function $p_i := \mathbb{P}(x = x_i)$, the entropy is defined as

$$H(x) := - \sum_{i=1}^n p_i \log p_i.$$

Entropy is maximized when all outcomes are equally likely, and minimized (i.e., zero) when the outcome is deterministic. A generalization of the concept of entropy for continuous variabili aleatorie is differential entropy.

See also: uncertainty, modello probabilistico.

epigraph The epigraph of a real-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is the set of points lying on or above its graph:

$$\text{epi}(f) = \{(\mathbf{x}, t) \in \mathbb{R}^n \times \mathbb{R} \mid f(\mathbf{x}) \leq t\}.$$

A function is convex if and only if its epigraph is a convex set [?], [?].

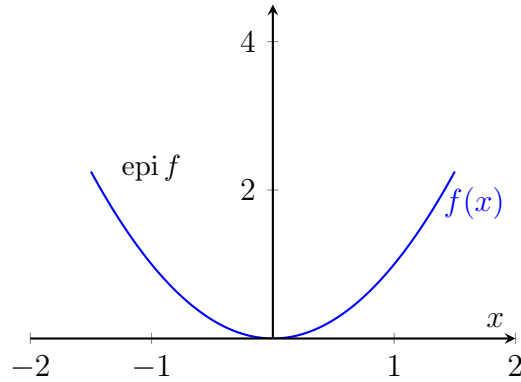


Figura 9: Epigraph of the function $f(x) = x^2$ (i.e., shaded area).

See also: function, graph, convex.

Erdős-Rényi graph (ER graph) An ER graph is a modello probabilistico for graphs defined over a given node set $i = 1, \dots, n$. One way to define the ER graph is via the collection of i.i.d. binary variabili aleatorie $b(\{i, i'\}) \in \{0, 1\}$, for each pair of different nodes i, i' . A specific realizzazione of an ER graph contains an edge $\{i, i'\}$ if and only if $b(\{i, i'\}) = 1$. The ER graph is parametrized by the number n of nodes and the probability $\mathbb{P}(b(\{i, i'\}) = 1)$.

See also: graph, modello probabilistico, i.i.d., variabile aleatoria, realizzazione, probability.

errore di addestramento La perdita media di una ipotesi nel predire le etichette di un punti dati in un insieme di addestramento. Talvolta, con il termine errore di addestramento ci si riferisce anche alla minima perdita media ottenuta da una soluzione del problema di ERM.

Si veda anche: perdita, ipotesi, etichetta, punto dati, insieme di addestramento, ERM.

errore di validazione Si consideri una ipotesi \hat{h} ottenuta mediante un qualche metodo di ML, ad esempio utilizzando la ERM su un insieme di addestramento. La perdita media di \hat{h} su un insieme di validazione, distinto dal insieme di addestramento, è detta errore di validazione. Si veda anche: ipotesi, ML, ERM, insieme di addestramento, perdita, insieme di validazione, validazione.

estimation error Consider punto datis, each with vettore delle caratteristiche \mathbf{x} and etichetta y . In some applications, we can model the relation between the vettore delle caratteristiche and the etichetta of a punto dati as $y = \bar{h}(\mathbf{x}) + \varepsilon$. Here, we use some true underlying ipotesi \bar{h} and a noise term ε which summarizes any modeling or labeling errors. The estimation error incurred by an ML method that learns a ipotesi \hat{h} , e.g., using ERM, is defined as $\hat{h}(\mathbf{x}) - \bar{h}(\mathbf{x})$, for some vettore delle caratteristiche. For a parametric spazio delle ipotesi, which consists of ipotesi maps determined by model parameters \mathbf{w} , we can define the estimation error as $\Delta\mathbf{w} = \hat{\mathbf{w}} - \bar{\mathbf{w}}$ [?, ?].

estremo superiore (o minimo dei maggioranti) L'estremo superiore di un insieme di numeri reali è il più piccolo tra i numeri che sono maggiori o uguali ad ogni elemento dell'insieme. In termini più rigorosi, un numero reale a è l'estremo superiore di un insieme $\mathcal{A} \subseteq \mathbb{R}$ se: 1) a è un maggiorante di \mathcal{A} ; e 2) nessun numero strettamente minore di a è un maggiorante di \mathcal{A} . Ogni insieme non vuoto di numeri reali che

sia superiormente limitato ammette un supremo, anche se questo non appartiene necessariamente all'insieme. [2, Sec. 1.4].

etichetta Un'informazione o una quantità di livello superiore associata a un punto dati. Ad esempio, se il punto dati è un'immagine, l'etichetta potrebbe indicare se l'immagine contiene o meno un gatto. Sinonimi di etichetta, comunemente utilizzati in ambiti specifici, includono variabile di risposta, variabile di output, label e target. [?], [?], [?].

Euclidean space The Euclidean space \mathbb{R}^d of dimension $d \in \mathbb{N}$ consists of vectors $\mathbf{x} = (x_1, \dots, x_d)$, with d real-valued entries $x_1, \dots, x_d \in \mathbb{R}$. Such an Euclidean space is equipped with a geometric structure defined by the inner product $\mathbf{x}^T \mathbf{x}' = \sum_{j=1}^d x_j x'_j$ between any two vectors $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ [2].

expectation-maximization (EM) Consider a modello probabilistico $\mathbb{P}(\mathbf{z}; \mathbf{w})$ for the punto datis \mathcal{D} generated in some ML application. The maximum likelihood estimator for the model parameters \mathbf{w} is obtained by maximizing $\mathbb{P}(\mathcal{D}; \mathbf{w})$. However, the resulting optimization problem might be computationally challenging. Valore atteso-maximization approximates the maximum likelihood estimator by introducing a latent variabile aleatoria \mathbf{z} such that maximizing $\mathbb{P}(\mathcal{D}, \mathbf{z}; \mathbf{w})$ would be easier [?, ?, ?]. Since we do not observe \mathbf{z} , we need to estimate it from the observed dataset \mathcal{D} using a conditional valore atteso. The resulting estimate $\hat{\mathbf{z}}$ is then used to compute a new estimate $\hat{\mathbf{w}}$ by solving $\max_{\mathbf{w}} \mathbb{P}(\mathcal{D}, \hat{\mathbf{z}}; \mathbf{w})$. The crux is that the conditional valore atteso $\hat{\mathbf{z}}$ depends on the model parameters $\hat{\mathbf{w}}$, which we have updated based on $\hat{\mathbf{z}}$. Thus, we have to re-calculate $\hat{\mathbf{z}}$, which, in turn, results in a new choice $\hat{\mathbf{w}}$ for the model

parameters. In practice, we repeat the computation of the conditional valore atteso (i.e., the E-step) and the update of the model parameters (i.e., the M-step) until some stopping criterion is met.

expert ML aims to learn a ipotesi h that accurately predicts the etichetta of a punto dati based on its caratteristicas. We measure the previsionone error using some loss function. Ideally, we want to find a ipotesi that incurs minimal perdita on any punto dati. We can make this informal goal precise via the independent and identically distributed assumption (i.i.d. assumption) and by using the Bayes risk as the baseline for the (average) perdita of a ipotesi. An alternative approach to obtaining a baseline is to use the ipotesi h' learned by an existing ML method. We refer to this ipotesi h' as an expert [?]. Regret minimization methods learn a ipotesi that incurs a perdita comparable to the best expert [?, ?].

explainability We define the (subjective) explainability of an ML method as the level of simulatability [?] of the previsionone delivered by an ML system to a human user. Quantitative measures for the (subjective) explainability of a trained modello can be constructed by comparing its previsionone with the previsionone provided by a user on a test set [?, ?]. Alternatively, we can use modello probabilisticos for dati and measure the explainability of a trained ML modello via the conditional (differential) entropy of its previsionone, given the user previsionone [?, ?].

explainable empirical risk minimization (EERM) Explainable ERM is an instance of SRM that adds a regularization term to the average perdita in the objective function of ERM. The regularization term

is chosen to favor ipotesi maps that are intrinsically explainable for a specific user. This user is characterized by their previsioni provided for the punto dati in a insieme di addestramento [?].

explainable machine learning (explainable ML) Explainable ML methods aim at complementing each previsione with an explanation of how the previsione has been obtained. The construction of an explicit explanation might not be necessary if the ML method uses a sufficiently simple (or interpretable) modello [?].

explanation One approach to make ML methods transparent is to provide an explanation along with the previsione delivered by an ML method. Explanations can take on many different forms. An explanation could be some natural text or some quantitative measure for the importance of individual caratteristiche of a punto dati [?]. We can also use visual forms of explanations, such as intensity plots for image classificazione [?].

feature learning Consider an ML application with punto dati characterized by raw caratteristiche $\mathbf{x} \in \mathcal{X}$. Caratteristica learning refers to the task of learning a map

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}' : \mathbf{x} \mapsto \mathbf{x}',$$

that reads in raw caratteristiche $\mathbf{x} \in \mathcal{X}$ of a punto dati and delivers new caratteristiche $\mathbf{x}' \in \mathcal{X}'$ from a new feature space \mathcal{X}' . Different caratteristica learning methods are obtained for different design choices of $\mathcal{X}, \mathcal{X}'$, for a spazio delle ipotesi \mathcal{H} of potential maps Φ , and for a quantitative measure of the usefulness of a specific $\Phi \in \mathcal{H}$. For example,

principal component analysis (PCA) uses $\mathcal{X} := \mathbb{R}^d$, $\mathcal{X}' := \mathbb{R}^{d'}$ with $d' < d$, and a spazio delle ipotesi

$$\mathcal{H} := \{ \Phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'} : \mathbf{x}' := \mathbf{F}\mathbf{x} \text{ with some } \mathbf{F} \in \mathbb{R}^{d' \times d} \}.$$

PCA measures the usefulness of a specific map $\Phi(\mathbf{x}) = \mathbf{F}\mathbf{x}$ by the minimo linear reconstruction error incurred on a dataset,

$$\min_{\mathbf{G} \in \mathbb{R}^{d' \times d}} \sum_{r=1}^m \left\| \mathbf{G}\mathbf{F}\mathbf{x}^{(r)} - \mathbf{x}^{(r)} \right\|_2^2.$$

feature map Caratteristica map refers to a map that transforms the original caratteristiche of a punto dati into new caratteristiche. The so-obtained new caratteristiche might be preferable over the original caratteristiche for several reasons. For example, the arrangement of punto datis might become simpler (or more linear) in the new feature space, allowing the use of linear models in the new caratteristiche. This idea is a main driver for the development of kernel methods [?]. Moreover, the hidden layers of a deep net can be interpreted as a trainable caratteristica map followed by a linear model in the form of the output layer. Another reason for learning a caratteristica map could be that learning a small number of new caratteristiche helps to avoid overfitting and ensures interpretability [?]. The special case of a caratteristica map delivering two numeric caratteristiche is particularly useful for dati visualization. Indeed, we can depict punto datis in a diagramma a dispersione by using two caratteristiche as the coordinates of a punto dati.

feature matrix Consider a dataset \mathcal{D} with m punto datis with vettore delle caratteristiche $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$. It is convenient to collect

the individual vettore delle caratteristiche into a caratteristica matrix $\mathbf{X} := (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})^T$ of size $m \times d$.

feature space The caratteristica space of a given ML application or method is constituted by all potential values that the vettore delle caratteristiche of a punto dati can take on. A widely used choice for the caratteristica space is the Euclidean space \mathbb{R}^d , with the dimension d being the number of individual characteristics of a punto dati.

federated averaging (FedAvg) FedAvg refers to an iterative FL algorithm that alternates between separately training local models and combining the updated local model parameters. The training of local models is implemented via several SGD steps [?].

federated learning (FL) FL is an umbrella term for ML methods that train modelli in a collaborative fashion using decentralized dati and computation.

federated learning network (FL network) A federated network is an undirected weighted graph whose nodes represent dati generators that aim to train a local (or personalized) modello. Each node in a federated network represents some device capable of collecting a local dataset and, in turn, train a local model. FL methods learn a local ipotesi $h^{(i)}$, for each node $i \in \mathcal{V}$, such that it incurs small perdita on the local datasets.

federated learning orizzontale (FL orizzontale) FL Orizzontale utilizza local dataset costituiti da punto dati differenti, ma impiega lo stesso insieme di caratteristica per descriverli [?]. Ad esempio, nella

previsione meteorologica si utilizza una rete di stazioni meteorologiche distribuite geograficamente, ciascuna delle quali misura le medesime variabili, come la temperatura giornaliera, la pressione atmosferica e le precipitazioni, ma in regioni spazio-temporali differenti. Ogni regione rappresenta un punto dati distinto, descritto attraverso lo stesso insieme di caratteristica. (ad esempio, la temperatura giornaliera o la pressione dell'aria).

FedProx FedProx refers to an iterative FL algorithm that alternates between separately training local models and combining the updated local model parameters. In contrast to FedAvg, which uses SGD to train local models, FedProx uses a proximal operator for the training [?].

Finnish Meteorological Institute (FMI) The FMI is a government agency responsible for gathering and reporting weather data in Finland.

fixed-point iteration A fixed-point iteration is an iterative method for solving a given optimization problem. It constructs a sequence $\mathbf{w}^{(0)}, \mathbf{w}^{(1)}, \dots$ by repeatedly applying an operator \mathcal{F} , i.e.,

$$\mathbf{w}^{(k+1)} = \mathcal{F}\mathbf{w}^{(k)}, \text{ for } k = 0, 1, \dots \quad (2)$$

The operator \mathcal{F} is chosen such that any of its fixed points is a solution $\hat{\mathbf{w}}$ to the given optimization problem. For example, given a differentiable and convex function $f(\mathbf{w})$, the fixed points of the operator $\mathcal{F} : \mathbf{w} \mapsto \mathbf{w} - \nabla f(\mathbf{w})$ coincide with the minimizers of $f(\mathbf{w})$. In general, for a given optimization problem with solution $\hat{\mathbf{w}}$, there are many different operators \mathcal{F} whose fixed points are $\hat{\mathbf{w}}$. Clearly, we should use an operator

\mathcal{F} in (??) that reduces the distance to a solution such that

$$\underbrace{\|\mathbf{w}^{(k+1)} - \widehat{\mathbf{w}}\|_2}_{\stackrel{(\text{??})}{=} \|\mathcal{F}\mathbf{w}^{(k)} - \widehat{\mathbf{w}}\|_2} \leq \|\mathbf{w}^{(k)} - \widehat{\mathbf{w}}\|_2.$$

Thus, we require \mathcal{F} to be at least non-expansive, i.e., the iteration (??) should not result in worse model parameters that have a larger distance to a solution $\widehat{\mathbf{w}}$. What is more, each iteration (??) should also make some progress, i.e., reduce the distance to a solution $\widehat{\mathbf{w}}$. This requirement can be made precise using the notion of a ?? [?], [?]. The operator \mathcal{F} is a ?? if, for some $\kappa \in [0, 1)$,

$$\|\mathcal{F}\mathbf{w} - \mathcal{F}\mathbf{w}'\|_2 \leq \kappa \|\mathbf{w} - \mathbf{w}'\|_2 \text{ holds for any } \mathbf{w}, \mathbf{w}'.$$

For a ?? \mathcal{F} , the fixed-point iteration (??) generates a sequence $\mathbf{w}^{(k)}$ that converges quite rapidly. In particular [2, Th. 9.23],

$$\|\mathbf{w}^{(k)} - \widehat{\mathbf{w}}\|_2 \leq \kappa^k \|\mathbf{w}^{(0)} - \widehat{\mathbf{w}}\|_2.$$

Here, $\|\mathbf{w}^{(0)} - \widehat{\mathbf{w}}\|_2$ is the distance between the initialization $\mathbf{w}^{(0)}$ and the solution $\widehat{\mathbf{w}}$. It turns out that a fixed-point iteration (??) with a firmly non-expansive operator \mathcal{F} is guaranteed to converge to a fixed-point of \mathcal{F} [?, Cor. 5.16]. Figure ?? depicts examples of a firmly non-expansive operator, a non-expansive operator, and a ??. All these operators are defined on the one-dimensional space \mathbb{R} . Another example of a firmly non-expansive operator is the proximal operator of a convex function [?], [?].

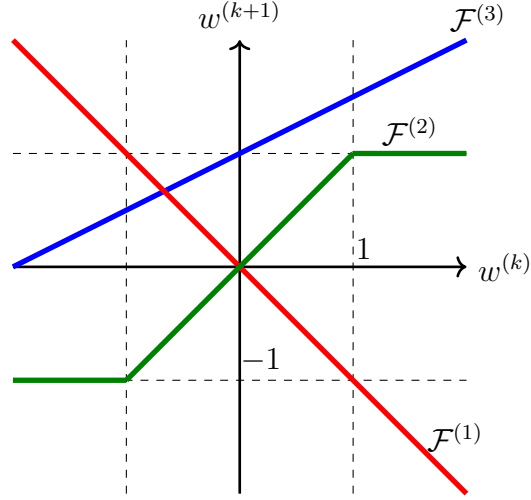


Figura 10: Example of a non-expansive operator $\mathcal{F}^{(1)}$, a firmly non-expansive operator $\mathcal{F}^{(2)}$, and a ?? $\mathcal{F}^{(3)}$.

See also: optimization problem, differenziabile, convex function, model parameters, ??, proximal operator.

flow-based clustering Flow-based clustering groups the nodes of an undirected graph by applying k -means clustering to node-wise vettore delle caratteristiche. These vettore delle caratteristiche are built from network flows between carefully selected sources and destination nodes [?].

function A function is a mathematical rule that assigns each element $u \in \mathcal{U}$ exactly one element $v \in \mathcal{V}$ [2]. We write this as $f : \mathcal{U} \rightarrow \mathcal{V}$, where \mathcal{U} is the domain and \mathcal{V} the co-domain of f . That is, a function f defines a unique output $f(u) \in \mathcal{V}$ for every input $u \in \mathcal{U}$.

Gaussian mixture model (GMM) A GMM is a particular type of modello probabilistico for a numeric vector \mathbf{x} (e.g., the caratteristiche of a punto dati). Within a GMM, the vector \mathbf{x} is drawn from a randomly selected multivariate normal distribution $p^{(c)} = \mathcal{N}(\boldsymbol{\mu}^{(c)}, \mathbf{C}^{(c)})$ with $c = I$. The index $I \in \{1, \dots, k\}$ is an variabile aleatoria with probabilities $\mathbb{P}(I = c) = p_c$. Note that a GMM is parametrized by the probability p_c , the media vector $\boldsymbol{\mu}^{(c)}$, and the matrice di covarianza $\boldsymbol{\Sigma}^{(c)}$ for each $c = 1, \dots, k$. GMMs are widely used for clustering, density estimation, and as a generative modello.

Gaussian process (GP) A GP is a collection of variabili aleatorie $\{f(\mathbf{x})\}_{\mathbf{x} \in \mathcal{X}}$ indexed by input values \mathbf{x} from some input space \mathcal{X} , such that, for any finite subset $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathcal{X}$, the corresponding variabili aleatorie $f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(m)})$ have a joint multivariate Gaussian distribution:

$$(f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(m)})) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$

For a fixed input space \mathcal{X} , a GP is fully specified (or parametrized) by

- a media function $\mu(\mathbf{x}) = \mathbb{E}\{f(\mathbf{x})\}$
- and a covariance function $K(\mathbf{x}, \mathbf{x}') = \mathbb{E}\{(f(\mathbf{x}) - \mu(\mathbf{x}))(f(\mathbf{x}') - \mu(\mathbf{x}'))\}$.

Example: We can interpret the temperature distribution across Finland (at a specific point in time) as the realizzazione of a GP $f(\mathbf{x})$, where each input $\mathbf{x} = (\text{lat}, \text{lon})$ denotes a geographic location. Temperature observations from FMI weather stations provide samples of $f(\mathbf{x})$ at

specific locations (see Figure ??). A GP allows us to predict the temperature nearby FMI weather stations and to quantify the uncertainty of these predictions.

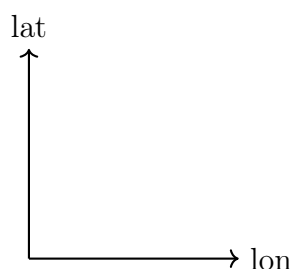


Figura 11: We can interpret the temperature distribution over Finland as a realizzazione of a GP indexed by geographic coordinates and sampled at FMI weather stations (indicated by blue dots).

See also: variabile aleatoria, media, function, realizzazione, FMI, sample, uncertainty.

general data protection regulation (GDPR) The GDPR was enacted by the European Union (EU), effective from May 25, 2018 [?]. It safeguards the privacy and dati rights of individuals in the EU. The GDPR has significant implications for how dati is collected, stored, and used in ML applications. Key provisions include the following:

- Data minimization principle: ML systems should only use the necessary amount of personal dati for their purpose.

- Transparency and explainability: ML systems should enable their users to understand how the systems make decisions that impact the users.
- Data subject rights: Users should get an opportunity to access, rectify, and delete their personal data, as well as to object to automated decision-making and profiling.
- Accountability: Organizations must ensure robust data security and demonstrate compliance through documentation and regular audits.

generalization Many current ML (and AI) systems are based on ERM:

At their core, they train a modello (i.e., learn a ipotesi $\hat{h} \in \mathcal{H}$) by minimizing the average perdita (or empirical risk) on some punto d'isidato $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$, which serve as a insieme di addestramento $\mathcal{D}^{(\text{train})}$. Generalization refers to an ML method's ability to perform well outside the insieme di addestramento. Any mathematical theory of generalization needs some mathematical concept for the "outside the insieme di addestramento." For example, statistical learning theory uses a modello probabilistico such as the i.i.d. assumption for data generation: the punto d'isidato in the insieme di addestramento are i.i.d. realizzazioni of some underlying distribuzione di probabilità $p(\mathbf{z})$. A modello probabilistico allows us to explore the outside of the insieme di addestramento by drawing additional i.i.d. realizzazioni from $p(\mathbf{z})$. Moreover, using the i.i.d. assumption allows us to define the risk of a trained modello $\hat{h} \in \mathcal{H}$ as the expected perdita $\bar{L}(\hat{h})$. What is more, we can use concentration bounds or convergence results for sequences of i.i.d. variabile aleatorias

to bound the deviation between the empirical risk $\hat{L}(\hat{h}|\mathcal{D}^{(\text{train})})$ of a trained modello and its risk [?]. It is possible to study generalization also without using modello probabilisticos. For example, we could use (deterministic) perturbations of the punto datis in the insieme di addestramento to study its outside. In general, we would like the trained modello to be robust, i.e., its previsiones should not change too much for small perturbations of a punto dati. Consider a trained modello for detecting an object in a smartphone snapshot. The detection result should not change if we mask a small number of randomly chosen pixels in the image [?].

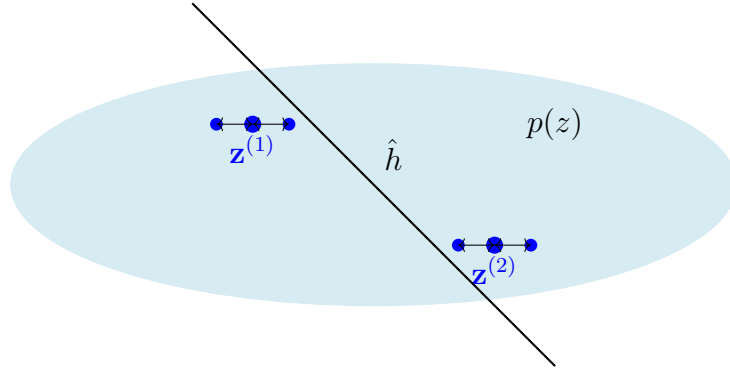


Figura 12: Two punto datis $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}$ that are used as a insieme di addestramento to learn a ipotesi \hat{h} via ERM. We can evaluate \hat{h} outside $\mathcal{D}^{(\text{train})}$ either by an i.i.d. assumption with some underlying distribuzione di probabilità $p(\mathbf{z})$ or by perturbing the punto datis.

generalized total variation (GTV) GTV is a measure of the variation of trained local models $h^{(i)}$ (or their model parameters $\mathbf{w}^{(i)}$) assigned to the nodes $i = 1, \dots, n$ of an undirected weighted graph \mathcal{G} with edges \mathcal{E} . Given a measure $d^{(h,h')}$ for the discrepanza between ipotesi maps h, h' , the GTV is

$$\sum_{\{i,i'\} \in \mathcal{E}} A_{i,i'} d^{(h^{(i)}, h^{(i')})}.$$

Here, $A_{i,i'} > 0$ denotes the weight of the undirected edge $\{i, i'\} \in \mathcal{E}$.

generalized total variation minimization (GTVMin) GTVMin is an instance of regularized empirical risk minimization (RERM) using the GTV of local model parameters as a regularizer [?].

See also: RERM, GTV, regularizer.

gradient descent (GD) Gradiente descent is an iterative method for finding the minimo of a differenziabile function $f(\mathbf{w})$ of a vector-valued argument $\mathbf{w} \in \mathbb{R}^d$. Consider a current guess or approximation $\mathbf{w}^{(k)}$ for the minimo of the function $f(\mathbf{w})$. We would like to find a new (better) vector $\mathbf{w}^{(k+1)}$ that has a smaller objective value $f(\mathbf{w}^{(k+1)}) < f(\mathbf{w}^{(k)})$ than the current guess $\mathbf{w}^{(k)}$. We can achieve this typically by using a gradient step

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \eta \nabla f(\mathbf{w}^{(k)}) \quad (3)$$

with a sufficiently small step size $\eta > 0$. Figure 8 illustrates the effect of a single gradiente descent step (2).

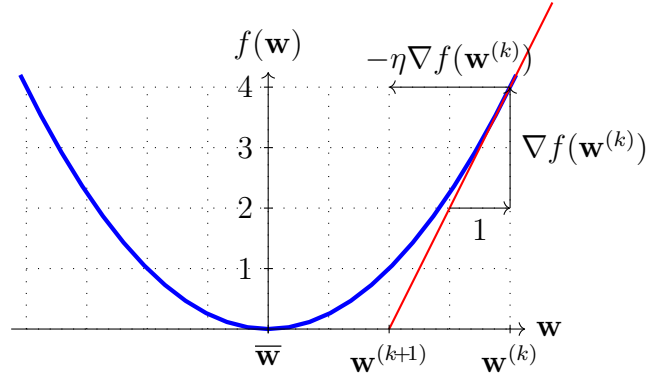


Figura 13: A single gradient step (2) towards the minimizer $\bar{\mathbf{w}}$ of $f(\mathbf{w})$.

gradient step Given a differenziabile real-valued function $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ and a vector $\mathbf{w} \in \mathbb{R}^d$, the gradiente step updates \mathbf{w} by adding the scaled negative gradiente $\nabla f(\mathbf{w})$ to obtain the new vector (see Figure 9)

$$\hat{\mathbf{w}} := \mathbf{w} - \eta \nabla f(\mathbf{w}). \quad (4)$$

Mathematically, the gradiente step is a (typically non-linear) operator $\mathcal{T}^{(f,\eta)}$ that is parametrized by the function f and the step size η .

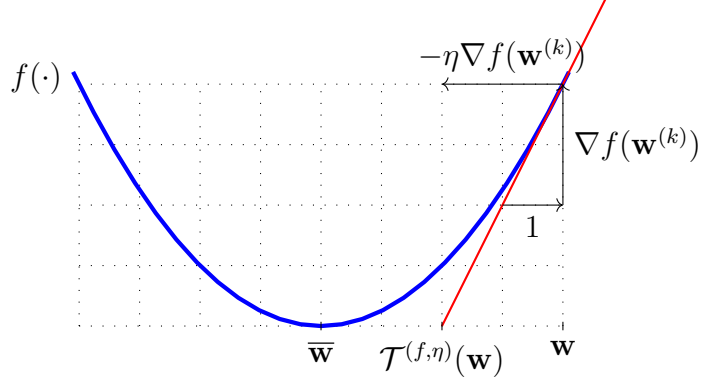


Figure 14: The basic gradient step (3) maps a given vector \mathbf{w} to the updated vector \mathbf{w}' . It defines an operator $\mathcal{T}^{(f,\eta)}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d : \mathbf{w} \mapsto \hat{\mathbf{w}}$.

Note that the gradient step (3) optimizes locally - in a neighborhood whose size is determined by the step size η - a linear approximation to the function $f(\cdot)$. A natural generalization of (3) is to locally optimize the function itself - instead of its linear approximation - such that

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}' \in \mathbb{R}^d} f(\mathbf{w}') + (1/\eta) \|\mathbf{w} - \mathbf{w}'\|_2^2. \quad (5)$$

We intentionally use the same symbol η for the parameter in (4) as we used for the step size in (3). The larger the η we choose in (4), the more progress the update will make towards reducing the function value $f(\hat{\mathbf{w}})$. Note that, much like the gradient step (3), also the update (4) defines a (typically non-linear) operator that is parametrized by the function $f(\cdot)$ and the parameter η . For a convex function $f(\cdot)$, this operator is known as the proximal operator of $f(\cdot)$ [?].

gradient-based methods Gradiente-based methods are iterative techniques for finding the minimo (or massimo) of a differenziabile objective function of the model parameters. These methods construct a sequence of approximations to an optimal choice for model parameters that results in a minimo (or massimo) value of the objective function. As their name indicates, gradiente-based methods use the gradients of the objective function evaluated during previous iterations to construct new, (hopefully) improved model parameters. One important example of a gradiente-based method is GD.

gradiente Per funzione a valori reali $f : \mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{w} \mapsto f(\mathbf{w})$, se esiste un vettore \mathbf{g} tale che $\lim_{\mathbf{w} \rightarrow \mathbf{w}'} \frac{f(\mathbf{w}) - (f(\mathbf{w}') + \mathbf{g}^T(\mathbf{w} - \mathbf{w}'))}{\|\mathbf{w} - \mathbf{w}'\|} = 0$ esso viene detto gradiente di f in \mathbf{w}' . Se esiste, il gradiente è unico e viene denotato con $\nabla f(\mathbf{w}')$ oppure con $\nabla f(\mathbf{w})|_{\mathbf{w}'}$ [2].

graph A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a pair that consists of a node set \mathcal{V} and an edge set \mathcal{E} . In its most general form, a graph is specified by a map that assigns each edge $e \in \mathcal{E}$ a pair of nodes [?]. One important family of graphs is simple undirected graphs. A simple undirected graph is obtained by identifying each edge $e \in \mathcal{E}$ with two different nodes $\{i, i'\}$. Weighted graphs also specify numeric weights A_e for each edge $e \in \mathcal{E}$.

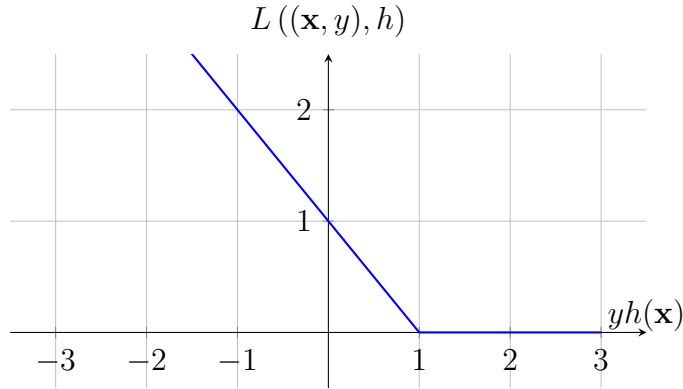
graph clustering Graph clustering aims at clustering punto dats that are represented as the nodes of a graph \mathcal{G} . The edges of \mathcal{G} represent pairwise similarities between punto dats. Sometimes we can quantify the extend of these similarities by an edge weight [?, ?].

hard clustering Hard clustering refers to the task of partitioning a given set of punto datis into (a few) non-overlapping clusters. The most widely used hard clustering method is k -means.

Hilbert space A Hilbert space is a linear vector space equipped with an inner product between pairs of vectors. One important example of a Hilbert space is the Euclidean space \mathbb{R}^d , for some dimension d , which consists of Euclidean vectors $\mathbf{u} = (u_1, \dots, u_d)^T$ along with the inner product $\mathbf{u}^T \mathbf{v}$.

hinge loss Consider a punto dati characterized by a vettore delle caratteristiche $\mathbf{x} \in \mathbb{R}^d$ and a binary etichetta $y \in \{-1, 1\}$. The hinge perdita incurred by a real-valued ipotesi map $h(\mathbf{x})$ is defined as

$$L((\mathbf{x}, y), h) := \max\{0, 1 - yh(\mathbf{x})\}. \quad (6)$$



A regularized variant of the hinge perdita is used by the support vector machine (SVM) [?].

histogram Consider a dataset \mathcal{D} that consists of m punti dati $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$, each of them belonging to some cell $[-U, U] \times \dots \times [-U, U] \subseteq \mathbb{R}^d$ with

side length U . We partition this cell evenly into smaller elementary cells with side length Δ . The histogram of \mathcal{D} assigns each elementary cell to the corresponding fraction of punti dati in \mathcal{D} that fall into this elementary cell. A visual example of such a histogram is provided in Fig. 10.

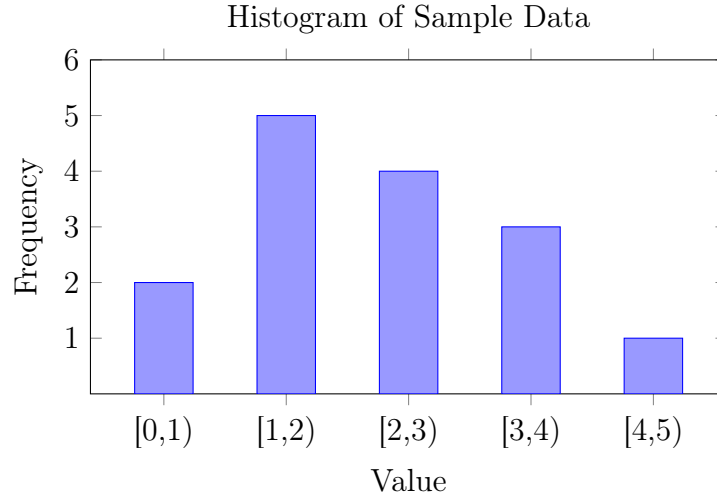


Figura 15: A histogram representing the frequency of punti dati falling within discrete value ranges (i.e., bins). Each bar height shows the count of samples in the corresponding interval.

See also: dataset, punto dati, sample.

Huber loss The Huber perdita unifies the squared error loss and the absolute error loss.

Huber regression Huber regressione refers to ERM-based methods that use the Huber loss as a measure of the previsione error. Two important special cases of Huber regressione are least absolute deviation regression

and linear regression. Tuning the threshold parameter of the Huber loss allows the user to trade the robustness of the absolute error loss against the computational benefits of the smooth squared error loss.

independent and identically distributed (i.i.d.) It can be useful to interpret punto dati $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$ as realizzazioni of i.i.d. variabile aleatorias with a common distribuzione di probabilità. If these variabile aleatorias are continuous-valued, their joint pdf is $p(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = \prod_{r=1}^m p(\mathbf{z}^{(r)})$, with $p(\mathbf{z})$ being the common marginal pdf of the underlying variabile aleatorias.

independent and identically distributed assumption (i.i.d. assumption)

The i.i.d. assumption interprets punti dati of a dataset as the realizzazioni of i.i.d. variabili aleatorie.

See also: i.i.d., punto dati, dataset, realizzazione, variabile aleatoria.

insieme di addestramento Un insieme di addestramento è un dataset \mathcal{D} che consiste di alcuni punti dati utilizzati nell'ambito della ERM per apprendere un'ipotesi \hat{h} . La perdita media di \hat{h} sull'insieme di addestramento è indicata come errore di addestramento. Il confronto tra l'errore di addestramento e l'errore di validazione di \hat{h} consente di diagnosticare il metodo di ML e fornisce indicazioni su come migliorare l'errore di validazione (ad esempio, usando uno spazio delle ipotesi diverso o raccogliendo un numero maggiore di punti dati) [6, Sec. 6.6]. Si veda anche: dataset, punto dati, ERM, ipotesi, perdita, errore di addestramento, errore di validazione, ML, spazio delle ipotesi.

insieme di validazione Un insieme di punti dati utilizzato per stimare il risk di una ipotesi \hat{h} appresa mediante un metodo di ML (ad esempio, risolvendo il problema di ERM). La perdita media di \hat{h} sull'insieme di validazione è indicata come errore di validazione e può essere utilizzata per diagnosticare un metodo di ML (si veda [6, Sec. 6.6]). Il confronto tra errore di addestramento e errore di validazione può fornire indicazioni utili per il miglioramento del metodo di ML (come ad esempio l'impiego di un diverso spazio delle ipotesi).

Si veda anche: punto dati, risk, ipotesi, ML, ERM, perdita, validazione, errore di validazione, errore di addestramento, spazio delle ipotesi.

interpretability An ML method is interpretable for a specific user if they can well anticipate the previsiones delivered by the method. The notion of interpretability can be made precise using quantitative measures of the uncertainty about the previsiones [?].

inverse matrix An inverse matrix \mathbf{A}^{-1} is defined for a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ that is of full rank, meaning its columns are linearly independent. In this case, \mathbf{A} is said to be invertible, and its inverse satisfies

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}.$$

A square matrix is invertible if and only if its determinant is non-zero. Inverse matrices are fundamental in solving systems of linear equations and in the closed-form solution of linear regression [?, ?]. The concept of an inverse matrix can be extended to matrices that are not square or not full rank. One may define a “left inverse” \mathbf{B} satisfying $\mathbf{B}\mathbf{A} = \mathbf{I}$, or a “right inverse” \mathbf{C} satisfying $\mathbf{A}\mathbf{C} = \mathbf{I}$. For general rectangular or

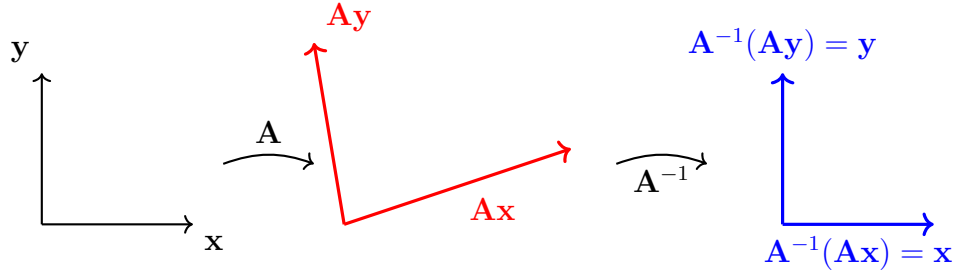


Figura 16: A matrix \mathbf{A} represents a linear transformation of \mathbb{R}^2 . The inverse matrix \mathbf{A}^{-1} represents the inverse transformation.

singular matrices, the Moore–Penrose pseudoinverse \mathbf{A}^+ provides a unified concept of generalized inverse matrix [3].

See also: determinant, linear regression.

ipotesi Una ipotesi è una mappa (o funzione) $h : \mathcal{X} \rightarrow \mathcal{Y}$ dallo feature space \mathcal{X} allo spazio delle etichette \mathcal{Y} . Dato un punto dati con caratteristiche \mathbf{x} , utilizziamo una funzione ipotesi h per stimare (o approssimare) l'etichetta y mediante la previsione $\hat{y} = h(\mathbf{x})$. Il ML si occupa essenzialmente di apprendere (o individuare) una funzione ipotesi h tale che $y \approx h(\mathbf{x})$ per qualsiasi punto dati (avente caratteristiche \mathbf{x} ed etichetta y).

kernel Consider a point characterized by a vector of features $\mathbf{x} \in \mathcal{X}$ with a generic feature space \mathcal{X} . A (real-valued) kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ assigns each pair of vectors of features $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ a real number $K(\mathbf{x}, \mathbf{x}')$. The value $K(\mathbf{x}, \mathbf{x}')$ is often interpreted as a measure for the similarity between \mathbf{x} and \mathbf{x}' . Kernel methods use a kernel to transform the vector of features \mathbf{x} to a new vector

delle caratteristiche $\mathbf{z} = K(\mathbf{x}, \cdot)$. This new vettore delle caratteristiche belongs to a linear feature space \mathcal{X}' which is (in general) different from the original feature space \mathcal{X} . The feature space \mathcal{X}' has a specific mathematical structure, i.e., it is a reproducing kernel Hilbert space [?, ?].

kernel method A kernel method is an ML method that uses a kernel K to map the original (raw) vettore delle caratteristiche \mathbf{x} of a punto dati to a new (transformed) vettore delle caratteristiche $\mathbf{z} = K(\mathbf{x}, \cdot)$ [?, ?]. The motivation for transforming the vettore delle caratteristiche is that, by using a suitable kernel, the punto datis have a "more pleasant" geometry in the transformed feature space. For example, in a binary classificazione problem, using transformed vettore delle caratteristiche \mathbf{z} might allow us to use linear models, even if the punto datis are not linearly separable in the original feature space (see Figure 11).

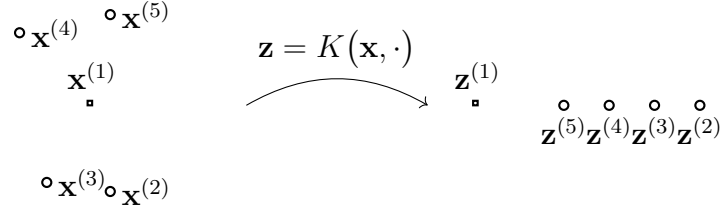


Figura 17: Five punto datis characterized by vettore delle caratteristiche $\mathbf{x}^{(r)}$ and etichettas $y^{(r)} \in \{\circ, \square\}$, for $r = 1, \dots, 5$. With these vettore delle caratteristiche, there is no way to separate the two classes by a straight line (representing the decision boundary of a linear classifier). In contrast, the transformed vettore delle caratteristiche $\mathbf{z}^{(r)} = K(\mathbf{x}^{(r)}, \cdot)$ allow us to separate the punto datis using a linear classifier.

Kullback-Leibler divergence (KL divergence) The KL divergence is a quantitative measure of how much one distribuzione di probabilità is different from another distribuzione di probabilità [?].

labeled datapoint A punto dati whose etichetta is known or has been determined by some means which might require human labor.

Laplacian matrix The structure of a graph \mathcal{G} , with nodes $i = 1, \dots, n$, can be analyzed using the properties of special matrices that are associated with \mathcal{G} . One such matrix is the graph Laplacian matrix $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{n \times n}$, which is defined for an undirected and weighted graph [?, ?]. It is defined

element-wise as (see Figure 12)

$$L_{i,i'}^{(\mathcal{G})} := \begin{cases} -A_{i,i'} & \text{for } i \neq i', \{i, i'\} \in \mathcal{E}, \\ \sum_{i'' \neq i} A_{i,i''} & \text{for } i = i', \\ 0 & \text{else.} \end{cases} \quad (7)$$

Here, $A_{i,i'}$ denotes the edge weight of an edge $\{i, i'\} \in \mathcal{E}$.

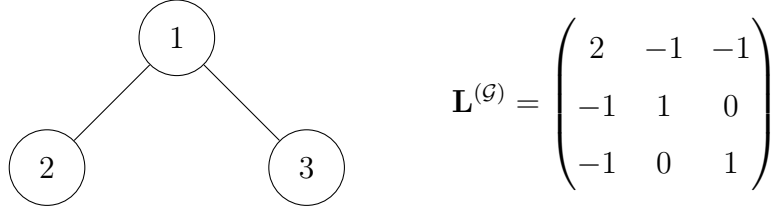


Figura 18: Left: Some undirected graph \mathcal{G} with three nodes $i = 1, 2, 3$. Right: The Laplacian matrix $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{3 \times 3}$ of \mathcal{G} .

large language model (LLM) Large language modellos is an umbrella term for ML methods that process and generate human-like text. These methods typically use deep nets with billions (or even trillions) of ???. A widely used choice for the network architecture is referred to as Transformers [?]. The training of large language modellos is often based on the task of predicting a few words that are intentionally removed from a large text corpus. Thus, we can construct labeled datapoints simply by selecting some words of a text as etichettas and the remaining words as caratteristiche di punto d'is. This construction requires very little human supervision and allows for generating sufficiently large insieme di addestramentos for large language modellos.

law of large numbers The law of large numbers refers to the convergence of the average of an increasing (large) number of i.i.d. variabile aleatorias to the media of their common distribuzione di probabilità. Different instances of the law of large numbers are obtained by using different notions of convergence [?].

learning rate Consider an iterative ML method for finding or learning a useful ipotesi $h \in \mathcal{H}$. Such an iterative method repeats similar computational (update) steps that adjust or modify the current ipotesi to obtain an improved ipotesi. One well-known example of such an iterative learning method is GD and its variants, SGD and projected gradient descent (projected GD). A key parameter of an iterative method is the learning rate. The learning rate controls the extent to which the current ipotesi can be modified during a single iteration. A well-known example of such a parameter is the step size used in GD [6, Ch. 5].

learning task Consider a dataset \mathcal{D} constituted by several punto datis, each of them characterized by caratteristiche \mathbf{x} . For example, the dataset \mathcal{D} might be constituted by the images of a particular database. Sometimes it might be useful to represent a dataset \mathcal{D} , along with the choice of caratteristiche, by a distribuzione di probabilità $p(\mathbf{x})$. A learning task associated with \mathcal{D} consists of a specific choice for the etichetta of a punto dati and the corresponding spazio delle etichette. Given a choice for the loss function and modello, a learning task gives rise to an instance of ERM. Thus, we could define a learning task also via an instance of ERM, i.e., via an objective function. Note that, for the same dataset,

we obtain different learning tasks by using different choices for the characteristics and etichetta of a punto dati. These learning tasks are related, as they are based on the same dataset, and solving them jointly (via multitask learning methods) is typically preferable over solving them separately [?], [?], [?].

least absolute deviation regression Least absolute deviation regression is an instance of ERM using the absolute error loss. It is a special case of Huber regression.

least absolute shrinkage and selection operator (Lasso) The Lasso is an instance of SRM. It learns the weights \mathbf{w} of a linear map $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ based on a insieme di addestramento. Lasso is obtained from linear regression by adding the scaled ℓ_1 -norma $\alpha \|\mathbf{w}\|_1$ to the average squared error loss incurred on the insieme di addestramento.

linear classifier Consider punto datis characterized by numeric caratteristicas $\mathbf{x} \in \mathbb{R}^d$ and a etichetta $y \in \mathcal{Y}$ from some finite spazio delle etichette \mathcal{Y} . A linear classifier is characterized by having decision regions that are separated by hyperplanes in \mathbb{R}^d [6, Ch. 2].

linear map A linear map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a function that satisfies additivity : $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$, and homogeneity : $f(c\mathbf{x}) = cf(\mathbf{x})$ for all vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and scalars $c \in \mathbb{R}$. In particular, $f(\mathbf{0}) = \mathbf{0}$. Any linear map can be represented as a matrix multiplication $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ for some matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$. The collection of real-valued linear maps for a given dimension n constitute a linear model which is used in many ML methods.

See also: linear model, linear regression, PCA, vettore delle caratteristiche.

linear model Consider *punto dati*, each characterized by a numeric vettore delle caratteristiche $\mathbf{x} \in \mathbb{R}^d$. A linear modello is a spazio delle ipotesi which consists of all linear maps,

$$\mathcal{H}^{(d)} := \{h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\}. \quad (8)$$

Note that (7) defines an entire family of spazio delle ipotesi, which is parametrized by the number d of caratteristiche that are linearly combined to form the previsione $h(\mathbf{x})$. The design choice of d is guided by computational aspects (e.g., reducing d means less computation), statistical aspects (e.g., increasing d might reduce previsione error), and interpretability. A linear modello using few carefully chosen caratteristiche tends to be considered more interpretable [?, ?].

linear regression Linear regression aims to learn a linear ipotesi map to predict a numeric etichetta based on the numeric caratteristiche of a *punto dati*. The quality of a linear ipotesi map is measured using the average squared error loss incurred on a set of labeled datapoints, which we refer to as the insieme di addestramento.

local dataset The concept of a local dataset is in between the concept of a *punto dati* and a dataset. A local dataset consists of several individual *punto dati*, which are characterized by caratteristiche and etichettas. In contrast to a single dataset used in basic ML methods, a local dataset is also related to other local datasets via different notions of

similarity. These similarities might arise from modello probabilisticos or communication infrastructure and are encoded in the edges of an FL network.

Local Interpretable Model-agnostic Explanations (LIME) Consider a trained modello (or learnt ipotesi) $\hat{h} \in \mathcal{H}$, which maps the vettore delle caratteristiche of a punto dati to the previsione $\hat{y} = \hat{h}$. Local Interpretable Model-agnostic Explanations (LIME) is a technique for explaining the behaviour of \hat{h} , locally around a punto dati with vettore delle caratteristiche $\mathbf{x}^{(0)}$ [?]. The explanation is given in the form of a local approximation $g \in \mathcal{H}'$ of \hat{h} (see Fig.). This approximation can be obtained by an instance of ERM with carefully designed insieme di addestramento. In particular, the insieme di addestramento consists of punto datis with vettore delle caratteristiche \mathbf{x} close to $\mathbf{x}^{(0)}$ and the (pseudo-)label $\hat{h}(\mathbf{x})$. Note that we can use a different modello \mathcal{H}' for the approximation than the original modello \mathcal{H} . For example, we can use a decision tree to approximate (locally) a deep net. Another widely-used choice for \mathcal{H}' is the linear model.

local model Consider a collection of local datasets that are assigned to the nodes of an FL network. A local modello $\mathcal{H}^{(i)}$ is a spazio delle ipotesi assigned to a node $i \in \mathcal{V}$. Different nodes might be assigned different spazio delle ipotesi, i.e., in general $\mathcal{H}^{(i)} \neq \mathcal{H}^{(i')}$ for different nodes $i, i' \in \mathcal{V}$.

logistic loss Consider a punto dati characterized by the caratteristics \mathbf{x} and a binary etichetta $y \in \{-1, 1\}$. We use a real-valued ipotesi h to

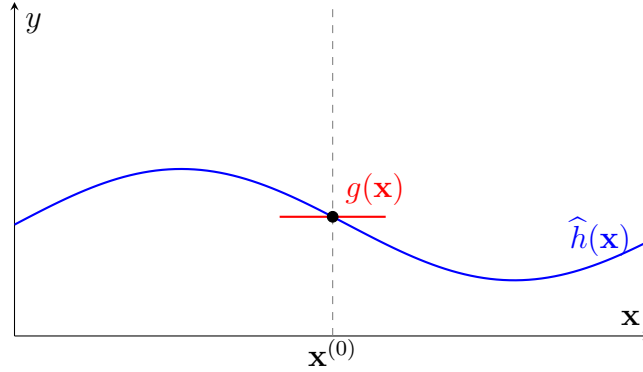


Figura 19: To explain a trained modello $\hat{h} \in \mathcal{H}$, around a given vettore delle caratteristiche $\mathbf{x}^{(0)}$, we can use a local approximation $g \in \mathcal{H}'$.

predict the etichetta y from the caratteristicas \mathbf{x} . The logistic perdita incurred by this previsione is defined as

$$L((\mathbf{x}, y), h) := \log(1 + \exp(-yh(\mathbf{x}))). \quad (9)$$

Carefully note that the expression (8) for the logistic perdita applies only for the spazio delle etichette $\mathcal{Y} = \{-1, 1\}$ and when using the thresholding rule (1).

logistic regression Logistic regressione learns a linear ipotesi map (or classifier) $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ to predict a binary etichetta y based on the numeric vettore delle caratteristiche \mathbf{x} of a punto dati. The quality of a linear ipotesi map is measured by the average logistic loss on some labeled datapoints (i.e., the insieme di addestramento).

loss function A perdita function is a map

$$L : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathbb{R}_+ : ((\mathbf{x}, y), h) \mapsto L((\mathbf{x}, y), h).$$

It assigns a non-negative real number (i.e., the perdita) $L((\mathbf{x}, y), h)$ to a pair that consists of a punto dati, with carateristicas \mathbf{x} and etichetta y , and a ipotesis $h \in \mathcal{H}$. The value $L((\mathbf{x}, y), h)$ quantifies the discrepancy between the true etichetta y and the previsione $h(\mathbf{x})$. Lower (closer to zero) values $L((\mathbf{x}, y), h)$ indicate a smaller discrepancy between previsione $h(\mathbf{x})$ and etichetta y . Figure 14 depicts a perdita function for a given punto dati, with carateristicas \mathbf{x} and etichetta y , as a function of the ipotesis $h \in \mathcal{H}$.

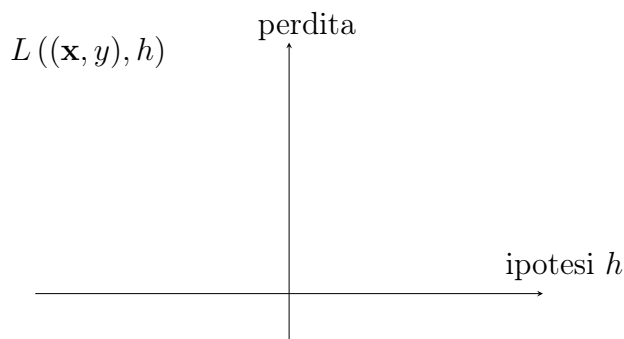


Figura 20: Some perdita function $L((\mathbf{x}, y), h)$ for a fixed punto dati, with vettore delle caratteristiche \mathbf{x} and etichetta y , and a varying ipotesis h . ML methods try to find (or learn) a ipotesis that incurs minimal perdita.

machine learning (ML) Il ML ha come scopo la previsione una etichetta a partire dalle carateristicas di un punto dati. Le tecniche di ML raggiungono quest'obiettivo apprendendo una ipotesis da uno spazio delle ipotesis (o modello) mediante la minimizzazione di una loss function [?, 6]. Una precisa formulazione di questo principio è ERM. I metodi di ML si differenziano in base a differenti scelte inerenti alla definizione dei

punto datis (in termini di caratteristica e etichetta), alla struttura del modello, e alla specificazione della loss function [6, Ch. 3].

map We use the term map as a synonym for function.

See also: function.

massimo Il massimo di un insieme $\mathcal{A} \subseteq \mathbb{R}$ di numeri reali, qualora esista, è l'elemento più grande appartenente a tale insieme. Un insieme \mathcal{A} ammette un massimo se è limitato superiormente e il suo estremo superiore (o minimo dei maggioranti) appartiene ad \mathcal{A} [2, Sec. 1.4].

matrice di covarianza La matrice di covarianza di una variabile aleatoria $\mathbf{x} \in \mathbb{R}^d$ è definita come $\mathbb{E}\left\{(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T\right\}$.

Si veda anche: variabile aleatoria.

maximum likelihood Consider punto datis $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ that are interpreted as the realizzaciones of i.i.d. variabile aleatorias with a common distribuzione di probabilità $\mathbb{P}(\mathbf{z}; \mathbf{w})$ which depends on the model parameters $\mathbf{w} \in \mathcal{W} \subseteq \mathbb{R}^n$. Massimo likelihood methods learn model parameters \mathbf{w} by maximizing the probability (density) $\mathbb{P}(\mathcal{D}; \mathbf{w}) = \prod_{r=1}^m \mathbb{P}(\mathbf{z}^{(r)}; \mathbf{w})$ of the observed dataset. Thus, the massimo likelihood estimator is a solution to the optimization problem $\max_{\mathbf{w} \in \mathcal{W}} \mathbb{P}(\mathcal{D}; \mathbf{w})$.

mean squared estimation error (MSEE) Consider an ML method that learns model parameters $\hat{\mathbf{w}}$ based on some dataset \mathcal{D} . If we interpret the punto datis in \mathcal{D} as i.i.d. realizzaciones of an variabile aleatoria \mathbf{z} , we define the estimation error $\Delta \mathbf{w} := \hat{\mathbf{w}} - \bar{\mathbf{w}}$. Here, $\bar{\mathbf{w}}$ denotes the true model parameters of the distribuzione di probabilità of \mathbf{z} . The media

squared estimation error is defined as the valore atteso $\mathbb{E}\{\|\Delta \mathbf{w}\|^2\}$ of the squared Euclidean norma of the estimation error [?, ?].

media La media di una variabile aleatoria \mathbf{x} , a valori in uno Euclidean space \mathbb{R}^d , è il suo valore atteso $\mathbb{E}\{\mathbf{x}\}$. Essa è definita come l'integrale di Lebesgue di \mathbf{x} rispetto alla distribuzione di probabilità P sottostante (si veda, ad esempio, [2] o [?]), ovvero,

$$\mathbb{E}\{\mathbf{x}\} = \int_{\mathbb{R}^d} \mathbf{x} dP(\mathbf{x}).$$

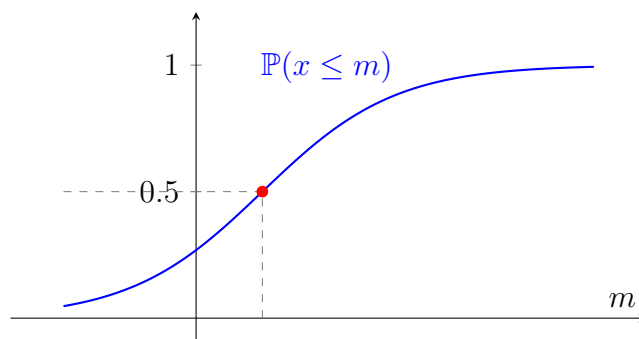
Talvolta è utile considerare la media come la soluzione del seguente problema di minimizzazione del risk [5]:

$$\mathbb{E}\{\mathbf{x}\} = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^d} \mathbb{E}\{\|\mathbf{x} - \mathbf{c}\|_2^2\}.$$

Usiamo inoltre il termine per riferirci alla media di una sequenza finita $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$. Tuttavia, queste due definizioni sono sostanzialmente equivalenti. Infatti, è possibile utilizzare la sequenza $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ per costruire una variabile aleatoria discreta $\tilde{\mathbf{x}} = \mathbf{x}^{(I)}$, dove l'indice I è scelto in modo uniforme a caso dall'insieme $\{1, \dots, m\}$. La media di $\tilde{\mathbf{x}}$ coincide esattamente con la media $\frac{1}{m} \sum_{r=1}^m \mathbf{x}^{(r)}$.

Si veda anche: variabile aleatoria, valore atteso, distribuzione di probabilità.

mediana Una mediana $\operatorname{med}(x)$ di una variabile aleatoria a valori reali x è un qualunque numero $m \in \mathbb{R}$ tale che $\mathbb{P}(x \leq m) \geq 1/2$ e $\mathbb{P}(x \geq m) \geq 1/2$ [?]. Possiamo definire la mediana $\operatorname{med}(\mathcal{D})$ di undataset $\mathcal{D} = \{x^{(1)}, \dots, x^{(m)} \in \mathbb{R}\}$ attraverso una specifica variabile aleatoria \tilde{x} naturalmente associata al \mathcal{D} . In particolare, questa variabile aleatoria è



costruita come $\tilde{x} = x^{(I)}$, dove l'indice I è scelto in modo uniforme a caso dall'insieme $\{1, \dots, m\}$, ossia $\mathbb{P}(I = r) = 1/m$ per ogni $r = 1, \dots, m$. Se la variabile aleatoria x è integrabile, una mediana di x è la soluzione del seguente problema di ottimizzazione:

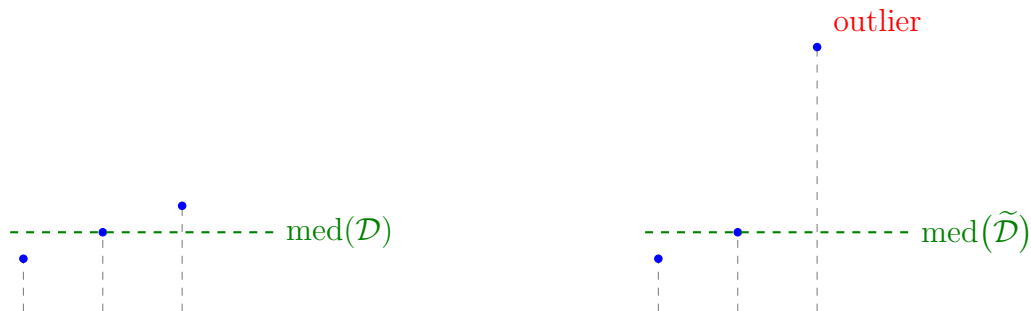
$$\min_{x' \in \mathbb{R}} \mathbb{E}|x - x'|.$$

Analogamente alla media, anche la mediana di un dataset \mathcal{D} può essere utilizzata per stimare i parametri di un modello probabilistico sottostante. Rispetto alla media, la mediana è meno influenzata dagli outliers. Ad esempio, la median di un dataset \mathcal{D} con più di un punto dati non varia anche se si aumenta arbitrariamente il valore del massimo elemento del \mathcal{D} . Al contrario, la media aumenterà anch'essa arbitrariamente.

Si veda anche: media, ??, outlier.

minimo Dato un insieme di numeri reali, il minimo è il più piccolo tra questi.

missing data Consider a dataset constituted by punto datis collected via some physical device. Due to imperfections and failures, some of the caratteristica or etichetta values of punto datis might be corrupted



(a) Dataset originale \mathcal{D} .

(b) Dataset rumoroso $\tilde{\mathcal{D}}$ in cui è presente un outlier.

Figura 21: La median è resistente alla presenza di outlier.

or simply missing. Data imputation aims at estimating these missing values [?]. We can interpret data imputation as an ML problem where the etichetta of a punto dati is the value of the corrupted caratteristica.

model parameters Modello parametri are quantities that are used to select a specific ipotesi map from a modello. We can think of a list of modello parametri as a unique identifier for a ipotesi map, similar to how a social security number identifies a person in Finland.

See also: modello, parametro, ipotesi, map.

model selection In ML, modello selection refers to the process of choosing between different candidate modelli. In its most basic form, modello selection amounts to: 1) training each candidate modello; 2) computing the errore di validazione for each trained modello; and 3) choosing the modello with the smallest errore di validazione [6, Ch. 6].

modello Nel contesto del ML, il termine modello si riferisce tipicamente allo

spazio delle ipotesi sottostante ad un metodo di ML [6], [?]. Tuttavia, il termine è utilizzato anche in altri ambiti, ma con un significato differente. Ad esempio, un modello probabilistico si riferisce a un insieme parametrizzato di distribuzioni di probabilità.

Si veda anche: ML, spazio delle ipotesi, modello probabilistico, distribuzione di probabilità.

modello probabilistico Un modello probabilistico interpreta i punti dati come realizzazioni di variabili aleatorie con una distribuzione di probabilità congiunta. Tale distribuzione di probabilità congiunta coinvolge tipicamente dei parametri che devono essere scelti manualmente oppure appresi mediante metodi di inferenza statistica, come la stima di maximum likelihood [?].

Si veda anche: modello, punto dati, realizzazione, variabile aleatoria, distribuzione di probabilità, parametro, maximum likelihood.

multi-armed bandit A multi-armed bandit (MAB) problem models a repeated decision-making scenario in which, at each time step k , a learner must choose one out of several possible actions, often referred to as arms, from a finite set \mathcal{A} . Each arm $a \in \mathcal{A}$ yields a stochastic reward $r^{(a)}$ drawn from an unknown distribuzione di probabilità with media $\mu^{(a)}$. The learner's goal is to maximize the cumulative reward over time by strategically balancing exploration (gathering information about uncertain arms) and exploitation (selecting arms known to perform well). This balance is quantified by the notion of regret, which measures the performance gap between the learner's strategy and the optimal

strategy that always selects the best arm. MAB problems form a foundational model in online learning, reinforcement learning, and sequential experimental design [?].

multi-label classification Multi-etichetta classificazione problems and methods use punto d'is that are characterized by several etichettas. As an example, consider a punto dati representing a picture with two etichettas. One etichetta indicates the presence of a human in this picture and another etichetta indicates the presence of a car.

multitask learning Multitask learning aims at leveraging relations between different learning tasks. Consider two learning tasks obtained from the same dataset of webcam snapshots. The first task is to predict the presence of a human, while the second task is to predict the presence of a car. It might be useful to use the same deep net structure for both tasks and only allow the weights of the final output layer to be different.

multivariate normal distribution The multivariate normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ is an important family of distribuzione di probabilità for a continuous variabile aleatoria $\mathbf{x} \in \mathbb{R}^d$ [?, ?, 5]. This family is parametrized by the media \mathbf{m} and the matrice di covarianza \mathbf{C} of \mathbf{x} . If the matrice di covarianza is invertible, the distribuzione di probabilità of \mathbf{x} is

$$p(\mathbf{x}) \propto \exp \left(- (1/2)(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}) \right).$$

mutual information (MI) The MI $I(\mathbf{x}; y)$ between two variabile aleatorias \mathbf{x}, y defined on the same spazio di probabilità is given by [?]

$$I(\mathbf{x}; y) := \mathbb{E} \left\{ \log \frac{p(\mathbf{x}, y)}{p(\mathbf{x})p(y)} \right\}.$$

It is a measure of how well we can estimate y based solely on \mathbf{x} . A large value of $I(\mathbf{x}; y)$ indicates that y can be well predicted solely from \mathbf{x} . This prevision could be obtained by a ipotesi learned by an ERM-based ML method.

nearest neighbor (NN) NN methods learn a ipotesi $h : \mathcal{X} \rightarrow \mathcal{Y}$ whose function value $h(\mathbf{x})$ is solely determined by the nearest neighbors within a given dataset. Different methods use different metrics for determining the nearest neighbors. If punto dati are characterized by numeric vettore delle caratteristiche, we can use their Euclidean distances as the metric.

neighborhood The neighborhood of a node $i \in \mathcal{V}$ is the subset of nodes constituted by the neighbors of i .

neighbors The neighbors of a node $i \in \mathcal{V}$ within an FL network are those nodes $i' \in \mathcal{V} \setminus \{i\}$ that are connected (via an edge) to node i .

networked data Networked dati consists of local datasets that are related by some notion of pairwise similarity. We can represent networked dati using a graph whose nodes carry local datasets and edges encode pairwise similarities. One example of networked dati arises in FL applications where local datasets are generated by spatially distributed devices. See also: dati, local dataset, graph, FL, device.

networked exponential families (nExpFam) A collection of exponential families, each of them assigned to a node of an FL network. The

model parameters are coupled via the network structure by requiring them to have a small GTV [?].

networked federated learning (NFL) Networked FL refers to methods that learn personalized modellos in a distributed fashion. These methods learn from local datasets that are related by an intrinsic network structure.

networked model A networked modello over an FL network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ assigns a local model (i.e., a spazio delle ipotesi) to each node $i \in \mathcal{V}$ of the FL network \mathcal{G} .

node degree The degree $d^{(i)}$ of a node $i \in \mathcal{V}$ in an undirected graph is the number of its neighbors, i.e., $d^{(i)} := |\mathcal{N}^{(i)}|$.

non-smooth We refer to a function as non-smooth if it is not smooth [?].

norma A norm è una funzione che associa a ciascun elemento (vettore) di uno spazio vettoriale lineare un numero reale non negativo. Tale funzione deve essere omogenea e definita, e deve soddisfare la disuguaglianza triangolare [?].

objective function An objective function is a map that assigns each value of an optimization variable, such as the model parameters \mathbf{w} of a ipotesi $h^{(\mathbf{w})}$, to an objective value $f(\mathbf{w})$. The objective value $f(\mathbf{w})$ could be the risk or the empirical risk of a ipotesi $h^{(\mathbf{w})}$.

online algorithm An online algorithm processes input dati incrementally, receiving dati items sequentially and making decisions or producing

outputs (or decisions) immediately without having access to the entire input in advance [?, ?]. Unlike an offline algorithm, which has the entire input available from the start, an online algorithm must handle uncertainty about future inputs and cannot revise past decisions. Similar to an offline algorithm, an online algorithm can be modeled formally as a collection of possible executions. However, the execution sequence for an online algorithm has a distinct structure:

$$\text{init}, s_1, \text{out}_1, \text{in}_2, s_2, \text{out}_2, \dots, \text{in}_T, s_T, \text{out}_T.$$

Each execution begins from an initial state (init) and proceeds through alternating computational steps, outputs (or decisions), and inputs. Specifically, at step k , the algorithm performs a computational step s_k , generates an output out_k , and then subsequently receives the next input in_{k+1} . A notable example of an online algorithm in ML is online gradient descent (online GD) (online gradient descent), which incrementally updates model parameters as new punto datis arrive.

online gradient descent (online GD) Consider an ML method that learns model parameters \mathbf{w} from some parameter space $\mathcal{W} \subseteq \mathbb{R}^d$. The learning process uses punto datis $\mathbf{z}^{(t)}$ that arrive at consecutive time-instants $t = 1, 2, \dots$. Let us interpret the punto datis $\mathbf{z}^{(t)}$ as i.i.d. copies of an variabile aleatoria \mathbf{z} . The risk $\mathbb{E}\{L(\mathbf{z}, \mathbf{w})\}$ of a ipotesi $h^{(\mathbf{w})}$ can then (under mild conditions) be obtained as the limit $\lim_{T \rightarrow \infty} (1/T) \sum_{t=1}^T L(\mathbf{z}^{(t)}, \mathbf{w})$. We might use this limit as the objective function for learning the model parameters \mathbf{w} . Unfortunately, this limit can only be evaluated if we wait infinitely long in order to collect all punto datis. Some ML applications

require methods that learn online: as soon as a new punto dati $\mathbf{z}^{(t)}$ arrives at time t , we update the current model parameters $\mathbf{w}^{(t)}$. Note that the new punto dati $\mathbf{z}^{(t)}$ contributes the component $L(\mathbf{z}^{(t)}, \mathbf{w})$ to the risk. As its name suggests, online GD updates $\mathbf{w}^{(t)}$ via a (projected) gradient step

$$\mathbf{w}^{(t+1)} := P_{\mathcal{W}}(\mathbf{w}^{(t)} - \eta_t \nabla_{\mathbf{w}} L(\mathbf{z}^{(t)}, \mathbf{w})). \quad (10)$$

Note that (9) is a gradient step for the current component $L(\mathbf{z}^{(t)}, \cdot)$ of the risk. The update (9) ignores all the previous components $L(\mathbf{z}^{(t')}, \cdot)$, for $t' < t$. It might therefore happen that, compared to $\mathbf{w}^{(t)}$, the updated model parameters $\mathbf{w}^{(t+1)}$ increase the retrospective average perdita $\sum_{t'=1}^{t-1} L(\mathbf{z}^{(t')}, \cdot)$. However, for a suitably chosen learning rate η_t , online GD can be shown to be optimal in practically relevant settings. By optimal, we mean that the model parameters $\mathbf{w}^{(T+1)}$ delivered by online GD after observing T punto datis $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(T)}$ are at least as good as those delivered by any other learning method [?, ?].

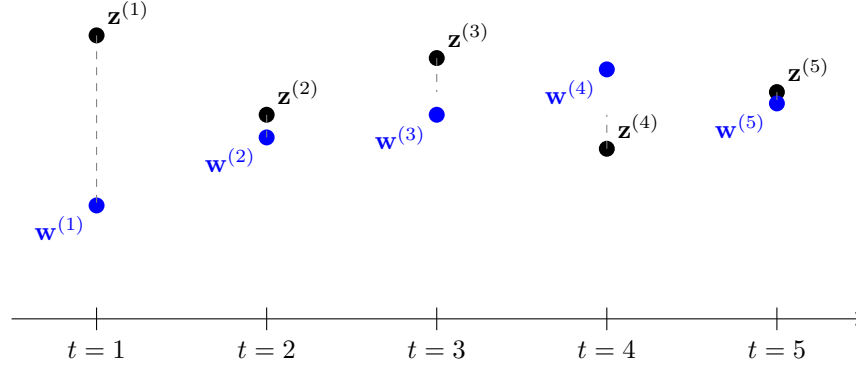


Figura 22: An instance of online GD that updates the model parameters $\mathbf{w}^{(t)}$ using the punto dati $\mathbf{z}^{(t)} = x^{(t)}$ arriving at time t . This instance uses the squared error loss $L(\mathbf{z}^{(t)}, w) = (x^{(t)} - w)^2$.

optimism in the face of uncertainty ML methods learn model parameters \mathbf{w} according to some performance criterion $\bar{f}(\mathbf{w})$. However, they usually cannot access $\bar{f}(\mathbf{w})$ directly but rely on an estimate (or approximation) of $f(\mathbf{w})$. As a case in point, ERM-based methods use the average perdita on a given dataset (i.e., the insieme di addestramento) as an estimate for the risk of a ipotesi. Using a modello probabilistico, one can construct a confidence interval $[l^{(\mathbf{w})}, u^{(\mathbf{w})}]$ for each choice \mathbf{w} for the model parameters. One simple construction is $l^{(\mathbf{w})} := f(\mathbf{w}) - \sigma/2$, $u^{(\mathbf{w})} := f(\mathbf{w}) + \sigma/2$, with σ being a measure of the (expected) deviation of $f(\mathbf{w})$ from $\bar{f}(\mathbf{w})$. We can also use other constructions for this interval as long as they ensure that $\bar{f}(\mathbf{w}) \in [l^{(\mathbf{w})}, u^{(\mathbf{w})}]$ with a sufficiently high probability. As optimists, we choose the model parameters according to the most favorable - yet still plausible - value $\tilde{f}(\mathbf{w}) := l^{(\mathbf{w})}$ of the

performance criterion. Two examples of ML methods that use such an optimistic construction of an objective function are SRM [?, Ch. 11] and upper confidence bound (UCB) methods for sequential decision making [?, Sec. 2.2].

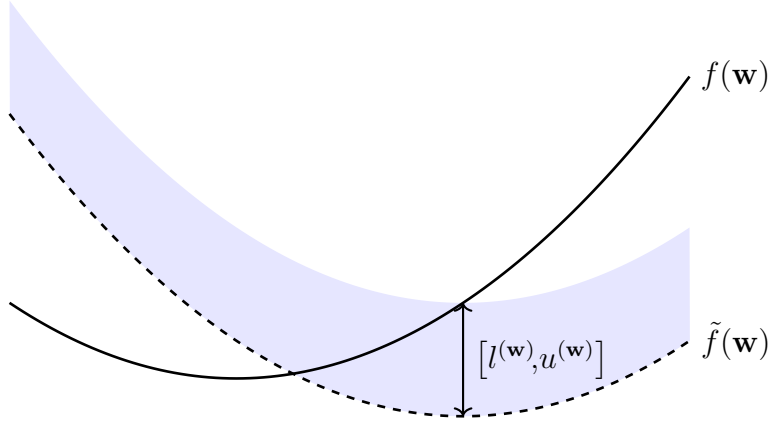


Figura 23: ML methods learn model parameters \mathbf{w} by using some estimate of $f(\mathbf{w})$ for the ultimate performance criterion $\tilde{f}(\mathbf{w})$. Using a modello probabilistico, one can use $f(\mathbf{w})$ to construct confidence intervals $[l^{(\mathbf{w})}, u^{(\mathbf{w})}]$ which contain $\tilde{f}(\mathbf{w})$ with high probability. The best plausible performance measure for a specific choice \mathbf{w} of model parameters is $\tilde{f}(\mathbf{w}) := l^{(\mathbf{w})}$.

optimization method An optimization method is an algorithm that reads in a representation of an optimization problem and delivers an (approximate) solution as its output [?], [?], [?].

See also: algorithm, optimization problem.

optimization problem An optimization problem is a mathematical structure consisting of an objective function $f : \mathcal{U} \rightarrow \mathcal{V}$ defined over an

optimization variable $\mathbf{w} \in \mathcal{U}$, together with a feasible set $\mathcal{W} \subseteq \mathcal{U}$. The co-domain \mathcal{V} is assumed to be ordered, meaning that for any two elements $\mathbf{a}, \mathbf{b} \in \mathcal{V}$, we can determine whether $\mathbf{a} < \mathbf{b}$, $\mathbf{a} = \mathbf{b}$, or $\mathbf{a} > \mathbf{b}$. The goal of optimization is to find those values $\mathbf{w} \in \mathcal{W}$ for which the objective $f(\mathbf{w})$ is extremal—i.e., minimal or maximal [?], [?], [?]. See also: objective function.

outlier Many ML methods are motivated by the i.i.d. assumption, which interprets *punto datis* as realizzazioni of i.i.d. *variabile aleatorias* with a common *distribuzione di probabilità*. The i.i.d. assumption is useful for applications where the statistical properties of the *dati* generation process are stationary (or time-invariant) [?]. However, in some applications the *dati* consists of a majority of regular *punto datis* that conform with an i.i.d. assumption as well as a small number of *punto datis* that have fundamentally different statistical properties compared to the regular *punto datis*. We refer to a *punto dati* that substantially deviates from the statistical properties of most *punto datis* as an outlier. Different methods for outlier detection use different measures for this deviation. Statistical learning theory studies fundamental limits on the ability to mitigate outliers reliably [?, ?].

overfitting Consider an ML method that uses ERM to learn a *ipotesi* with the minimo empirical risk on a given *insieme di addestramento*. Such a method is overfitting the *insieme di addestramento* if it learns a *ipotesi* with a small empirical risk on the *insieme di addestramento* but a significantly larger *perdita* outside the *insieme di addestramento*.

parameter space The parameter space \mathcal{W} of an ML modello \mathcal{H} is the set of all feasible choices for the model parameters (see Figure 17). Many important ML methods use a modello that is parametrized by vectors of the Euclidean space \mathbb{R}^d . Two widely used examples of parametrized modelli are linear models and deep nets. The parameter space is then often a subset $\mathcal{W} \subseteq \mathbb{R}^d$, e.g., all vectors $\mathbf{w} \in \mathbb{R}^d$ with a norma smaller than one.

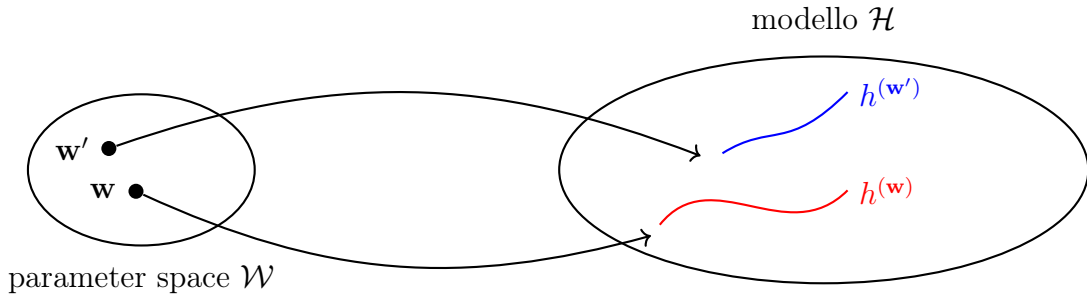


Figura 24: The parameter space \mathcal{W} of an ML modello \mathcal{H} consists of all feasible choices for the model parameters. Each choice \mathbf{w} for the model parameters selects a ipotesi map $h^{(\mathbf{w})} \in \mathcal{H}$.

parametro Il parametro di un modello di ML modello è una quantità regolabile (cioè apprendibile o adattabile) che ci consente di selezionare tra diverse maps ipotesi. Ad esempio, il linear model $\mathcal{H} := \{h^{(\mathbf{w})} : h^{(\mathbf{w})}(x) = w_1x + w_2\}$ consiste di tutte maps ipotesi $h^{(\mathbf{w})}(x) = w_1x + w_2$ con una particolare scelta dei parametri $\mathbf{w} = (w_1, w_2)^T \in \mathbb{R}^2$. Un altro esempio di parametro di un modello è dato dai weights assegnati a una connessione tra due neuroni di una ANN.

Si veda anche: ML, modello, ipotesi, map, linear model, weights, ANN.

perdita metodi di glsml utilizzano una loss function $L(\mathbf{z}, h)$ per misurare l'errore commesso applicando una specifica ipotesi ad un determinato punto dati. Con un lieve abuso di notazione, utilizziamo il termine perdita sia per indicare la loss function L in sé, sia per il valore specifico $L(\mathbf{z}, h)$, relativo a un punto dati \mathbf{z} e ad una ipotesi h .
Si ved anche: ML, loss function, ipotesi, punto dati.

polynomial regression Polynomial regression aims at learning a polynomial ipotesi map to predict a numeric etichetta based on the numeric characteristics of a punto dati. For punto datis characterized by a single numeric caratteristica, polynomial regression uses the spazio delle ipotesi $\mathcal{H}_d^{(\text{poly})} := \{h(x) = \sum_{j=0}^{d-1} x^j w_j\}$. The quality of a polynomial ipotesi map is measured using the average squared error loss incurred on a set of labeled datapoints (which we refer to as the insieme di addestramento).

positive semi-definite (psd) A (real-valued) symmetric matrix $\mathbf{Q} = \mathbf{Q}^T \in \mathbb{R}^{d \times d}$ is referred to as psd if $\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq 0$ for every vector $\mathbf{x} \in \mathbb{R}^d$. The property of being psd can be extended from matrices to (real-valued) symmetric kernel maps $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ (with $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$) as follows: For any finite set of vettore delle caratteristiche $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$, the resulting matrix $\mathbf{Q} \in \mathbb{R}^{m \times m}$ with entries $Q_{r,r'} = K(\mathbf{x}^{(r)}, \mathbf{x}^{(r')})$ is psd [?].

predictor A predictor is a real-valued ipotesi map. Given a punto dati with characteristics \mathbf{x} , the value $h(\mathbf{x}) \in \mathbb{R}$ is used as a previsione for the true numeric etichetta $y \in \mathbb{R}$ of the punto dati.

previsione

principal component analysis (PCA) PCA determines a linear feature map such that the new characteristics allow us to reconstruct the original characteristics with the minimum reconstruction error [6].

privacy attack A privacy attack on an ML system aims to infer sensitive attributes of individuals by exploiting partial access to a trained ML model. One form of a privacy attack is ϵ -DP.

See also: attack, sensitive attribute, ϵ -DP, trustworthy AI, general data protection regulation (GDPR).

privacy funnel The privacy funnel is a method for learning privacy-friendly characteristics of a point \mathbf{x} [?].

privacy leakage Consider an ML application that processes a dataset \mathcal{D} and delivers some output, such as the predictions obtained for new point \mathbf{x} . Privacy leakage arises if the output carries information about a private (or sensitive) characteristic of a point \mathbf{x} (which might be a human) of \mathcal{D} . Based on a probabilistic model for the data generation, we can measure the privacy leakage via the MI between the output and the sensitive characteristic. Another quantitative measure of privacy leakage is DP. The relations between different measures of privacy leakage have been studied in the literature (see [?]).

privacy protection Consider some ML method \mathcal{A} that reads in a dataset \mathcal{D} and delivers some output $\mathcal{A}(\mathcal{D})$. The output could be the learned model parameters $\hat{\mathbf{w}}$ or the prediction $\hat{h}(\mathbf{x})$ obtained for a specific point \mathbf{x} with characteristics \mathbf{x} . Many important ML applications involve point \mathbf{x} representing humans. Each point \mathbf{x} is characterized by

caratteristiche \mathbf{x} , potenzialmente un'etichetta y , e un attributo sensibile s (e.g., una recente diagnosi medica). Roughly speaking, privacy protection means that it should be impossible to infer, from the output $\mathcal{A}(\mathcal{D})$, any of the sensitive attributes of punti dati in \mathcal{D} . Matematicamente, privacy protection requires non-invertibility of the map $\mathcal{A}(\mathcal{D})$. In general, just making $\mathcal{A}(\mathcal{D})$ non-invertible is typically insufficient for privacy protection. We need to make $\mathcal{A}(\mathcal{D})$ sufficiently non-invertible.

probabilistic principal component analysis (PPCA) Probabilistic PCA extends basic PCA by using a modello probabilistico for punti dati. The modello probabilistico of probabilistic PCA reduces the task of dimensionality reduction to an estimation problem that can be solved using EM methods.

probability We assign a probability value, typically chosen in the interval $[0, 1]$, to each event that might occur in a random experiment [?, ?, ?, 5].

probability density function (pdf) The pdf $p(x)$ of a real-valued variable aleatoria $x \in \mathbb{R}$ is a particolare rappresentazione della sua distribuzione di probabilità. If the pdf exists, it can be used to compute the probability that x takes on a value from a measurable set $\mathcal{B} \subseteq \mathbb{R}$ via $\mathbb{P}(x \in \mathcal{B}) = \int_{\mathcal{B}} p(x') dx'$ [5, Ch. 3]. If the pdf of a vector-valued variable aleatoria $\mathbf{x} \in \mathbb{R}^d$ exists, it allows us to compute the probability of \mathbf{x} belonging to a measurable region \mathcal{R} via $\mathbb{P}(\mathbf{x} \in \mathcal{R}) = \int_{\mathcal{R}} p(\mathbf{x}') dx'_1 \dots dx'_d$ [5, Ch. 3].

See also: variabile aleatoria, distribuzione di probabilità, probability.

projected gradient descent (projected GD) Consider an ERM-based method that uses a parametrized model with parameter space $\mathcal{W} \subseteq \mathbb{R}^d$. Even if the objective function of ERM is smooth, we cannot use basic GD, as it does not take into account constraints on the optimization variable (i.e., the model parameters). Projected GD extends basic GD to handle constraints on the optimization variable (i.e., the model parameters). A single iteration of projected GD consists of first taking a gradient step and then projecting the result back onto the parameter space.

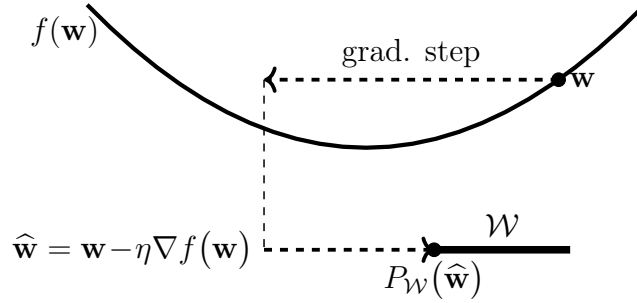


Figure 25: Projected GD augments a basic gradient step with a projection back onto the constraint set \mathcal{W} .

projection Consider a subset $\mathcal{W} \subseteq \mathbb{R}^d$ of the d -dimensional Euclidean space.

We define the projection $P_{\mathcal{W}}(\mathbf{w})$ of a vector $\mathbf{w} \in \mathbb{R}^d$ onto \mathcal{W} as

$$P_{\mathcal{W}}(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}' \in \mathcal{W}} \|\mathbf{w} - \mathbf{w}'\|_2. \quad (11)$$

In other words, $P_{\mathcal{W}}(\mathbf{w})$ is the vector in \mathcal{W} which is closest to \mathbf{w} . The projection is only well-defined for subsets \mathcal{W} for which the above minimum exists [?].

proximable A convex function for which the proximal operator can be computed efficiently is sometimes referred to as proximable or simple [?].

proximal operator Given a convex function $f(\mathbf{w}')$, we define its proximal operator as [?, ?]

$$\mathbf{prox}_{f(\cdot), \rho}(\mathbf{w}) := \underset{\mathbf{w}' \in \mathbb{R}^d}{\operatorname{argmin}} \left[f(\mathbf{w}') + (\rho/2) \|\mathbf{w} - \mathbf{w}'\|_2^2 \right] \text{ with } \rho > 0.$$

As illustrated in Figure 19, evaluating the proximal operator amounts to minimizing a penalized variant of $f(\mathbf{w}')$. The penalty term is the scaled squared Euclidean distance to a given vector \mathbf{w} (which is the input to the proximal operator). The proximal operator can be interpreted as a generalization of the gradient step, which is defined for a smooth convex function $f(\mathbf{w}')$. Indeed, taking a gradient step with step size η at the current vector \mathbf{w} is the same as applying the proximal operator of the function $\tilde{f}(\mathbf{w}') = (\nabla f(\mathbf{w}))^T (\mathbf{w}' - \mathbf{w})$ and using $\rho = 1/\eta$.

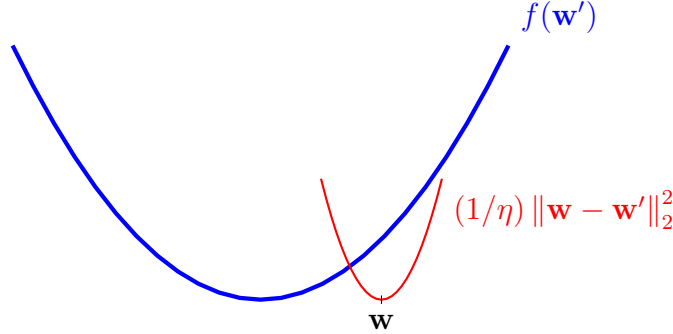


Figure 26: A generalized gradient step updates a vector \mathbf{w} by minimizing a penalized version of the function $f(\cdot)$. The penalty term is the scaled squared Euclidean distance between the optimization variable \mathbf{w}' and the given vector \mathbf{w} .

pseudoinverse The Moore–Penrose pseudoinverse \mathbf{A}^+ of a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ generalizes the notion of an inverse matrix [3]. The pseudoinverse arises naturally within ridge regression when applied to a dataset with arbitrary etichette \mathbf{y} and feature matrix $\mathbf{X} = \mathbf{A}$ [?, Ch. 3]. The model parameters learned by ridge regression are given by

$$\hat{\mathbf{w}}^{(\alpha)} = (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}, \quad \alpha > 0.$$

We can then define the pseudoinverse $\mathbf{A}^+ \in \mathbb{R}^{d \times m}$ via the limit [?, Ch. 3]

$$\lim_{\alpha \rightarrow 0^+} \hat{\mathbf{w}}^{(\alpha)} = \mathbf{A}^+ \mathbf{y}.$$

punto dati Si definisce punto dati qualsiasi oggetto che contiene o veicola informazione [?]. Possono essere punti dati studenti, segnali radio, alberi, foreste, immagini, variabile aleatorias, numeri reali o proteine. I punti dati si caratterizzano utilizzando due tipologie di proprietà. Un tipo di proprietà sono le caratteristiche. Caratteristiche sono proprietà di un punto dati che possono essere misurate o calcolate in modo automatico. Un secondo tipo di proprietà viene definita etichetta. L'etichetta di un punto dati rappresenta un'informazione di livello superiore (o una quantità di interesse). A differenza delle caratteristiche, la determinazione dell'etichetta di un punto dati richiede tipicamente l'intervento di esperti umani (esperti del dominio). In termini generali, l'obiettivo del ML è quello di prevedere l'etichetta di un punto dati basandosi esclusivamente sulle sue caratteristiche.

See also: dati, variabile aleatoria, caratteristica, etichetta, ML.

quadratic function A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ of the form

$$f(\mathbf{w}) = \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{q}^T \mathbf{w} + a,$$

with some matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$, vector $\mathbf{q} \in \mathbb{R}^d$, and scalar $a \in \mathbb{R}$.

See also: function.

random forest A random forest is a set (or ensemble) of different decision trees. Each of these decision trees is obtained by fitting a perturbed copy of the original dataset.

realizzazione Si consideri una variabile aleatoria x che associa a ciascun elemento (ossia, esito o evento elementare) $\omega \in \mathcal{P}$ di uno spazio di probabilità \mathcal{P} ad un elemento a di uno spazio misurabile \mathcal{N} [2], [?], [?]. Una realizzazione di x è un qualunque elemento $a' \in \mathcal{N}$ per il quale si ha $x(\omega') = a'$ per qualche $\omega' \in \mathcal{P}$.

Si veda anche: variabile aleatoria, spazio di probabilità.

rectified linear unit (ReLU) The ReLU is a popular choice for the activation function of a neuron within an ANN. It is defined as $\sigma(z) = \max\{0, z\}$, with z being the weighted input of the artificial neuron.

regime ad alta dimensionalità Il regime ad alta dimensionalità nella ERM si manifesta quando la dimensione effettiva del modello risulta superiore alla sample size, ovvero al numero di punti dati etichettati presenti nel insieme di addestramento. Ad esempio, i metodi di linear regression operano in tale regime ogni volta che il numero d di caratteristiche impiegate per descrivere i punti dati supera il numero di punti dati

presenti nel insieme di addestramento. Un ulteriore esempio di metodi di ML che operano in tale regime è costituito dalle grandi ANNs, le quali presentano un numero di weights (e termini di bias) molto più grande del numero totale di punti dati nel insieme di addestramento. La statistica ad alta dimensionalità costituisce un filone recente della teoria della probability, che studia il comportamento dei metodi di ML nel regime ad alta dimensionalità [?], [?].

Si veda anche: ERM, dimensione effettiva, overfitting, regularization.

regressione I problemi di regressione riguardano la previsione di un'etichetta numerica esclusivamente a partire dalle caratteristiche di un punto dati [6, Ch. 2].

Si veda anche: previsione, etichetta, caratteristica, punto dati.

regret The regret of a ipotesi h relative to another ipotesi h' , which serves as a baseline, is the difference between the perdita incurred by h and the perdita incurred by h' [?]. The baseline ipotesi h' is also referred to as an expert.

regularization A key challenge of modern ML applications is that they often use large modelli, which have an dimensione effettiva in the order of billions. Training a high-dimensional modello using basic ERM-based methods is prone to overfitting: the learned ipotesi performs well on the insieme di addestramento but poorly outside the insieme di addestramento. Regularization refers to modifications of a given instance of ERM in order to avoid overfitting, i.e., to ensure that the learned

ipotesi performs not much worse outside the insieme di addestramento. There are three routes for implementing regularization:

- 1) Modello pruning: We prune the original modello \mathcal{H} to obtain a smaller modello \mathcal{H}' . For a parametric modello, the pruning can be implemented via constraints on the model parameters (such as $w_1 \in [0.4, 0.6]$ for the weight of caratteristica x_1 in linear regression).
- 2) Perdita penalization: We modify the objective function of ERM by adding a penalty term to the errore di addestramento. The penalty term estimates how much larger the expected perdita (or risk) is compared to the average perdita on the insieme di addestramento.
- 3) Data augmentation: We can enlarge the insieme di addestramento \mathcal{D} by adding perturbed copies of the original punto dats in \mathcal{D} . One example for such a perturbation is to add the realizzazione of an variabile aleatoria to the vettore delle caratteristiche of a punto dati.

Figure 20 illustrates the above three routes to regularization. These routes are closely related and sometimes fully equivalent: data augmentation using variabile aleatoria Gaussianas to perturb the vettore delle caratteristiche in the insieme di addestramento of linear regression has the same effect as adding the penalty $\lambda \|\mathbf{w}\|_2^2$ to the errore di addestramento (which is nothing but ridge regression). The decision on which route to use for regularization can be based on the available

computational infrastructure. For example, it might be much easier to implement data augmentation than modello pruning.

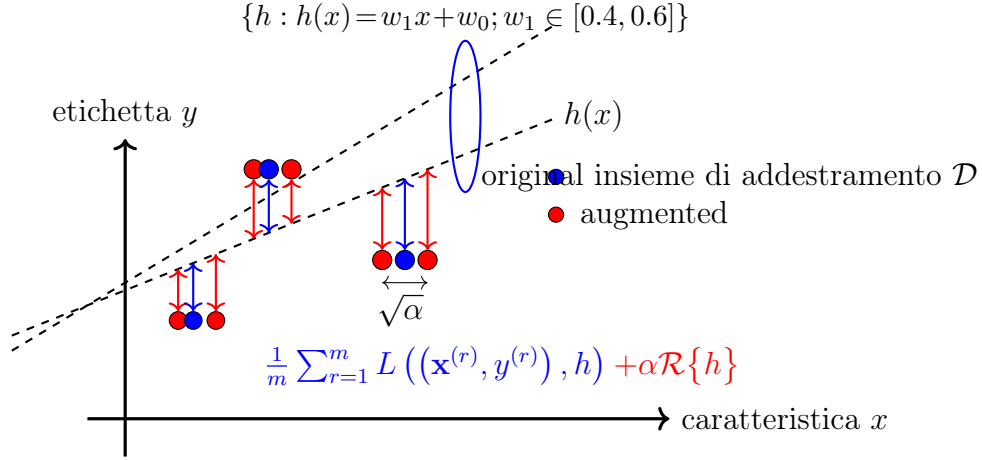


Figura 27: Three approaches to regularization: 1) data augmentation; 2) perdita penalization; and 3) modello pruning (via constraints on model parameters).

regularized empirical risk minimization (RERM) Basic ERM learns a ipotesi (or trains a modello) $h \in \mathcal{H}$ based solely on the empirical risk $\widehat{L}(h|\mathcal{D})$ incurred on a insieme di addestramento \mathcal{D} . To make ERM less prone to overfitting, we can implement regularization by including a (scaled) regularizer $\mathcal{R}\{h\}$ in the learning objective. This leads to regularized empirical risk minimization (RERM),

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}} \widehat{L}(h|\mathcal{D}) + \alpha \mathcal{R}\{h\}. \quad (12)$$

The parameter $\alpha \geq 0$ controls the regularization strength. For $\alpha = 0$, we recover standard ERM without regularization. As α increases, the learned ipotesi is increasingly biased toward small values of $\mathcal{R}\{h\}$. The component $\alpha \mathcal{R}\{h\}$ in the objective function of (11) can be intuitively understood as a surrogate for the increased average perdita that may occur when predicting etichettas for punto datis outside the insieme di addestramento. This intuition can be made precise in various ways. For example, consider a linear model trained using squared error loss and the regularizer $\mathcal{R}\{h\} = \|\mathbf{w}\|_2^2$. In this setting, $\alpha \mathcal{R}\{h\}$ corresponds to the expected increase in perdita caused by adding variabile aleatoria Gaussianas to the vettore delle caratteristiche in the insieme di addestramento [6, Ch. 3]. A principled construction for the regularizer $\mathcal{R}\{h\}$ arises from approximate upper bounds on the generalization error. The resulting RERM instance is known as SRM [?, Sec. 7.2].

regularized loss minimization (RLM) See RERM.

regularizer A regularizer assigns each ipotesi h from a spazio delle ipotesi \mathcal{H} a quantitative measure $\mathcal{R}\{h\}$ for how much its previsione error on a

insieme di addestramento might differ from its prevision errors on punto datis outside the insieme di addestramento. Ridge regression uses the regularizer $\mathcal{R}\{h\} := \|\mathbf{w}\|_2^2$ for linear ipotesi maps $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$ [6, Ch. 3]. Lasso uses the regularizer $\mathcal{R}\{h\} := \|\mathbf{w}\|_1$ for linear ipotesi maps $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$ [6, Ch. 3].

Rényi divergence The Rényi divergence measures the (dis)similarity between two distribuzione di probabilitàs [?].

reward A reward refers to some observed (or measured) quantity that allows us to estimate the perdita incurred by the prevision (or decision) of a ipotesi $h(\mathbf{x})$. For example, in an ML application to self-driving vehicles, $h(\mathbf{x})$ could represent the current steering direction of a vehicle. We could construct a reward from the measurements of a collision sensor that indicate if the vehicle is moving towards an obstacle. We define a low reward for the steering direction $h(\mathbf{x})$ if the vehicle moves dangerously towards an obstacle.

ridge regression Ridge regressione learns the weights \mathbf{w} of a linear ipotesi map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. The quality of a particular choice for the model parameters \mathbf{w} is measured by the sum of two components. The first component is the average squared error loss incurred by $h^{(\mathbf{w})}$ on a set of labeled datapoints (i.e., the insieme di addestramento). The second component is the scaled squared Euclidean norma $\alpha \|\mathbf{w}\|_2^2$ with a regularization parameter $\alpha > 0$. Adding $\alpha \|\mathbf{w}\|_2^2$ to the average squared error loss is equivalent to replacing each original punto datis by the

realizzazione of (infinitely many) i.i.d. variabile aleatorias centered around these punto datis (see regularization).

riduzione della dimensionalità Con riduzione della dimensionalità ci si riferisce a metodi che apprendono una trasformazione $h : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ di un insieme (tipicamente ampio) di caratteristiche grezze x_1, \dots, x_d in un insieme più contenuto di caratteristiche significative $z_1, \dots, z_{d'}$. L'utilizzo di un numero ridotto di caratteristiche comporta molteplici vantaggi:

- Vantaggio statistico: in genere si riduce il rischio di overfitting, poiché la riduzione del numero di caratteristiche comporta spesso una diminuzione della dimensione effettiva di un modello.
- Vantaggio computazionale: un numero inferiore di caratteristiche implica un minor carico computazionale durante l'addestramento dei modelli di ML. Ad esempio, i metodi di linear regression richiedono l'inversione di una matrice la cui dimensione dipende dal numero di caratteristiche.
- **Visualizzazione:** la riduzione della dimensionalità riveste inoltre un ruolo essenziale nella visualizzazione dei dati. Si può, ad esempio, apprendere una trasformazione che restituisce due caratteristiche z_1, z_2 utilizzabili come coordinate in uno diagramma a dispersione. La Fig. 21 mostra lo diagramma a dispersione di cifre scritte a mano posizionate secondo le caratteristiche trasformate. In questo caso, i punti dati sono originariamente rappresentati da

un ampio numero di valori in una scala di grigi (uno per ciascun pixel).

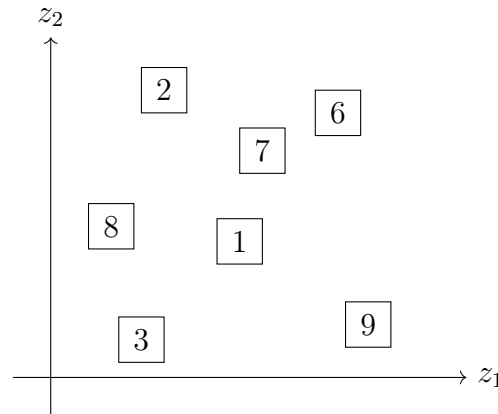


Figura 28: Esempio di riduzione della dimensionalità: dati di immagini ad alta dimensionalità (ad esempio, immagini ad alta risoluzione di cifre scritte a mano) proiettati in 2D utilizzando caratteristiche apprese (z_1, z_2) e visualizzati in un diagramma a dispersione .

Si veda anche: caratteristica, overfitting, dimensione effettiva, modello, ML, linear regression, dati, diagramma a dispersione, punto dati.

risk Consider a ipotesi h used to predict the etichetta y of a punto dati based on its caratteristicas \mathbf{x} . We measure the quality of a particular prevision using a loss function $L((\mathbf{x}, y), h)$. If we interpret punto datis as the realizzaciones of i.i.d. variabile aleatorias, also the $L((\mathbf{x}, y), h)$ becomes the realizzazione of an variabile aleatoria. The i.i.d. assumption allows us to define the risk of a ipotesi as the expected perdita $\mathbb{E}\{L((\mathbf{x}, y), h)\}$. Note that the risk of h depends on both the specific choice for the loss function and the distribuzione di probabilità of the punto datis.

sample A finite sequence (or list) of punto datis $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$ that is obtained or interpreted as the realizzazione of m i.i.d. variabile aleatorias with a common distribuzione di probabilità $p(\mathbf{z})$. The length m of the sequence is referred to as the sample size.

sample covariance matrix The sample matrice di covarianza $\hat{\Sigma} \in \mathbb{R}^{d \times d}$ for a given set of vettore delle caratteristiche $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ is defined as

$$\hat{\Sigma} = (1/m) \sum_{r=1}^m (\mathbf{x}^{(r)} - \hat{\mathbf{m}})(\mathbf{x}^{(r)} - \hat{\mathbf{m}})^T.$$

Here, we use the sample mean $\hat{\mathbf{m}}$.

sample mean The sample media $\mathbf{m} \in \mathbb{R}^d$ for a given dataset, with vettore delle caratteristiche $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$, is defined as

$$\mathbf{m} = (1/m) \sum_{r=1}^m \mathbf{x}^{(r)}.$$

sample size The number of individual punto datis contained in a dataset.

semi-supervised learning (SSL) SSL methods use unlabeled datapoints to support the learning of a ipotesi from labeled datapoints [?]. This approach is particularly useful for ML applications that offer a large amount of unlabeled datapoints, but only a limited number of labeled datapoints.

sensitive attribute ML revolves around learning a ipotesi map that allows us to predict the etichetta of a punto dati from its caratteristiche. In some applications, we must ensure that the output delivered by an ML system does not allow us to infer sensitive attributes of a punto dati.

Which part of a punto dati is considered a sensitive attribute is a design choice that varies across different application domains.

similarity graph Some ML applications generate punto datis that are related by a domain-specific notion of similarity. These similarities can be represented conveniently using a similarity graph $\mathcal{G} = (\mathcal{V} := \{1, \dots, m\}, \mathcal{E})$. The node $r \in \mathcal{V}$ represents the r -th punto dati. Two nodes are connected by an undirected edge if the corresponding punto datis are similar.

singular value decomposition (SVD) The SVD for a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ is a factorization of the form

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{U}^T,$$

with orthonormal matrices $\mathbf{V} \in \mathbb{R}^{m \times m}$ and $\mathbf{U} \in \mathbb{R}^{d \times d}$ [3]. The matrix $\mathbf{\Lambda} \in \mathbb{R}^{m \times d}$ is only non-zero along the main diagonal, whose entries $\Lambda_{j,j}$ are non-negative and referred to as singular values.

smooth A real-valued function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is smooth if it is differenziabile and its gradiente $\nabla f(\mathbf{w})$ is continuous at all $\mathbf{w} \in \mathbb{R}^d$ [?, ?]. A smooth function f is referred to as β -smooth if the gradiente $\nabla f(\mathbf{w})$ is Lipschitz continuous with Lipschitz constant β , i.e.,

$$\|\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}')\| \leq \beta \|\mathbf{w} - \mathbf{w}'\|, \text{ for any } \mathbf{w}, \mathbf{w}' \in \mathbb{R}^d.$$

The constant β quantifies the amount of smoothness of the function f : the smaller the β , the smoother f is. Optimization problems with a smooth objective function can be solved effectively by gradient-based

methods. Indeed, gradient-based methods approximate the objective function locally around a current choice \mathbf{w} using its gradient. This approximation works well if the gradient does not change too rapidly. We can make this informal claim precise by studying the effect of a single gradient step with step size $\eta = 1/\beta$ (see Figure 23).

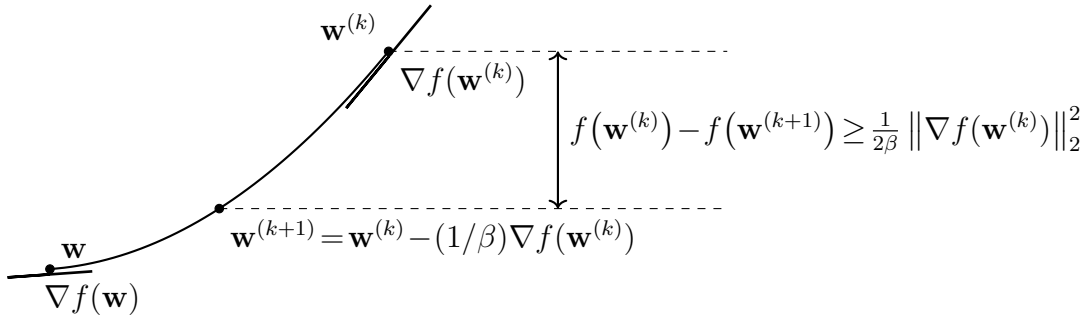


Figure 29: Consider an objective function $f(\mathbf{w})$ that is β -smooth. Taking a gradient step, with step size $\eta = 1/\beta$, decreases the objective by at least $\frac{1}{2\beta} \|\nabla f(\mathbf{w}^{(k)})\|_2^2$ [?, ?, ?]. Note that the step size $\eta = 1/\beta$ becomes larger for smaller β . Thus, for smoother objective functions (i.e., those with smaller β), we can take larger steps.

soft clustering Soft clustering refers to the task of partitioning a given set of punto datis into (a few) overlapping clusters. Each punto dati is assigned to several different clusters with varying degrees of belonging. Soft clustering methods determine the degree of belonging (or soft cluster assignment) for each punto dati and each cluster. A principled approach to soft clustering is by interpreting punto datis as i.i.d. realizzaciones of a GMM. We then obtain a natural choice for the degree of belonging

as the conditional probability of a punto dati belonging to a specific mixture component.

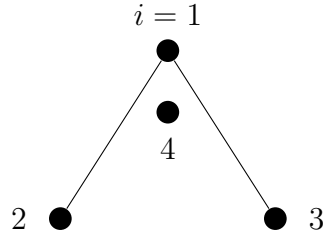
spazio delle etichette Si consideri un'applicazione di ML che coinvolge punti dati caratterizzati da caratteristiche e etichette. Lo spazio delle etichette è costituito da tutti i possibili valori che l'etichetta di un punto dati può assumere. I metodi di regressione, il cui obiettivo è predire etichette numeriche, utilizzano spesso lo spazio delle etichette $\mathcal{Y} = \mathbb{R}$. I metodi di classificazione binaria usano uno spazio delle etichette costituito da due elementi distinti, ad esempio $\mathcal{Y} = \{-1, 1\}$, $\mathcal{Y} = \{0, 1\}$, oppure $\mathcal{Y} = \{\text{"immagine di un gatto"}, \text{"non un immagine di un gatto"}\}$. Si veda anche: ML, punto dati, caratteristica, etichetta, regressione, classificazione.

spazio delle ipotesi Ogni metodo pratico di ML utilizza uno spazio delle ipotesi (o modello) \mathcal{H} . Lo spazio delle ipotesi di un metodo di ML è un sottoinsieme di tutte le possibili maps dallo feature space allo spazio delle etichette. La scelta nella definizione dello spazio delle ipotesi dovrebbe tenere in considerazione le risorse computazionali disponibili e lo statistical aspects. Qualora l'infrastruttura computazionale consenta operazioni matriciali efficienti e sussista una relazione (approssimativamente) lineare tra un insieme di caratteristiche e una etichetta, una scelta opportuna per lo spazio delle ipotesi potrebbe essere il linear model.

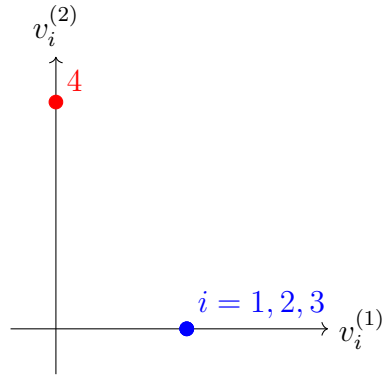
Si veda anche: ML, ipotesi, modello, map, feature space, spazio delle etichette, statistical aspects, caratteristica, etichetta, linear model.

spazio di probabilità

spectral clustering Spectral clustering is a particular instance of graph clustering, i.e., it clusters punto datis represented as the nodes $i = 1, \dots, n$ of a graph \mathcal{G} . Spectral clustering uses the eigenvectors of the Laplacian matrix $\mathbf{L}^{(\mathcal{G})}$ to construct vettore delle caratteristiche $\mathbf{x}^{(i)} \in \mathbb{R}^d$ for each node (i.e., for each punto dati) $i = 1, \dots, n$. We can feed these vettore delle caratteristiche into Euclidean space-based clustering methods, such as k -means or soft clustering via GMM. Roughly speaking, the vettore delle caratteristiche of nodes belonging to a well-connected subset (or cluster) of nodes in \mathcal{G} are located nearby in the Euclidean space \mathbb{R}^d (see Figure 24).



$$\mathbf{L}^{(\mathcal{G})} = \begin{pmatrix} 2 & -1 & -1 & 0 \\ -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$$



$$\mathbf{V} = (\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \mathbf{v}^{(3)}, \mathbf{v}^{(4)})$$

$$\mathbf{v}^{(1)} = \frac{1}{\sqrt{3}} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \end{pmatrix}, \quad \mathbf{v}^{(2)} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

Figura 30: **Top.** Left: An undirected graph \mathcal{G} with four nodes $i = 1, 2, 3, 4$, each representing a punto dati. Right: The Laplacian matrix $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{4 \times 4}$ and its EVD. **Bottom.** Left: A diagramma a dispersione of punto datis using the vettore delle caratteristiche $\mathbf{x}^{(i)} = (v_i^{(1)}, v_i^{(2)})^T$. Right: Two eigenvectors $\mathbf{v}^{(1)}, \mathbf{v}^{(2)} \in \mathbb{R}^d$ corresponding to the eigenvalue $\lambda = 0$ of the Laplacian matrix $\mathbf{L}^{(\mathcal{G})}$.

spectrogram A spectrogram represents the time-frequency distribution of the energy of a time signal $x(t)$. Intuitively, it quantifies the amount of signal energy present within a specific time segment $[t_1, t_2] \subseteq \mathbb{R}$ and frequency interval $[f_1, f_2] \subseteq \mathbb{R}$. Formally, the spectrogram of a signal is defined as the squared magnitude of its short-time Fourier transform (STFT) [?]. Figure 25 depicts a time signal along with its spectrogram.



Figura 31: Left: A time signal consisting of two modulated Gaussian pulses. Right: An intensity plot of the spectrogram.

The intensity plot of its spectrogram can serve as an image of a signal. A

simple recipe for audio signal classification is to feed this signal image into deep nets originally developed for image classification and object detection [?]. It is worth noting that, beyond the spectrogram, several alternative representations exist for the time-frequency distribution of signal energy [?, ?].

squared error loss The squared error perdita measures the prevision error of a ipotesi h when predicting a numeric etichetta $y \in \mathbb{R}$ from the characteristics \mathbf{x} of a punto dati. It is defined as

$$L((\mathbf{x}, y), h) := \left(y - \underbrace{h(\mathbf{x})}_{=\hat{y}} \right)^2.$$

stability Stability is a desirable property of a ML method \mathcal{A} that maps a dataset \mathcal{D} (e.g., a insieme di addestramento) to an output $\mathcal{A}(\mathcal{D})$, such as learned model parameters or the prevision for a specific punto dati. Intuitively, \mathcal{A} is stable if small changes in the input dataset \mathcal{D} lead to small changes in the output $\mathcal{A}(\mathcal{D})$. Several formal notions of stability exist that enable bounds on the generalization error or risk of the method; see [?, Ch. 13]. To build intuition, consider the three datasets depicted in Fig. 26, each of which is equally likely under the same dati-generating distribuzione di probabilità. Since the optimal model parameters are determined by this underlying distribuzione di probabilità, an accurate ML method \mathcal{A} should return the same (or very similar) output $\mathcal{A}(\mathcal{D})$ for all three dataset. In other words, any useful \mathcal{A} must be robust to variability in sample realizzazioni from the same distribuzione di probabilità, i.e., it must be stable.

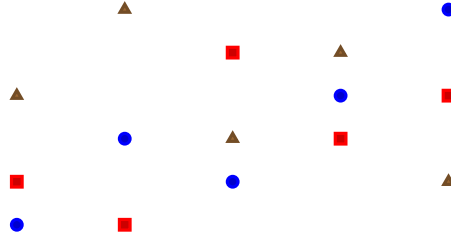


Figura 32: Three datasets $\mathcal{D}^{(*)}$, $\mathcal{D}^{(\square)}$, and $\mathcal{D}^{(\triangle)}$, each sampled independently from the same data-generating distribution. A stable ML method should return similar outputs when trained on any of these datasets.

statistical aspects By statistical aspects of an ML method, we refer to (properties of) the distribution of its output under a probabilistic model for the data fed into the method.

step size See learning rate.

stochastic A process or method is called stochastic if it involves a random component or is governed by probabilistic laws. In ML, stochastic methods often incorporate randomness for reasons such as optimization (e.g., SGD) or uncertainty modeling (e.g., probabilistic models). A stochastic process is a collection of random variables indexed by time or space, which are used to model random phenomena evolving over time (e.g., noise in sensors or financial time series).

See also: ML, SGD, uncertainty, probabilistic model, variable

aleatoria, stochastic block model (SBM).

stochastic block model (SBM) The stochastic block modello is a probabilistic generative modello for an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with a given set of nodes \mathcal{V} [?]. In its most basic variant, the stochastic block modello generates a graph by first randomly assigning each node $i \in \mathcal{V}$ to a cluster index $c_i \in \{1, \dots, k\}$. A pair of different nodes in the graph is connected by an edge with probability $p_{i,i'}$ that depends solely on the etichettas $c_i, c_{i'}$. The presence of edges between different pairs of nodes is statistically independent.

stochastic gradient descent (SGD) Stochastic GD is obtained from GD by replacing the gradiente of the objective function with a stochastic approximation. A main application of stochastic GD is to train a parametrized modello via ERM on a insieme di addestramento \mathcal{D} that is either very large or not readily available (e.g., when punto datis are stored in a database distributed all over the planet). To evaluate the gradiente of the empirical risk (as a function of the model parameters \mathbf{w}), we need to compute a sum $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ over all punto datis in the insieme di addestramento. We obtain a stochastic approximation to the gradiente by replacing the sum $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ with a sum $\sum_{r \in \mathcal{B}} \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ over a randomly chosen subset $\mathcal{B} \subseteq \{1, \dots, m\}$ (see Figure 27). We often refer to these randomly chosen punto datis as a batch. The batch size $|\mathcal{B}|$ is an important parameter of stochastic GD. Stochastic GD with $|\mathcal{B}| > 1$ is referred to as mini-batch stochastic GD [?].

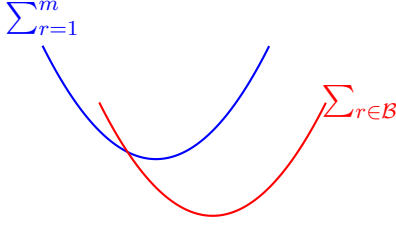


Figura 33: Stochastic GD for ERM approximates the gradiente $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ by replacing the sum over all punto datis in the insieme di addestramento (indexed by $r = 1, \dots, m$) with a sum over a randomly chosen subset $\mathcal{B} \subseteq \{1, \dots, m\}$.

stopping criterion Many ML methods use iterative algorithms that construct a sequence of model parameters in order to minimize the errore di addestramento. For example, gradient-based methods iteratively update the parametri of a parametric modello, such as a linear model or a deep net. Given a finite amount of computational resources, we need to stop updating the parametri after a finite number of iterations. A stopping criterion is any well-defined condition for deciding when to stop updating.

See also: algorithm, gradient-based methods.

strongly convex A continuously differenziabile real-valued function $f(\mathbf{x})$ is strongly convex with coefficient σ if $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + (\sigma/2) \|\mathbf{y} - \mathbf{x}\|_2^2$ [?], [?, Sec. B.1.1].

structural risk minimization (SRM) Structural risk minimization (SRM) is an instance of RERM, which which the modello \mathcal{H} can be expressed as a countable union of sub-models: $\mathcal{H} = \bigcup_{n=1}^{\infty} \mathcal{H}^{(n)}$. Each sub-model

$\mathcal{H}^{(n)}$ permits the derivation of an approximate upper bound on the generalization error incurred when applying ERM to train $\mathcal{H}^{(n)}$. These individual bounds—one for each sub-model—are then combined to form a regularizer used in the RERM objective. These approximate upper bounds (one for each $\mathcal{H}^{(n)}$) are then combined to construct a regularizer for RERM [?, Sec. 7.2].

subgradient For a real-valued function $f : \mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{w} \mapsto f(\mathbf{w})$, a vector \mathbf{a} such that $f(\mathbf{w}) \geq f(\mathbf{w}') + (\mathbf{w} - \mathbf{w}')^T \mathbf{a}$ is referred to as a subgradient of f at \mathbf{w}' [?, ?].

subgradient descent Subgradient descent is a generalization of GD that does not require differentiability of the function to be minimized. This generalization is obtained by replacing the concept of a gradient with that of a subgradient. Similar to gradients, also subgradients allow us to construct local approximations of an objective function. The objective function might be the empirical risk $\widehat{L}(h^{(\mathbf{w})}|\mathcal{D})$ viewed as a function of the model parameters \mathbf{w} that select a ipotesi $h^{(\mathbf{w})} \in \mathcal{H}$.

support vector machine (SVM) The SVM is a binary classificazione method that learns a linear ipotesi map. Thus, like linear regression and logistic regression, it is also an instance of ERM for the linear model. However, the SVM uses a different loss function from the one used in those methods. As illustrated in Figure 28, it aims to maximally separate punto d'is from the two different classes in the feature space (i.e., massimo margin principle). Maximizing this separation is equivalent to minimizing a regularized variant of the hinge loss (5) [?, ?, ?].

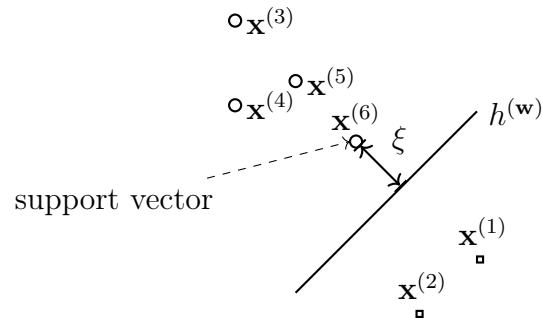


Figura 34: The SVM learns a ipotesi (or classifier) $h^{(\mathbf{w})}$ with minimal average soft-margin hinge loss. Minimizing this perdita is equivalent to maximizing the margin ξ between the decision boundary of $h^{(\mathbf{w})}$ and each class of the insieme di addestramento.

The above basic variant of SVM is only useful if the punto datis from different categories can be (approximately) linearly separated. For an ML application where the categories are not derived from a kernel.

test set A set of punti dati that have been used neither to train a modello (e.g., via ERM) nor in a insieme di validazione to choose between different modelli.

See also: punto dati, modello, ERM, insieme di validazione.

total variation See GTV.

transparency Transparency is a fundamental requirement for trustworthy AI [?]. In the context of ML methods, transparency is often used interchangeably with explainability [?, ?]. However, in the broader scope of AI systems, transparency extends beyond explainability and

includes providing information about the system’s limitations, reliability, and intended use. In medical diagnosis systems, transparency requires disclosing the confidence level for the previsions delivered by a trained modello. In credit scoring, AI-based lending decisions should be accompanied by explanations of contributing factors, such as income level or credit history. These explanations allow humans (e.g., a loan applicant) to understand and contest automated decisions. Some ML methods inherently offer transparency. For example, logistic regression provides a quantitative measure of classificazione reliability through the value $|h(\mathbf{x})|$. Decision trees are another example, as they allow human-readable decision rules [?]. Transparency also requires a clear indication when a user is engaging with an AI system. For example, AI-powered chatbots should notify users that they are interacting with an automated system rather than a human. Furthermore, transparency encompasses comprehensive documentation detailing the purpose and design choices underlying the AI system. For instance, modello data-sheets [?] and AI system cards [?] help practitioners understand the intended use cases and limitations of an AI system [?].

trustworthy artificial intelligence (trustworthy AI) Besides the computational aspects and statistical aspects, a third main design aspect of ML methods is their trustworthiness [?]. The EU has put forward seven key requirements (KRs) for trustworthy AI (that typically build on ML methods) [?]:

- 1) KR1 - Human agency and oversight;

- 2) KR2 - Technical robustness and safety;
- 3) KR3 - Privacy and data governance;
- 4) KR4 - Transparency;
- 5) KR5 - Diversity, non-discrimination and fairness;
- 6) KR6 - Societal and environmental well-being;
- 7) KR7 - Accountability.

uncertainty Uncertainty refers to the degree of confidence—or lack thereof—associated with a quantity such as a model prediction, parameter estimate, or observed data point. In ML, uncertainty arises from various sources, including noisy data, limited training samples, or ambiguity in model assumptions. Probability theory offers a principled framework for representing and quantifying such uncertainty.

underfitting Consider an ML method that uses ERM to learn a hypothesis with the minimum empirical risk on a given dataset. Such a method is underfitting the dataset if it is not able to learn a hypothesis with a sufficiently small empirical risk on the dataset. If a method is underfitting, it will typically also not be able to learn a hypothesis with a small risk.

upper confidence bound (UCB) Consider a ML application that requires selecting, at each time step k , an action a_k from a finite set of alternatives \mathcal{A} . The utility of selecting action a_k is quantified by a numeric reward signal $r^{(a_k)}$. A widely used probabilistic model for this type of

sequential decision-making problem is the stochastic multi-armed bandit setting [?]. In this model, the reward $r^{(a)}$ is viewed as the realizzazione of a variabile aleatoria with unknown media $\mu^{(a)}$. Ideally, we would always choose the action with the largest expected reward $\mu^{(a)}$, but these means are unknown and must be estimated from observed dati. Simply choosing the action with the largest estimate $\hat{\mu}^{(a)}$ can lead to suboptimal outcomes due to estimation uncertainty. The UCB strategy addresses this by selecting actions not only based on their estimated means but also by incorporating a term that reflects the uncertainty in these estimates—favouring actions with high potential reward and high uncertainty. Theoretical guarantees for the performance of UCB strategies, including logarithmic regret bounds, are established in [?].

validazione Si consideri una ipotesi \hat{h} appresa mediante qualche metodo di ML, ad esempio risolvendo un problema di ERM su un insieme di addestramento \mathcal{D} . La validazione si riferisce alla pratica di valutare la perdita associata alla ipotesi \hat{h} su un insieme di punti dati che non appartengono all'insieme di addestramento \mathcal{D} .

Si veda anche: ipotesi, ML, ERM, insieme di addestramento, perdita, punto dati.

valore atteso Si consideri un vettore delle caratteristiche numeriche $\mathbf{x} \in \mathbb{R}^d$ che interpretiamo come una realizzazione di una variabile aleatoria con una distribuzione di probabilità $p(\mathbf{x})$. Il valore atteso di \mathbf{x} è definito come l'integrale $\mathbb{E}\{\mathbf{x}\} := \int \mathbf{x}p(\mathbf{x})$. Si noti che il valore atteso è definito solo se tale integrale esiste, ovvero se la variabile aleatoria è integrabile [?, ?, 2]

Vapnik–Chervonenkis dimension (VC dimension) The VC dimension of an infinite spazio delle ipotesi is a widely-used measure for its size. We refer to the literature (see [?]) for a precise definition of VC dimension as well as a discussion of its basic properties and use in ML.

variabile aleatoria Una variabile aleatoria è una function che mappa da uno spazio di probabilità \mathcal{P} ad uno spazio dei valori [?], [?]. Lo spazio di probabilità è costituito da eventi elementari ed è dotato di una misura di probability che assegna probability a sottoinsiemi di \mathcal{P} . Diverse tipologie di variabili aleatorie includono:

- variabili aleatorie binarie, che mappano ciascun evento elementare ad un elemento di un insieme binario (ad esempio, $\{-1, 1\}$ oppure $\{\text{cat}, \text{no cat}\}$);
- variabili aleatorie a valori reali, che assumono valori nei numeri reali \mathbb{R} ;
- variabili aleatorie a valori vettoriali, che associano eventi elementari allo Euclidean space \mathbb{R}^d .

La teoria della probability utilizza il concetto di spazi misurabili per definire in modo rigoroso e studiare le proprietà di collezioni (anche numerose) di variabili aleatorie [?].

Si veda anche: function, spazio di probabilità, probability, Euclidean space.

variabile aleatoria Gaussiana Una variabile aleatoria Gaussiana standard

è una variabile aleatoria a valori reali x con pdf [?, ?, 5]

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp^{-x^2/2}.$$

Data una variabile aleatoria Gaussiana standard x , è possibile costruire una generica variabile aleatoria Gaussiana x' con media μ e varianza σ^2 tramite $x' := \sigma(x + \mu)$. La distribuzione di probabilità di una variabile aleatoria Gaussiana è detta distribuzione normale, indicata con $\mathcal{N}(\mu, \sigma)$.

Un vettore aleatorio Gaussiano $\mathbf{x} \in \mathbb{R}^d$ con matrice di covarianza \mathbf{C} e media $\boldsymbol{\mu}$ può essere costruito come [?], [?], [?] $\mathbf{x} := \mathbf{A}(\mathbf{z} + \boldsymbol{\mu})$, dove $\mathbf{z} := (z_1, \dots, z_d)^T$ è un vettore di variabili aleatorie gaussiane standard i.i.d., e $\mathbf{A} \in \mathbb{R}^{d \times d}$ è una matrice qualsiasi tale che $\mathbf{A}\mathbf{A}^T = \mathbf{C}$. La distribuzione di probabilità di un vettore aleatorio gaussiano è detta multivariate normal distribution, ed è indicata con $\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$.

I vettori casuali gaussiani sorgono come marginali di dimensione finita dei GPs, i quali definiscono distribuzioni gaussiane congiunte consistenti su insiemi di indici arbitrari (potenzialmente infiniti) [?].

Le variabili aleatorie Gaussiane sono ampiamente utilizzate come modelli probabilistici nell'analisi statistica dei metodi di ML. La loro rilevanza deriva in parte dal central limit theorem (CLT), che rappresenta una formulazione matematica rigorosa della seguente regola empirica: la media di un grande numero di variabili aleatorie indipendenti (non necessariamente Gaussiane) tende a una variabile aleatoria Gaussiana [?].

Rispetto ad altre distribuzioni di probabilità, la multivariate normal distribution si distingue anche per il fatto che— in un senso matematicamente preciso— rappresenta l'incertezza massima. Tra tutti i vettori aleatori continui con una data matrice di covarianza \mathbf{C} , il vettore alea-

torio Gaussiano $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ massimizza l'entropia differenziale [?, Th. 8.6.5]. Questo rende le distribuzioni Gaussianhe una scelta naturale per rappresentare l'incertezza (o la mancanza di conoscenza) in assenza di ulteriori informazioni strutturali.

Si veda anche: multivariate normal distribution, GP, modello probabilistico, CLT, entropy.

varianza La varianza di una variabile aleatoria a valori reali x è definita come il valore atteso $\mathbb{E}\{(x - \mathbb{E}\{x\})^2\}$ del quadrato della differenza tra x ed il suo valore atteso $\mathbb{E}\{x\}$. Estendiamo questa definizione alle variabili aleatorie vettoriali \mathbf{x} come $\mathbb{E}\{\|\mathbf{x} - \mathbb{E}\{\mathbf{x}\}\|_2^2\}$.

Si veda anche: variabile aleatoria, valore atteso.

vector space A vector space (also called linear space) is a collection of elements (called vectors) closed under vector addition and scalar multiplication, i.e.,

- If $\mathbf{x}, \mathbf{y} \in \mathcal{V}$, then $\mathbf{x} + \mathbf{y} \in \mathcal{V}$.
- If $\mathbf{x} \in \mathcal{V}$ and $c \in \mathbb{R}$, then $c\mathbf{x} \in \mathcal{V}$.
- In particular, $\mathbf{0} \in \mathcal{V}$.

The Euclidean space \mathbb{R}^n is a vector space. Linear models and linear maps operate within such spaces.

See also: Euclidean space, linear model, linear map.

vertical federated learning (vertical FL) Vertical FL uses local datasets that are constituted by the same punto d'analisi but characterizing them with different characteristics [?]. For example, different healthcare

providers might all contain information about the same population of patients. However, different healthcare providers collect different measurements (e.g., blood values, electrocardiography, lung X-ray) for the same patients.

vettore delle caratteristiche Il vettore delle caratteristiche si riferisce a un vettore $\mathbf{x} = (x_1, \dots, x_d)^T$ i cui elementi sono singole caratteristiche x_1, \dots, x_d . Molti metodi di ML utilizzano vettori delle caratteristiche appartenenti ad uno Euclidean space di dimensione finita \mathbb{R}^d . Tuttavia, per alcuni metodi di ML, risulta preferibile utilizzare vettori delle caratteristiche che risiedono in uno spazio vettoriale di dimensione infinita (si veda, ad esempio, kernel method).

weights Consider a parametrized spazio delle ipotesi \mathcal{H} . We use the term weights for numeric model parameters that are used to scale characteristics or their transformations in order to compute $h^{(\mathbf{w})} \in \mathcal{H}$. A linear model uses weights $\mathbf{w} = (w_1, \dots, w_d)^T$ to compute the linear combination $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. Weights are also used in ANNs to form linear combinations of characteristics or the outputs of neurons in hidden layers.

zero-gradient condition Consider the unconstrained optimization problem $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ with a smooth and convex objective function $f(\mathbf{w})$. A necessary and sufficient condition for a vector $\hat{\mathbf{w}} \in \mathbb{R}^d$ to solve this problem is that the gradiente $\nabla f(\hat{\mathbf{w}})$ is the zero vector,

$$\nabla f(\hat{\mathbf{w}}) = \mathbf{0} \Leftrightarrow f(\hat{\mathbf{w}}) = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}).$$

0/1 loss The 0/1 perdita $L^{(0/1)}((\mathbf{x}, y), h)$ measures the quality of a classifier $h(\mathbf{x})$ that delivers a previsione \hat{y} (e.g., via thresholding (1)) for the etichetta y of a punto dati with caratteristiche \mathbf{x} . It is equal to 0 if the previsione is correct, i.e., $L^{(0/1)}((\mathbf{x}, y), h) = 0$ when $\hat{y} = y$. It is equal to 1 if the previsione is wrong, i.e., $L^{(0/1)}((\mathbf{x}, y), h) = 1$ when $\hat{y} \neq y$.

Indice analitico

- 0/1 loss, 123
- k -fold cross-validation (k -fold CV), 15
- k -means, 15
- absolute error loss, 15
- accuracy, 15
- activation function, 16
- algorithm, 16
- application programming interface (API), 17
- artificial intelligence (AI), 17
- artificial neural network (ANN), 18
- attack, 18
- autoencoder, 18
- backdoor, 19
- bagging, 19
- baseline, 19
- batch, 21
- Bayes estimator, 22
- Bayes risk, 22
- bias, 22
- bootstrap, 23
- caratteristica, 23
- central limit theorem (CLT), 23
- classificazione, 23
- classifier, 24
- cluster, 24
- clustered federated learning (CFL), 25
- clustering, 26
- clustering assumption, 26
- computational aspects, 26
- condition number, 27
- confusion matrix, 27
- connected graph, 27
- convex, 27
- convex clustering, 28
- Courant–Fischer–Weyl min-max characterization, 28
- covarianza, 29
- data augmentation, 29
- data minimization principle, 30
- data normalization, 31
- data poisoning, 31
- dataset, 31
- dati, 33

decision boundary, 33
 decision region, 34
 decision tree, 34
 deep net, 35
 degree of belonging, 35
 denial-of-service attack, 35
 density-based spatial clustering of
 applications with noise
 (DBSCAN), 36
 determinant, 36
 device, 36
 diagramma a dispersion, 37
 differential entropy, 38
 differential privacy (DP), 38
 differenziabile, 39
 dimensione effettiva, 39
 discrepancy, 39
 distributed algorithm, 39
 distribuzione di probabilità, 40

 edge weight, 41
 eigenvalue, 41
 eigenvalue decomposition (EVD),
 41
 eigenvector, 41
 empirical risk, 41
 empirical risk minimization
 (ERM), 41
 epigraph, 42
 errore di addestramento, 43
 errore di validazione, 44
 estimation error, 44
 estremo superiore (o minimo dei
 maggioranti), 44
 etichetta, 45
 Euclidean space, 45
 expectation-maximization (EM),
 45
 expert, 46
 explainability, 46
 explainable empirical risk
 minimization (EERM), 46
 explainable machine learning
 (explainable ML), 47
 explanation, 47

 feature learning, 47
 feature map, 48
 feature matrix, 48
 feature space, 49
 federated averaging (FedAvg), 49
 federated learning (FL), 49

- federated learning network (FL network), 49
- federated learning orizzontale (FL orizzontale), 49
- FedProx, 50
- Finnish Meteorological Institute (FMI), 50
- fixed-point iteration, 50
- flow-based clustering, 52
- function, 52
- Gaussian mixture model (GMM), 53
- Gaussian Process (GP), 53
- general data protection regulation (GDPR), 54
- generalization, 55
- generalized total variation (GTV), 57
- generalized total variation minimization (GTVMin), 57
- gradient descent (GD), 57
- gradient step, 58
- gradient-based methods, 60
- gradiente, 60
- graph, 60
- graph clustering, 60
- hard clustering, 61
- Hilbert space, 61
- hinge loss, 61
- histogram, 61
- Huber loss, 62
- Huber regression, 62
- hypothesis, 65
- independent and identically distributed (i.i.d.), 63
- independent and identically distributed assumption (i.i.d. assumption), 63
- insieme di addestramento, 63
- insieme di validazione, 64
- interpretability, 64
- inverse matrix, 64
- kernel, 65
- kernel method, 66
- Kullback-Leibler divergence (KL divergence), 67
- label space, 106
- labeled datapoint, 67

- Laplacian matrix, 67
- large language model (LLM), 68
- law of large numbers, 69
- learning rate, 69
- learning task, 69
- least absolute deviation regression, 70
- least absolute shrinkage and selection operator (Lasso), 70
- linear classifier, 70
- linear map, 70
- linear model, 71
- linear regression, 71
- local dataset, 71
- Local Interpretable Model-agnostic Explanations (LIME), 72
- local model, 72
- logistic loss, 72
- logistic regression, 73
- loss function, 73
- machine learning (ML), 74
- map, 75
- massimo, 75
- matrice di covarianza, 75
- maximum likelihood, 75
- mean squared estimation error (MSEE), 75
- media, 76
- mediana, 76
- minimo, 77
- missing data, 77
- model parameters, 78
- model selection, 78
- modello, 78
- multi-armed bandit (MAB), 79
- multi-label classification, 80
- multitask learning, 80
- multivariate normal distribution, 80
- mutual information (MI), 80
- nearest neighbor (NN), 81
- neighborhood, 81
- neighbors, 81
- networked data, 81
- networked exponential families (nExpFam), 81
- networked federated learning (NFL), 82
- networked model, 82

- node degree, 82
- non-smooth, 82
- norm, 82
- objective function, 82
- online algorithm, 82
- online gradient descent (online GD), 83
- optimism in the face of uncertainty, 85
- optimization method, 86
- optimization problem, 86
- outlier, 87
- overfitting, 87
- parameter, 88
- parameter space, 88
- perdita, 89
- polynomial regression, 89
- positive semi-definite (psd), 89
- predictor, 89
- principal component analysis (PCA), 90
- privacy attack, 90
- privacy funnel, 90
- privacy leakage, 90
- privacy protection, 90
- probabilistic model, 79
- probabilistic principal component analysis (PPCA), 91
- probability, 91
- probability density function (pdf), 91
- projected gradient descent (projected GD), 92
- projection, 92
- proximable, 93
- proximal operator, 93
- pseudoinverse, 94
- punto dati, 94
- quadratic function, 95
- Rényi divergence, 100
- random forest, 95
- realizzazione, 95
- rectified linear unit (ReLU), 95
- regime ad alta dimensionalità, 95
- regressione, 96
- regret, 96
- regularization, 96
- regularized empirical risk minimization (RERM), 99

regularized loss minimization	step size, 111
(RLM), 99	stochastic, 111
regularizer, 99	stochastic block model (SBM), 112
reward, 100	stochastic gradient descent (SGD),
ridge regression, 100	112
riduzione della dimensionalità, 101	stopping criterion, 113
risk, 102	strongly convex, 113
sample, 103	structural risk minimization
sample covariance matrix, 103	(SRM), 113
sample mean, 103	subgradient, 114
sample size, 103	subgradient descent, 114
semi-supervised learning (SSL),	support vector machine (SVM),
103	114
sensitive attribute, 103	test set, 115
similarity graph, 104	total variation, 115
singular value decomposition	transparency, 115
(SVD), 104	trustworthy artificial intelligence
smooth, 104	(trustworthy AI), 116
soft clustering, 105	uncertainty, 117
spazio delle ipotesi, 106	underfitting, 117
spectral clustering, 107	upper confidence bound (UCB),
spectrogram, 109	117
squared error loss, 110	validation, 118
stability, 110	valore atteso, 118
statistical aspects, 111	

Vapnik–Chervonenkis dimension

(VC dimension), 119

variabile aleatoria, 119

variabile aleatoria Gaussiana, 119

varianza, 121

vector space, 121

vertical federated learning (vertical

FL), 121

vettore delle caratteristiche, 122

weights, 122

zero-gradient condition, 122

Riferimenti bibliografici

- [1] W. Rudin, *Real and Complex Analysis*, 3rd ed. New York, NY, USA: McGraw-Hill, 1987.
- [2] W. Rudin, *Principles of Mathematical Analysis*, 3rd ed. New York, NY, USA: McGraw-Hill, 1976.
- [3] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 4th ed. Baltimore, MD, USA: The Johns Hopkins Univ. Press, 2013.
- [4] G. H. Golub and C. F. Van Loan, “An analysis of the total least squares problem,” *SIAM J. Numer. Anal.*, vol. 17, no. 6, pp. 883–893, Dec. 1980, doi: 10.1137/0717073.
- [5] D. P. Bertsekas and J. N. Tsitsiklis, *Introduction to Probability*, 2nd ed. Belmont, MA, USA: Athena Scientific, 2008.
- [6] A. Jung, *Machine Learning: The Basics*. Singapore, Singapore: Springer Nature, 2022.