

Machine Learning per la fisica applicata e fisica delle alte energie

Lezione 1: Uno sguardo d'insieme al Machine Learning.

1.1 GENERALITA' SUL CORSO

48 h di lezioni (6 CFU) di cui

24 Dr. Nocera (12 + 12)

24 Prof. Badger (24)

Orario: Lun. 9-11; Mar. 14-16 sempre in AULA D

Obiettivi: familiarizzare con gli algoritmi e i software di Machine Learning in vista sia di una formazione di terzo livello sia di un successivo impiego nel settore ricerca e sviluppo.

Prerequisiti: familiarità con un moderno linguaggio di programmazione (Python, C++) come dal corso di Algoritmi Numerici per la Fisica.

Risultati dell'apprendimento:

- descrivere le tecniche di Machine Learning utilizzate nella risoluzione di problemi fisici complessi;
- individuare le caratteristiche di ciascuna tecnica e metterle in relazione con un problema fisico;
- applicare software specializzato agli scopi del punto precedente.

Programma: come da sito CAMPUSNET; flessibile; aggiornabile.

Modalità di insegnamento: lezioni frontali + set di esercizi proposti a cadenza bi-settimanale + discussione in classe di un paper a scelta tra quelli proposti in un database. Le attività non sono né obbligatorie, né sottoposte a giudizio, MA sono CALDAMENTE CONSIGLIATE.

Esercizio: ORALE. Colloquio di circa 30 minuti in cui lo studente espone la realizzazione di un progetto computazionale a sua scelta, ma concordato con i docenti a lezione, in cui una o più delle tecniche di Machine Learning sono applicate alla risoluzione di un problema fisico. Il problema fisico non dev'essere necessariamente originale e dev'essere possibilmente circoscritto. Le modalità di presentazione sono a scelta dello studente. La lingua (italiano o inglese) pure.

Bibliografia: sul sito CAMPUSNET

- living review
- Hastie, Tibshirani, Friedman "The elements of Statistical Learning" Data Mining, Inference and Prediction", Springer (2013) 2nd Ed.
- James, Witten, Hastie, Tibshirani, "An Introduction to Statistical Learning", Springer (2021) 2nd Ed.
- L. P. Murphy, "Probabilistic Machine Learning: An Introduction" MIT Press (2022).

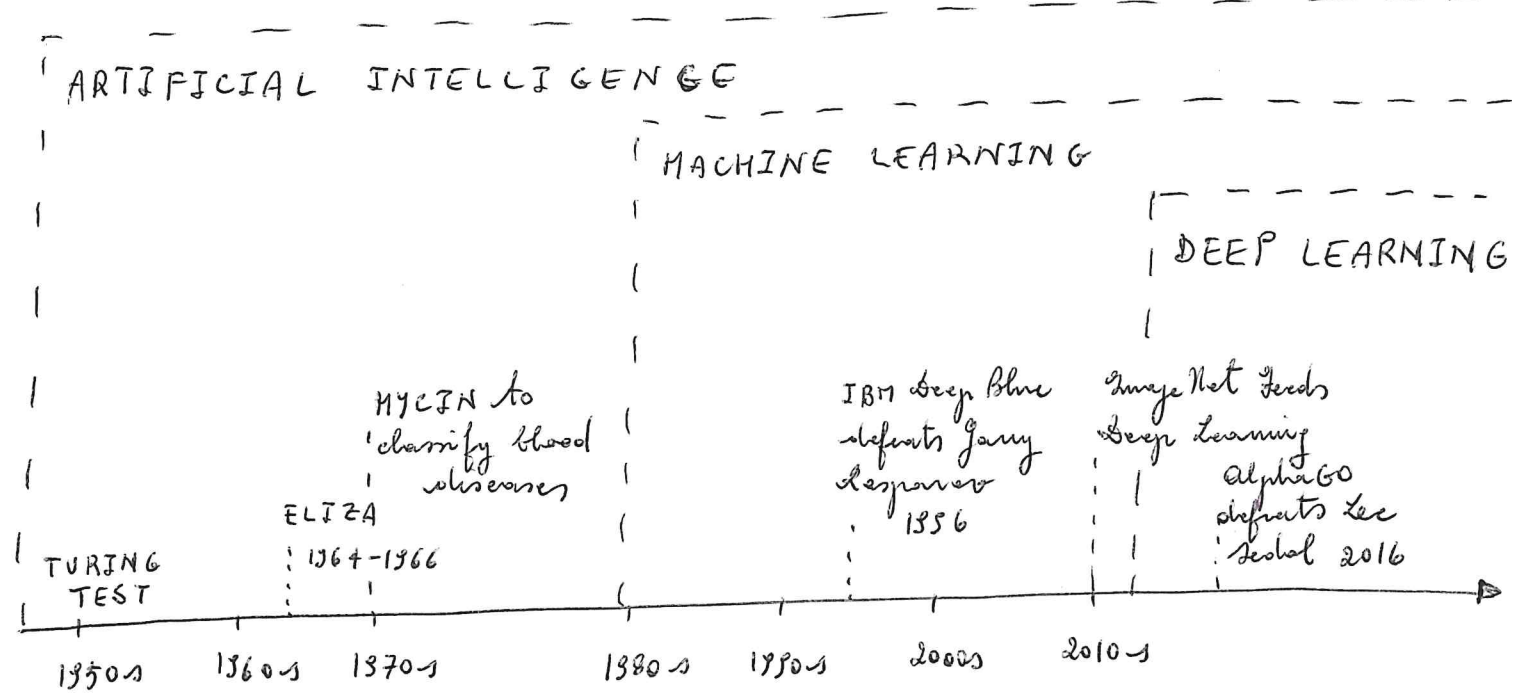
DOMANDE?

1.2 Che cos'è il Machine Learning?

2. Mitchell (1997)

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

Il termine è molto più vecchio e fu introdotto da Arthur Samuel nel contesto del computer gaming nel 1959.



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

MACHINE LEARNING

- A) SUPERVISED LEARNING
 B) UNSUPERVISED LEARNING
 C) REINFORCEMENT LEARNING

Il corso svilupperà essenzialmente A) e B).

1.3 SUPERVISED LEARNING

Il compito T consiste nell'imparare una mappa f dagli inputs $x \in X$ agli outputs $y \in Y$; gli input x sono chiamati FEATURES (o COVARIATES o PREDICTORS) e sono generalmente costituiti da un vettore reale di dimensione finita: $X \equiv \mathbb{R}^D$

(per esempio l'altezza ed il peso di una persona). Gli output y sono noti come LABEL (o TARGET & RESPONSE). L'esperienza E è data come set di N coppie input-output: $D = \{(x_n, y_n)\}_{n=1}^N$ e si chiama TRAINING SET. N si chiama SAMPLE SIZE. La performance P dipende dal task T .

1.3.1. CLASSIFICAZIONE

In un problema di classificazione lo spazio degli output C è un set NON ordinato ed esclusivo di LABELS y , dette CLASSES, $Y = \{1, 2, \dots, C\}$. Il problema di predire una classe dato un input è detto PATTERN RECOGNITION. Se gli output sono solo due, il problema è detto CLASSIFICAZIONE BINARIA.

Esempio: Classificare la specie di Iris (Setosa, Versicolor, Virginica).

In generale, in IMAGE CLASSIFICATION, gli input sono immagini

$$X = \mathbb{R}^D, \quad D = C \times D_1 \times D_2 \quad C = 3 \text{ channels (RGB)}$$

$$f: X \rightarrow Y \quad (\text{è l'immagine un cane? un gatto? oppure?})$$

Nel caso delle specie di Iris i botanisti hanno individuato 4 caratteristiche numeriche: lunghezza e larghezza del sepalo e del petalo. Pertanto $X = \mathbb{R}^4$. Il training set è una collezione di 150 esempi delle 3 specie, 50 per ogni specie.

index	sl[cm]	sw[cm]	pl[cm]	pw[cm]	label	} TABULAR DATA (DESIGN MATRIX)
0	5,1	3,5	1,4	0,2	Setosa	
1	4,9	3,0	1,4	0,2	Setosa	
⋮						
50	7,0	3,2	4,7	1,4	Versicolor	
⋮						
150	5,9	3,0	5,1	1,8	Virginica	

D

BIG DATA: N grande ($N \gg D$)

WIDE DATA: D grande ($D \gg N$)

E' una buona idea fare un'esplorazione dei dati (EXPLORATORY DATA ANALYSIS) per vedere se ci siano dei pattern ovvi, per esempio mediante pair plots. Per dati di stati grandi, si può procedere mediante DIMENSIONALITY REDUCTION.

$$f(\vec{x}, \vec{\vartheta}) = \begin{cases} \text{Setore se p.l.} < 2,45 \text{ cm} \\ \text{Bevercolor o Origine altrimenti} \end{cases} \left\{ \begin{array}{l} \text{DECISION} \\ \text{TREE} \end{array} \right.$$

↑
THRESHOLD
PARAMETER

$$\begin{cases} \text{p.l.} < 2,75 \text{ cm} & \text{Bevercolor} \\ \text{Origine altrimenti} \end{cases}$$

Lo scopo del supervised learning è quello di ottenere modelli di classificazione (tra le altre cose). La PERFORMANCE può essere misurata come:

$$\mathcal{L}(\vartheta) \triangleq \frac{1}{N} \sum_{n=1}^N \mathbb{I}(y_n \neq f(\vec{x}_n, \vec{\vartheta})) \quad \text{MISCLASSIFICATION RATE}$$

dove $\mathbb{I}(e)$ è l'indicatore binario $\mathbb{I}(e) = \begin{cases} 1 & \text{se } e \text{ è vero} \\ 0 & \text{se } e \text{ è falso} \end{cases}$

Nel caso in cui alcuni errori di classificazione siano più dannosi di altri, posso definire una loss function $l(y, \hat{y})$ e

$$\mathcal{L}(\vartheta) \triangleq \frac{1}{N} \sum_{n=1}^N l(y_n, f(\vec{x}_n, \vec{\vartheta})) \quad \text{EMPIRICAL RISK}$$

Una maniera per definire il TRAINING (o MODEL FITTING) è minimizzare il rischio empirico

$$\hat{\vartheta} = \underset{\vartheta}{\operatorname{argmin}} \mathcal{L}(\vartheta) = \underset{\vartheta}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N l(y_n, f(\vec{x}_n, \vec{\vartheta}))$$

sebbene lo scopo sia di considerare dati che non conosciamo.

Nella determinazione di modelli parametrici è comune avere la probabilità logaritmica negativa come loss function

$$\ell(y, f(\bar{x}; \bar{\theta})) = -\ln p(y/f(\bar{x}; \bar{\theta}))$$

RAGIONE: un buon modello (con un loss basso) assegna un'alta probabilità all'output vero y \forall input x . La media è

$$NLL(\bar{\theta}) = -\frac{1}{N} \sum_{n=1}^N \ln p(y_n/f(\bar{x}_n; \bar{\theta})) \quad \text{NEGATIVE LOG-LIKELIHOOD}$$

Minimizzando NLL , si calcola

$$\hat{\theta}_{MLE} = \underset{\bar{\theta}}{\operatorname{argmin}} NLL(\bar{\theta}) \quad \text{MAXIMUM LIKELIHOOD ESTIMATE}$$

1.3.2 REGRESSIONE

In un problema di regressione lo spazio degli output è $y \in \mathbb{R}$. Esempio: y potrebbe essere il grado di tossicità di un feto quando mangiate da una vacca al pascolo.

In tal caso, occorre definire una loss function differente. Una scelta naturale è il QUADRATIC LOSS (o ℓ_2):

$$\ell_2(y, \hat{y}) = (y - \hat{y})^2$$

Segue che

$$MSE(\bar{\theta}) = \frac{1}{N} \sum_{n=1}^N (y_n - f(\bar{x}_n; \bar{\theta}))^2$$

Supponendo che la distribuzione di probabilità degli output sia Gaussiana

$$\mathcal{N}(y/\mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$$

allora:

$$p(y/\bar{x}; \bar{\theta}) = \mathcal{N}(y/f(\bar{x}; \bar{\theta}), \sigma^2)$$

e pertanto

$$MLL(\vartheta) = -\frac{1}{N} \sum_{n=1}^{N_D} \ln \left[\left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left(-\frac{1}{2\sigma^2} (y_n - f(x_n; \vartheta))^2 \right) \right]$$
$$= \frac{1}{2\sigma^2} MSE(\bar{\vartheta}^?) + \text{const}$$

Pertanto MLE sui parametri minimizza ℓ_2 .

1.3.2.1 REGRESSIONE LINEARE

$$f(x, \bar{\vartheta}^?) = w_0 x + b$$

↑ ↑
slope offset

$$\bar{\vartheta}^? = (w_0, b)$$

$$\hat{\vartheta} = \underset{\bar{\vartheta}^?}{\operatorname{argmin}} MSE(\bar{\vartheta}^?) \quad (\text{LEAST SQUARE SOLUTION})$$

In generale

$$f(\vec{x}, \bar{\vartheta}^?) = b + w_0 x_1 + w_1 x_2 + \dots + w_D x_D = b + \vec{w}^? \cdot \vec{x}^?$$

(MULTIPLE LINEAR REGRESSION)

1.3.2.2 REGRESSIONE POLINOMIALE

$$f(x; \vec{w}) = \vec{w}^T \Phi(x) \quad \text{con} \quad \Phi(x) = [1, x, x^2, \dots, x^D]$$

Se $D = N - 1$ allora si parla di INTERPOLATION ($MSE=0$)

$$f(\vec{x}; \vec{w}) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + \dots$$

1.3.2.3 RETI NEURALI

A due esempi precedenti prevedono che f sia lineare in $\vec{w}^?$

$$f(\vec{x}, \vec{w}, \vec{V}) = \vec{w}^T \Phi(\vec{x}, \vec{V}) \quad \vec{V} \text{ è un set di parametri per } \Phi(x)$$

$$f(\vec{x}, \vec{\vartheta}) = f_L(f_{L-1}(\dots(f_1(\vec{x}'))\dots)) \quad \text{con} \quad f_\ell(\vec{x}) = \vec{w}_\ell^T f_{1:L-1}(\vec{x})$$

NEURAL NETWORK

Overfitting and generalisation. Musi però descrivere

$$\mathcal{L}(\vec{\theta}, \mathcal{D}_{\text{train}}) \triangleq \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(\vec{x}, \vec{y}) \in \mathcal{D}_{\text{train}}} \ell(\vec{y}, f(\vec{x}, \vec{\theta}))$$

rendendolo esplicito il set di training $\mathcal{D}_{\text{train}}$. Se $\mathcal{L}(\vec{\theta}, \mathcal{D}_{\text{train}})$ è tale per cui $\text{MSE} = 0$, allora si parla di OVERFITTING.

Per stabilire se un modello soffre di OVERFITTING, supponiamo di conoscere la "vera" distribuzione $p^*(\vec{x}, \vec{y})$ usata per generare il training set. Allora, posso calcolare il POPULATION RISK

$$\mathcal{L}(\vec{\theta}, p^*) \triangleq \mathbb{E}_{p^*(x, y)} [\ell(\vec{y}, f(\vec{x}, \vec{\theta}))]$$

La differenza $\mathcal{L}(\vec{\theta}, p^*) - \mathcal{L}(\vec{\theta}, \mathcal{D}_{\text{train}})$ si chiama GENERALISATION GAP: GG large \Rightarrow OVERFITTING. In practice, partition the data in $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$

$$\mathcal{L}(\vec{\theta}, \mathcal{D}_{\text{test}}) \triangleq \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(\vec{x}, \vec{y}) \in \mathcal{D}_{\text{test}}} \ell(\vec{y}, f(\vec{x}, \vec{\theta})) \approx \text{POPULATION RISK}$$

How do we choose our model? Hyperoptimisation.

No free lunch theorem. "All models are wrong, but some are useful" (George Box, 1987). There is no single best model that works optimally for all kinds of problems.

1.4. UNSUPERVISED LEARNING

Il compito τ consiste nel trovare certe caratteristiche degli inputs $x \in X$ senza conoscere gli output y . Se un punto di vista probabilistico, si tratta di fittare un modello $p(\vec{x})$ non condizionato (per contro, il supervised learning fitte $p(\vec{y}/\vec{x})$).

Differenze tra SL e US.

US: non c'è necessità di collezionare grandi training sets

- US: non c'è necessità di partizionare X
- US: obbliga il modello a "spiegare" gli inputs

1.4.1 CLUSTERING and LATENT FACTORS

Il problema consiste nel trovare partizioni in X che corrispondono a x simili. Esempio: classificazione degli Iris. (CLUSTERING).

Il problema consiste nel proiettare un data set "high-dimensional" in caratteristiche "low dimensional". Esempio: suggerimento che

$x_n \in \mathbb{R}^D$ sia generato da low-dimensional LATENT FACTORS $z_n \in \mathbb{R}^K$: $z_n \rightarrow x_n$. Assumiamo un prior Gaussiano $p(z_n)$

$$p(\bar{x}_n / \bar{z}_n; \bar{\theta}) = \mathcal{N}(\bar{x}_n / \bar{W} \bar{z}_n + \bar{\mu}, \bar{\Sigma}) \quad \begin{matrix} \text{FACTOR ANALYSIS} \\ \text{PRINCIPAL COMPONENT ANALYSIS} \end{matrix}$$

(e $\bar{\Sigma} = \sigma^2 \mathbb{I}$)

se il modello è non-lineare, allora si parla di VARIATIONAL AUTOENCODER

Come valutare la bontà dell'un-supervised learning?

Minimare la probabilità assegnata dal modello a test samples non viste dal modello

$$\mathcal{L}(\bar{\theta}; \mathcal{D}) = - \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \ln p(\bar{x} / \bar{\theta})$$

cioè un "buon modello" non si farà sorprendere da data samples che obbediscono al modello.

1.5 REINFORCEMENT LEARNING

In questa classe di problemi, il sistema o AGENTE, deve imparare ad interagire con l'ambiente. Questo può essere realizzato mediante una POLICY $\bar{a} = \pi(\bar{x})$ che specifichi che cosa l'agente deve fare \forall possibile input \bar{x} dell'ambiente. L'agente riceve una ricompensa occasionale (o una punizione) per le sue azioni. È simile ad imparare con un critico (anziché con un insegnante)

