

ALMA MATER STUDIORUM · UNIVERSITÀ DI BOLOGNA

---

SCUOLA DI SCIENZE  
Corso di Laurea in Matematica

ELEMENTI DI PROBABILITÀ  
PER IL  
MACHINE LEARNING

Relatore:  
Chiar.mo Prof.  
Andrea Pascucci

Presentata da:  
Elia Saltarelli

Sessione I  
Anno Accademico 2019/2020



*A me stesso,  
alla mia famiglia,  
ai miei amici piú cari.*



# Indice

<b>Introduzione</b>	<b>5</b>
<b>1 Preliminari</b>	<b>7</b>
1.1 Nozioni probabilistiche . . . . .	7
1.2 Nozioni di Teoria dei Grafi . . . . .	13
<b>2 Macchine di Boltzmann</b>	<b>14</b>
2.1 Neurone formale . . . . .	14
2.2 Rete neurale stocastica . . . . .	15
2.3 Macchine di Boltzmann ristrette . . . . .	16
2.4 Metodo della massima verosimiglianza . . . . .	18
2.5 Addestramento di una RBM . . . . .	22
<b>3 Algoritmo</b>	<b>25</b>
<b>Bibliografia</b>	<b>30</b>

# Introduzione

Il machine learning o apprendimento automatico é quella branca dell'intelligenza artificiale che si pone l'obiettivo di far apprendere in modo automatico a delle macchine che funzionano imitando il comportamento del cervello umano.

Fu proprio la possibilitá che le macchine imparassero dai dati ad interessare negli anni '40 il matematico Marvin Minsky e i ricercatori Arthur Samuel e Frank Rosenblatt, sfruttando teorie e concetti probabilistici.

Fu lo stesso Minsky a costruire la prima rete neurale stocastica, lo *SNARC* (Stochastic neural analog reinforcement calculator) che comprendeva 40 neuroni artificiali interconnessi.

Nel '58, Frank Rosenblatt pubblicó il primo modello di perceptrone, una rete neurale che in un numero finito di cicli di addestramento poteva apprendere qualsiasi tipo di compito, nei limiti delle sue dimensioni. Per ulteriori informazioni si veda [6].

Oggi, le tecniche di machine learning alimentano molti settori della societá moderna:

**Ingegneria e automazione:** realizzare sistemi automatizzati in grado di apprendere dal mutare delle condizioni, di ridurre gli sprechi e di imparare dai propri errori.

**Informatica:** le ricerche sul web utilizzano motori di ricerca, i quali attraverso una o piú parole chiave, restituiscono liste di risultati attinenti alla ricerca effettuata. Nello stesso modo avviene la selezione dei contenuti sui social network per abbinare notizie, post e pubblicitá agli interessi degli utenti [1].

**Finanza:** si utilizzano le reti neurali come strumento di predizione dei movimenti futuri di un dato titolo.

**Medicina:** la possibilitá di usare macchine di Boltzmann per l'analisi di immagini mediche, per la refertazione in radiologia ecc. Per esempio si veda [4].

**Genetica:** nel sequenziamento, per identificare specifici pattern, all'interno di una grande quantitá di dati genetici, con l'obiettivo di aiutare a predire le probabilitá di un paziente di sviluppare una data patologia o di progettarne un potenziale percorso di cura.

**Sicurezza:** tramite sistemi di prevenzione dalle frodi, ad esempio i filtri anti-spam delle e-mail basati su sistemi di machine learning in grado di riconoscere messaggi di posta fraudolenti. Per la clonazione di carte di credito, gli algoritmi mettono in relazione eventi ed abitudini del soggetto in modo da individuare comportamenti anomali. Fino ad arri-

vare al riconoscimento facciale e a quello vocale per poter accedere ai nostri smartphone.

**Industria automobilistica:** il machine learning é alla base delle auto a guida autonoma, consentendo loro di riconoscere l'ambiente circostante, distinguere un pedone da un lampione oppure rilevare segnali, strisce pedonali e corsie. Si veda per esempio [5].

L'obiettivo di questo elaborato é quello di definire le reti neurali e studiare alcune tecniche di machine learning da un punto di vista matematico.

# Capitolo 1

## Preliminari

### 1.1 Nozioni probabilistiche

Per le definizioni e i teoremi di questo capitolo si veda [11].

#### **Definizione 1.1. Spazio Misurabile**

Uno spazio misurabile é una coppia  $(\Omega, F)$  dove:

- $\Omega$  é un insieme non vuoto;
- $F$  é una  $\sigma$ -algebra su  $\Omega$  cioè una famiglia non vuota di sottoinsiemi di  $\Omega$  che soddisfa le seguenti proprietà:
  1. se  $A \in F$  allora  $A^c \in F$ ;
  2. l'unione numerabile di elementi di  $F$  appartiene ad  $F$ .

#### **Definizione 1.2. Misura**

Una misura sullo spazio misurabile  $(\Omega, F)$  é una funzione

$$\mu : F \rightarrow [0, +\infty] \tag{1.1}$$

tale che:

1.  $\mu(\emptyset) = 0$
2.  $\mu$  é  $\sigma$ -additiva su  $F$ , ossia per ogni successione  $(A_n)_{n \in \mathbb{N}}$  di elementi disgiunti di  $F$  vale:

$$\mu \left( \bigsqcup_{n=1}^{+\infty} A_n \right) = \sum_{n=1}^{+\infty} \mu(A_n) \tag{1.2}$$

**Definizione 1.3. Spazio di probabilità**

Uno spazio con misura  $(\Omega, F, \mu)$  in cui  $\mu(\Omega) = 1$  è detto spazio di probabilità. In particolare  $\mu = P$  è chiamata *misura di probabilità* o semplicemente *probabilità*. Ogni elemento  $\omega \in \Omega$  è detto *esito*; ogni  $A \in F$  è detto *evento* e il numero  $P(A)$  è la *probabilità di A*.

Può essere utile, come di seguito, introdurre il concetto di probabilità condizionata per studiare la dipendenza tra eventi.

**Definizione 1.4. Probabilità condizionata**

In uno spazio di probabilità  $(\Omega, F, P)$  sia B un evento *non trascurabile*, cioè  $P(B) > 0$ , la probabilità di A condizionata a B è definita da:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad A \in F \quad (1.3)$$

**Teorema 1.1. Formula della probabilità totale**

Per ogni evento B tale che  $0 < P(B) < 1$  vale

$$P(A) = P(A|B)P(B) + P(A|B^C)(1 - P(B))$$

Piú in generale, se  $(B_i)_{i \in I}$  è una partizione finita o numerabile di  $\Omega$ , con  $P(B_i) > 0$  per ogni  $i \in I$ , vale

$$P(A) = \sum_{i \in I} P(A \cap B_i) = \sum_{i \in I} P(A|B_i)P(B_i) \quad (1.4)$$

*dimostrazione.*

La dimostrazione deriva direttamente dalla  $\sigma$ -additività della misura P.

**Teorema 1.2. Formula di Bayes**

Siano A, B eventi non trascurabili. Allora vale:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1.5)$$

*dimostrazione.*

Il risultato viene direttamente dalla definizione di probabilità condizionata, considerando

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(B|A) = \frac{P(B \cap A)}{P(A)}$$

Notando che  $P(A \cap B)$  e  $P(B \cap A)$  sono equivalenti, otteniamo

$$P(A|B)P(B) = P(B|A)P(A).$$

Da qui il risultato.

**Definizione 1.5. Indipendenza di eventi**

In uno spazio di probabilità  $(\Omega, F, P)$  diciamo che due eventi A e B sono *indipendenti in P* se

$$P(A \cap B) = P(A)P(B) \tag{1.6}$$

Questo significa che la misura dell'intersezione tra i due eventi é il prodotto delle misure dei due eventi.

Il concetto di indipendenza tra eventi é relativo alla misura di probabilità che sto considerando. Essere eventi indipendenti in P sta a significare che l'accadere dell'evento B non influenza la probabilità di A.

Le nozioni e gli esempi svolti in questo elaborato si riferiscono a spazi di probabilità discreti, ossia quando  $\Omega$  é finito o al piú numerabile, in quanto é sufficiente rivolgersi al caso discreto per introdurci allo studio del machine learning.

Capiremo, in seguito, come ci si puó occupare di apprendimento automatico utilizzando spazi di probabilità continui.

Consideriamo  $(\Omega, F, P)$  e fissiamo  $d \in \mathbb{N}$ . dati  $H \subseteq \mathbb{R}^d$  e una funzione  $X : \Omega \rightarrow \mathbb{R}^d$ , indichiamo con

$$(X \in H) := \{\omega \in \Omega | X(\omega) \in H\} = X^{-1}(H)$$

la contro-immagine di H mediante X. Intuitivamente  $(X \in H)$  rappresenta l'insieme degli esiti  $\omega$ , ossia i possibili stati del fenomeno aleatorio, tali che  $X(\omega) \in H$ .

**Definizione 1.6. Variabile aleatoria**

Una *variabile aleatoria* su  $(\Omega, F, P)$  a valori in  $\mathbb{R}^d$  é una funzione

$$X : \Omega \rightarrow \mathbb{R}^d$$

tale che  $(X \in H) \in F$  per ogni  $H \in \mathbb{B}_d$ : scriviamo che  $X \in mF$  e diciamo che X é F-misurabile.

### Esempio 1.1. Distribuzione di Bernoulli

Sia  $p \in [0, 1]$ . La *distribuzione di Bernoulli di parametro  $p$*  si indica con  $Be_p$  ed è definita come combinazione di due delta di Dirac:

$$Be_p \approx p\delta_1 + (1-p)\delta_0$$

dove

$$\delta_a(H) = \begin{cases} 1 & \text{se } a \in H \\ 0 & \text{se } a \notin H \end{cases}$$

In particolare, per  $H \in B_d$  si ha

$$Be_p(H) = \begin{cases} 0 & \text{se } 0, 1 \notin H \\ 1 & \text{se } 0, 1 \in H \\ p & \text{se } 1 \in H \text{ e } 0 \notin H \\ 1-p & \text{se } 0 \in H \text{ e } 1 \notin H \end{cases}$$

Diciamo che una variabile aleatoria  $T$  ha distribuzione di Bernoulli,  $T \approx Be_p$  con  $p \in [0, 1]$ , se

$$P(X = i) = p^i(1-p)^{1-i} \quad \text{per } i = 0, 1.$$

In teoria della probabilità uno dei concetti più importanti relativi ad una variabile aleatoria discreta è il suo *valore atteso*, altro non è che la media dei possibili stati della variabile aleatoria pesata secondo le rispettive probabilità.

### Definizione 1.7. Valore atteso

In uno spazio di probabilità  $(\Omega, F, P)$  discreto, sia  $X$  una variabile aleatoria con distribuzione discreta finita

$$X \approx \sum_{k=1}^n p_k \delta_{x_k}$$

dove  $p_k = P(X = x_k)$  per  $k = 1, \dots, n$ , il valore atteso di  $X$  è

$$E[X] := \sum_{k=1}^n X_k p_k \tag{1.7}$$

$$\tag{1.8}$$

Nel caso in cui  $X$  sia una v.a. con distribuzione discreta non finita è sufficiente fare delle ipotesi sulla convergenza della serie.

Un risultato importante è

### Teorema 1.3. Del calcolo della media per v.a. discrete

Siano

$$X : \Omega \rightarrow \mathbb{R}^d \quad e \quad f : \mathbb{R}^d \rightarrow \mathbb{R}^N$$

rispettivamente una v.a. discreta su  $(\Omega, F, P)$  con legge  $\mu_X = \sum_{k=1}^{+\infty} p_k \delta_{X_k}$  e  $f$  una funzione  $\mathbb{B}_d$ -misurabile allora se  $\sum_{k=1}^{+\infty} |x_k(t)| < +\infty$  vale

$$E[f(X)] = \sum_{k=1}^{+\infty} f(X_k) p_k \quad (1.9)$$

$$(1.10)$$

### Definizione 1.8. Varianza

Data una v.a. discreta reale tale che  $\sum_{k=1}^{+\infty} |x_k|^2 < +\infty$ . Si definisce *varianza di X* il numero reale non negativo

$$\text{var}(X) := E[(X - E[X])^2] = E[X^2] - E[X]^2 \quad (1.11)$$

La radice quadrata della varianza  $\sqrt{\text{var}(X)}$  è detta *deviazione standard*.

Per assimilare in modo corretto le definizioni date possiamo pensare al valore atteso come la media scolastica di uno studente, e la deviazione standard come una misura della costanza dello studente, con una relazione proporzionalmente inversa, all'aumentare del valore di  $\sqrt{\text{var}(X)}$  diminuisce la perseveranza dello studente.

La positività di  $\text{var}(X)$  ci viene garantita dal seguente risultato

### Teorema 1.4. Disuguaglianza di Jensen

Siano  $-\infty < a < b < +\infty$  e

$$X : \Omega \rightarrow ]a, b[ \quad e \quad f : ]a, b[ \rightarrow \mathbb{R}$$

rispettivamente una v.a. sullo spazio  $(\Omega, F, P)$  e una funzione convessa. Se  $\sum_{k=1}^{+\infty} |X| < +\infty$  e  $\sum_{k=1}^{+\infty} |f(X)| < +\infty$  allora si ha:

$$f(E[X]) \leq E[f(X)] \quad (1.12)$$

*dimostrazione.*

Se  $f$  è convessa vale che per ogni  $z \in ]a, b[$  esiste  $m \in \mathbb{R}$  tale che

$$f(w) \geq f(z) + m(w - z), \quad w \in ]a, b[$$

Posto ora  $z = E[X]$  ( $E[X] \in ]a, b[$  per ipotesi) si ha

$$f(X(w)) \geq f(E[X]) + m(X(w) - E[X]), \quad w \in \Omega$$

da cui, prendendone il valore atteso e utilizzando la sua proprietà di monotonia ottengo

$$E[f(X)] \geq E[f(E[X]) + m(X(w) - E[X])] =$$

(utilizzando la proprietà di linearità del valore atteso e il fatto che  $E[c] = c$ )

$$= f(E[X]) + m[X - E[X]] = f(E[X]).$$

### Definizione 1.9. Covarianza

La covarianza di due variabili aleatorie è il numero reale

$$\text{cov}(X, Y) := E[(X - E[X])(Y - E[Y])]$$

Per la disuguaglianza di Cauchy-Schwartz si ha che

$$|\text{cov}(X, Y)| \leq \sqrt{\text{var}(X)\text{var}(Y)} \quad (1.13)$$

### Definizione 1.10. Correlazione

Siano  $X, Y$  due variabili aleatorie reali tali che  $\text{var}(X)$  e  $\text{var}(Y)$  siano  $> 0$ . Il coefficiente di correlazione tra  $X$  e  $Y$  è

$$\rho_{XY} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} \quad (1.14)$$

per la (1.14) abbiamo che  $\rho \in [-1, 1]$  e in particolare se  $\rho$  è negativo le v.a. si dicono *negativamente correlate*, se  $\rho$  è positivo le v.a. si dicono *positivamente correlate*. Inoltre  $|\rho| = 1$  se e soltanto se c'è una dipendenza lineare del tipo  $Y = aX + b$ .

## 1.2 Nozioni di Teoria dei Grafi

Per le nozioni di questo paragrafo si veda [7]

### Definizione 1.11. Grafo

Un grafo  $G$  è una coppia  $(S, T)$  composta da un insieme finito di elementi detti nodi (o vertici)  $S$  e da un insieme di *connessioni* (o archi)  $T$  che connettono coppie di nodi.

Due nodi connessi da un arco sono detti *adiacenti*.

### Definizione 1.12. Grafo orientato

Un grafo  $G = (S, T)$  è un grafo orientato se  $T \subset S \times S$  è formato da coppie ordinate di nodi. L'arco  $(s_i, s_j)$  è un arco da  $s_i$  a  $s_j$ , si definisce  $s_j$  come la *testa* dell'arco e  $s_i$  come la *coda* dell'arco.

Se invece la relazione tra i nodi fosse simmetrica, il grafo si direbbe *non orientato*.

### Definizione 1.13. Grafo pesato

Dati  $S, T$  e una funzione  $\omega : T \rightarrow \mathbb{R}$  che associa un valore (o *peso*) ad ogni arco, il grafo  $G = (S, T, \omega)$  è un grafo pesato.

Un grafo non pesato è un grafo pesato con  $\omega$  costante uguale ad 1.

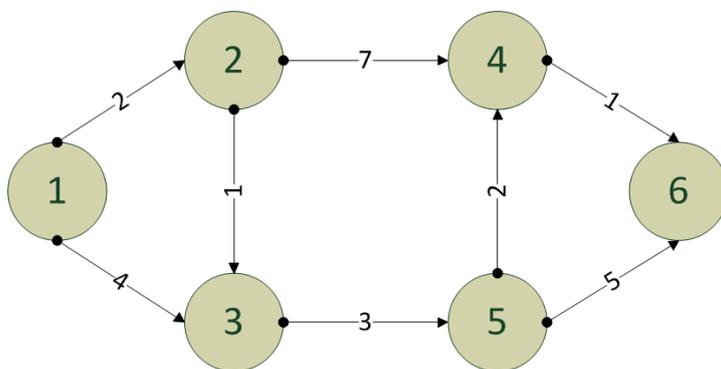


Figura 1.1: Esempio di grafo orientato pesato

### Definizione 1.14. Grafo multipartito

Un grafo  $G$  è multipartito se l'insieme dei nodi  $S$  può essere partizionato in  $k$  sottoinsiemi  $S_1, \dots, S_k$  tali che ciascun arco del grafo ha una delle due estremità in  $S_l$  e l'altra in  $S_h$  con  $l \neq h$ , cioè se  $(s_i, s_j)$  è un arco del grafo allora  $s_i$  e  $s_j$  non possono appartenere allo stesso sottoinsieme di vertici.

# Capitolo 2

## Macchine di Boltzmann

Le macchine di Boltzmann vennero introdotte per la prima volta nel 1985 da Geoffrey Hinton e da Terry Sejnowski [10]. Vogliamo arrivare a questa struttura costruendola in modo rigoroso utilizzando dei concetti matematici.

### 2.1 Neurone formale

Nel 1943 il neurofisiologo Mc Cullog e il matematico Pitts introdussero il concetto di neurone formale.

Consideriamo il seguente spazio campionario

$$\Omega = \mathbb{R}$$

con la  $\sigma$ -algebra degli eventi data da

$$F = \{\emptyset, \Omega, A, A^C\}, \quad A = \left( f\left(\sum_{i=1}^N \omega_{i,1} s_i + b\right) > m \right)$$

dove:

- $b \in \mathbb{R}$  é un valore numerico detto *bias*
- $\{s_i\}_{i=1}^N$  é una famiglia di v.a. (*ingressi*) indipendenti con distribuzione di Bernoulli che possono assumere i valori 0 o 1
- $\omega_{i,1} \in \mathbb{R}$   $i = 1, \dots, N$  é detto il *peso* relativo ad ogni  $s_i$
- $f : \mathbb{R} \rightarrow \mathbb{R}$  é chiamata *funzione di attivazione*

- $m \in \mathbb{R}$  é un *valore di soglia* fissato

Una misura di probabilità  $P$  su  $(\Omega, F)$  può essere introdotta assegnando una probabilità all'evento  $A$ .

**Definizione 2.1. Neurone formale**

Un neurone formale é una variabile aleatoria dello spazio  $(\Omega, F, P)$  tale che

$$X = \mathbb{1}_A = \begin{cases} 1 & \text{se } A \\ 0 & \text{altrimenti} \end{cases} \quad (2.1)$$

$X$  ha distribuzione di Bernoulli,  $X \approx Be_p$ .

**Osservazione 2.1.** La definizione di neurone formale prevede l'assunzione di due soli possibili stati (uscite del neurone), uno attivo e l'altro inattivo, determinati dal valore di soglia prescelto.

La probabilità di attivazione di un neurone formale cioè  $P(X = 1) = p$  é determinata dalla funzione di attivazione considerata.

## 2.2 Rete neurale stocastica

In seguito, un neurone verrà sempre inteso come neurone formale secondo la definizione data.

**Definizione 2.2. Rete neurale stocastica**

Sia  $(S, T, \omega)$  un grafo orientato pesato in cui  $S$  é un insieme di neuroni formali, sia  $b : S \rightarrow \mathbb{R}$  la funzione che associa ad ogni  $s_i$  il proprio bias  $b_i$ . Allora  $(S, T, \omega, b)$  é una rete neurale stocastica e la funzione

$$\begin{aligned} \omega : T &\rightarrow \mathbb{R} \\ (s_i, s_j) &\mapsto \omega_{ij} \end{aligned}$$

é detta *peso delle connessioni*.

Inoltre, se  $\omega_{ij} > 0$  la connessione é *eccitatoria*, se  $\omega_{ij} < 0$  la connessione é *inibitoria*.

Perció una rete neurale stocastica é costituita da un insieme di neuroni formali, ognuno con un certo numero di connessioni in ingresso e/o un certo numero di connessioni in uscita. Ciascuna di queste unità é una variabile aleatoria che dipende dalle proprie connessioni in ingresso e influenza tutte le altre unità connesse tramite le connessioni in uscita.

**Osservazione 2.2.** In una rete neurale stocastica  $(S, T, \omega)$ , le connessioni  $(s_i, s_j) \in T$  non sono neutre, infatti se consideriamo il neurone  $s_j$ , la sua funzione di attivazione riceverá come input il valore delle v.a.  $s_i \forall s_i | (s_i, s_j) \in T$  moltiplicato per il peso relativo alla connessione  $\omega_{ij}$ .

Il peso  $\omega_{ij}$  si riferisce alla connessione che va dal neurone  $s_i$  al neurone  $s_j$  e non viceversa.

### Definizione 2.3. Macchina di Boltzmann

Una rete neurale stocastica  $(S, T, \omega, b)$  che soddisfa la seguente proprietá di simmetria

$$\text{Se } (s_i, s_j) \in E \text{ e } (s_j, s_i) \in E \Rightarrow \omega_{ij} = \omega_{ji}$$

é una macchina di Boltzmann.

In una macchina di Boltzmann definiamo  $\omega_{ij}$  la covarianza tra le due variabili aleatorie  $s_i$  e  $s_j$ .

**Osservazione 2.3.** Allo stesso modo potremmo definire una macchina di Boltzmann con  $\omega : T \rightarrow [-1, 1]$  in modo che  $\omega_{ij}$  rappresenti la correlazione tra  $s_i$  e  $s_j$ .

## 2.3 Macchine di Boltzmann ristrette

### Definizione 2.4. Macchine di Boltzmann ristrette

Una macchina di Boltzmann ristretta o RBM é una macchina di Boltzmann  $(S, T, \omega, b)$  multipartita in  $k$  insiemi di neuroni  $S_1, \dots, S_k$  detti *strati* con la proprietá che  $\exists (s_i, s_j) \in T \forall s_i \in S_l \text{ e } s_j \in S_{l+1} \forall l = 1, \dots, k - 1$ . Questo vale a dire che le connessioni tra i neuroni sono consentite con tutti i neuroni degli strati adiacenti.

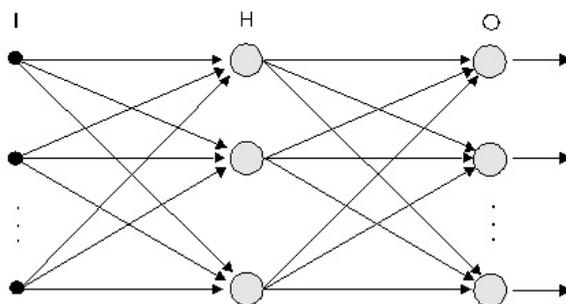


Figura 2.1: Rappresentazione di una RBM

**Osservazione 2.4.** Una macchina di Boltzmann non ha vincoli su numero di connessioni o numero di neuroni perciò la quantità di parametri da determinare durante un processo

di addestramento (definito in seguito) é molto elevata, questa puó essere ridotta imponendo dei limiti sulle connessioni (si veda ad esempio [9]), ed é per questo che abbiamo introdotto il concetto di RBM.

**Osservazione 2.5.** In seguito considereremo soltanto RBM a due strati (bipartite) in cui

- $V = \{0, 1\}^D$  é lo spazio dei *vettori visibili* detto *strato visibile*
- $H = \{0, 1\}^K$  é lo spazio dei *vettori nascosti* detto *strato nascosto*

**Notazione 2.1.** Indicheremo con  $v \in V$  il vettore che ha come componenti gli stati dei neuroni visibili (v.a. dello strato visibile) e con  $h \in H$  il vettore che ha come componenti gli stati dei neuroni nascosti (v.a. dello strato nascosto).

**Definizione 2.5. Energia di una configurazione**

Sia  $(S, T, \omega, b)$  una RBM a due strati, l'energia di una configurazione congiunta dei vettori  $v$  ed  $h$  é definita come

$$E(v, h) = v^T b + h^T c + v^T W h$$

dove

- $b \in \mathbb{R}^D$  é il vettore dei bias relativi ai neuroni dello strato visibile
- $c \in \mathbb{R}^K$  é il vettore dei bias relativi ai neuroni dello strato nascosto
- $W \in \mathbb{R}^{D \times K}$  é la matrice dei pesi  $\omega_{ij}$  che mettono in relazione l' $i$ -esimo neurone dello strato visibile con il  $j$ -esimo neurone dello strato nascosto.

**Osservazione 2.6.** Una RBM definisce la *distribuzione di Boltzmann*, in cui la probabilità di una certa configurazione congiunta del vettore visibile e di quello nascosto é determinata solamente dall'energia della configurazione considerata rispetto all'energia di tutte le configurazioni possibili

$$p(v, h) = \frac{e^{E(v, h)}}{\sum_{v \in V} \sum_{h \in H} e^{E(v, h)}}$$

**Prop 2.1.** Per una RBM a due strati valgono le seguenti relazioni

$$P(h|v) = \prod_{i=1}^K P(h_i|v) \quad P(v|h) = \prod_{j=1}^D P(v_j|h)$$

*dimostrazione.*

É possibile dimostrare che se  $X, Y$  sono due v.a. con distribuzione rispettivamente  $\mu_X$  e  $\mu_Y$ , il vettore aleatorio  $(X, Y)$  ha distribuzione  $\mu_{(X, Y)} = \mu_X \mu_Y \Leftrightarrow X$  e  $Y$  sono indipendenti.

In una RBM i neuroni dello stesso strato non sono connessi tra loro cioè sono indipendenti. Infatti se  $s_i, s_j \in V$  si ha  $\omega_{ij} = cov(s_i, s_j) = 0 \forall i \neq j$ . Vista l'indipendenza

$$P(h|v) = \frac{P(h \cap v)}{P(v)} = \prod_{i=1}^K \frac{P(h_i \cap v)}{P(v)} = \prod_{i=1}^K P(h_i|v)$$

e analogamente

$$P(v|h) = \frac{P(v \cap h)}{P(h)} = \prod_{i=1}^K \frac{P(v_i \cap h)}{P(h)} = \prod_{i=1}^K P(v_i|h)$$

## 2.4 Metodo della massima verosimiglianza

### Definizione 2.6. Processo di addestramento

Per processo di addestramento di una rete neurale stocastica  $(S, T, \omega, b)$  si intende un procedimento per la determinazione dei pesi  $\omega_{ij}$  e dei bias relativi ai neuroni in modo da rendere noti i parametri della distribuzione di probabilità che la caratterizza.

Sia  $x = (x_1, \dots, x_n)$  un *campione* costituito da  $n$  eventi indipendenti, la cui probabilità dipende da un vettore di parametri  $\theta$ .

### Definizione 2.7. Funzione di verosimiglianza

Dato  $x$  campione, la funzione

$$L(x, \theta) = \prod_{i=1}^n P_{\theta}(x_i)$$

è la funzione di verosimiglianza associata al campione.

**Osservazione 2.7.** La funzione di verosimiglianza associata ad un campione rappresenta la funzione di probabilità  $P_{\theta}(x)$  del campione stesso, vista l'indipendenza degli eventi che lo costituiscono.

Il metodo della massima verosimiglianza come processo di addestramento di una rete neurale stocastica consiste nella determinazione del vettore dei parametri  $\theta$  relativi ad una distribuzione di probabilità che massimizzi la funzione di verosimiglianza.

### Definizione 2.8. Stimatore di massima verosimiglianza

Data una funzione di verosimiglianza  $L(x, \theta)$

$$\theta^* = \arg \max L(x, \theta)$$

è detto stimatore di massima verosimiglianza.

**Osservazione 2.8.** Per facilitarne il calcolo si introduce la funzione di *log-verosimiglianza*

$$l(x, \theta) := \log L(x, \theta)$$

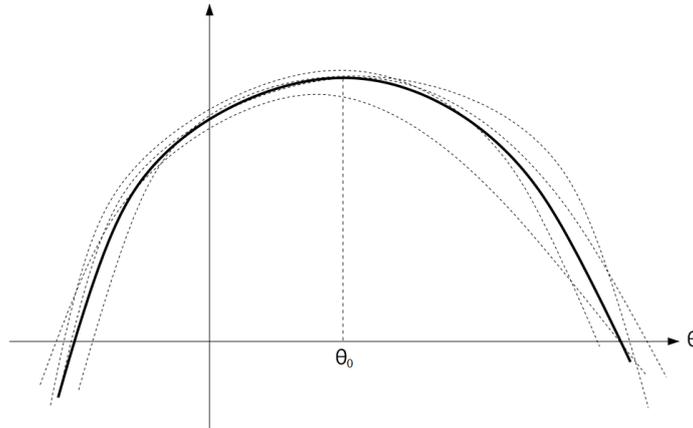
Il logaritmo é una funzione monotona crescente che non va a modificare le caratteristiche della funzione  $L(x, \theta)$  in termini di crescita e decrescenza, ma soprattutto si ottiene una forma piú semplice da trattare. In effetti, nel caso di variabili i.i.d, come nel nostro campione, si ottiene

$$l(x, \theta) = \log \prod_{i=1}^n P_{\theta}(x_i) = \sum_{i=1}^n \log P_{\theta}(x_i)$$

**Teorema 2.1. del valore atteso della log-verosimiglianza**

La funzione  $E[l(x, \theta)]$  assume un massimo in corrispondenza di  $\theta_0$  cioè del vero valore della funzione di probabilità di  $x$ .

La figura seguente rappresenta l'andamento della funzione di log-verosimiglianza dei campioni (linee tratteggiate) relativamente a quello di  $E[l(x, \theta)]$  (linea continua).



*dimostrazione.*

Consideriamo  $A[\theta] = E \left[ \frac{L(x, \theta)}{L(x, \theta_0)} \right]$

Ricordiamo che stiamo considerando il caso discreto e che  $L(x, \theta)$  é una funzione di probabilità di  $x$  per ogni valore di  $\theta$  vale

$$A[\theta] = \sum_x \frac{L(x, \theta)}{L(x, \theta_0)} L(x, \theta_0) = \sum_x L(x, \theta) = 1$$

Da ciò segue che  $\log A[\theta] = 0$ . Siccome la funzione  $\log$  é concava, possiamo applicare la disuguaglianza di Jensen (1.14)

$$E \left[ \log \frac{L(x, \theta)}{L(x, \theta_0)} \right] \leq \log E \left[ \frac{L(x, \theta)}{L(x, \theta_0)} \right] = 0$$

$$E[l(x, \theta) - l(x, \theta_0)] \leq 0$$

Utilizzando la linearità del valore atteso

$$E[l(x, \theta)] \leq E[l(x, \theta_0)]$$

Questa relazione è vera per ogni  $\theta$ !

**Notazione 2.2.** In seguito, per provare la consistenza del metodo della massima verosimiglianza, indicheremo la funzione di log-verosimiglianza con  $l_n(x, \theta)$  e il suo stimatore di massima verosimiglianza con  $\theta_n^*$  per esplicitare la dimensione  $n$  del campione.

**Osservazione 2.9.** La relazione derivata nel teorema precedente ci mette in evidenza che se  $\theta \neq \theta_0$ , per ogni  $\epsilon > 0$ , per ogni  $\delta \in ]0, 1[ \exists n_0 \in \mathbb{N}$  tale che:

$$P(l_n(x, \theta_0) - l_n(x, \theta_j) > \epsilon \text{ per ogni } n > n_0) \geq 1 - \delta$$

perché quando  $n$  diverge la funzione di verosimiglianza in  $\theta_0$  tende a diventare molto più grande di ogni altro valore di  $\theta$ .

**Definizione 2.9. Stimatore consistente**

Sia  $\theta_n^*$  uno stimatore per la funzione di log-verosimiglianza e  $\theta_0$  il valore di  $\theta$  corrispondente alla vera distribuzione di probabilità del campione  $x$ . Diremo che  $\theta_n^*$  è uno stimatore consistente del parametro se

$$\lim_{n \rightarrow \infty} P(|\theta_n^* - \theta_0| > \epsilon) = 0 \quad \text{per ogni } \epsilon > 0$$

**Teorema 2.2. Consistenza dello stimatore di massima verosimiglianza**

Data  $l_n(x, \theta)$  funzione di log-verosimiglianza del campione  $x$ , lo stimatore di massima verosimiglianza  $\theta^*$  è uno stimatore consistente del parametro  $\theta$ .

*dimostrazione.*

Per la dimostrazione ci restringiamo al caso in cui  $\theta$  possa assumere un numero finito di valori  $\theta_j$  con  $j = 0, \dots, m$  (vale analogamente anche nel caso in cui  $\theta \in \mathbb{R}$ ).

Definiamo

$$A_j = \{l_n(x, \theta_0) - l_n(x, \theta_j) > \epsilon \text{ per ogni } n > n_0\} \quad \text{per } j = 1, \dots, m.$$

la cui probabilità, per quanto visto nel teorema (2.1), può essere resa maggiore di  $1 - \delta$ , scegliendo un certo  $n_0$  sufficientemente elevato. Perciò

$$\begin{aligned}
 P\left(\bigcap_{j=1}^m A_j\right) &= 1 - P\left(\bigcap_{j=1}^m A_j^C\right) \\
 &\geq 1 - \sum_{j=1}^m P(A_j^C) \\
 &\geq 1 - m(1 - (1 - \delta)) \\
 &= 1 - m\delta
 \end{aligned}
 \tag{2.2}$$

e così abbiamo mostrato che

$$P(l_n(x, \theta_0) - l_n(x, \theta_j) > \epsilon \text{ per ogni } j \neq 0 \text{ per ogni } n > n_0) \geq 1 - m\delta$$

con  $\delta$  arbitrario. Ciò è un risultato ancora più forte della consistenza definita per uno stimatore!

**Osservazione 2.10.** Per massimizzare la funzione di log-verosimiglianza è necessario porre

$$\frac{\partial l(x, \theta)}{\partial \theta} = 0$$

Naturalmente, se  $\theta$  fosse un vettore di parametri, la formula qui sopra verrà intesa come derivata rispetto ad ogni componente del vettore considerando costanti le altre.

Se il campione è composto da  $n$  eventi indipendenti possiamo porre

$$\frac{\partial l(x, \theta)}{\partial \theta} = \sum_{i=1}^n \frac{\partial \log P_\theta(x_i)}{\partial \theta} = 0$$

### Esempio 2.1. Stimatore di una distribuzione di Bernoulli

Diciamo che una variabile aleatoria  $X$  ha distribuzione di Bernoulli se

$$P(X = i) = p^i (1 - p)^{1-i} \quad \text{per } i = 0, 1.$$

Dato un campione costituito da una famiglia  $\{x_i\}_{i=1}^n$  di eventi indipendenti, uno stimatore per il parametro  $p$  incognito (in questo caso  $\theta = p$ ) possiamo trovarlo applicando il metodo della massima verosimiglianza. La funzione di verosimiglianza è

$$L(x, p) = \prod_{i=1}^n p^{x_i} (1 - p)^{1-x_i} = p^\xi (1 - p)^{n-\xi}$$

dove  $\xi = \sum_{i=1}^n x_i$

La funzione di log-verosimiglianza é data da

$$l(x, p) = \log p \sum_{i=1}^n x_i + \log(1-p) \sum_{i=1}^n (1-x_i)$$

$$l(x, p) = \xi \log p + (n - \xi) \log 1 - p$$

Derivando la funzione di log-verosimiglianza

$$\frac{\partial l(x, p)}{\partial p} = \frac{\xi}{p} - \frac{n - \xi}{1 - p}$$

e ponendo

$$\frac{\partial l(x, p)}{\partial p} = \frac{\xi}{p} - \frac{n - \xi}{1 - p} = 0$$

otteniamo uno stimatore per p

$$p^* = \frac{1}{n} \xi = \frac{1}{n} \sum_{i=1}^n x_i$$

ció lo stimatore p per una variabile con distribuzione di Bernoulli é la media dei dati del campione.

## 2.5 Addestramento di una RBM

Il caso piú semplice che possiamo trattare é quello in cui la nostra RBM bipartita sia costituita da due soli neuroni che hanno dimensione 1 cioé due variabili aleatorie che possono assumere il valore 0 o 1. Indichiamo con b il bias relativo al neurone visibile, con c il bias relativo al neurone nascosto e siano inoltre  $V = H = \{0, 1\}$ .

In questo caso il campione é rappresentato da

$$\phi_1 = (v^{(1)}, h^{(1)})$$

$$\phi_2 = (v^{(2)}, h^{(2)})$$

...

$$\phi_N = (v^{(N)}, h^{(N)})$$

La distribuzione di probabilitá é

$$P_{\theta}(v, h) = \frac{e^{bv+ch+\omega vh}}{\sum_{v \in V} \sum_{h \in H} e^{bv+ch+\omega vh}}$$

dove  $\omega$  é il paramentro di interazione tra i due neuroni e  $\theta = (b, c, \omega)$  é il vettore dei parametri da determinare durante il processo di addestramento.

Costruiamo la funzione di log-verosimiglianza per calcolare b, c e W

$$l(b, c, \omega) = \log \prod_{i=1}^N P_{\theta}(v^{(i)}, h^{(i)})$$

Per comoditá pongo  $Z = \sum_{v \in V} \sum_{h \in H} e^{bv+ch+\omega vh}$

$$\begin{aligned} l(b, c, \omega) &= \log \prod_{i=1}^N \frac{e^{bv^{(i)}+ch^{(i)}+\omega v^{(i)}h^{(i)}}}{Z} \\ &= \log \frac{e^{\sum_{i=1}^N bv^{(i)}+ch^{(i)}+\omega v^{(i)}h^{(i)}}}{Z^N} \\ &= \log e^{\sum_{i=1}^N bv^{(i)}+ch^{(i)}+\omega v^{(i)}h^{(i)}} - \log Z^N \\ &= \sum_{i=1}^N bv^{(i)} + ch^{(i)} + \omega v^{(i)}h^{(i)} - N \log Z \end{aligned}$$

(2.3)

Cerchiamo b, c,  $\omega$  che massimizzino la funzione di log-verosimiglianza, perció poniamo

$$\frac{\partial l}{\partial b} = 0, \quad \frac{\partial l}{\partial c} = 0, \quad \frac{\partial l}{\partial \omega} = 0.$$

Svolgendo i calcoli

$$\frac{\partial l}{\partial b} = \sum_{i=1}^N v^{(i)} - N \frac{\sum_{v \in V} \sum_{h \in H} v e^{bv+ch+\omega vh}}{Z} = 0$$

$$\frac{\sum_{v \in V} \sum_{h \in H} v e^{bv+ch+\omega vh}}{Z} = \frac{\sum_{i=1}^N v^{(i)}}{N}$$

$$\sum_{v \in V} \sum_{h \in H} p(v, h) v = \frac{\sum_{i=1}^N v^{(i)}}{N} \quad (2.4)$$

Analogamente rispetto c e  $\omega$  si ha

$$\sum_{v \in V} \sum_{h \in H} p(v, h) h = \frac{\sum_{i=1}^N h^{(i)}}{N} \quad (2.5)$$

$$\sum_{v \in V} \sum_{h \in H} p(v, h)vh = \frac{\sum_{i=1}^N v^{(i)}h^{(i)}}{N} \quad (2.6)$$

Si noti che le espressioni in (2.3), (2.4) e (2.5) sono equivalenze tra il valore atteso della distribuzione congiunta e la media aritmetica dei dati del campione.

# Capitolo 3

## Algoritmo

In questo capitolo verrà presentato un algoritmo di machine learning in Python per la classificazione di oggetti. Dato un fiore di iris in input vogliamo ottenere la sua categoria di appartenenza in output.

**Primo passo:** è necessario scaricare tutto quello di cui avremo bisogno per eseguire l'algoritmo.

Partiamo con il download di *Anaconda*, una piattaforma per il Data Science con Python (linguaggio di programmazione), una volta completato il download dobbiamo installare nel seguente ordine:

- *Numpy*, una libreria opensource che contiene funzioni matematiche utili per lavorare con matrici e vettori multidimensionali;
- *Sklearn*, una libreria opensource di apprendimento automatico che contiene algoritmi di classificazione e machine learning;
- *Tensorflow*, una libreria opensource che fornisce modelli sperimentati e ottimizzati per il machine learning.

**Secondo passo:** per il processo di addestramento di un modello occorre fornire all'algoritmo di apprendimento i dati di addestramento da cui possa apprendere.

---

```
1     from sklearn import datasets
2     iris_dataset = datasets.load_iris()
3     print(iris_dataset['DESCR'])
```

---



Il dataset di iris é un insieme di addestramento composto da esempi che misurano la lunghezza e la larghezza del petalo e del sepal di Iris-Setosa, Iris-Versicolor e Iris-Virginica. L'ultimo comando ci fornisce la descrizione del dataset di Iris:

```
:Number of Instances: 150 (50 in each of three classes)
:Number of Attributes: 4 numeric, predictive attributes and the class
:Attribute Information:
- sepal length in cm
- sepal width in cm
- petal length in cm
- petal width in cm
- class:
  - Iris-Setosa
  - Iris-Versicolour
  - Iris-Virginica
```

```
:Summary Statistics:
```

```
=====  =====  =====  =====  =====
                Min  Max   Mean   SD   Class Correlation
=====  =====  =====  =====  =====
sepal length:   4.3  7.9   5.84   0.83   0.7826
sepal width:    2.0  4.4   3.05   0.43  -0.4194
petal length:   1.0  6.9   3.76   1.76   0.9490 (high!)
petal width:    0.1  2.5   1.20   0.76   0.9565 (high!)
=====  =====  =====  =====  =====
```

```
:Class Distribution: 33.3\% for each of 3 classes.
:Creator: R.A. Fisher
:Donor: Michael Marshall (July, 1988)
```



Il parametro `test_size` in `train_test_split` indica la percentuale di esempi per l'insieme di test (30%).

**Passo 3:** consiste nel scegliere e addestrare il modello sui dati di training:

---

```
8     from sklearn.linear_model import Perceptron
9     model= Perceptron(max_iter=100)
10    model.fit(X_train, y_train)
```

---

Perceptron (perceptrone) é uno dei primi algoritmi di classificazione della storia, a livello strutturale é una macchina di Boltzmann ristretta.

L'attributo `max_iter` fissa il numero massimo di iterazioni, altri parametri sono predefiniti.

Arrivati a questo punto il modello é addestrato, dobbiamo capire se i dati scelti sono consistenti con il metodo (ovvero permettono "di generalizzare").

**Passo 4:** facciamo predire al modello sui dati di training che ha già appreso (per vedere se restituisce gli stessi output che gli abbiamo dato) e sui dati di test (per verificare la capacità di classificare un oggetto del quale non ha precedentemente appreso la classe di appartenenza).

---

```
11    predicted_train = model.predict(X_train)
12    predicted_test  = model.predict(X_test)
```

---

Siamo in grado di ottenere l'accuratezza del modello, sui dati di training e sui dati di test:

---

```
13    from sklearn.metrics import accuracy_score
14    print('Train accuracy')
15    print( accuracy_score(y_train, predicted_train) )
16    print('Test accuracy')
17    print( accuracy_score(y_test, predicted_test) )
```

---

Facendo un primo tentativo:

```
Train accuracy
0.9047619047619048
Test accuracy
0.8444444444444444
```

cioé il modello ha un'accuratezza del 90,5% sui dati di training e del 84,4% sui dati di test.

Esistono molti altri algoritmi per il machine learning per la classificazione, anche molto piú efficienti di quello visto.

Nello specifico con la libreria tensorflow é possibile addestrare modelli RBM scegliendo il numero di strati nascosti e per ognuno di questi il numero di neuroni:

---

```
1 import tensorflow as tf
2 columns = [tf.contrib.layers.real_valued_column("", dimension=4)]
3 classifier = tf.contrib.learn.DNNClassifier(feature_columns=columns,
4       hidden_units=[10, 20, 10], n_classes=3,
5       model_dir=iris_dataset)
6 classifier.fit(iris_dataset.data, iris_dataset.target)
```

---

# Bibliografia

- [1] Alpaydin, E. (2020). *Introduction to machine learning*. MIT press.
- [2] Domingos, P. (2012). *A few useful things to know about machine learning*. Communications of the ACM, 55(10), 78-87.
- [3] Benjamin, A., Chartrand, G., Zhang, P. (2017). *The fascinating world of graph theory*. Princeton University Press.
- [4] Rajkumar, A., Dean, J., Kohane, I. (2019). *Machine learning in medicine*. New England Journal of Medicine, 380(14), 1347-1358.
- [5] Li, Q., Zheng, N., Cheng, H. (2004). Springrobot: A prototype autonomous vehicle and its algorithms for lane detection. *IEEE Transactions on Intelligent Transportation Systems*, 5(4), 300-308.
- [6] Minsky, M., Papert, S. A. (2017). *Perceptrons: An introduction to computational geometry*. MIT press.
- [7] West, D. B. (2001). *Introduction to graph theory* (Vol. 2). Upper Saddle River: Prentice hall.
- [8] Riedmiller, M. (1994). Advanced supervised learning in multi-layer perceptrons—from backpropagation to adaptive learning algorithms. *Computer Standards Interfaces*, 16(3), 265-278.
- [9] Larochelle, H., Mandel, M., Pascanu, R., Bengio, Y. (2012). *Learning algorithms for the classification restricted boltzmann machine*. Journal of Machine Learning Research, 13(Mar), 643-669.
- [10] Hinton, G. E., Sejnowski, T. J., Ackley, D. H. (1984). *Boltzmann machines: Constraint satisfaction networks that learn*. Pittsburgh: Carnegie-Mellon University, Department of Computer Science.
- [11] Pascucci, Andrea. (2020). *Teoria della probabilità*. Springer-Verlag Mailand.

- [12] Minini, Andrea. <http://www.andreaminini.com/python/sklearn/come-fare-machine-learning-con-scikit-learn-di-python>.
- [13] Minini, Andrea. <http://www.andreaminini.com/ai/tensorflow/esempio-tutorial-tensorflow>.