

Indice

PICCIALLI	4
1 Classificazione supervisionata	4
2 Classificazione non supervisionata	4
3 Regressione	4
4 Problemi di ottimizzazione continua definiti su \mathbb{R}^n	4
5 Condizioni di Minimo	5
5.1 Punto di minimo globale	5
5.2 Punto di minimo locale	5
5.3 Punto di minimo locale isolato	5
5.4 Convessità	5
Funzione convessa	5
Insieme convessi	6
Corollario	6
Corollario	7
5.5 Concavità	7
Funzione concava	7
5.6 Applicazioni convessità	7
5.7 Derivata direzionale di f in x lungo direzione d	9
5.7.1 Sviluppo di Taylor I ordine	10
5.7.2 Sviluppo Taylor II ordine	10
5.8 Segno di una matrice	10
5.8.1 Matrice Hessiana	11
5.8.2 Criterio dei minori di Nord-Ovest	11
5.8.3 Segno di una matrice - Autovalori/Autovettori	12
6 Condizioni di esistenza delle soluzioni	13
6.1 Definizioni preliminari	13
6.2 Esistenza del minimo	13
7 Condizioni di ottimo per problemi non vincolati	14
7.0.1 Direzione di discesa	15
7.0.2 Direzioni a curvatura negativa in x	16
7.0.3 Direzioni a curvatura negativa in x	18
8 Condizioni di ottimo per problemi vincolati con insieme ammissibile	18
8.1 Condizioni di ottimo per problemi vincolati con insieme ammissibile convesso	20
8.2 Problemi con vincoli lineari	21
8.3 Condizioni di ottimalità per problemi con vincoli di box	22
9 Condizioni di Ottimo per problemi vincolati	24
9.1 Condizioni di ottimo analitiche	24
9.1.1 Lagrangiana generalizzata	24
9.1.2 Condizioni necessarie di Fritz John	24
9.2 Condizioni di Qualificazione dei Vincoli o Condizioni di Regolarità	24
9.2.1 Condizioni di qualificazione dei vincoli	25
9.2.2 Condizioni necessarie di KKT (Karush-Kuhn-Tucker)	26
9.2.3 Moltiplicatori di Lagrange	26
9.3 Problemi con vincoli di box	27
9.4 Condizioni Sufficienti KKT	28

10 Cenni sulla Teoria Statistica	30
10.1 Teoria di Vapinil-Chervonenkis	31
11 Iperpiani Orientati	32
11.1 Iperpiani con gap di tolleranza	33
11.1.1 Variazione di VC dimension	33
12 Teoria della Dualità	41
12.1 Introduzione	41
12.2 Dualità di Wolfe	42
12.3 Programmazione Quadratica - Duale	44
13 SVM Lineari	46
14 SVM Non Lineari	50
14.1 Kernel trick	50
14.2 Funzioni Kernel	51
14.3 Kernel Polinomiale	52
14.4 Kernel Gaussiano	52
15 Working Set	52
16 SMO	58
16.1 Criterio di Arresto	60
17 Riepilogo SVM	61
18 ν-Classification	62
19 Problemi di Regressione con SVM	63
19.1 Regressione Logistica	63
19.2 Come si addestra la regressione logistica?	64
19.3 Valutazione di un classificatore in probabilità	65
19.4 ROC CURVE (receiver operating characteristic)	66
20 Clustering	66
20.1 Algoritmo di K-Means	70
20.2 Scelta del numero di Cluster	71
PACIFICI	72
1 Alberi decisionali per la classificazione	72
1 Introduzione	72
Apprendimento statistico	72
Apprendimento supervisionato	73
Apprendimento non supervisionato	73
Response	73
2 Alberi Decisionali	74
Pro vs Contro alberi di decisione	74
TDITD Family	75
Algoritmo di Hunt	75
Algoritmo ID3	75

3 Errori negli alberi di classificazione	76
Indici di (in)purity	76
Test	76
4 Albero Ottimo di Classificazione	77
Introduzione	77
Decisioni del problema	78
Struttura del problema	78
Variabili	79
Vincoli	79
Allocazione	79
Variabili	79
Vincoli	79
Scelta M_1 e M_2	80
Funzione Obiettivo	80
Dati	81
Task - obiettivo	81
Riepilogo	82
Warm Starts	83
5 Addestramento di un OCT	83
6 OCT: Modelli Multivariati	85
6.1 Introduzione	85
6.2 Formulazione di OCT-H	85
6.3 Warm Start OCT-H	86
6.4 Tuning degli Iperparametri	87
7 Alberi di Decisione Pro e Contro	87
8 Bagging, Random Forests, Boosting	87
8.1 Bagging	87
8.2 Random Forests	88
8.3 Boosting	88
8.4 Adaboost ST	88
8.5 Errore di classificazione al passo m-esimo	89
8.6 Scelta del Weak Learner al passo m	89
8.7 Scelta del parametro α_m	89
2 Classificazione Robusta	90
1 Introduzione	90
2 Robust Classification	90
3 Approccio MaxMin	90
4 Uncertainty Set	91
5 Modello MIP (Programmazione Intera Mista)	91
5.1 Vincoli Strutturali	92
5.2 Allocazione Punti nelle Foglie	92
5.3 Funzione Obiettivo	92
6 Robustezza vs. Incertezza delle feature	93

7 Robustezza vs. Incertezza nelle label	93
7.1 Dualita PL	93
7.2 Robustezza nelle label	94
7.3 Proprietà	95

PICCIALLI

1 Classificazione supervisionata

Nella **classificazione supervisionata** sono noti a priori, tra i dati di training e di set i pattern che rappresentano le diverse classi e si vuole determinare un modello matematica che definisca la classe di appartenenza.

2 Classificazione non supervisionata

Nella **classificazione non supervisionata** non sono noti a priori i pattern rappresentativi delle classi. Si vuole determinare il numero di classi di similitudine e un modello matematico che dato un pattern dello spazio delle caratteristiche, definisca la corrispondente classe di appartenenza.

3 Regressione

Nella **regressione** sono note coppie di pattern/target rappresentative di una funzione incognita a valori reali. Si vuole determinare una funzione analitica che approssima la funzione incognita.

4 Problemi di ottimizzazione continua definiti su \mathbb{R}^n

Nei problemi continui a dimensione finita, lo spazio delle variabile o **spazio delle feature**, è uno spazio lineare in \mathbb{R}^n . In particolare si ha un problema del tipo:

$$\min f(x)$$

$$x \in S \subseteq \mathbb{R}^n$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R} := \text{funzione obiettivo}$$

$$S := \text{insieme ammissibile}$$

Questi problemi possono essere classificate in base alla struttura dell'insieme ammissibile e alle ipotesi sulla funzione obiettivo:

- Se $S = \mathbb{R}^n$ si parla di **ottimizzazione non vincolata**;
- Se $S \subset \mathbb{R}^n$ si parla di **ottimizzazione vincolata**. Il problema più comune di questa categoria è quello in cui S è descritto attraverso un **insieme finito di vincoli di uguaglianza e disequaglianza**:

$$S = \{x \in \mathbb{R}^n: g(x) \leq 0, h(x) = 0\}$$

In cui: $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$ e $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$ vettori di funzioni assegnate.

5 Condizioni di Minimo

5.1 Punto di minimo globale

Un punto $x^* \in S$ è detto punto di **minimo globale** di f su S se:

$$f(x^*) \leq f(x) \quad \forall x \in S$$

e il valore $f(x^*)$ è il minimo globale di f su S .

Si dice che $x^* \in S$ è un punto di **minimo globale stretto** di f su S se:

$$f(x^*) < f(x) \quad \forall x \in S, x \neq x^*$$

5.2 Punto di minimo locale

Un punto $x^* \in S$ è detto di **minimo locale** di f su S se esiste un intorno $B(x^*, \rho)$ con $\rho > 0$ tale che:

$$f(x^*) \leq f(x) \quad \forall x \in S \cap B(x^*, \rho)$$

Un punto $x^* \in S$ è detto di **minimo locale stretto** se esiste un intorno $B(x^*, \rho)$ con $\rho > 0$ tale che:

$$f(x^*) < f(x) \quad \forall x \in S \cap B(x^*, \rho)$$

5.3 Punto di minimo locale isolato

Dato un punto x^* è un punto di **minimo locale isolato** quando è un punto di minimo locale, ma l'intorno non è totalmente contenuto in S .

5.4 Convessità

Funzione convessa

Un funzione $f(x)$ è **convessa** nell'insieme S se:

$$\forall x, y \in S \quad \forall \lambda \in [0, 1] \rightarrow f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

In cui:

- $f(\lambda x + (1 - \lambda)y) :=$ segmento che connette x, y ;

- $\lambda f(x) + (1 - \lambda)f(y) := \text{segmento che connette } f(x) \text{ e } f(y)$

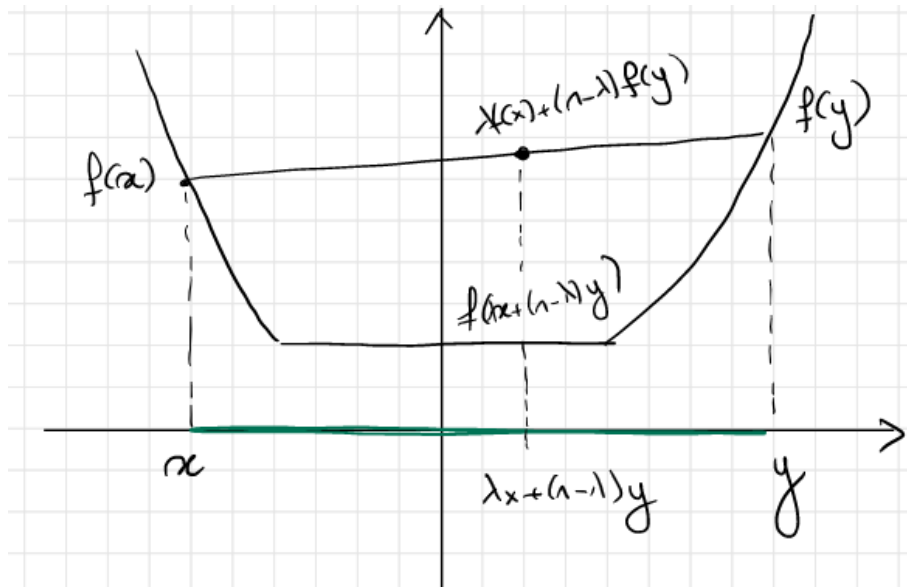


Figura 1.

Inoltre, $f(x)$ è **strettamente convessa** su S se $\forall x, y \in S, \forall \lambda \in [0, 1]$:

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$$

Oss.:

Convessità \nRightarrow Esistenza di minimo

es. $f(x) = e^x$, è strettamente convessa su \mathbb{R} , ma non ammette minimo globale.

Insieme convessi

S è un **insieme convesso** se $\forall x, y \in S \rightarrow \lambda x + (1 - \lambda)y \in S \quad \forall \lambda \in [0, 1]$

Corollario

Se $g(x)$ è convesso $\rightarrow S: \{x \in \mathbb{R}: g(x) \leq 0\}$ è convesso

Dimostrazione

Occorre dimostrare che $\forall x, y \in S \quad \lambda x + (1 - \lambda)y \in S \quad \forall \lambda \in [0, 1]$

Quindi:

$$x, y \in S \rightarrow g(x) \leq 0, g(y) \leq 0 \rightarrow \lambda x + (1 - \lambda)y \in S \wedge g(\lambda x + (1 - \lambda)y) \leq 0$$

Poiché g è convessa $\rightarrow g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y)$:

- $(1 - \lambda) \geq 0$ poiché $\lambda \in [0, 1]$
- $g(y), g(x) \leq 0$ per ipotesi;

Quindi:

$$g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y) \leq 0 \rightarrow \lambda x + (1 - \lambda)y \in S \rightarrow \textbf{Insieme convesso}$$

Corollario

Se S_1, \dots, S_n sono convessi $\rightarrow S_1 \cap \dots \cap S_n$ è convesso.

Dimostrazione

$$x, y \in S_1 \iff x, y \in S$$

Se S_i convesso $\rightarrow \lambda x + (1 - \lambda)y \in S_i \quad \forall \lambda \in [0, 1] \quad \forall i \rightarrow \lambda x + (1 - \lambda)y \in S_i \quad i = 1, \dots, n$
 $\rightarrow \lambda x + (1 - \lambda)y \in S \rightarrow S$ è convesso

Inoltre:

- Se $S = \{x \in \mathbb{R}^n: g_i(x) \leq 0 \quad i = 1, 2, \dots, n\}$ e g_i convessa $i = 1, \dots, n \Rightarrow S$ è convessa
- Se $S = \{h(x) = 0\}$ è convesso $\iff h(x)$ è lineare

5.5 Concavità

Funzione concava

Una funzione $f: \mathbb{R}^n \rightarrow \mathbb{R}$ è **concava** su S se $-f(x)$ è convessa su S . In altre parole, se e solo se:

$$f(\lambda x + (1 - \lambda)y) \geq \lambda f(x) + (1 - \lambda)f(y) \quad \forall x, y \in S \quad \forall \lambda \in [0, 1]$$

Una funzione $f: \mathbb{R}^n \rightarrow \mathbb{R}$ è **strettamente concava** se $-f(x)$ è strettamente convessa. In altre parole, se e solo se:

$$f(\lambda x + (1 - \lambda)y) > \lambda f(x) + (1 - \lambda)f(y) \quad \forall x, y \in S \quad \forall \lambda \in [0, 1]$$

Una funzione lineare è sia concava che convessa.

5.6 Applicazioni convessità

Ricordiamo che il problema che vogliamo risolvere è del tipo:

$$\begin{cases} \min f(x) \\ x \in S \end{cases}$$

Ipotesi:

1. f convessa su S : $\forall x, y \in S \quad f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad \forall \lambda \in [0, 1]$
(strettamente convesso $\rightarrow f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$)
2. S convesso $\forall x, y \in S \quad \lambda x + (1 - \lambda)y \in S \quad \forall \lambda \in [0, 1]$

Proposizione 1. se $\begin{cases} \min_{x \in S} f(x) \end{cases}$, f convessa su S e S convessa. Allora valgono le seguenti proposizioni:

1. Non \exists minimi locali non globali; [Non si assicura l' \exists della soluzione];
2. L'insieme delle soluzioni ottime è un insieme convesso;

Dimostrazione. (1)

Supponiamo che $x^* \in S$ è un minimo locale $\rightarrow \exists \rho > 0$ t.c. $f(x^*) \leq f(x) \quad \forall x \in S \cap B(\rho, x^*)$

Considero un punto $x \in S, x \notin B(x^*, \rho)$ e quest'ultimo non deve essere un punto di minimo globale:

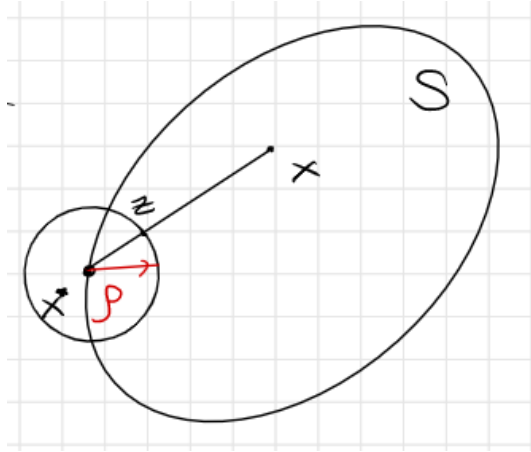


Figura 2.

In questa figura abbiamo rappresentato il nostro insieme S convesso e il punto di minimo globale x^* con il suo corrispettivo intorno. Inoltre, punto x soddisfa l'ipotesi cui sopra.

Poiché S è convesso \rightarrow il segmento $[x^*, x]$ è contenuto in S . In particolare è tale che:

$$\lambda x + (1 - \lambda)x^* \in S \quad \forall \lambda \in [0, 1]$$

Il punto z è quel punto che si trova nell'intorno di x^* e sul segmento che connette $[x^*, x]$. Allora:

$$z = \bar{\lambda}x + (1 - \bar{\lambda})x^* \in S \cap B(\rho, x^*)$$

Inoltre poiché x^* è un punto di minimo:

$$f(x^*) \leq f(z) \quad \text{dato che } z \in B(\rho, x^*)$$

$$f(x^*) \leq f(\bar{\lambda}x + (1 - \bar{\lambda})x^*) \leq \bar{\lambda}f(x) + (1 - \bar{\lambda})f(x^*) \text{ poiché } f \text{ convessa}$$

Quindi:

$$f(x^*) \leq \bar{\lambda}f(x) + f(x^*) - \bar{\lambda}f(x^*)$$

$$f(x^*) \leq f(x) \quad \forall x \in S$$

Quindi abbiamo dimostrato che non se il problema è di minimizzazione, la funzione e l'insieme sono convessi, allora non esiste un minimo locale che non è globale.

Inoltre, se considerassimo la convessità stretta otterremmo che x^* è l'**unico** punto di minimo globale. \square

Dimostrazione. (2)

Se l'insieme S è convesso, definiamo $X^* := \{\text{insieme soluzioni globali}\}$. I casi di interesse sono i seguenti:

- $X^* = \emptyset \rightarrow$ Insieme convesso;
- $X^* = \{x^*\} \rightarrow$ Insieme convesso;
- $|X^*| \geq 2$ ci sono almeno due soluzioni globali;

\square

Nell'ultimo caso, supponiamo che $x_1^* \in X^*, x_2^* \in X^*$ se sono soluzioni globali allora:

$$x_1^* \in S \rightarrow f(x_1^*) \leq f(x) \quad \forall x \in S$$

$$x_2^* \in S \rightarrow f(x_2^*) \leq f(x) \quad \forall x \in S$$

Necessariamente se sono **punti di minimo globale** deve valere:

$$f(x_1^*) = f(x_2^*) = f^*$$

Inoltre S è convesso $\rightarrow \forall \lambda \in [0, 1] \quad \lambda x_1^* + (1 - \lambda)x_2^* \in S$

Poiché f convessa: $f(\lambda x_1^* + (1 - \lambda)x_2^*) \leq \lambda f(x_1^*) + (1 - \lambda)f(x_2^*) = \lambda f^* + f^* - \lambda f^*$

$$f(\lambda x_1^* + (1 - \lambda)x_2^*) \leq f^*$$

Poiché f^* è ottimo globale per ipotesi $\rightarrow \lambda x_1^* + (1 - \lambda)x_2^* \in X^* \rightarrow X^*$ è un insieme convesso.

5.7 Derivata direzionale di f in x lungo direzione d

In questo corso consideriamo funzioni $f \in C^1$ che sono **continuamente differenziabili** e le sue derivate sono **funzioni continue**.

Definizione 2. Una **derivata direzionale** $Df(x, d)$ se esiste:

$$\lim_{t \rightarrow 0} \frac{f(x + t d) - f(x)}{t} = Df(x, d)$$

Definizione 3. Definiamo la **derivata parziale i -esima in d** :

$$Df(x, e_i) \quad d = e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i\text{-esima componente}$$

$$\frac{\partial f(x)}{\partial x_i} = \lim_{t \rightarrow 0} \frac{f(x_1, \dots, x_i + t, \dots, x_n) - f(x_i)}{t}$$

Definizione 4. Gradiente di $f(x)$

Il gradiente di $f(x)$ è un vettore colonna nella forma:

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{pmatrix}$$

Corollario 5.

$$\text{Se } f \in C^1 \text{ su } S \rightarrow Df(x, d) = \nabla f(x)^T d = \begin{pmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \end{pmatrix} \begin{pmatrix} d_1 \\ \vdots \\ d_n \end{pmatrix} = \sum_{i=1}^n \frac{\partial f(x)}{\partial x_i} d_i$$

Nota 6.

$$x^T y = \sum_{i=1}^n x_i y_i = \|x\| \|y\| \cos(\theta) \rightarrow x^T y \leq \|x\| \|y\|$$

5.7.1 Sviluppo di Taylor I ordine

$$f(x+d) = f(x) + \nabla f(x)^T d + \alpha(x, d)$$

Il resto $\alpha(x, d)$ è **trascurabile** poiché:

$$\lim_{\|d\| \downarrow 0} \frac{\alpha(x, d)}{\|d\|} = 0$$

Quindi se $\|d\|$ è sufficientemente piccolo $\rightarrow f(x+d) \approx f(x) + \nabla f(x)^T d$ è **lineare**

Proposizione 7. Convessità e derivata del primo ordine

$$f \in C^1 \text{ su } S \text{ e } f \text{ è convessa su } S \iff \forall x, y \in S \quad f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

Inoltre:

$$f \in C^1 \text{ su } S \text{ e } f \text{ strettamente convessa su } S \iff \forall x, y \in S \quad f(y) > f(x) + \nabla f(x)^T (y - x)$$

5.7.2 Sviluppo Taylor II ordine

Definizione 8. Matrice Hessiana

Una matrice $\nabla^2 f(x)$ è una **matrice hessiana** se è simmetrica ed è nella seguente forma:

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f(x)}{\partial^2 x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_1} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial^2 x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} & \cdots & \cdots & \frac{\partial^2 f(x)}{\partial^2 x_n} \end{pmatrix}$$

Per esistere, $f \in C^2$

A questo punto possiamo introdurre lo **sviluppo al II ordine di Taylor**:

$$f(x+d) = f(x) + \nabla f(x)^T d + \frac{1}{2} d^T \nabla^2 f(x) d + \beta(x, d)$$

Questo sviluppo è migliore dello sviluppo del primo ordine poiché tende a zero più velocemente. Infatti:

$$\lim_{\|d\| \downarrow 0} \frac{\beta(x, d)}{\|d\|^2} = 0$$

In pratica è più accurato in uno spazio più grande.

5.8 Segno di una matrice

Definizione 9. (Semi)definita positiva/negativa

Data una matrice $Q \in \mathbb{R}^{n \times n}$, $Q = Q^T$. Definiamo:

$$- \quad Q \text{ semidefinita positiva: } Q \succcurlyeq 0 \iff y^T Q y \geq 0 \quad \forall y \in \mathbb{R}^n$$

- Q **definita positiva**: $Q \succ 0 \iff y^T Q y > 0 \quad \forall y \in \mathbb{R}^n$
- Q **definita negativa**: $Q \prec 0 \iff y^T Q y < 0 \quad \forall y \in \mathbb{R}^n$
- Q **semidefinita negativa**: $Q \preceq 0 \iff y^T Q y \leq 0 \quad \forall y \in \mathbb{R}^n$
- Q **indefinita**: $\iff \exists y_1 \in \mathbb{R}^n: y_1^T Q y_1 > 0 \wedge \exists y_2 \in \mathbb{R}^n: y_2^T Q y_2 < 0$

5.8.1 Matrice Hessiana

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f(x)}{\partial^2 x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_1} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial^2 x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} & \cdots & \cdots & \frac{\partial^2 f(x)}{\partial^2 x_n} \end{pmatrix}$$

Proposizione 10. $f(x)$ convessa su $\forall x, y \in S \quad \lambda \in [0, 1] \quad f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$

Proposizione 11. $f \in C^2$ e f convessa su $S \iff \nabla^2 f(x) \succeq 0 \forall x \in S$

Nota 12.

- Se si trova 1 punti per cui non vale 11, allora non si ha la convessità

5.8.2 Criterio dei minori di Nord-Ovest

Ipotizziamo $Q \in \mathbb{R}^{n \times n}$, $Q = Q^T$

Definizione 13. *Minori di Nord-Ovest di ordine k*

I minori di Nord-Ovest di ordine k sono quei minori ottenuti cancellando le ultime $n - k$ righe e colonne.

Quindi se :

$$Q \succ 0 \iff \text{tutti } i \text{ minori di nord - ovest hanno determinante } > 0 \text{ di ordine } k = 1, \dots, n$$

Definizione 14. *Matrice diagonale dominante*

La matrice diagonale dominante è una condizione **sufficiente** di definita positività. In particolare date le ipotesi 5.8.2, si ha una matrice diagonale dominante se:

$$q_{ii} > \sum_{j \neq i} |q_{ij}| \quad \forall i$$

Nel caso in cui non è possibile determinare la semidefinita/definita positività occorre considerare i **minori principali**.

Definizione 15. *Minori principali di ordine k*

I minori principali di ordine k sono tutti i minori che si appoggiano sulla diagonale eliminando $n-k$ righe e colonne in tutti i modi possibili.

Importante 16. Algoritmo segno di una matrice

1. Applico criterio di nord-ovest:

A) tutti i minori hanno determinante $>0 \rightarrow \text{STOP}$ $Q \succ 0$

B) Si ha un minore con determinante $=0 \rightarrow$ **Criterio minori principali:**

i. tutti i minori principali hanno determinante $\geq 0 \rightarrow \text{STOP} \rightarrow Q \succeq 0$

ii. trovo minore con determinante <0 :

a) se i minori principali hanno determinante >0 e un determinante <0

$\rightarrow \text{STOP} \Rightarrow Q$ indefinita

b) Ritorno al passo 1 considerando $-Q$ per dedurre la (semi)definita negatività.

5.8.3 Segno di una matrice - Autovalori/Autovettori

Sia data una matrice $Q = Q^T, Q \in \mathbb{R}^{n \times n}$:

Definizione 17. *Autovalore, Autovettore*

λ_i è un **autovalore** di Q con corrispondente **autovettore** u_i se:

$$Q u_i = \lambda_i u_i \quad i = 1, \dots, n \quad \|u_i\| = 1$$

Nota 18. $u_i^T u_j = 0$ gli autovettori sono ortogonali per autovalori diversi $\lambda_i \neq \lambda_j$

Quindi:

$$Q = \sum_{i=1}^n \lambda_i u_i u_i^T$$

Si ha che:

- Q semidefinita positiva $\iff \lambda_1, \dots, \lambda_n \geq 0$ ($y^T Q y \geq 0 \forall y \in \mathbb{R}^n \iff \lambda_{\min}(Q) \geq 0$)
- Q definita positiva $\iff \lambda_1, \dots, \lambda_n > 0$ ($y^T Q y > 0 \forall y \in \mathbb{R}^n \iff \lambda_{\min}(Q) > 0$)
- Q semidefinita negativa $\iff \lambda_1, \dots, \lambda_n \leq 0$ ($y^T Q y \leq 0 \forall y \in \mathbb{R}^n \iff \lambda_{\max}(Q) \leq 0$)
- Q definita negativo $\iff \lambda_1, \dots, \lambda_n < 0$ ($y^T Q y < 0 \forall y \in \mathbb{R}^n \iff \lambda_{\max}(Q) < 0$)
- Q indefinita $\iff \exists \lambda_i > 0, \lambda_j < 0$

$$\lambda_{\min}(Q) = \min_{x \in \mathbb{R}^n} \frac{x^T Q x}{\|x\|^2} := \text{quoziente di Rayleigh}$$

La seguente formula permette di valutare la coercività:

$$x^T Q x \geq \lambda_{\min}(Q) \|x\|^2$$

6 Condizioni di esistenza delle soluzioni

Consideriamo il problema nella forma:

$$\begin{cases} \min f(x) \\ x \in S \end{cases}$$

6.1 Definizioni preliminari

Definizione 19. *Insieme aperto*

Un insieme S è **aperto** se $\forall x \in S, \exists B(x, \rho)$ aperta tutta contenuta in S .

Definizione 20. *Insieme chiuso*

Un insieme S è **chiuso** se il suo complemento \mathbb{R}^n / S è aperto.

Definizione 21. *Insieme limitato*

Un insieme S è **limitato** se $\exists M > 0: \forall x \in S \quad \|x\| \leq M$.

Definizione 22. *Insieme compatto*

Un insieme S è **compatto** se è **chiuso** e **limitato**.

6.2 Esistenza del minimo

Teorema 23. Teorema di Weierstrass

Se vogliamo $\begin{cases} \min_{x \in S} f(x) \\ f \text{ continua e } S \text{ compatta} \end{cases} \longrightarrow \exists \text{ minimo (massimo) globale.}$

Importante 24. Il teorema di Weierstrass fornisce condizioni solo sufficienti: se non è verificata non si può concludere nulla.

L'ipotesi di compattezza di S può essere rilassata, imponendo una condizione più stringente sulla f .

Definizione 25. *Insieme di Livello vincolato di f*

Dato il problema di minimizzazione:

$$\begin{cases} \min f(x) \\ x \in S \end{cases}$$

Si definisce **insieme di livello vincolato di f** :

$$\mathcal{L}(\alpha) = \{x \in S: f(x) \leq \alpha\} \quad \alpha = f(x_0), x_0 \in S \text{ iniziale}$$

Teorema 26.

Se $\exists x_0: \mathcal{L}(f(x_0))$ compatto e f continua $\implies \exists$ una soluzione globale

Dimostrazione.

Se si considera il problema di minimizzazione $\left\{ \min_{x \in \mathcal{L}(f(x_0))} f(x) \right\}$, vale il teorema di Weierstrass.

Quindi esiste una soluzione globale x^* tale che $f(x^*) \leq f(x) \quad \forall x \in \mathcal{L}(f(x_0))$.

Se:

$x \notin \mathcal{L}(f(x_0))$, ma $x \in S \rightarrow f(x) > f(x_0) \geq f(x^*) \rightarrow f(x^*) \leq f(x) \quad \forall x \in \mathcal{L}(f(x_0)), \forall x \in S: x \notin \mathcal{L}(f(x_0)) \implies f(x^*) \leq f(x) \quad \forall x \in S \rightarrow x^*$ è un minimo globale. \square

Definizione 27. Coercività di f

La funzione f si dice **coerciva** su S se comunque scelgo $\{x_k\} \in S$ tale che:

$$\lim_{k \rightarrow \infty} \|x_k\| = \infty \text{ e } \lim_{k \rightarrow \infty} f(x_k) = \infty$$

Teorema 28.

f è **coerciva** su $S \iff$ tutti i suoi insiemi di livello sono compatti

Teorema 29. Esistenza del minimo (Sufficiente)

Se si considera il problema di minimizzazione $\left\{ \min_{x \in S} f(x) \right\}$ e la funzione f è continua e coerciva su S . Allora esiste la soluzione di minimo globale.

Importante 30.

Se $Q > 0$ ed $f(x)$ è strettamente convessa allora la soluzione esiste ed è unica su qualunque S chiuso e non vuoto.

Definizione 31. Termine di regolarizzazione

Sia dato il problema di minimizzazione $\left\{ \min_{x \in S} f(x) \right\} \longrightarrow \min f(x) + \tau \|x - x_k\|^2 \quad \tau > 0$

Nota 32.

- Il termine $\tau \|x - x_k\|^2$ rende coerciva la funzione e la convessifica rendendo il problema più semplice. Poiché la matrice Hessiana diventa costante.

7 Condizioni di ottimo per problemi non vincolati

Con riferimento ad un problema di minimizzazione del tipo:

$$\min f(x), \quad x \in S$$

una **condizione di ottimalità** è una condizione (necessaria, sufficiente, necessaria e sufficiente) perché un punto x^* risulti una soluzione ottima locale o globale del problema.

Importante 33. Condizioni necessarie

Una **condizione necessaria** può servire a restringere l'insieme in cui ricercare le soluzioni del problema necessario e a costruire algoritmi che soddisfino tali condizioni:

$$x^* \text{ minimo locale} \in S \longrightarrow \text{vale la condizione}$$

Quindi forniscono **criteri di arresto** per gli algoritmi di risoluzioni e una indicazione su dove cercare l'ottimo.

Importante 34. Condizioni sufficienti

Una **condizione sufficiente** può servire a dimostrare che un punto ottenuto sia una soluzione ottima, locale o globale, del problema.

$$\text{vale la condizione} \longrightarrow x^* \text{ minimo locale} \in S$$

Importante 35. Condizioni necessarie e sufficienti

Questo tipo di condizioni si hanno solamente nelle tipologie di problemi convessi.

$$x^* \text{ minimo globale} \longleftrightarrow \text{vale condizione}$$

Da queste considerazioni possiamo trarne una rappresentazione insiemistica:

Da cui possiamo evincere che, se la condizione necessaria è violata, il punto può essere escluso e che le condizioni sufficienti sono utilissime solo se soddisfatte; in alternativa non si può concludere nulla dal problema.

7.0.1 Direzione di discesa

L'idea è quella di definire una **direzione di discesa** in un punto x come un vettore $d \in \mathbb{R}^n$ tale che, per tutti gli spostamenti sufficientemente piccoli lungo d , si ha una diminuzione stretta del valore di f rispetto al valore di x

Definizione 36. Direzione di discesa

Sia $f: \mathbb{R}^n \rightarrow \mathbb{R}$ e $x \in \mathbb{R}^n$. Si dice che un vettore $d \in \mathbb{R}^n, d \neq 0$ è una **direzione di discesa** per f in x se esiste $\bar{t} > 0$ tale che:

$$f(x + t d) < f(x) \quad \forall t \in (0, \bar{t}]$$

In modo analogo si potrebbe definire una **direzione di salita** nel punto x se f aumenta strettamente per tutti gli spostamenti abbastanza piccoli lungo d

La direzione di discesa la si può caratterizzare tramite le derivate direzionali della funzione obiettivo. Infatti ricordiamo che la **derivata direzionale** è definita nel seguente modo:

$$Df(x, d) = \lim_{t \rightarrow 0} \frac{f(x + t d) - f(x)}{t}$$

Proposizione 37. Derivata direzionale e direzione di discesa

Se $f \in C^1$ continua e derivabile almeno una volta $\rightarrow Df(x, d) = \nabla f(x)^T d$. Quindi, per t sufficientemente piccolo, otteniamo:

1. $\nabla f(x)^T d < 0 \rightarrow f(x + td) - f(x) < 0 \rightarrow d$ è di **discesa**;
2. $\nabla f(x)^T d > 0 \rightarrow f(x + td) - f(x) > 0 \rightarrow d$ è di **salita**;
3. $\nabla f(x)^T d = 0 \rightarrow f(x + td) - f(x) = 0 \rightarrow$????

Teorema 38. Condizione di discesa del primo ordine (necessaria e sufficiente)

Se f convessa e C^1 , d discesa $\iff \nabla f(x)^T d < 0$

Dimostrazione.

Prima implicazione: (\Leftarrow)

Se $\nabla f(x)^T d < 0$, la direzione d è di discesa in quanto, per la definizione di derivata direzionale, per valori positivi sufficientemente piccoli di t , deve valere $f(x + td) - f(x) < 0$.

Seconda Implicazione: (\Rightarrow)

Poniamo $y = x + td$ ed assumiamo che f sia convessa. Dalla definizione di convessità:

$$\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T (y - x) \\ f(x + td) &\geq f(x) + t \nabla f(x)^T (y - x) \\ &\geq f(x) + t \nabla f(x)^T d \geq f(x) \end{aligned}$$

per cui $\nabla f(x)^T d \geq 0$ implica $f(x + td) \geq f(x)$. Tuttavia, se d è una direzione di discesa, deve necessariamente essere $\nabla f(x)^T d < 0$. \square

Definizione 39. Segno derivata direzionale e antigradiente

Dal punto di vista geometrico, ricordando il concetto di angolo fra due vettori possiamo studiare il **segno derivata direzionale**:

$$\nabla f(x)^T d = \|\nabla f(x)\| \|d\| \cos(\theta)$$

Se:

- $\theta > 90^\circ \rightarrow \cos(\theta) < 0 \rightarrow \nabla f(x)^T d < 0 \rightarrow d$ è di **discesa**;
- $\theta < 90^\circ \rightarrow \cos(\theta) > 0 \rightarrow \nabla f(x)^T d > 0 \rightarrow d$ è di **salita**;
- Il caso “limite”, cioè di **massima discesa** si ottiene con il cosiddetto **antigradiente**:

$$d = -\nabla f(x) \rightarrow \theta = 180^\circ \rightarrow \cos(\theta) = -1 \text{ è sempre di discesa}$$

$$\nabla f(x)^T d = -\|\nabla f(x)\|^2 < 0$$

7.0.2 Direzioni a curvatura negativa in x

Se f è differenziabile due volte è possibile caratterizzare l'andamento di f lungo una direzione assegnata utilizzando anche le derivate seconde per stabilire le condizioni di ottimo del secondo ordine.

Definizione 40. Direzione a curvatura negativa

Sia $f: \mathbb{R}^n \rightarrow \mathbb{R}$ due volte continuamente differenziabile nell'intorno di un punto $x \in \mathbb{R}^n$. Si dice che un vettore $d \in \mathbb{R}^n, d \neq 0$ è una **direzione a curvatura negativa** per f in x se risulta:

$$d^T \nabla^2 f(x) d < 0$$

Nota 41. Una direzione a curvatura negativa è quindi tale che la derivata direzionale seconda è negativa in x , per cui diminuisce localmente la derivata direzionale del primo ordine.

Teorema 42. Condizione di discesa del secondo ordine

Sia $f: \mathbb{R}^n \rightarrow \mathbb{R}$ due volte continuamente differenziabile nell'intorno di un punto $x \in \mathbb{R}^n$ e dia $d \in \mathbb{R}^n \neq 0$. Supponiamo che risulti $\nabla f(x)^T d = 0$, e che d sia una direzione a curvatura negativa in x , ossia tale che $d^T \nabla^2 f(x) d < 0$. Allora d è una **direzione di discesa** per f in x .

Dimostrazione.

Poiché f è differenziabile due volte, si ha:

$$f(x + t d) = f(x) + t \nabla f(x)^T d + \frac{1}{2} t^2 d^T \nabla^2 f(x) d + \beta(x, t d)$$

In cui β è il resto:

$$\lim_{t \rightarrow 0} \frac{\beta(x, t d)}{t^2} = 0$$

Essendo per ipotesi $\nabla f(x)^T d = 0$, si può scrivere:

$$\frac{f(x + t d) - f(x)}{t^2} = \frac{1}{2} d^T \nabla^2 f(x) d + \frac{\beta(x, t d)}{t^2}$$

Poiché $d^T \nabla^2 f(x) d < 0$ per definizione di curvatura negativa e $\frac{\beta(x, t d)}{t^2} \rightarrow 0$ per $t \rightarrow 0$. Per valori sufficientemente piccoli di t , si ha che $f(x + t d) - f(x) < 0$ e quindi d è una direzione di discesa. \square

Definizione 43. Segno derivata direzionale e antigradiente

Dal punto di vista geometrico, ricordando il concetto di angolo fra due vettori possiamo studiare il **segno derivata direzionale**:

$$\nabla f(x)^T d = \|\nabla f(x)\| \|d\| \cos(\theta)$$

Se:

- $\theta > 90^\circ \rightarrow \cos(\theta) < 0 \rightarrow \nabla f(x)^T d < 0 \rightarrow d$ è di **discesa**;
- $\theta < 90^\circ \rightarrow \cos(\theta) > 0 \rightarrow \nabla f(x)^T d > 0 \rightarrow d$ è di **salita**;
- Il caso “limite”, cioè di **massima discesa** si ottiene con il cosiddetto **antigradiente**:

$$d = -\nabla f(x) \rightarrow \theta = 180^\circ \rightarrow \cos(\theta) = -1 \text{ è sempre di } \mathbf{discesa}$$

$$\nabla f(x)^T d = -\|\nabla f(x)\|^2 < 0$$

7.0.3 Direzioni a curvatura negativa in x

Se f è differenziabile due volte è possibile caratterizzare l'andamento di f lungo una direzione assegnata utilizzando anche le derivate seconde per stabilire le condizioni di ottimo del secondo ordine.

Definizione 44. Direzione a curvatura negativa

Sia $f: \mathbb{R}^n \rightarrow \mathbb{R}$ due volte continuamente differenziabile nell'intorno di un punto $x \in \mathbb{R}^n$. Si dice che un vettore $d \in \mathbb{R}^n, d \neq 0$ è una **direzione a curvatura negativa** per f in x se risulta:

$$d^T \nabla^2 f(x) d < 0$$

Nota 45. Una direzione a curvatura negativa è quindi tale che la derivata direzionale seconda è negativa in x , per cui diminuisce localmente la derivata direzionale del primo ordine.

Teorema 46. Condizione di discesa del secondo ordine

Sia $f: \mathbb{R}^n \rightarrow \mathbb{R}$ due volte continuamente differenziabile nell'intorno di un punto $x \in \mathbb{R}^n$ e sia $d \in \mathbb{R}^n \neq 0$. Supponiamo che risulti $\nabla f(x)^T d = 0$, e che d sia una direzione a curvatura negativa in x , ossia tale che $d^T \nabla^2 f(x) d < 0$. Allora d è una **direzione di discesa** per f in x .

Dimostrazione.

Poiché f è differenziabile due volte, si ha:

$$f(x + t d) = f(x) + t \nabla f(x)^T d + \frac{1}{2} t^2 d^T \nabla^2 f(x) d + \beta(x, t d)$$

In cui β è il resto:

$$\lim_{t \rightarrow 0} \frac{\beta(x, t d)}{t^2} = 0$$

Essendo per ipotesi $\nabla f(x)^T d = 0$, si può scrivere:

$$\frac{f(x + t d) - f(x)}{t^2} = \frac{1}{2} d^T \nabla^2 f(x) d + \frac{\beta(x, t d)}{t^2}$$

Poiché $d^T \nabla^2 f(x) d < 0$ per definizione di curvatura negativa e $\frac{\beta(x, t d)}{t^2} \rightarrow 0$ per $t \rightarrow 0$. Per valori sufficientemente piccoli di t , si ha che $f(x + t d) - f(x) < 0$ e quindi d è una direzione di discesa. \square

8 Condizioni di ottimo per problemi vincolati con insieme ammissibile

Il problema ce vogliamo risolvere è del tipo:

$$\min_{x \in S} f(x)$$

cercando le direzioni di discesa ammissibile del problema, cioè che non violano le condizioni dettate da $x \in S$.

Definizione 47. Direzione ammissibile

Sia S un sottoinsieme di \mathbb{R}^n e $x \in S$. Si dice che un vettore $d \in \mathbb{R}^n, d \neq 0$ è una **direzione ammissibile** per S in s se:

$$d \in \mathbb{R}^n: \exists \bar{t} > 0 \quad x + t d \in S \quad \forall t \in (0, \bar{t}]$$

Proposizione 48. *Condizione di ottimo necessaria di minimo locale*

Sia $x^* \in S$ un punto di minimo locale del problema:

$$\min f(x), x \in S;$$

allora non può esistere una direzione ammissibile in x^* che sia anche di discesa.

Dimostrazione.

Per assurdo, se esistesse una direzione d che sia ammissibile e di discesa in x^* , allora in ogni intorno di x^* sarebbe possibile trovare un punto $x^* + t d \in B(x^*, \rho)$ tale che:

$$f(x^* + t d) \leq f(x^*)$$

il che contraddice l'ipotesi che x^* sia un punto di minimo globale. □

Se f è differenziabile la condizione:

$$\nabla f(x)^T d < 0$$

è una condizione sufficiente perché d sia una direzione di discesa per f in x , otteniamo la condizione encessaria seguente.

Proposizione 49. *Condizione necessaria del primo ordine*

Supponiamo che $f: \mathbb{R}^n \rightarrow \mathbb{R}$ sia continuamente differenziabile nell'intorno in un punto di minimo locale $x^* \in S$ del problema:

$$\min f(x) \quad x \in S$$

Allora non può essitere una direzione ammissibile in d in x^* tale che:

$$\nabla f(x^*)^T d < 0$$

e quindi si ha:

$$\nabla f(x^*)^T d \geq 0 \quad \forall d \in \mathbb{R}^n \text{ ammissibile in } x^*$$

Proposizione 50. *Condizione necessaria del secondo ordine*

Supponiamo che $f: \mathbb{R}^n \rightarrow \mathbb{R}$ sia due volte continuamente differenziabile nell'intorno di un punto di minimo locale $x^* \in S$ del problema:

$$\min f(x) \quad x \in S$$

Allora non può esistere una direzione ammissibile d in x^* tale che:

$$\nabla f(x^*)^T d = 0 \quad d^T \nabla^2 f(x) d < 0$$

e quindi si ha:

$$\forall d \text{ ammissibile tale che } \nabla f(x^*)^T d = 0 \quad d^T \nabla^2 f(x) d > 0$$

Nota 51.

Non si ha l'uguaglianza poiché avremmo una direzione a curvatura negativa e quindi ammissibile nell'insieme.

8.1 Condizioni di ottimo per problemi vincolati con insieme ammissibile convesso

Se l'insieme ammissibile S è convesso, vogliamo caratterizzare le condizioni di ottimo. A tale scopo occorre caratterizzare le direzioni ammissibili.

Proposizione 52. Direzione ammissibile (S convesso)

$$\begin{aligned} S \text{ convesso} \\ \iff d = y - x \quad \text{per qualche } y \in S \\ d \text{ ammissibile in } x \end{aligned}$$

Dimostrazione.

Sia $x \in S$ allora comunque si fissi $y \in S$ tale che $x \neq y$, per la convessità di S si ha che:

$$\lambda y + (1 - \lambda)x \in S \quad \forall \lambda \in [0, 1]$$

$$x + \lambda(y - x) \in S \longrightarrow d = y - x \in S \text{ ammissibile}$$

Al contrario se $d \in \mathbb{R}^n \neq 0$ è una direzione ammissibile per S , in x esiste un punto $y \in S$ ed uno scalare $\lambda > 0$ tali che $d = \lambda(y - x)$. \square

Proposizione 53. Condizione necessaria di minimo locale del primo ordine (S convesso)

Sia $x^* \in S$ un punto di minimo locale del problema:

$$\min f(x) \quad x \in S$$

in cui $S \subseteq \mathbb{R}^n$ è un insieme convesso e supponiamo f sia continuamente differenziabile in un intorno di x^* . Allora si ha necessariamente:

$$\nabla f(x^*)^T (y - x^*) \geq 0 \quad \forall y \in S$$

Proposizione 54. Condizione necessaria di minimo locale del secondo ordine (S convesso)

Sia $x^* \in S$ un punto di minimo locale del problema:

$$\min f(x) \quad x \in S$$

in cui $S \subseteq \mathbb{R}^n$ è un insieme convesso e supponiamo f sia due volte continuamente differenziabile in un intorno di x^* . Allora si ha necessariamente:

$$\forall y \in S \quad \nabla f(x^*)^T (y - x^*) = 0 \longrightarrow (y - x^*)^T \nabla^2 f(x^*) (y - x^*) \geq 0$$

Questo si riconduce alla ricerca dei **punti di stazionarietà** della funzione. Infatti:

- Se il punto è interno all'insieme S :

$$\nabla f(x^*)^T (y - x^*) \geq 0 \quad \forall y \in S$$



$$\nabla f(x^*) = 0$$

Questa condizione è detta di **stazionarietà non vincolata**.

- Se ci troviamo sul bordo dell'insieme S :

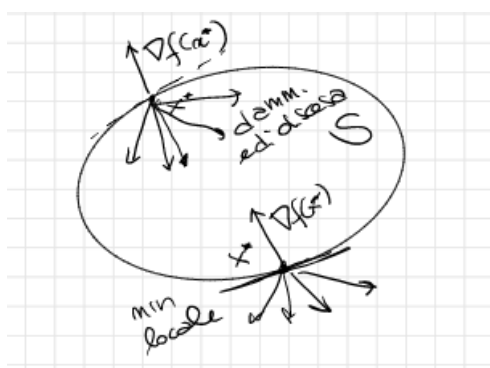


Figura 3.

Dipende poiché tutte le direzioni di discesa sono non ammissibili. In alto, se il $\cos > 90$ otteniamo direzioni ammissibili di discesa e non abbiamo minimo; in basso, se il $\cos > 90$ le direzioni di discesa non sono ammissibili ed otteniamo un minimo locale.

Dimostrazione.

\Rightarrow) x^* minimo globale $\rightarrow x^*$ minimo locale \rightarrow vale la C.N. del primo ordine;

\Leftarrow) Poiché f convessa su $S \quad \forall y \in S$

$$f(y) \geq f(x^*) + \nabla f(x^*)^T (y - x^*) \geq f(x^*)$$

In cui $\nabla f(x^*)^T (y - x^*) \geq 0$ per ipotesi. Allora:

$$f(x^*) \leq f(y) \quad \forall y \in S \rightarrow x^* \text{ minimo globale}$$

□

Importante 55.

Se f strettamente convessa \rightarrow minimo globale diventa unico

8.2 Problemi con vincoli lineari

$$\begin{aligned} \min \quad & f(x) \quad A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m \\ & Ax \leq b \end{aligned}$$

Inoltre:

$$a_i^T x \leq b_i \quad i = 1, \dots, m; a^T \text{ righe di } A$$

Supponiamo $\bar{x} \in S \iff a_i^T \bar{x} \leq b_i \quad i = 1, \dots, m$

Definizione 56. Vincoli attivi nel punto

$$I(x) := \{i = 1, \dots, m: a_i^T x = b_i\}$$

Nota 57.

I vincoli attivi influenzano le direzioni ammissibili

Proposizione 58. Direzioni ammissibili: disuguaglianze lineari

Sia x^* tale che $a_i^T x^* \leq b_i$ con $i = 1, \dots, m$. Allora un vettore non nullo $d \in \mathbb{R}^n$ è una **direzione ammissibile** per S se e solo se:

$$a_i^T d \leq 0 \quad \forall i \in I(x^*)$$

Dimostrazione.

\implies) Dati un vettore non nullo $d \in \mathbb{R}^n$ e numero $t \geq 0$ qualsiasi sufficientemente piccolo, il punto $x^* + td$ è ammissibile relativamente ai vincoli attivi in x^* se e solo se:

$$a_i^T (x^* + td) = b_i + t a_i^T d \leq b_i \quad \forall i \in I(x^*)$$

Quindi se e solo se vale la proposizione 58.

\Leftarrow) Se vale $a_i^T d \leq 0 \quad \forall i \in I(x^*)$, d è una direzione ammissibile anche per i vincoli non attivi. Infatti, essendo $a_i^T x^* \leq b_i$ per $i \notin I(x^*)$ si può assumere t sufficientemente piccolo da avere:

$$a_i^T (x^* + td) \leq b_i \quad \forall i \notin I(x^*)$$

Quindi si può concludere che esiste un numero $\bar{t} > 0$ tale che $A(x^* + td) \leq b \quad \forall t \in [0, \bar{t}]$, ossia che d è una direzione ammissibile per S in x^* , \square

Importante 59. L'insieme delle direzioni ammissibili è un **poliedro**.

Importante 60. Se ho vincoli soddisfatti all'uguaglianza la disequazione 58 viene soddisfatta all'uguaglianza.

8.3 Condizioni di ottimalità per problemi con vincoli di box

Sia $x^* \in \mathbb{R}^n$ un punto di minimo locale del problema

$$\begin{cases} \min f(x) \\ l \leq x \leq u \end{cases}$$

In cui $l_i < u_i \quad i = 1, \dots, n$. Allora, la condizione di stazionarietà del primo ordine è utile per calcolare il minimo locale x^* ; quindi, risulta:

$$\frac{\partial f}{\partial x_i} = \begin{cases} \geq 0 & \text{se } x_i^* = l_i \\ = 0 & \text{se } l_i < x_i^* < u_i \\ \leq 0 & \text{se } x_i^* = u_i \end{cases}$$

Nota 61. L'insieme appena descritto è un insieme chiuso e limitato, quindi **compatto**.

Grazie al **teorema di Weierstrass** 23 possiamo dedurre l'esistenza della soluzione.

A questo punto, dobbiamo discutere tre casistiche:

- A) Se $x_i^* = l_i$, consideriamo la direzione $d = e_i$, in cui e_i è l' i -esimo asse coordinato. Ricordiamo che:

$$\text{Il vincolo è attivo } x_i^* = l_i \iff a_i^T x \leq b_i \longrightarrow d \text{ ammissibile} \iff a_i^T d \leq 0$$

Quindi segue necessariamente che l'unica direzione ammissibile è data da:

$$\nabla f(x^*)^T d = \frac{\partial f(x^*)}{\partial x_i} \geq 0$$



Figura 4. direzioni ammissibili

- B) Se $l_i < x_i^* < u_i$ le direzioni $d = e_i, d = -e_i$ sono ammissibili e dalla proposizione si ottiene:

$$\nabla f(x^*)^T d = \frac{\partial f(x^*)}{\partial x_i} \geq 0; \nabla f(x^*)^T d = -\frac{\partial f(x^*)}{\partial x_i} \geq 0$$

Da cui:

$$\nabla f(x^*)^T d = \frac{\partial f(x^*)}{\partial x_i} = 0$$

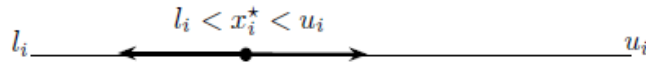


Figura 5. Direzioni ammissibili

- C) Se $x_i^* = u_i$, si ha la direzione $d = -e_i$ ammissibile per cui otteniamo che:

$$\nabla f(x^*)^T d = \frac{\partial f(x^*)}{\partial x_i} \leq 0$$



Figura 6.

9 Condizioni di Ottimo per problemi vincolati

Consideriamo i problemi con vincoli di eguaglianza e disequaglianza del tipo:

$$\begin{cases} \min f(x) \\ g_i(x) \leq 0 & i = 1, \dots, n \\ h_j(x) = 0 & j = 1, \dots, p \end{cases}$$

9.1 Condizioni di ottimo analitiche

9.1.1 Lagrangiana generalizzata

Prima di tutto vogliamo esprimere le condizioni di ottimo come estensione della *regole dei moltiplicatori di Lagrange*. A tale scopo, introduciamo i seguenti moltiplicatori:

1. λ_0 associato alla funzione obiettivo;
2. λ_i associato ai vincoli di disequaglianza;
3. μ_j associato ai vincoli di eguaglianza;

Così facendo, possiamo definire la **Lagrangiana generalizzata**:

$$L(x, \lambda_0, \lambda, \mu) = \lambda_0 f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x)$$

9.1.2 Condizioni necessarie di Fritz John

Teorema 62. Condizioni necessarie di Fritz John

Sia x^* un punto di minimo locale del problema:

$$\begin{cases} \min f(x) \\ g(x) \leq 0 \\ h(x) = 0 \end{cases}$$

e supponiamo che f, g, h siano continuamente differenziabili in un intorno di x^* . Allora esistono moltiplicatori $\lambda_0^* \in \mathbb{R}, \lambda^* \in \mathbb{R}^m, \mu^* \in \mathbb{R}^p$ tali che:

$$\nabla_x L(x^*, \lambda_0^*, \lambda^*, \mu^*) = 0 \quad = \quad \lambda_0^* f(x^*) + \sum_{i=1}^m \lambda_i^* g_i(x^*) + \sum_{j=1}^p \mu_j^* h_j(x^*) := \text{Annullamento del gradiente}$$

$$\lambda_i^* g_i(x^*) = 0 \quad i = 1, \dots, m \quad := \text{Condizione di complementarietà}$$

$$\lambda_0^*, \lambda^* \geq 0; (\lambda_0^*, \lambda^*, \mu^*) \neq 0 := \text{Condizione di non negatività}$$

$$g_i(x^*) \leq 0, h_j(x^*) = 0 := \text{Condizione di ammissibilità}$$

9.2 Condizioni di Qualificazione dei Vincoli o Condizioni di Regolarità

Dalla condizione di negatività si hanno due possibilità:

$$\lambda_0^* \geq 0 \rightarrow \begin{cases} \lambda_0^* = 0 \\ \lambda_0^* > 0 \end{cases}$$

Nel caso in cui $\lambda_0^* = 0 \longrightarrow \exists \lambda^*, \mu^*$ t.c.:

$$\sum_{i=1}^m \lambda_i^* g_i(x^*) + \sum_{j=1}^p \mu_j^* h_j(x^*) = 0$$

$$\lambda^* \geq 0, (\lambda^*, \mu^*) \neq 0$$

$$g_i(x^*) \leq 0 \quad h_j(x^*) = 0$$

Quindi la funzione obiettivo non è considerata: un qualunque punto x^* che soddisfa le condizioni di Fritz John è un candidato all'ottimo indipendentemente da $f \longrightarrow$ è un punto particolare dell'insieme ammissibile: un **punto non regolare**.

9.2.1 Condizioni di qualificazione dei vincoli

In questo ambito siamo interessati solo ai candidati per cui $\lambda_0^* \neq 0$ poiché coinvolgono la nostra funzione obiettivo. Le **condizioni di qualificazione dei vincoli** sono delle condizioni sui vincoli che definiscono l'insieme ammissibile in x^* che se sono soddisfatte, allora $\lambda_0^* > 0$:

a) LICQ (Linear Independence Constraint Qualification)

x^* , candidato di minimo locale, soddisfa **LICQ** se per i vettori $\nabla g_i(x^*)$ con $i \in I(x^*) = \{i = 1, \dots, n: g_i(x^*) = 0\}$ (insieme dei vincoli attivi) e $\nabla h_j(x^*) \quad j = 1, \dots, p$ sono **linearmente indipendenti** allora $\lambda_0^* \neq 0$;

Dimostrazione. Si suppone per assurdo che $\lambda_0^* = 0 \rightarrow$ la **condizione di annullamento** della Lagrangiana diventa:

$$\sum_{i=1}^m \lambda_i^* g_i(x^*) + \sum_{j=1}^p \mu_j^* h_j(x^*) = 0$$

Dalla **complementarietà**:

$$\lambda_i^* g_i(x^*) = 0 \quad i = 1, \dots, m$$

Quindi, se $i \notin I(x^*) \iff g_i(x^*) < 0 \longrightarrow \lambda_i^* = 0, \mu_j^* = 0$

Non viene soddisfatta la complementarietà.

Inoltre, se vale LICQ $\longrightarrow \nabla g_i(x^*) \quad i \in I(x^*)$ e $\nabla h_j(x^*)$ sono linearmente indipendenti. Quindi l'unica soluzione è:

$$d_i = 0, i \in I(x^*) \quad \mu_j^* = 0 \quad j = 1, \dots, p \quad \text{ASSURDO}$$

b) Linearità dei vincoli di eguaglianza e concavità dei vincoli attivi in x^*

Supponiamo che i vincoli di eguaglianza siano *lineari* e che i vincoli di disequaglianza attivi siano concavi nel punto x^* . In tali ipotesi è possibile trovare un intorno $B(x^*, \rho)$ di x^* tale che, per ogni $x \in B(x^*, \rho)$ si abbia:

$$\begin{aligned} h_i(x) &= h_i(x^*) + \nabla h_i(x^*)^T (x - x^*) \quad i = 1, \dots, p \\ g_i(x) &\leq g_i(x^*) + \nabla g_i(x^*)^T (x - x^*) \quad i \in I(x^*) \end{aligned}$$

Un caso particolare è quando tutti i vincoli attivi in x^* siano lineari. □

c) Condizioni di qualificazione dei vincoli di Mangasarian-Fromovitz

Sia $x^* \in S$ e supponiamo che g, h siano continuamente differenziabili in un intorno di x^* . Si dice che è soddisfatta in x^* la condizione di qualificazione dei vincoli di **Mangasarian-Fromovitz** se:

- i. I gradienti dei vincoli di eguaglianza $\{\nabla h_j(x^*), j = 1, \dots, p\}$ sono linearmente indipendenti;
- ii. $\nabla g_i(x^*)^T d < 0 \quad \forall i \in I(x^*) \quad \nabla h_i(x^*)^T d = 0 \quad \forall i = 1, \dots, p;$

d) Condizioni di Slater

Nelle **condizioni di Slater** consideriamo il caso in cui l'insieme ammissibile è definito attraverso vincoli convessi di disuguaglianza. Quindi:

Supponiamo che le funzioni g_i siano convesso e continuamente differenziabili su un insieme aperto convesso contenente l'insieme ammissibile:

$$S = \{x \in \mathbb{R}^n : g(x) \leq 0\}$$

Si dice che la condizione di Slater è soddisfatta se esiste $\hat{x} \in S$ tale che:

$$g(\hat{x}) < 0$$

Cioè che il punto \hat{x} sia interno all'insieme.

Se vale una di queste condizioni di qualificazione dei vincoli, allora $\lambda_0^* > 0$ e il punto è un **punto regolare**.

Importante 63.

Se $\lambda_0^* > 0$, si può assumere $\lambda_0^* = 1$ poiché equivale a $(\lambda_0, \lambda, \mu) \rightarrow \left(1, \frac{\lambda}{\lambda_0}, \frac{\mu}{\lambda_0}\right)$, lasciando invariato il problema. Inoltre, la Lagrangiana generalizzata diventa **Lagrangiana**:

$$L(x, \lambda, \mu) = 0 = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x)$$

9.2.2 Condizioni necessarie di KKT (Karush-Kuhn-Tucker)

Sia $x^* \in S$ un punto di minimo locale e supponiamo che f, g, h siano continuamente differenziabili in un intorno di x^* . Supponiamo inoltre che ne punto x^* sia soddisfatta una delle condizioni di qualificazione dei vincoli. Allora esistono $\lambda^* \in \mathbb{R}^m$ e $\mu^* \in \mathbb{R}^p$ tale che:

$$\nabla f(x^*) + \nabla g(x^*) \lambda^* + \nabla h(x^*) \mu^* = 0$$

$$g(x^*) \leq 0, h(x^*) = 0$$

$$\lambda^{*T} g(x^*) = 0$$

$$\lambda^* \geq 0$$

9.2.3 Moltiplicatori di Lagrange

Teorema 64. Teorema di Lagrange

Sia $x^* \in S$ un punto di minimo locale del problema:

$$\min f(x)$$

$$h(x) = 0$$

e supponiamo che f, h siano continuamente differenziabili in un intorno di x^* . Supponiamo inoltre che nel punto x^* sia soddisfatta una delle condizioni di qualificazione dei vincoli.

Allora esiste $\mu^* \in \mathbb{R}^p$ tale che:

$$\nabla f(x^*) + \nabla h(x^*) \mu^* = 0$$

$$h(x^*) = 0$$

Se vale il teorema, la ricerca dei punti che soddisfano le condizioni necessarie si riconduce alla ricerca delle soluzioni del sistema di $n + p$ equazioni nelle $n + p$ incognite (x, μ) .

9.3 Problemi con vincoli di box

$$\begin{cases} \min f(x) \\ l_i \leq x_i \leq u_i \end{cases} \quad i = 1, \dots, n$$

$$l_i - x_i \leq 0 \quad i = 1, \dots, n$$

$$x_i - u_i \leq 0 \quad i = 1, \dots, n$$

Definiamo, inoltre, i moltiplicatori $\hat{\lambda}_i$ e λ_i .

Essendo vincoli lineari, possiamo applicare le condizioni di KKT per ottenere condizioni necessarie di minimo locale. A tale scopo definiamo la Lagrangiana:

$$L(x, \hat{\lambda}, \lambda) = f(x) + \sum_{i=1}^m \hat{\lambda}_i (l_i - x_i) + \sum_{j=1}^n \lambda_j (x_j - u_j)$$

Scriviamo le condizioni di KKT:

(a)

$$\exists (\hat{\lambda}^*, \lambda^*) \text{ t.c. } \nabla L(x^*, \hat{\lambda}^*, \lambda^*) = 0 = \nabla f(x^*) - \hat{\lambda}^* + \lambda^*$$

$$\frac{\partial f(x^*)}{\partial x_i} - \hat{\lambda}_i^* + \lambda_i^* = 0 \quad i = 1, \dots, n$$

(b)

$$\hat{\lambda}_i^* (l_i - x_i^*) = 0 \quad i = 1, \dots, n$$

$$\lambda_i^* (x_i^* - u_i) = 0 \quad i = 1, \dots, n$$

(c)

$$\hat{\lambda}_i^*, \lambda_i^* \geq 0 \quad i = 1, \dots, n$$

(d)

$$l_i \leq x_i^* \leq u_i$$

La condizione (d) si divide in:

$$\begin{cases} x_i^* = l_i \textbf{(1)} \\ l_i < x_i^* < u_i \textbf{(2)} \\ x_i^* = u_i \textbf{(3)} \end{cases}$$

(1) $x_i^* = l_i$

Se sostituiamo (1) in (b) si ottiene che: $\lambda_i^* = 0$. Inoltre in (a):

$$\frac{\partial f(x^*)}{\partial x} - \hat{\lambda}_i^* + \lambda_i^* = 0 \quad i = 1, \dots, n$$

Si ottiene:

$$\frac{\partial f(x^*)}{\partial x} = \hat{\lambda}_i^* \geq 0$$

(2) $l_i < x_i^* < u_i$

Se sostituiamo (2) in (b):

$$\hat{\lambda}_i^* (l_i - x_i^*) = 0 \quad i = 1, \dots, n$$

$$\lambda_i^* (x_i^* - u_i) = 0 \quad i = 1, \dots, n$$

$$(l_i - x_i^*) < 0 \longrightarrow \hat{\lambda}_i^* = 0$$

$$(x_i^* - u_i) < 0 \longrightarrow \lambda_i^* = 0$$

Inoltre in (a):

$$\frac{\partial f(x^*)}{\partial x} - \hat{\lambda}_i^* + \lambda_i^* = 0 \quad i = 1, \dots, n$$

$$\frac{\partial f(x^*)}{\partial x} = 0$$

(3) $x_i^* = u_i$

Se sostituiamo (3) in (b) si ottiene che: $\lambda_i^* = 0$. Inoltre in (a):

$$\frac{\partial f(x^*)}{\partial x} - \hat{\lambda}_i^* + \lambda_i^* = 0 \quad i = 1, \dots, n$$

Si ottiene:

$$\frac{\partial f(x^*)}{\partial x} = \hat{\lambda}_i^* \leq 0$$

$$\frac{\partial f(x^*)}{\partial x} \leq 0$$

Ricapitolando siamo tornati al risultato precedente 8.3:

$$\frac{\partial f}{\partial x_i} = \begin{cases} \geq 0 & \text{se } x_i^* = l_i \\ = 0 & \text{se } l_i < x_i^* < u_i \\ \leq 0 & \text{se } x_i^* = u_i \end{cases}$$

9.4 Condizioni Sufficienti KKT

Le **condizioni di KKT** divengono *condizioni sufficienti* (e necessarie se il punto p regolare) di minimo globale nell'ipotesi che f e g , siano funzioni convesse e che i vincoli di eguaglianza siano lineari.

Teorema 65. *Condizioni sufficienti(e necessarie) di KKT*

Supponiamo che la funzione obiettivo f sia convessa, che le funzioni g_i , per $i = 1, \dots, m$ siano convesse e che i vincoli di uguaglianza siano lineari, ossia: $h(x) = Ax - b$. Supponiamo inoltre che f, g, h siano continuamente differenziabili su un insieme aperto contenente l'insieme ammissibile. Allora, se esistono moltiplicatori λ^* e μ^* tali che valgono le seguenti condizioni:

$$\begin{aligned}\nabla f(x^*) + \nabla g(x^*) \lambda^* + \nabla h(x^*) \mu^* &= 0 \\ g(x^*) &\leq 0, h(x^*) = 0 \\ \lambda^{*T} g(x^*) &= 0 \\ \lambda^* &\geq 0\end{aligned}$$

il punto x^* è un punto di minimo globale vincolato. Se inoltre f è strettamente convessa, allora x^* è l'unico punto di minimo globale vincolato.

Dimostrazione.

Osserviamo innanzitutto che l'insieme ammissibile è convesso e che f e g sono convesse e appartengono a C^1 .

Considerato un qualunque punto x ammissibile, essendo $\lambda^* \geq 0$, si ha:

$$\begin{aligned}f(x) &\geq f(x) + \lambda^{*T} g(x) + \mu^{*T} h(x) \\ h(x) &= h(x^*) + \nabla h(x^*)^T (x - x^*) \\ g_i(x) &\geq g_i(x^*) + \nabla g_i(x^*)^T (x - x^*)\end{aligned}$$

Da cui, essendo $\lambda^* \geq 0$:

$$\lambda^{*T} g_i(x) \geq \lambda^{*T} g_i(x^*) + \lambda^{*T} \nabla g_i(x^*)^T (x - x^*)$$

per la convessità di f :

$$f(x) \geq f(x^*) + \nabla f(x^*)^T (x - x^*)$$

Per dimostrare che x^* è minimo globale occorre dimostrare che $f(x^*) \leq f(x) \quad \forall x$ ammissibile tale che $g_i(x) \leq 0, h_j(x) = 0$

$$\begin{aligned}f(x) &\geq f(x) + \lambda^{*T} g(x) + \mu^{*T} h(x) \\ f(x) &\geq f(x^*) + \nabla f(x^*)^T (x - x^*) + \lambda^{*T} g_i(x^*) + \lambda^{*T} \nabla g_i(x^*)^T (x - x^*) + \mu^{*T} h(x^*) + \mu^{*T} \nabla h(x^*)^T (x - x^*) \\ &\geq f(x^*) + (\nabla f(x^*) + \lambda^{*T} \nabla g_i(x^*)^T + \nabla h(x^*) \mu^*)(x - x^*) = f(x^*)\end{aligned}$$

Quindi abbiamo dimostrato che:

$$f(x) \geq f(x^*) \quad \forall x \text{ ammissibile}$$

□

L'iter da seguire è il seguente:

- KKT se soddisfatta allora x^* è un candidato all'ottimo. Altrimenti:
 - x^* è regolare:
 - Si, il punto viene scartato;

– Verifico le condizioni di Fritz John con $\lambda_0 = 0$:

- Sì, allora è un candidato;
- No, si scarta il punto

10 Cenni sulla Teoria Statistica

Consideriamo un problema di **classificazione binario supervisionato** in cui, noto un data set di cui conosco le etichette, vogliamo trarne dei pattern “*nascosti*” all’interno dei dati. Quindi:

$$T = \{(x^i, y^i), x^i \in \mathbb{R}^n, y^i \in \{-1, 1\}, i = 1, \dots, l\}$$

Ipotesi:

\exists una distribuzione di probabilità $P(x, y)$ **non nota** che lega x e y .

Inoltre, le osservazioni del training set siano state generate in maniera indipendente e identicamente distribuita secondo $P(x, y)$

Tesi:

Vogliamo implementare una macchina di apprendimento che impari la relazione $x \rightarrow y$. Ciò equivale ad ottenere un modello $f(\alpha): \mathbb{R}^n \rightarrow \{-1, 1\}$ con $\alpha :=$ **parametro modificabile**.

La $f(\alpha)$ è **deterministica** poiché, una volta fissato α , ad ogni ingresso x fornisce un output y . Quello che dobbiamo fare è addestrare la macchina di apprendimento scegliendo il parametro α .

A tale scopo occorre introdurre una misura che mi quantifica l’errore effettuato nella classificazione del test e minimizzarla. Questa misura viene detta **rischio effettivo** (coincide con il valore atteso \mathbb{E} data la probabilità):

$$R(\alpha) = \int \frac{1}{2} |y - f(x, \alpha)| dP(x, y)$$

In cui:

- y è l’uscita reale;
- $f(x, \alpha)$ è l’uscita ottenuta dalla macchina;

Tuttavia il **rischio effettivo** non è calcolabile dato che abbiamo supposto che la distribuzione di probabilità non è nota. Un’alternativa viene data dal **rischio empirico**:

$$R_{\text{emp}}(\alpha) = \frac{1}{l} \sum_{i=1}^l |y^i - f(x^i, \alpha)|$$

Occorre notare che abbiamo scorrelato il rischio dalla conoscenza della distribuzione di probabilità, ma l’abbiamo legata ai dati del nostro training set e quindi è una misura che possiamo calcolare.

Nota 66.

In generale non vi è una relazione tra rischio empirico e rischio effettivo.

Inoltre, se confrontassimo due rischi empirici di due macchine diverse anche se possono avere lo stesso valore si potrebbero comportare in maniera diversa in presenza di dati del test set.

10.1 Teoria di Vapinil-Chervonenkis

Sia $\eta \in [0, 1]$ con probabilità $1 - \eta$ vale la seguente diseguaglianza:

$$R(\alpha) \leq R_{\text{emp}} + \sqrt{\frac{h \log\left(\frac{2l}{h}\right) - \log\left(\frac{\eta}{4}\right)}{l}}$$

In cui definiamo la **VC Confidence**:

$$\text{VCC} = \sqrt{\frac{h \log\left(\frac{2l}{h}\right) - \log\left(\frac{\eta}{4}\right)}{l}}$$

essa dipende da η, l, h . L'ultimo parametro è detto **VC Dimension** e misura la capacità di classificazione di una macchina $f(\alpha)$.

Quindi il termine a sinistra è composto da una media pesata degli errori del training set più la VC Confidence che dipende dalla complessità della macchina, dal numero di feature e da η , costituendo così un upper-bound per l'errore reale.

$$R(\alpha) \leq R_{\text{emp}}(\alpha) + \text{VC}(h, \eta, l) = U(\alpha, h)$$

Nella progettazione della nostra macchina vogliamo un **trade-off** fra complessità della macchina e il rischio reale. A tale scopo, occorre minimizzare l'upper bound $U(\alpha, h)$.

Definizione 67. VC dimension

h è il massimo numero di punti che possono essere **frammentati**, ovvero che possono essere separati tramite la macchina di apprendimento $f(\alpha)$ in due classi quando etichettati in tutti i modi possibili.

La nuova funzione obiettivo da massimizzare è il **rischio strutturale**:

$$U(\alpha, h) = R_{\text{emp}}(\alpha) + \text{VCC}(h, \alpha)$$

Si definiscono delle classi di funzioni di apprendimento:

$$F_1 \subseteq F_2 \subseteq \dots \subseteq F_p$$

$$h_1 \leq h_2 \leq \dots \leq h_p$$

Procedura euristica per minimizzare il rischio strutturale

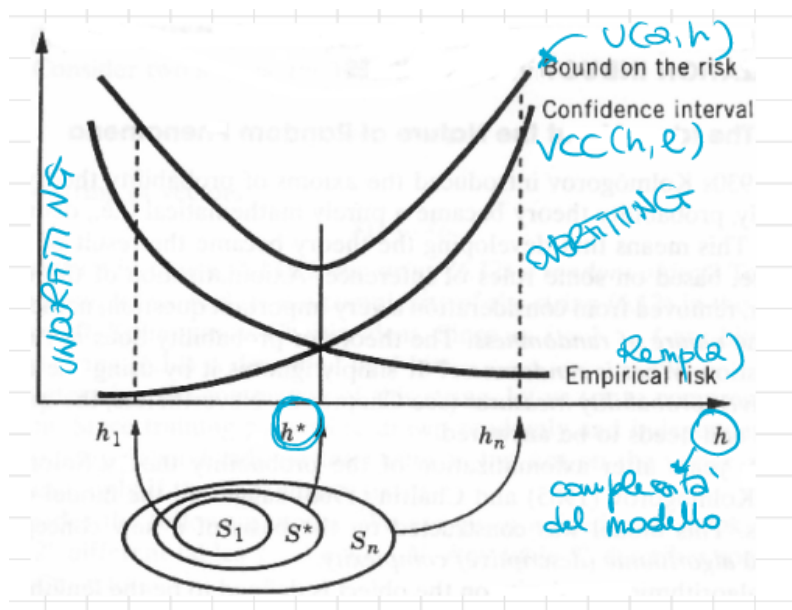
Per ogni classi F_j con VC dimension h_j , minimizza il rischio empirico $R_{\text{emp}}(\alpha)$:

$$\alpha_j = \arg \min_{F_j} R_{\text{emp}}(\alpha)$$

Inoltre, calcola:

$$U^j = R_{\text{emp}}(\alpha_j) + \text{VCC}(h_j, l) \quad j = 1, \dots, p$$

Scelgo la classe di macchine che mi restituisce il valore più basso di U^j .



Nota 68. h non è necessariamente proporzionale alla dimensione di α .

Si hanno due approcci possibili per minimizzare $U(\alpha, w)$:

1. Fisso h e minimizzo $R_{\text{emp}}(\alpha) \rightarrow$ **reti neurali**;
2. Fisso $R_{\text{emp}}(\alpha) = 0$ e minimizzo $VCC(h, l) \rightarrow$ **SVM** (Super Vector Machine);

11 Iperpiani Orientati

Un **iperpiano orientato** è definito come:

$$f(\alpha) = \text{sgn}(w^T x + b) \quad \alpha = \begin{pmatrix} w \\ b \end{pmatrix} \quad x \in \text{TS}$$

Teorema 69.

Un insieme di punti $\{x^1, \dots, x^m\}$ è frammentabile dalla famiglia degli iperpiani orientati $f(\alpha) = \text{sgn}(w^T x + b) \iff \{x^1, \dots, x^m\}$ sono **affinemente indipendenti** cioè $x_2 - x_1, \dots, x_m - x_1$ sono linearmente indipendenti.

Importante 70.

h è il massimo numero di punti affinemente indipendenti in \mathbb{R}^n , quindi $h = n + 1$.

Definizione 71. Margine

Il **margin** è la distanza minima dell'iperpiano orientato dai punti del training set.

Intuizione

Quanto più l'iperpiano è lontano dai dati di training, tanto più sarà migliore la sua capacità di generalizzazione.

Idea

Ci si restringe ad un sottocaso degli iperpiani detti **iperpiani orientati con gap di tolleranza** (margin).

Ipotesi: Supponiamo che i miei dati appartengano ad una sfera di diametro D .

$$y = H(w, b) = \begin{cases} 1 & \text{se } \frac{w^T x + b}{\|w\|} \geq \rho \\ -1 & \text{se } \frac{w^T x + b}{\|w\|} \leq -\rho \end{cases}$$

Nota 72.

- $\frac{w^T x + b}{\|w\|} :=$ è la distanza di x dall'iperpiano;
- $\rho = 0 \longrightarrow$ iperpiani orientati.

Ora devo calcolare la VC dimensione, essa dipende sia da ρ che da D .

11.1 Iperpiani con gap di tolleranza

Supponiamo che i dati siano generati all'interno di una sfera di diametro D e che si assegnino le uscite:

$$y = \begin{cases} 1 & \text{se } \frac{w^T x + b}{\|w\|} \geq \rho \\ -1 & \text{se } \frac{w^T x + b}{\|w\|} \leq -\rho \end{cases}$$

Definizione 73. *Insiemi del data set*

$$A = \{x^i \in T: y^i = 1\}$$

$$B = \{x^j \in T: y^j = -1\}$$

11.1.1 Variazione di VC dimension

L'intuizione è che all'aumentare di ρ , h diminuisce.

Teorema 74. *VC dimension*

Per gli iperpiani con margine si ha:

$$h \leq \min \left\{ \left\lceil \frac{D^2}{\rho^2} \right\rceil, n \right\} + 1$$

Nella minimizzazione del *rischio strutturale* si prende:

$$F_i = \{f(x) = w^T x + b, \rho(w, b) \geq \rho_i\}$$

Definizione 75. *Margine Iperpiano con gap di tolleranza*

Definiamo il **margine dell'iperpiano** come:

$$\rho(w, b) = \min_{x^i \in T} \left\{ \frac{|w^T x^i + b|}{\|w\|} \right\}$$

Data una famiglia di iperpiani F_i vogliamo determinare l'iperpiano che minimizza l'upper-bound così definito:

$$R(\alpha) \leq R_{\text{emp}}(\alpha) + \text{VCC}(h, l, \eta) = U(\alpha, h)$$

così facendo otteniamo l'**iperpiano ottimo** che ha un margine massimo ed ha un livello di generalizzazione maggiore su nuovi dati appartenenti al test set.

Ipotesi: A, B sono linealmente separabili e quindi:

$$\exists(\bar{w}, \bar{b}) \text{ che separa correttamente } T \text{ in due insiemi } A \text{ e } B$$

In particolare, dato un *training set*:

$$TS = \{(x^i, y^i) \mid i = 1, \dots, l, x^i \in \mathbb{R}^n, y^i \in \{-1, 1\}\}$$

$$A = \{x^i \in T : y^i = 1\}$$

$$B = \{x^j \in T : y^j = -1\}$$

I due insiemi sono **linearmente separabili** se $\exists(w, b) : y(x) = \text{sgn}(w^T x + b)$ classifica correttamente il TS , ovvero:

$$w^T x^i + b > 0 \quad x^i \in A$$

$$w^T x^j + b < 0 \quad x^j \in B$$

Poiché i risolutori si comportano male in presenza di disequazioni non strette, riscriviamo le relazioni superiori nel seguente modo:

$$w^T x^i + b \geq \varepsilon \quad x^i \in A$$

$$w^T x^j + b \leq -\varepsilon \quad x^j \in B$$

Inoltre, dato che ε è uno scalare >0 , possiamo rinominare $w = \frac{w}{\varepsilon}$ e $b = \frac{b}{\varepsilon}$ ottenendo:

$$w^T x^i + b \geq 1 \quad x^i \in A$$

$$w^T x^j + b \leq -1 \quad x^j \in B$$

Che corrisponde ad avere un $R_{\text{emp}}(w, b) = 0$.

Fra tutti questi iperpiano, come da ipotesi, vogliamo determinare quello ottimo e con il margine massimo:

$$\begin{aligned} \max_{w, b} \rho(w, b) &= \min_{x^i \in A \cup B} \left\{ \frac{|w^T x^i + b|}{\|w\|} \right\} \\ w^T x^i + b &\geq 1 \quad \forall x^i \in A \\ w^T x^j + b &\leq -1 \quad \forall x^j \in B \end{aligned}$$

Questo problema combina un problema di massimo e di minimo, il che può essere riscritto come:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 \\ w^T x^i + b &\geq 1 \quad \forall x^i \in A \\ w^T x^j + b &\leq -1 \quad \forall x^j \in B \end{aligned}$$

Nota 76.

Questo problema è convesso poiché, i vincoli sono lineari e la funzione obiettivo è quadratica e anch'essa convessa.

Lemma 77.

Per ogni iperpiano ammissibile (\hat{w}, \hat{b}) si ha:

$$\rho(\hat{w}, \hat{b}) \geq \frac{1}{\|\hat{w}\|}$$

Dimostrazione.

$$\rho(w, b) = \min_{x^i \in A \cup B} \left\{ \frac{|w^T x^i + b|}{\|w\|} \right\}$$

Poiché è ammissibile:

$$w^T x^i + b \geq 1 \quad \forall x^i \in A \rightarrow \forall x^i \in A \quad |w^T x^i + b| \geq 1 = w^T \hat{x}^i + \hat{b}$$

$$w^T x^j + b \leq -1 \quad \forall x^j \in B \rightarrow |w^T x^j + b| = -(w^T x^j + b) \geq 1$$

Quindi:

$$|w^T x^i + b| \geq 1 \quad \forall x^i \in A \cup B$$

$$\min_{x^i \in A \cup B} \{|w^T x^i + b|\} \geq 1 \rightarrow \rho(\hat{w}, \hat{b}) \geq \frac{1}{\|\hat{w}\|} \quad \square$$

Teorema 78.

Per ogni iperpiano ammissibile (\hat{w}, \hat{b}) , esiste un altro iperpiano ammissibile (\bar{w}, \bar{b}) tale che:

$$\rho(\bar{w}, \bar{b}) = \frac{1}{\|\bar{w}\|} \geq \rho(\hat{w}, \hat{b})$$

Inoltre, esistono due punti $x^+ \in A$ e $x^- \in B$ tali che:

$$w^T x^+ + b = 1 \quad w^T x^- + b = -1$$

Dimostrazione.

Prendo $x^i \in A$ e $x^j \in B$ più vicini, per calcolare la distanza definisco:

$$\hat{d}_i = \min_{x^i \in A} \left\{ \frac{|\hat{w}^T x^i + \hat{b}|}{\|\hat{w}\|} \right\}$$

$$\hat{d}_j = \min_{x^j \in B} \left\{ \frac{|\hat{w}^T x^j + \hat{b}|}{\|\hat{w}\|} \right\}$$

Per l'insieme A :

$$\hat{d}_i = \frac{|\hat{w}^T \hat{x}^i + \hat{b}|}{\|\hat{w}\|} \leq \frac{|\hat{w}^T x^i + \hat{b}|}{\|\hat{w}\|} \quad \forall x^i \in A \rightarrow \hat{w}^T \hat{x}^i \leq \hat{w}^T x^i$$

Per l'insieme B :

$$\hat{d}_j = \frac{|\hat{w}^T \hat{x}^j + \hat{b}|}{\|\hat{w}\|} \leq \frac{|\hat{w}^T x^j + \hat{b}|}{\|\hat{w}\|} \quad \forall x^j \in B$$

$$\frac{-\hat{w}^T \hat{x}^j - \hat{b}}{\|\hat{w}\|} \leq \frac{-\hat{w}^T x^j - \hat{b}}{\|\hat{w}\|} \rightarrow -\hat{w}^T \hat{x}^j \leq -\hat{w}^T x^j$$

Definiamo $\rho(\hat{w}, \hat{b}) = \min \{\hat{d}_i, \hat{d}_j\}$

Corollario 79. *Proprietà del minimo*

$$\min \{a, b\} \leq \frac{1}{2}a + \frac{1}{2}b$$

Dimostrazione.

$$\min \{a, b\} \leq a \quad \min \{a, b\} \leq b \longrightarrow \frac{1}{2} \min \{a, b\} \leq \frac{1}{2}a \quad \frac{1}{2} \min \{a, b\} \leq \frac{1}{2}b$$

Sommo entrambe le relazioni:

$$\min \{a, b\} \leq \frac{1}{2}a + \frac{1}{2}b \quad \square$$

Quindi, ritornando alla prima dimostrazione:

$$\begin{aligned} \rho(\hat{w}, \hat{b}) = \min \{\hat{d}_i, \hat{d}_j\} &\leq \frac{1}{2}\hat{d}_i + \frac{1}{2}\hat{d}_j = \frac{1}{2} \frac{|\hat{w}^T \hat{x}^i + \hat{b}|}{\|\hat{w}\|} + \frac{1}{2} \frac{|\hat{w}^T \hat{x}^j + \hat{b}|}{\|\hat{w}\|} \\ &= \frac{1}{2} \frac{\hat{w}^T \hat{x}^i + \hat{b}}{\|\hat{w}\|} - \frac{1}{2} \frac{\hat{w}^T \hat{x}^j + \hat{b}}{\|\hat{w}\|} \\ &= \frac{\hat{w}^T (\hat{x}^i - \hat{x}^j)}{2\|\hat{w}\|} \end{aligned}$$

Modifico (\hat{w}, \hat{b}) introducendo due scalari α , coefficiente di rotazione, e β , coefficiente di traslazione, $\longrightarrow (\bar{w}, \bar{b})$:

$$\bar{w} = \alpha \hat{w}$$

$$\bar{b} = \beta$$

Questi due parametri devono soddisfare le seguenti relazioni:

$$\begin{cases} \alpha \hat{w}^T \hat{x}^i + \beta = 1 \\ \alpha \hat{w}^T \hat{x}^j + \beta = -1 \end{cases}$$

Che rappresentano i punti sul margine dell'iperpiano. Equivalgono a vincoli attivi cioè soddisfatti all'uguaglianza. Risolviamo ora il sistema:

$$\alpha \hat{w}^T (\hat{x}^i - \hat{x}^j) = 2 \rightarrow \alpha = \frac{2}{\hat{w}^T (\hat{x}^i - \hat{x}^j)} \geq 2$$

$$\beta = 1 - \alpha \hat{w}^T \hat{x}^i = 1 - \frac{2\hat{w}^T \hat{x}^i}{\hat{w}^T (\hat{x}^i - \hat{x}^j)} = \frac{\hat{w}^T \hat{x}^i - \hat{w}^T \hat{x}^j - 2\hat{w}^T \hat{x}^i}{\hat{w}^T (\hat{x}^i - \hat{x}^j)} = \frac{\hat{w}^T (\hat{x}^i + \hat{x}^j)}{\hat{w}^T (\hat{x}^i - \hat{x}^j)}$$

Nota 80.

$$0 < \alpha \leq 1$$

A questo punto dobbiamo verificare che il nuovo iperpiano (\bar{w}, \bar{b}) sia ammissibile. A tale scopo, so che:

$$1 \leq \hat{w}^T \hat{x}^i + \hat{b} \leq \hat{w}^T \hat{x}^j + \hat{b} \quad \forall x^i \in A$$

$$\begin{aligned} |\hat{w}^T \hat{x}^i + \hat{b}| &\leq |\hat{w}^T \hat{x}^i + \hat{b}| & \forall x^i \in A \\ |\hat{w}^T \hat{x}^j + \hat{b}| &\leq |\hat{w}^T \hat{x}^j + \hat{b}| & \forall x^j \in B \longrightarrow |\hat{w}^T \hat{x}^j + \hat{b}| \geq 1 \end{aligned}$$

Quindi:

$$\begin{aligned}\bar{w}^T x^i + \bar{b} &= \alpha \hat{w}^T x^i + \beta \quad \text{Ricordando che } \hat{w}^T \hat{x}^i \leq \hat{w}^T x^i \\ &\geq \alpha \hat{w}^T \hat{x}^i + \beta = 1 \quad \forall x^i \in A\end{aligned}$$

e

$$\begin{aligned}\bar{w}^T x^j + \bar{b} &= \alpha \hat{w}^T x^j + \beta \quad \text{Ricordando che } -\hat{w}^T \hat{x}^j \leq -\hat{w}^T x^j \\ &\leq \alpha \hat{w}^T \hat{x}^j + \beta = -1 \quad \forall x^j \in B\end{aligned}$$

Quindi l'iperpiano ottenuto ruotando e traslando è ancora ammissibile.

In conclusione, (\bar{w}, \bar{b}) è l'iperpiano di separazione:

$$\rho(\bar{w}, \bar{b}) = \min_{x^i \in A \cup B} \left\{ \frac{|\bar{w}^T x^i + \bar{b}|}{\|\bar{w}\|} \right\} = \frac{1}{\|\bar{w}\|} = \frac{1}{\alpha \|\hat{w}\|}$$

Nota 81.

Il numeratore vale 1 per costruzione.

$$= \frac{1}{\alpha \|\hat{w}\|} = \frac{\hat{w}^T (\hat{x}^i - \hat{x}^j)}{2 \|\hat{w}\|} = \frac{1}{2} \hat{d}_i + \frac{1}{2} \hat{d}_j \geq \rho(\hat{w}, \hat{b})$$

Inoltre, so che:

$$\begin{cases} \bar{w}^T \hat{x}^i + \beta = 1 \\ \bar{w}^T \hat{x}^j + \beta = -1 \end{cases} \longrightarrow \begin{cases} x^+ = x^i \\ x^- = x^j \end{cases}$$

Che erano l'ipotesi del teorema. □

Ricapitolando, l'**iperpiano ottimo**:

$$\max_{w, \rho} \rho(w, b) = \min_{x^i \in A \cup B} \left\{ \frac{|w^T x^i + b|}{\|w\|} \right\}$$

$$w^T x^i + b \geq 1 \quad \forall x^i \in A$$

$$w^T x^j + b \leq -1 \quad \forall x^j \in B$$

Lemma 82.

Per ogni iperpiano ammissibile (\hat{w}, \hat{b}) si ha $\rho(\hat{w}, \hat{b}) \geq \frac{1}{\|\hat{w}\|}$

Teorema 83.

Per ogni iperpiano ammissibile (\hat{w}, \hat{b}) esiste un iperpiano (\bar{w}, \bar{b}) tale che:

$$\rho(\bar{w}, \bar{b}) = \frac{1}{\|\bar{w}\|} \geq \rho(\hat{w}, \hat{b})$$

Inoltre esistono $x^+ \in A, x^- \in B$ tali che:

$$\bar{w} x^+ + \bar{b} = 1$$

$$\bar{w} x^- + \bar{b} = -1$$

Vogliamo far vedere che il problema iniziale di max min è equivalente a :

$$\min \frac{1}{2} \|w\|^2 = f(w, b)$$

$$w^T x^i + b \geq 1 \quad \forall x^i \in A$$

$$w^T x^j + b \leq -1 \quad \forall x^j \in b$$

L'insieme ammissibile è chiuso ma potrebbe essere illimitato. Inoltre, è anche convesso e non vuoto:

$$\nabla f = \begin{pmatrix} w \\ 0 \end{pmatrix} \quad \nabla^2 f = \begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix} \succcurlyeq 0$$

Quindi f è convessa e coerciva solo rispetto a w . Tuttavia, non è né convessa né coerciva né vale Weierstrass-Ermann, occorre quindi determinare l'unicità della soluzione.

Teorema 84.

Il problema equivalente ammette una soluzione unica (w^, b^*) .*

Dimostrazione.

L'insieme ammissibile è chiuso, ma non limitato. Dobbiamo, quindi, far vedere l'insieme di livello è compatto, se lo è allora il problema ammette soluzione.

$$L_0 = \left\{ (w, b) \in F : \frac{1}{2} \|w\|^2 \leq \frac{1}{2} \|w_0\|^2 \right\}$$

$$F(w, b) : w^T x^i + b \geq 1 \quad \forall x^i \in A, w^T x^j + b \leq -1 \quad \forall x^j \in b$$

L_0 è chiuso poiché F è chiuso, devo far vedere che è limitato.

Suppongo per assurdo che $\exists \{w^k, b^k\}$ tale che $w^k, b^k \in L_0$ $\left\| \begin{pmatrix} w^k \\ b^k \end{pmatrix} \right\| \uparrow \infty$

Dato che l'insieme di livello è definito come $\|w^k\|^2 \leq \|w_0\|^2 \longrightarrow \|w^k\|$ è finita, suppongo che $\lim_{b \rightarrow \infty} b^k = \infty$:

$$\begin{pmatrix} w^k \\ b^k \end{pmatrix} \in F \longrightarrow w^{k^T} x^i + b^k \geq 1 \quad \forall x^i \in A$$

$$w^{k^T} x^j + b^k \leq -1 \quad \forall x^j \in B$$

Se $b^k \rightarrow \infty$ il primo gruppo di vincoli è soddisfatto. Ci rimane da soddisfare il secondo:

$$w^{k^T} x^j + b^k \leq -1 \quad \forall x^j \in B$$

Per essere soddisfatto deve valere $w^{k^T} x^j \rightarrow -\infty \longrightarrow \|w^k\| \rightarrow \infty$. Ciò comporta un assurdo poiché la successione non appartiene più a L_0 .

Allo stesso modo se $b^k \rightarrow -\infty$ devo soddisfare:

$$w^{k^T} x^i + b^k \geq 1 \quad \forall x^i \in A \longrightarrow w^{k^T} x^i \rightarrow \infty \longrightarrow \|w^k\| \rightarrow \infty$$

Ciò comporta un assurdo. Quindi, L_0 è un insieme compatto e quindi esiste soluzione.

L'unico aspetto che ci occorre dimostrare che la soluzione è unica: suppongo per assurdo che non lo sia. In particolare:

$$\exists(w^*, b^*), (\bar{w}, \bar{b}) \text{ ottime}$$

Quindi:

$$\frac{1}{2}\|w^*\|^2 = \frac{1}{2}\|\bar{w}\|^2 \leq \frac{1}{2}\|w\|^2 \quad \forall \begin{pmatrix} w \\ b \end{pmatrix} \text{ ammissibile}$$

Suppongo che $\bar{w} \neq w^*$. Siccome l'insieme è convesso, tramite la definizione di convessità:

$$\exists \lambda \in [0, 1] \quad \lambda \begin{pmatrix} w^* \\ b^* \end{pmatrix} + (1 - \lambda) \begin{pmatrix} \bar{w} \\ \bar{b} \end{pmatrix} \in F$$

Scegliamo $\lambda = \frac{1}{2}$:

$$\begin{aligned} \begin{pmatrix} \tilde{w} \\ \tilde{b} \end{pmatrix} &= \frac{1}{2} \begin{pmatrix} \bar{w} \\ \bar{b} \end{pmatrix} + \frac{1}{2} \begin{pmatrix} w^* \\ b^* \end{pmatrix} \\ f(\tilde{w}, \tilde{b}) &= \frac{1}{2} \left\| \frac{1}{2} \bar{w} + \frac{1}{2} w^* \right\|^2 \end{aligned}$$

Calcoliamo il valore della funzione obiettivo. La funzione obiettivo è strettamente convessa in w se e solo se:

$$f\left(\lambda \begin{pmatrix} w \\ b \end{pmatrix} + (1 - \lambda) \begin{pmatrix} w^* \\ b^* \end{pmatrix}\right) < \lambda f\left(\begin{pmatrix} \bar{w} \\ \bar{b} \end{pmatrix}\right) + (1 - \lambda) f\left(\begin{pmatrix} w^* \\ b^* \end{pmatrix}\right)$$

In \tilde{w} :

$$f(\tilde{w}, \tilde{b}) = \frac{1}{2} \|\tilde{w}\|^2 < \frac{1}{2} \frac{1}{2} \|w^*\|^2 + \frac{1}{2} \frac{1}{2} \|\bar{w}\|^2 = \frac{1}{2} \|w^*\|^2$$

Questo è un assurdo poiché avremmo trovato un iperpiano ammissibile migliore di quello ottimo e quindi deve valere necessariamente che:

$$\bar{w} = w^*$$

Ora ci rimane da controllare il valore del termine noto b : supponiamo $b^* > \bar{b}$.

So che $\exists x^+ \in A$:

$$w^{*T} x^+ + b^* = 1 = \bar{w}^T x^+ + b^* > \bar{w}^T x^+ + \bar{b} \longrightarrow \bar{w}^T x^+ + \bar{b} < 1 \longrightarrow \text{ASSURDO}$$

Questo è un assurdo poiché $x^+ \in A$ e $\begin{pmatrix} \bar{w} \\ \bar{b} \end{pmatrix}$ è ammissibile. Allo stesso modo, supponiamo che $b^* < \bar{b}$:

$$w^{*T} x^- + b^* = -1 = \bar{w}^T x^- + b^* < \bar{w}^T x^- + \bar{b} \longrightarrow \bar{w}^T x^- + \bar{b} > 1 \longrightarrow \text{ASSURDO}$$

Questo è un assurdo poiché $x^- \in B$ e $\begin{pmatrix} \bar{w} \\ \bar{b} \end{pmatrix}$ è ammissibile.

Quindi per valere entrambe le condizioni, dobbiamo avere che:

$$\bar{b} = b^*$$

La soluzione, infine, è unica. □

Teorema 85.

L'unica soluzione $\begin{pmatrix} w^ \\ b^* \end{pmatrix}$ del problema:*

$$\min \frac{1}{2} \|w\|^2$$

$$w^T x^i + b \geq 1 \quad \forall x^i \in A$$

$$w^T x^j + b \leq -1 \quad \forall x^j \in B$$

è l'unica soluzione del problema:

$$\max \rho(w, b)$$

$$w^T x^i + b \geq 1 \quad \forall x^i \in A$$

$$w^T x^j + b \leq -1 \quad \forall x^j \in B$$

Dimostrazione.

Dal lemma 82 so che per ogni iperpiano ammissibile:

$$\rho(\hat{w}, \hat{b}) \geq \frac{1}{\|\hat{w}\|}$$

Inoltre, so che per ogni iperpiano ammissibile $\begin{pmatrix} \hat{w} \\ \hat{b} \end{pmatrix}$ esiste un iperpiano $\begin{pmatrix} \bar{w} \\ \bar{b} \end{pmatrix}$ tale che:

$$\rho(\hat{w}, \hat{b}) \leq \rho(\bar{w}, \bar{b}) = \frac{1}{\|\bar{w}\|}$$

In aggiunta, so che:

$$\frac{1}{2} \|w^*\|^2 \leq \frac{1}{2} \|w\|^2 \quad \forall \begin{pmatrix} w \\ b \end{pmatrix} \text{ ammissibile}$$

$$\longrightarrow \|w^*\| \leq \|w\| \quad \forall \begin{pmatrix} w \\ b \end{pmatrix} \text{ ammissibile}$$

$$\longrightarrow \frac{1}{\|w^*\|} \geq \frac{1}{\|w\|} \quad \forall \begin{pmatrix} w \\ b \end{pmatrix} \text{ ammissibile}$$

Per qualunque $\begin{pmatrix} \hat{w} \\ \hat{b} \end{pmatrix}$ ammissibile:

$$\frac{1}{\|\hat{w}\|} \leq \rho(\hat{w}, \hat{b}) \leq \rho(\bar{w}, \bar{b}) \leq \frac{1}{\|w^*\|}$$

Allora:

$$\rho(\hat{w}, \hat{b}) \leq \rho(w^*, b^*) \quad \forall \begin{pmatrix} \hat{w} \\ \hat{b} \end{pmatrix}$$

Cioè è l'iperpiano a massimo margine.

Abbiamo dimostrato che è l'**ottimo globale**. □

Il problema della riformulazione del problema è dato dalla presenza di vincoli per ogni punto e quindi risulta oneroso dal punto di vista computazionale.

12 Teoria della Dualità

12.1 Introduzione

Il problema che vogliamo risolvere è:

$$\begin{cases} \min \frac{1}{2} \|w\|^2 \\ y^i(w^T x^i + b) \geq 1 \quad \forall i = 1, \dots, l \end{cases}$$

A tale scopo generiamo una formulazione alternativa del problema, detto problema duale, e, di conseguenza, otteniamo i seguenti vantaggi:

1. Teorici: otteniamo risultati e informazioni sul problema originario → condizioni di ottimo;
2. Pratici: il duale è più semplice da risolvere → il numero di vincoli diminuisce mentre il numero delle variabili aumentano.

Quindi la situazione che abbiamo davanti è del tipo:

Problema Primale

$$\begin{cases} \min f(x) \\ x \in S \end{cases}$$

Problema Duale

$$\begin{cases} \max \psi(u) \\ u \in U \end{cases}$$

Tipicamente vale la proprietà di **dualità debole**:

$$\sup \psi(u) \leq \inf f(x)$$

$$u \in U \quad x \in S$$

In alcuni casi, cioè quando il problema è convesso, vale la **dualità forte**:

$$\sup \psi(u) = \inf f(x)$$

$$u \in U \quad x \in S$$

La dualità ci permette di ottenere:

1. Stima del valore ottimo (lower-bound) del primale;
2. Condizioni di ottimo;
3. Algoritmi di soluzione.

12.2 Dualità di Wolfe

Dato un problema **primale** non lineare convesso, del tipo:

$$\begin{cases} \min f(x) \\ g_i(x) \leq 0 & i = 1, \dots, m \\ h_j(x) = 0 & j = 1, \dots, p \end{cases}$$

$$f, g_i, h_j \in C^1 \quad i = 1, \dots, m \quad j = 1, \dots, p$$

In particolare:

$$f, g_i \quad i = 1, \dots, m \quad \text{convesse}$$

$$h_j \quad j = 1, \dots, p \quad \text{lineare}$$

Quindi abbiamo che le condizioni di KKT risultano necessarie e sufficienti di ottimo globale:

$$L(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^p \mu_j h_j(x)$$

Definiamo a questo punto il **duale di Wolfe**:

$$\begin{cases} \max_{x, \lambda, \mu} L(x, \lambda, \mu) \\ \nabla_x L(x, \lambda, \mu) = 0 = \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla g_i(x) + \sum_{j=1}^p \mu_j \nabla h_j(x) \\ \lambda \geq 0 \end{cases}$$

Teorema 86.

Se il problema **primale** che ammette soluzione x^* con moltiplicatori (λ^*, μ^*) allora (x^*, λ^*, μ^*) è soluzione del duale e il **gap di dualità** è nullo, cioè:

$$f(x^*) = L(x^*, \lambda^*, \mu^*)$$

Dimostrazione.

(x^*, λ^*, μ^*) è soluzione del problema primale \longleftrightarrow le condizioni di KKT sono soddisfatte:

$$\nabla_x L(x^*, \lambda^*, \mu^*) = 0 \quad (1)$$

$$\lambda^* \geq 0 \quad (2)$$

$$\lambda_i g_i(x^*) = 0 \quad i = 1, \dots, m$$

$$g_i(x^*) \leq 0 \quad i = 1, \dots, m$$

$$h_j(x^*) = 0 \quad j = 1, \dots, p$$

La (1)+(2) forniscono la soluzione (x^*, λ^*, μ^*) ammissibile per il duale.

$$L(x^*, \lambda^*, \mu^*) = f(x^*) + \sum_{i=1}^m \lambda_i g_i(x^*) + \sum_{j=1}^p \mu_j h_j(x^*) = f(x^*)$$

Nota 87.

- $\sum_{i=1}^m \lambda_i g_i(x) = 0$ per complementarità;
- $\sum_{j=1}^p \mu_j h_j(x) = 0$ per ammissibilità

Quindi il gap di dualità è nullo.

A questo punto devo far vedere che la soluzione è **ottima** per il duale. Quindi devo far vedere che:

$$L(x^*, \lambda^*, \mu^*) \geq L(x, \lambda, \mu) \quad \forall (x, \lambda, \mu) \text{ ammissibile per il duale}$$

$$L(x^*, \lambda^*, \mu^*) = f(x^*) \geq f(x^*) + \sum_{i=1}^m \lambda_i g_i(x^*) + \sum_{j=1}^p \mu_j h_j(x^*) = L(x^*, \lambda, \mu) \quad (\alpha)$$

Nota 88.

- λ_i poiché λ è ammissibile per il duale;
- $g_i(x^*) \leq 0$ per l'ammissibilità di x^* nel problema primale;
- $h_j(x^*) = 0$ poiché x^* è ammissibile

Otteniamo quindi:

$$L(x^*, \lambda^*, \mu^*) \geq L(x^*, \lambda, \mu)$$

Inoltre sappiamo che $L(x, \lambda, \mu)$ è convessa e appartiene a C^1 . Allora:

$$L(x^*, \lambda, \mu) \geq L(x, \lambda, \mu) + \nabla_x L(x, \lambda, \mu)^T (x^* - x) \quad \forall (x, \lambda, \mu) \text{ ammissibile nel duale } (\beta)$$

Nota 89.

- $\nabla_x L(x, \lambda, \mu)^T (x^* - x) = 0$ per l'ammissibilità duale

Infine, unendo $(\alpha) + (\beta)$, si giunge a:

$$L(x^*, \lambda^*, \mu^*) \geq L(x, \lambda, \mu) \quad \forall (x, \lambda, \mu) \text{ ammissibile per il duale}$$

Il che implica che la terna (x^*, λ^*, μ^*) è ottima per il duale. □

Nota 90.

Dopo aver effettuato queste operazioni di trasformazione del problema primale, il duale non è più lineare e non vale neanche la convessità. Ciò corrisponde ad un aumento di complessità. Tuttavia nel nostro caso, con funzione obiettivo quadratica e vincoli lineari, il problema si semplifica.

12.3 Programmazione Quadratica - Duale

Dato il problema dell'iperpiano ottimo:

$$\begin{cases} \min \frac{1}{2} x^T Q x + c^T x \\ A x \geq b \\ A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m \end{cases} \quad Q \succcurlyeq 0$$

Definiamo la Lagrangiana.

$$L(x, \lambda) = \frac{1}{2} x^T Q x + c^T x + \lambda^T (A x - b)$$

Scriviamo il duale:

$$\max_{x, \lambda} \frac{1}{2} x^T Q x + c^T x + \lambda^T (A x - b)$$

$$Q x + c + A^T \lambda = 0$$

$$\lambda \geq 0$$

Moltiplicando il vincolo per x^T :

$$x^T Q x + c^T x - x^T A^T \lambda = 0 \longrightarrow c^T x + \lambda^T A x = -x^T Q x$$

e sostituisco nella funzione obiettivo:

$$\begin{aligned} \max \frac{1}{2} x^T Q x - x^T Q x - \lambda^T b &= -\frac{1}{2} x^T Q x - \lambda^T b \\ &\longrightarrow -\min \frac{1}{2} x^T Q x + b^T \lambda \end{aligned}$$

Quindi posso ora scrivere il duale:

$$\begin{cases} -\min_{x, \lambda} \frac{1}{2} x^T Q x + b^T \lambda = \theta(x, \lambda) \\ Q x + c + A^T \lambda = 0 \\ \lambda \geq 0 \end{cases} \quad \textbf{Problema Convesso}$$

$$\nabla^2 \theta(x, \lambda) = \begin{pmatrix} Q & 0 \\ 0 & 0 \end{pmatrix} \succcurlyeq 0 \quad \text{semidefinito positivo} \rightarrow \text{solo strettamente convessa}$$

Quindi le condizioni di KKT sono necessarie e sufficienti di ottimo globale.

Teorema 91.

Dato il problema primale e duale:

$$\begin{cases} \min \frac{1}{2} x^T Q x + c^T x \\ A x \leq b \end{cases} \quad (P) \quad \begin{cases} -\min \frac{1}{2} x^T Q x + b^T \lambda \\ Q x + c + A^T \lambda = 0 \\ \lambda \geq 0 \end{cases} \quad \begin{matrix} [v] \\ [z] \end{matrix} \quad (D)$$

La condizione necessaria è che $Q \succcurlyeq 0$.

Sia $(\bar{x}, \bar{\lambda})$ soluzione del duale, allora $\exists x^*$, non necessariamente uguale a \bar{x} , tale che:

1. $Q(\bar{x} - x^*) = 0$;
2. x^* soluzione del problema primale (P);
3. $(\bar{x}, \bar{\lambda})$ è una coppia di minimo globale-moltiplicatore del primale (P);

Nota 92.

La x , passando dal primale al duale può cambiare, ma la $\bar{\lambda}$ è lo stesso.

Dimostrazione.

$(\bar{x}, \bar{\lambda})$ è soluzione ottima del duale (D) \longrightarrow KKT è soddisfatta nel duale.

Scriviamo ora la Lagrangiana del duale W con le corrispettive condizioni di KKT:

$$\begin{aligned} W(x, \lambda, v, z) &= \frac{1}{2}x^T Qx + b^T \lambda - v^T(Qx + c + A^T \lambda) - z^T \lambda \\ \nabla_x W &= Qx - Qv = 0 \\ \nabla_\lambda W &= b - Av - z = 0 \\ z^T \lambda &= 0 \\ z &\geq 0 \\ Q\bar{x} + c + A^T \bar{\lambda} &= 0 \\ \lambda &\geq 0 \end{aligned} \tag{a}$$

$$\tag{\Delta}$$

Sono soddisfatte in $\bar{x}, \bar{\lambda}$ poiché per le ipotesi di KKT. Ora, devo trovare la soluzione x^* del primale. A tale scopo, scriviamo KKT del primale e si deve avere:

$$\begin{aligned} \nabla_x L &= 0 = -Qx + c + A^T \lambda \\ \lambda^T (Ax - b) &= 0 \\ \lambda &\geq 0 \\ Ax &\leq b \end{aligned}$$

Da (a) è possibile ricavarsi $z = b - Av$ che, sostituito nella complementarietà del duale:

$$\bar{\lambda}^T (b - Av) = 0 \longrightarrow z \geq 0 \text{ (ammissibilità)} \longrightarrow Av \leq b \text{ (ammissibilità primale di } v)$$

Nota 93.

Quindi abbiamo trovato l'ammissibilità tra primale e duale con lo stesso λ .

Inoltre, so che:

$$Q(\bar{x} - v) = 0 \longrightarrow Q\bar{x} = Qv$$

Sostituendo quest'ultimo risultato in (Δ):

$$Qv + c + A^T \bar{\lambda} = 0 \longrightarrow (x, v) \text{ soddisfa KKT}$$

A questo punto, abbiamo scritto KKT per il primale, in cui:

$$\begin{aligned}
 x^* = v &\longrightarrow \text{La coppia } (x^*, \bar{\lambda}) \text{ soddisfa KKT del primale} \\
 &\longrightarrow x^* \text{ è soluzione ottima del primale} \longrightarrow x^* = v \text{ è soluzione globale } Q \bar{x} - Q x^* = 0 \\
 &\longrightarrow Q(\bar{x} - x^*) = 0
 \end{aligned}$$

□

Nota 94.

- $\bar{\lambda}$ viene trovato risolvendo il problema duale: ci si ricava la soluzione duale con cui risolvere il primale;
- Se $Q \succ 0$, $\bar{x} = x^*$ poiché $Q(\bar{x} - x^*) = 0 \iff \bar{x} - x^* = 0$

13 SVM Lineari

Il nostro obiettivo è quello di applicare la teoria della dualità della programmazione quadratica per risolvere il problema dell'iperpiano ottimo seguente:

$$\begin{cases} \min \frac{1}{2} \|w\|^2 \\ y^i (w^T x^i + b) \geq 1 \quad i = 1, \dots, l \quad (\lambda_i) \end{cases}$$

In questo tipo di problema è possibile applicare la dualità di Wolfe poiché la funzione obiettivo è convessa ed in particolare:

$$Q = \begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix} \succcurlyeq 0 \quad \text{semidefinita positiva}$$

Inoltre, i vincoli sono lineari.

Quindi iniziamo con lo scrivere la Lagrangiana di questo problema:

$$L(w, b, \lambda) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^l \lambda_i (y^i (w^T x^i + b) - 1)$$

Il corrispondente duale di Wolfe, risulterà:

$$\max_{w, b, \lambda} \frac{1}{2} w^T w - \sum_{i=1}^l \lambda_i y^i w^T x^i - \sum_{i=1}^l \lambda_i y^i b + \sum_{i=1}^l \lambda_i$$

$$\nabla_w L(w, b, \lambda) = w - \sum_{i=1}^l \lambda_i y^i x^i = 0$$

$$\frac{\partial L(w, b, \lambda)}{\partial b} = - \sum_{i=1}^l \lambda_i y^i = 0$$

$$\lambda \geq 0$$

Dall'annullamento della Lagrangiana rispetto a b otteniamo che $\sum_{i=1}^l \lambda_i y^i b = 0$; mentre dall'annullamento della Lagrangiana rispetto a w otteniamo:

$$w = \sum_{i=1}^l \lambda_i y^i x^i$$

Sostituendo l'ultima espressione nella funzione obiettivo, il problema può essere riscritto come:

$$\begin{aligned} \max \frac{1}{2} \sum_{i=1}^l \lambda_i y^i (x^i)^T \sum_{j=1}^l \lambda_j y^j x^j - \sum_{i=1}^l \lambda_i y^i (x^i)^T \sum_{j=1}^l \lambda_j y^j x^j + \sum_{i=1}^l \lambda_i = \\ = \max -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y^i y^j (x^i)^T x^j - \sum_{i=1}^l \lambda_i \end{aligned}$$

Riscrivendo la funzione obiettivo in forma otteniamo il duale di Wolfe:

$$\begin{cases} -\min \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y^i y^j (x^i)^T x^j - \sum_{i=1}^l \lambda_i \\ \sum_{i=1}^l \lambda_i y^i = 0 \\ \lambda \geq 0 \quad i = 1, \dots, l \end{cases}$$

A questo punto definiscono la forma matriciale di questo problema introducendo la matrice:

$$X = [y^1 x^1 \dots y^l x^l] \quad X \in \mathbb{R}^{n \times l}$$

In cui l'elemento ij -esimo risulterà nella forma:

$$[X^T X]_{ij} = y^i y^j (x^i)^T x^j$$

Quindi otteniamo il **duale in forma compatta**:

$$\min \frac{1}{2} \lambda^T X^T X \lambda - e^T \lambda$$

$$y^T \lambda = 0$$

$$\lambda \geq 0$$

$$\text{in cui } e = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}.$$

Nota 95.

$$Q = X^T X \longrightarrow z^T Q z = z^T X^T X z = (X z)^T X z = \|X z\|^2 \geq 0 \longrightarrow Q \succcurlyeq 0$$

Il problema è convesso in λ .

Quello a cui vogliamo rispondere è che:

- Il problema duale ammette soluzione?

Per il teorema 86 so che il problema primale ammette soluzione unica, allora il problema duale ammette soluzione.

- Come si trovano w^* e b^* una volta che ho risolto il problema duale?

Sia λ^* la soluzione del problema duale:

$$w^* = \sum_{i=1}^l \lambda_i^* y^i x^i \quad \lambda^* \geq 0$$

Gli unici punti del Training Set che contribuiscono ad individuare w sono gli x^i per cui $\lambda_i^* > 0$. I punti che soddisfano tale condizioni vengono detti **vettori di supporto**.

Risolvere il problema duale mi fornisce l'informazione sulla ridondanza del TS. Nel caso in cui $\lambda_i^* > 0$, allora il vincolo i – esimo del primale è attivo all'uguaglianza e quindi il punto si trova sul margine.

- Come si ricava b^* ?

Per ricavare b^* si utilizza il teorema della dualità di Wolfe nel caso quadratico. In particolare λ^* è il moltiplicatore associato a w^*, b^* e quindi vale la complementarità:

$$\lambda_i^* (y^i (w^{*T} x^i + b^*) - 1) = 0 \quad i = 1, \dots, l$$

Se $\lambda_i^* > 0 \longrightarrow y^i (w^{*T} x^i + b^*) = 1$. Quindi:

$$b^* = \frac{1}{y^i} - w^{*T} x^i = \frac{1}{y^i} - \sum_{j=1}^l \lambda_j^* y^j (x^j)^T x^i$$

Quindi la **SVM** diventa:

$$y(x) = \text{sgn}(w^{*T} x + b^*) = \text{sgn}\left(\sum_{i=1}^l \lambda_i^* y^i (x^i)^T x + b^*\right)$$

Riassunto, per trovare l'iperpiano ottimo (cioè quello con migliore capacità di generalizzazione tra quelli che separano i due insiemi A e B) posso risolvere o il problema primale o il duale e sono entrambi problemi di programmazione quadratica convessa.

Tutta la teoria funziona se A e B sono linearmente separabili. In caso contrario, l'insieme ammissibile è un insieme **vuoto** e quindi occorre rilassare i vincoli tramite l'aggiunta di variabili ausiliarie dette **variabili di slack**:

$$y^i (w^T x^i + b) \geq 1 - \xi_i \quad i = 1, \dots, l$$

$$\xi_i \geq 0$$

Quindi un punto x^i è **mal classificato** se i segni di y^i e $w^T x^i + b$ sono discordi. Allora:

$$y^i (w^T x^i + b) < 0$$

$$y^i (w^T x^i + b) \geq 1 - \xi_i \longleftrightarrow \xi_i \geq 1 - y^i (w^T x^i + b) \longrightarrow \xi_i > 1 \text{ se } x^i \text{ è mal classificato}$$

In particolare, $\sum_{i=1}^l \xi_i$ è un **upperbound** sul numero di punti mal classificati.

A fronte di ciò modifichiamo la funzione obiettivo del problema primale nel seguente modo;

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i$$

$$\begin{aligned} y^i (w^T x^i + b) &\geq 1 - \xi_i & i = 1, \dots, l & \quad \lambda_i \\ \xi_i &\geq 0 & i = 1, \dots, l & \quad \mu_i \end{aligned}$$

Nota 96.

Il termine $C \sum_{i=1}^l \xi_i$ rappresenta una **penalità** sull'errore sul TS e $C > 0$ viene scelto in cross validation.

Scriviamo a questo punto la nuova Lagrangiana

$$L(w, b, \xi, \lambda, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \lambda_i (y^i (w^T x^i + b) - 1 + \xi_i) - \sum_{i=1}^l \mu_i \xi_i$$

e il suo duale:

$$\max \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \lambda_i (y^i (w^T x^i + b) - 1 + \xi_i) - \sum_{i=1}^l \mu_i \xi_i$$

$$\Delta_w L(w, b, \xi, \lambda, \mu) = w - \sum_{i=1}^l \lambda_i y^i x^i = 0$$

$$\frac{\partial L}{\partial b}(w, b, \xi, \lambda, \mu) = - \sum_{i=1}^l \lambda_i y^i = 0$$

$$\frac{\partial L}{\partial \xi_i}(w, b, \xi, \lambda, \mu) = C - \lambda_i - \mu_i \quad i = 1, \dots, l$$

$$\lambda_i \geq 0 \quad i = 1, \dots, l$$

$$\mu_i \geq 0 \quad i = 1, \dots, l$$

Dall'annullamento della Lagrangiana rispetto a w e ξ possiamo ricavare:

$$w = \sum_{i=1}^l \lambda_i y^i x^i \quad \mu_i = C - \lambda_i \quad i = 1, \dots, l$$

Sostituendo queste due ultime espressioni, il problema diventa:

$$\max \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y^i y^j (x^i)^T x^j + \textcolor{red}{C} \sum_{i=1}^l \textcolor{red}{\xi_i} - \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y^i y^j (x^i)^T x^j - \sum_{i=1}^l \textcolor{red}{\lambda_i y^i b} + \sum_{i=1}^l \lambda_i - \sum_{i=1}^l \textcolor{red}{\lambda^i \xi^i} - \sum_{i=1}^l \textcolor{red}{(C - \lambda_i) \xi_i}$$

Nota 97. I termini in rosso si semplificano o diventano nulli.

$$\max -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y^i y^j (x^i)^T x^j + \sum_{i=1}^l \lambda_i$$

Riscrivendolo in forma di minimo:

$$-\min \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y^i y^j (x^i)^T x^j - \sum_{i=1}^l \lambda_i$$

$$\sum_{i=1}^l \lambda^i y^i = 0$$

$$\lambda_i \geq 0$$

$$C - \lambda_i \geq 0 \quad i = 1, \dots, l$$

Gli ultimi due vincoli vengono riscritti come $0 \leq \lambda_i \leq C \quad i = 1, \dots, l$

Alla fine il duale diventa:

$$-\min \quad \frac{1}{2} \lambda^T X^T X \lambda - e^T \lambda$$

$$y^T \lambda = 0$$

$$0 \leq \lambda_i \leq C \quad i = 1, \dots, l$$

A questo punto supponiamo di risolvere il duale λ^* , vogliamo capire come si trovano w^*, b^* del primale. A tale scopo:

$$w^* = \sum_{i=1}^l \lambda_i^* y^i x^i$$

Importante 98. I vettori di supporto sono gli x^i per cui $\lambda_i^* > 0$.

In aggiunta, b^* si trova tramite la complementarità. In particolare, (λ^*, μ^*) dove $\mu_i^* = C - \lambda_i^*$ sono i moltiplicatori associati alla soluzione primale. Allora questi moltiplicatori devono soddisfare la complementarità:

$$\lambda_i^* (y^i (w^{*T} x^i + b^*) - 1 + \xi_i^*) = 0 \quad i = 1, \dots, l$$

$$\mu_i^* \xi_i^* = (C - \lambda_i^*) \xi_i^* = 0 \quad i = 1, \dots, l$$

$$\xi_i^* (C - \lambda_i^*) = 0 \longrightarrow \text{Se } \lambda_i^* < C \longrightarrow \xi_i^* = 0 \longrightarrow x^i \text{ è ben classificato}$$

Per trovare b^* scelgo $i: 0 < \lambda_i^* < C$. Allora $\xi_i^* = 0 \longrightarrow y^i (w^{*T} x^i + b^*) = 1$.

I punti x^i per cui $\lambda_i^* = C$ non posso calcolare ξ_i^* . Allora sono candidati ad essere malclassificati. Infine, la nostra macchina SVN sarà del tipo:

$$y(x) = \text{sgn}(w^{*T} x + b^*)$$

14 SVM Non Lineari

Per alcuni insiemi di dati una superficie di separazione lineare non cattura correttamente la classificazione. A tale scopo, passiamo dallo spazio di input a uno spazio di dimensione maggiore, detto **feature space**.

Quello che non era linearmente separabile nello spazio di input lo può diventare nello spazio delle feature. La superficie che si ottiene è lineare nello spazio delle feature, ma non lineare nello spazio di partenza. In pratica, dobbiamo utilizzare il **kernel trick**.

14.1 Kernel trick

La funzione kernel calcola il prodotto scalare di due trasformazioni dell'input. L'idea è quella di passare da:

$$x \text{ (spazio di input)} \longrightarrow \Phi(x) = \{\Phi(x_1), \Phi(x_2), \dots\} \text{ (spazio delle feature)}$$

L'iperpiano di separazione viene cercato nello spazio delle feature e la corrispondente macchina SVM diventa:

$$y(x) = \text{sgn}(w^T \Phi(x) + b)$$

Lo spazio di arrivo può avere dimensione infinita.

La funzione $\Phi(x)$ non viene calcolata esplicitamente, ma si sfrutta il fatto che x , nella SVM lineari non compare mai da sola, ma sempre come prodotto $x_i^T x_j$ oppure $x^T x^i$. Quindi, mi serve calcolare $\Phi(x_i)^T \Phi(x_j)$ o $\Phi(x)^T \Phi(x^i)$.

Di conseguenza, il problema di ottimizzazione nel duale diventa:

$$\min \frac{1}{2} \sum_i \sum_j y^i y^j \lambda_i \lambda_j \Phi(x^i)^T \Phi(x^j) - \sum_{i=1}^l \lambda_i$$

$$\sum_i \lambda_i y^i = 0$$

$$0 \leq \lambda_i \leq C$$

La b^* si calcola scegliendo $0 \leq \lambda_i^* \leq C \longrightarrow b^* = \frac{1}{y^i} - \sum_{j=1}^l y^j \lambda_j^* \Phi(x^j)^T \Phi(x^i)$

La SVM diventa quindi:

$$y(x) = \text{sgn} \left(\sum_{i=1}^l \lambda_i^* y^i \Phi(x^j)^T \Phi(x^i) + b^* \right)$$

14.2 Funzioni Kernel

Dato un insieme $X \subseteq \mathbb{R}^n$, $K: X \times X \rightarrow \mathbb{R}$ è una **funzione Kernel** se:

$$k(x, y) = \Phi(x)^T \Phi(y)$$

In cui $\Phi: \mathbb{R}^n \rightarrow H$ spazio Euclideo.

Nota 99.

- Se K è un Kernel allora $\forall x^1, \dots, x^l \in X$ definisco:

$$K = [K(x^i, x^j)] \quad i = 1, \dots, l; j = 1, \dots, l$$

è simmetrica e semidefinita positiva, quindi non si perdono le proprietà di convessità.

$$K = \begin{pmatrix} k(x^1, x^1) & \dots & k(x^1, x^l) \\ \vdots & \ddots & \vdots \\ k(x^l, x^1) & \dots & k(x^l, x^l) \end{pmatrix} \succcurlyeq 0$$

- Si usa la funzione Kernel nel duale:

$$\min \frac{1}{2} \sum_i \sum_j y^i y^j \lambda_i \lambda_j K(x^i, x^j) - \sum_{i=1}^l \lambda_i$$

$$\sum_i \lambda_i y^i = 0$$

$$0 \leq \lambda_i \leq C$$

- Il problema rimane convesso, si risolve il duale trovando λ^* e si trova b^* per complementarità:

$$b^* = \frac{1}{y^i} - \sum_{j=1}^l y^j \lambda_j^* k(x^i, x^j)$$

La SVM diventa:

$$y(x) = \text{sgn} \left(\sum_{i=1}^l \lambda_i y^i k(x^i, x) + b^* \right)$$

I Kernel più utilizzati sono:

1. **Lineare:** $K(x, y) = x^T y$
2. **Polinomiale:** $k(x, y) = (x^T y + \gamma)^p \quad \gamma \geq 0, p \geq 1$ (Scelti in cross validation)
3. **Gaussiano:** $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}} \quad \sigma > 0$
4. **Tangente Iperbolica:** $k(x, y) = \tanh(\beta x^T y + \gamma)$

14.3 Kernel Polinomiale

$$k(x, y) = (x^T y + \gamma)^p = \Phi(x)^T \Phi(y) \quad \Phi: \mathbb{R}^n \rightarrow H$$

La dimensione minima di H è $\binom{n+p-1}{n} = \frac{(n+p-1)!}{n!(p-1)!}$

14.4 Kernel Gaussiano

$$k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}} \quad \gamma = \frac{1}{2\sigma^2} \rightarrow k(x, y) = e^{-\gamma \|x-y\|^2}$$

Tramite un Kernel Gaussiano, nello spazio delle feature, tutto è linearmente separabile poiché lo spazio delle feature ha dimensione infinita. (Dato che con un semplice esempio possiamo esprimere l'esponenziale tramite una serie infinita)

Data questa proprietà del Kernel Gaussiano possiamo quindi affermare che la teoria di Vapnik rimane valida e quindi assicurare l'esistenza dell'iperpiano ottimo. Tuttavia, questa scelta non risulta sempre ottimale poiché il costo computazionale potrebbe risultare oneroso e si cade facilmente nel caso di overfitting dei dati.

15 Working Set

Dato un Training Set:

$$\text{TS} = \{(x^i, y^i), x^i \in \mathbb{R}^n, y^i \in \{-1, 1\} \mid i = 1, \dots, k\}$$

Trovare l'iperpiano a massimo margine vuol dire risolvere:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$y^i (w^T x^i + b) \geq 1 - \xi_i \quad i = 1, \dots, k \quad (\lambda_i)$$

$$\xi_i \geq 0$$

Equivalentemente possiamo risolvere il duale:

$$\begin{aligned} \min \frac{1}{2} \sum_i \sum_j y^i y^j (x^i)^T x^j - \sum_{i=1}^l \lambda_i \\ \sum_{i=1}^l \lambda_i y^i = 0 \\ 0 \leq \lambda_i \leq C \end{aligned}$$

Che si risolve ottenendo:

$$w^* = \sum_{i=1}^l \lambda_i^* y^i x^i$$

e per complementarità:

$$b^*$$

La SVM risultante sarà:

$$y(x) = \text{sgn}(w^{*T} x + b^*)$$

Così facendo si produce una SVM lineare.

Nel caso delle SVM non lineare si usa il kernel trick e si risolve il problema:

$$\begin{aligned} \min \frac{1}{2} \sum_i \sum_j y^i y^j K(x^i, x^j) - \sum_{i=1}^l \lambda_i \\ \sum_{i=1}^l \lambda_i y^i = 0 \\ 0 \leq \lambda_i \leq C \end{aligned}$$

La SVM risultante sarà del tipo:

$$y(x) = \text{sgn} \left(\sum_{i=1}^l \lambda_i^* y^i K(x^i, x) + b^* \right)$$

Indipendentemente dal Kernel che si utilizza, il problema duale è nella forma:

$$\begin{aligned} \min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - e^T \alpha \\ y^T \alpha = 0 \quad Q \succcurlyeq 0 \\ 0 \leq \alpha \leq C \end{aligned}$$

$$Q_{ij} = y^i y^j K(x^i, x^j) \text{ con } Q = Q^T \quad Q \in \mathbb{R}^{l \times l}$$

Talvolta, Q può essere talmente grande da avere problemi di memoria allora si trova in cache un sottoinsieme di colonne. Allora si utilizzano delle **tecniche di decomposizione**: si ottimizza ogni volta rispetto ad un sottoinsieme di variabili detto **working set** dipendente dal tipo di kernel utilizzato. In particolare:

Si partiziona $\alpha = \begin{pmatrix} \alpha_w \\ \alpha_{\bar{w}} \end{pmatrix}$ $\bar{w} = \{1, \dots, n\} / W$ con le seguenti proprietà:

- $W \cup \bar{W} = \{1, \dots, l\}$
- $W \cap \bar{W} = \emptyset$

Parto da α_k ammissibile: $y^T \alpha_k = 0$ con $0 \leq \alpha_k \leq C$ $i = 1, \dots, l$

Calcolo α_W^{k+1} fissando tutte le variabili tranne quelle in W :

$$\alpha_{\bar{W}} = \alpha_{\bar{W}}^k$$

$$Q = \begin{pmatrix} Q_{WW} & Q_{W\bar{W}} \\ Q_{\bar{W}W} & Q_{\bar{W}\bar{W}} \end{pmatrix}$$

Il problema da risolvere diventa:

$$\begin{aligned} \min_{\alpha_w} f(\alpha_w, \alpha_{\bar{W}}^k) &= \frac{1}{2} \begin{pmatrix} \alpha_w^T & \alpha_{\bar{W}}^{kT} \end{pmatrix} \begin{pmatrix} Q_{WW} & Q_{W\bar{W}} \\ Q_{\bar{W}W} & Q_{\bar{W}\bar{W}} \end{pmatrix} \begin{pmatrix} \alpha_w \\ \alpha_{\bar{W}}^k \end{pmatrix} - e^T \alpha_w - e^T \alpha_{\bar{W}} = \\ &= \frac{1}{2} \alpha_w^T Q_{WW} \alpha_w + \frac{1}{2} \alpha_{\bar{W}}^{kT} Q_{\bar{W}\bar{W}} \alpha_{\bar{W}}^k + \alpha_{\bar{W}}^{kT} Q_{\bar{W}W} \alpha_w - e^T \alpha_w = \\ &= \frac{1}{2} \alpha_w^T Q_{WW} \alpha_w + c^T \alpha_w \quad c = Q_{\bar{W}W} \alpha_{\bar{W}}^k - e \end{aligned}$$

$$\alpha_W^{k+1} = \arg \min_{\alpha_w} \frac{1}{2} \alpha_w^T Q_{WW} \alpha_w + c^T \alpha_w.$$

$$y_W^T \alpha_w + y_{\bar{W}}^T \alpha_{\bar{W}}^k = 0 \rightarrow y_W^T \alpha_w = -y_{\bar{W}}^T \alpha_{\bar{W}}^k$$

$$0 \leq \alpha_w \leq C$$

La struttura del problema rimane invariata ma la dimensione diventa $|W|$ più piccolo.

L'algoritmo implementativo è il seguente.

Algoritmo 1

DATI

$$\alpha^o (\alpha^o = 0)$$

INIT $k = 0$

WHILE criterio di arresto non è soddisfatto

1. Seleziono tramite WSS un insieme W^k

2. Poni $W = W^k$ e risolvi:

$$\alpha_W^* = \arg \min \frac{1}{2} \alpha_W^T Q_{WW} \alpha_W + c^T \alpha_w$$

$$y_W^T \alpha_w = -y_{\bar{W}}^T \alpha_{\bar{W}}^k$$

$$0 \leq \alpha_w \leq C$$

3. Poni:

$$\alpha_i^{k+1} = \begin{cases} \alpha_i^* & \text{se } i \in W \\ \alpha_i^k & \text{se } i \in \bar{W} \end{cases}$$

4. $k = k + 1$
END

Importante 100. La convergenza dipende dalla selezione del WS.

Devo scegliere la cardinalità di W , allora se è piccolo il sottoproblema facile, ma si necessitano di più iterazioni e quindi la convergenza risulta più lenta. La cardinalità minima è 2

Se $|W| = 2$ allora si hanno algoritmi di tipo SMO (Sequential Minimization Ottimization) e la soluzione è analitica, ma non si ha necessità di algoritmi e l'iterazione è molto veloce.

In generale non si supera $|W| \leq 10$, se no l'iterazione diventa molto costosa. Lo svantaggio degli algoritmi a di decomposizione è la lentezza dato il grande numero di iterazioni e non sono molto precisi.

Teorema 101. Riformulazione KKT

$$\alpha^* \text{ ammissibile ottimo globale} \iff \exists \lambda^*: \nabla_i f(\alpha^*) + \lambda^* y^i = \begin{cases} \geq 0 & \text{se } i \in L(\alpha^*) \\ = 0 & \text{se } i \notin L(\alpha^*) \cup U(\alpha^*) \\ \leq 0 & \text{se } i \in U(\alpha^*) \end{cases}$$

In cui:

$$L(\alpha^*) = \{i: \alpha_i^* = 0\}$$

$$U(\alpha^*) = \{i: \alpha_i^* = C\}$$

Dimostrazione.

\implies Se vale KKT allora vale il teorema

Supponiamo che $i \in L(\alpha^*)$ allora $\alpha_i^* = 0$. La sostuiamo nella complementarità $\hat{\xi}_i = 0$. A questo punto sostituisco nell'annullamento del gradiente del Lagrangiano e si ottiene:

$$\nabla_i f(\alpha^*) + \lambda^* y^i - \xi_i = 0$$

$$\nabla_i f(\alpha^*) + \lambda^* y^i = \xi_i = 0$$

Supponiamo che $i \in U(\alpha^*) \longrightarrow \alpha_i^* = C$. La sostituiamo nella complementarità $\xi_i = 0$. A questo punto sostituisco nell'annullamento del gradiente del Lagrangiano:

$$\nabla_i f(\alpha^*) + \lambda^* y^i + \hat{\xi}_i = 0$$

$$\nabla_i f(\alpha^*) + \lambda^* y^i = -\hat{\xi}_i \leq 0$$

Supponiamo infine che $i \notin L(\alpha^*) \cup U(\alpha^*) \longrightarrow 0 < \alpha_i^* < C$. Dalla complementarità $\xi_i = 0$ $\hat{\xi}_i = 0$. A questo punto sostituiamo nell'annullamento del gradiente della Lagrangiana:

$$\nabla_i f(\alpha^*) + \lambda^* y^i = 0$$

\longleftarrow Se vale il teorema allora vale KKT

Dobbiamo generare dei moltiplicatori che devono soddisfare KKT e ξ_i e $\hat{\xi}_i$. Nel caso in cui:

- $i \in L(\alpha^*) \longrightarrow \alpha_i^* = 0$ per soddisfare la complementarità prendo:

$$\hat{\xi}_i = 0 \longrightarrow \hat{\xi}_i(\alpha_i - C) = 0$$

Devo scegliere $\xi_i \geq 0$ per soddisfare l'annullamento del gradiente della lagrangiana ed in particolare scelgo:

$$\xi_i = \nabla_i f(\alpha^*) + \lambda^* y^i \geq 0 \text{ Per ipotesi}$$

- $i \in U(\alpha^*) \rightarrow \alpha_i = C$ per soddisfare la complementarità $\xi_i = 0$

$$(\alpha_i - C) \hat{\xi}_i = 0$$

è soddisfatta perché $\alpha_i^* = C$. Devo scegliere $\hat{\xi}_i \geq 0$ tale che sia soddisfatto l'annullamento del gradiente della lagrangiana e quindi si ottiene:

$$\hat{\xi}_i = -\nabla_i f(\alpha^*) + \lambda^* y^i \geq 0$$

- $i \notin L(\alpha^*) \cup U(\alpha^*) \rightarrow 0 < \alpha_i^* < C$. Quindi:

$$\xi_i = \hat{\xi}_i = 0 \rightarrow \xi_i \alpha_i = 0 \quad \hat{\xi}_i (\alpha_i - C) = 0 \text{ e per ipotesi } \nabla_i L = 0$$

□

Nota 102.

Dal teorema si possono definire i seguenti insiemi:

$$L(\alpha^*) = L^+(\alpha^*) \cup L^-(\alpha^*)$$

In cui:

$$L^+(\alpha^*) = \{i: \alpha_i^* = 0 \text{ e } y^i = 1\}$$

$$L^-(\alpha^*) = \{i: \alpha_i^* = 0 \text{ e } y^i = -1\}$$

Inoltre:

$$U(\alpha^*) = U^+(\alpha^*) \cup U^-(\alpha^*)$$

In cui:

$$U^+(\alpha^*) = \{i: \alpha_i^* = C \text{ e } y^i = 1\}$$

$$U^-(\alpha^*) = \{i: \alpha_i^* = C \text{ e } y^i = -1\}$$

Ponendo così riscrivere le condizioni del teorema come segue:

$$\lambda^* y^i = \begin{cases} \geq -\nabla_i f(\alpha^*) & i \in L(\alpha^*) \\ = -\nabla_i f(\alpha^*) & i \notin L(\alpha^*) \cup U(\alpha^*) \\ \leq -\nabla_i f(\alpha^*) & i \in U(\alpha^*) \end{cases}$$

$$\lambda^* = \begin{cases} \geq -\frac{\nabla_i f(\alpha^*)}{y^i} & i \in L^+(\alpha^*) \cup U^-(\alpha^*) \\ = -\frac{\nabla_i f(\alpha^*)}{y^i} & i \in L(\alpha^*) \cup U(\alpha^*) \\ \leq -\frac{\nabla_i f(\alpha^*)}{y^i} & i \in L^-(\alpha^*) \cup U^+(\alpha^*) \end{cases}$$

Oltre a questi parizionamenti possiamo definire i seguenti insiemi aggiuntivi:

$$R(\alpha) = L^+(\alpha^*) \cup U^-(\alpha^*) \cup \{i: 0 < \alpha_i < C\}$$

$$S(\alpha) = L^-(\alpha^*) \cup U^+(\alpha^*) \cup \{i: 0 < \alpha_i < C\}$$

di cui:

$$R(\alpha) \cup S(\alpha) = \{1, \dots, l\}$$

$$R(\alpha) \cap S(\alpha) = \{i: 0 < \alpha_i < C\}$$

Teorema 103.

$$\begin{aligned} \alpha^* \text{ ammissibile} \\ \implies \alpha^* \text{ è ottimo} \\ R(\alpha^*) \cap S(\alpha^*) \text{ non vuoto} \end{aligned}$$

Dimostrazione.

Supponiamo $R(\alpha^*) = \emptyset \rightarrow S(\alpha^*) = \{1, \dots, l\}$ e $R(\alpha^*) \cap S(\alpha^*) = \emptyset$,

Allora $\nexists i$ tale che; $0 < \alpha_i < C \rightarrow S(\alpha^*) = L^-(\alpha^*) \cup U^-(\alpha^*) = \{1, \dots, l\}$

La condizione di ottimo diventa:

$$\lambda^* \leq -\frac{\nabla_i f(\alpha^*)}{y^i} \quad i = 1, \dots, l$$

$$\rightarrow \lambda^* = \min_{i=1, \dots, l} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\}$$

Soddisfa le condizioni di ottimo globale.

Al contrario se $S(\alpha^*) = \emptyset$ $R(\alpha^*) = \{1, \dots, l\}$ e $R(\alpha^*) \cap S(\alpha^*) = \emptyset$

Allora $\nexists i: 0 < \alpha^* < C$:

$$R(\alpha^*) = L^+(\alpha^*) \cup U^-(\alpha^*)$$

La condizione di ottimo diventa:

$$\lambda^* \geq -\frac{\nabla_i f(\alpha^*)}{y^i} \quad i = 1, \dots, l$$

$$\lambda^* = \max_{i=1, \dots, l} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\}$$

Il che soddisfa le condizioni di ottimo. □

Teorema 104. *Condizione di ottimo*

$$\alpha^* \text{ ammissibile è ottima} \iff \max_{i=1, \dots, l} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\} \leq \min_{i=1, \dots, l} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\}$$

16 SMO

SMO sta per sequential minimization optimization ed in questi tipi di algoritmi sono interessato ad avere un Working Set di cardinalità 2 di cui muovo due variabili. In particolare, anche in questa tipologia di algoritmi effettuiamo decomposizione e vogliamo cercare di caratterizzare le direzioni ammissibili su cui devo decomporre. A tale scopo, definiamo l'insieme:

$$F = \{\alpha \in \mathbb{R}^l: y^T \alpha = 0, 0 \leq \alpha_i \leq C, i = 1, \dots, l\}$$

Ho vincoli lineari, uno di uguaglianza e 2*1 vincoli di box:

$$a_i^T \alpha_i = b_i \longrightarrow a_i^T (a + d) = b_i \longleftrightarrow a_i^T \alpha + a_i^T d = b_i \longrightarrow a_i^T d = 0$$

In questo caso:

$$y^T \alpha = 0 \longrightarrow d \text{ è ammissibile} \longleftrightarrow y^T d = 0$$

Se mi trovo in α ammissibile e $\alpha_i = 0 \longrightarrow d_i \geq 0$

Se mi trovo in α ammissibile e $\alpha_i = C \longrightarrow d_i \leq 0$

Ricapitolando in α ammissibile l'insieme delle direzioni ammissibile è definito come:

$$\{d \in \mathbb{R}^l: y^T d = 0, d_i \geq 0 \quad \forall i \in L(\bar{\alpha}), d_i \leq 0, \forall U(\bar{\alpha})\}$$

Teorema 105.

Sia $\bar{\alpha}$ ammissibile e sia $\{(i, j) \in \{1, \dots, l\} \mid i \neq j\}$. Definisco $d^{i,j}$ come:

$$d_h^{i,j} = \begin{cases} \frac{1}{y^i} & \text{se } h = i \\ 0 & \text{altrimenti} \\ -\frac{1}{y^j} & \text{se } h = j \end{cases}$$

Allora:

- d ammissibile $\iff i \in R(\bar{\alpha}) \quad e \quad j \in S(\bar{\alpha})$
- $d^{i,j}$ è di discesa per f in $\bar{\alpha}$ se:

$$\frac{\nabla_i f(\alpha^*)}{y^i} - \frac{\nabla_j f(\alpha^*)}{y^j} < 0$$

Dimostrazione.

- Se d ammissibile $\longrightarrow y^T d = 0$ banale.

Se $d^{i,j}$ ammissibile $\longrightarrow i \in R(\bar{\alpha}) \quad e \quad j \in S(\bar{\alpha})$. Suppongo per assurdo che non si vero:

$$j \in R(\bar{\alpha}) = L^+(\bar{\alpha}) \cup U^-(\bar{\alpha}) \cup \{i: 0 < \alpha_i < C\}$$

Poiché $j \in S(\bar{\alpha}) \longrightarrow j \notin \{i: 0 < \alpha_i < C\} \longrightarrow j \in L^+(\bar{\alpha})$.

Quindi se $j \in L^+(\bar{\alpha}) \rightarrow d^{i,j} = -\frac{1}{y^i} \leq 0$

Il che è un assurdo perché per essere ammissibile al lower bound deve essere ≥ 0 .

- $d^{i,j}$ è di discesa $\longleftrightarrow \nabla f(\bar{\alpha})^T d^{i,j} \leq 0$

Quindi:

$$\nabla f(\bar{\alpha})^T d^{i,j} = \frac{\nabla_i f(\alpha^*)}{y^i} - \frac{\nabla_j f(\alpha^*)}{y^j} < 0$$

□

Definisco:

- $I(\alpha) = \left\{ i: \arg \max_{i \in R(\alpha)} \left\{ -\frac{\nabla_i f(\alpha)}{y_i} \right\} \right\}$
- $J(\alpha) = \left\{ j: \arg \max_{j \in S(\alpha)} \left\{ -\frac{\nabla_j f(\alpha)}{y_j} \right\} \right\}$

Se α non è ottimo $\rightarrow I(\alpha)$ e $J(\alpha)$ sono ben definiti e non vuoti.

Se α non è ottimo, la condizione di ottimalità risulta violata:

$$\max_{i \in R(\alpha)} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\} \leq \min_{j \in J(\alpha)} \left\{ -\frac{\nabla_j f(\alpha^*)}{y^j} \right\}$$

Supponiamo di scegliere $i \in I(\alpha)$ e $j \in J(\alpha)$ e considerare una direzione d^{ij} . Allora:

1. d^{ij} è ammissibile poiché $i \in I(\alpha) \subseteq R(\alpha)$ e $j \in J(\alpha) \subseteq S(\alpha)$
2. d^{ij} è di discesa ed è quella con derivata direzionale più bassa.

$i \in I(\alpha)$

$$-\frac{\nabla_i f(\alpha^*)}{y^i} \geq -\frac{\nabla_t f(\alpha^*)}{y^t} \quad \forall t \in R(\alpha) \quad (1)$$

$j \in J(\alpha)$

$$-\frac{\nabla_j f(\alpha^*)}{y^j} \leq -\frac{\nabla_s f(\alpha^*)}{y^s} \quad \forall s \in S(\alpha) \quad (2)$$

Moltiplicando (1) per -1 :

$$\frac{\nabla_i f(\alpha^*)}{y^i} \leq \frac{\nabla_t f(\alpha^*)}{y^t} \quad (3)$$

Sommando (2) e (3):

$$\frac{\nabla_i f(\alpha^*)}{y^i} - \frac{\nabla_j f(\alpha^*)}{y^j} \leq \frac{\nabla_t f(\alpha^*)}{y^t} - \frac{\nabla_s f(\alpha^*)}{y^s} \quad \forall t \in R(\alpha) \quad \forall s \in S(\alpha)$$

Dalla condizione di ottimo violata, si ottiene:

$$-\frac{\nabla_i f(\alpha^*)}{y^i} > -\frac{\nabla_j f(\alpha^*)}{y^j} \rightarrow \frac{\nabla_i f(\alpha^*)}{y^i} - \frac{\nabla_j f(\alpha^*)}{y^j} < 0$$

è la massima direzione di discesa che viola maggiormente le condizioni di ottimo.

Algoritmo

Dati $\alpha^0 = 0, \nabla f(\alpha^0) = -e$

Init $k = 0$

While criterio di arresto non soddisfatto

1. Determina Working Set $WS: i \in I(\alpha^k), j \in J(\alpha^k)$ e poni $W_k = \{i, j\}$
2. Determina analiticamente una soluzione α_i^*, α_j^* del problema ristretto:

$$\begin{aligned} \min_{\alpha_i, \alpha_j} q(\alpha_i, \alpha_j) &= \frac{1}{2} \begin{pmatrix} \alpha_i & \alpha_j \end{pmatrix} \begin{pmatrix} q_{ii} & q_{ij} \\ q_{ji} & q_{jj} \end{pmatrix} \begin{pmatrix} \alpha_i \\ \alpha_j \end{pmatrix} + \sum_{h \neq i, j} (q_{ih} \alpha_i + q_{jh} \alpha_j) \alpha_h^k - \alpha_i - \alpha_j + \text{cost} \\ y^i \alpha_i + y^j \alpha_j &= - \sum_{h \neq i, j} y^h \alpha_h^k \\ 0 &\leq \alpha_i \leq C \quad 0 \leq \alpha_j \leq C \end{aligned}$$

3. Poni:

$$\alpha_h^{k+1} = \begin{cases} \alpha_i^* & h = i \\ \alpha_j^* & h = j \\ \alpha^* & h \neq \{i, j\} \end{cases}$$

4. Poni $\nabla f(\alpha^{k+1}) = Q \alpha^{k+1} - e = Q \alpha^{k+1} \pm Q \alpha^k - e = Q \alpha^k - e + Q(\alpha^{k+1} - \alpha^k) =$
 $= \nabla f \alpha^k + Q_i(\alpha^{k+1} - \alpha_i^k) + Q_j(\alpha^{k+1} - \alpha_j^k)$

5. $k = k + 1$

END

Teorema 106.

Se $Q \succcurlyeq 0 \quad Q = Q^T$ ogni punto limite delle sequenza $\{\alpha^k\}$ è soluzione del duale.

16.1 Criterio di Arresto

La condizione di ottimo è:

$$\max_{i \in R(\alpha)} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\} \leq \min_{i \in S(\alpha)} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\}$$

Per rilassare questa condizione imponiamo:

$$m(\alpha) = \begin{cases} \max_{i \in R(\alpha)} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\} & \text{se } R(\alpha) \neq \emptyset \\ -\infty & \text{altrimenti} \end{cases}$$

$$M(\alpha) = \begin{cases} \min_{i \in R(\alpha)} \left\{ -\frac{\nabla_i f(\alpha^*)}{y^i} \right\} & \text{se } S(\alpha) \neq \emptyset \\ \infty & \text{altrimenti} \end{cases}$$

Quindi, α è soluzione ottima $\iff m(\alpha) \leq M(\alpha)$.

Mentre la sequenza $\{\alpha^k\}$ convergente all'ottimo si potrebbe avere $m(\alpha) \geq M(\alpha)$.

Un criterio di arresto utilizzato è:

$$m(\alpha^k) \leq M(\alpha^k) + \varepsilon \quad \varepsilon > 0$$

Si potrebbe avere:

$$\lim_{k \rightarrow \infty} m(\alpha^k) \neq m(\alpha^*) \text{ di norma discontinua}$$

La convergenza viene garantita teoricamente aggiungendo ipotesi sulla struttura della matrice Q .

17 Riepilogo SVM

Nel problema delle SVM si vuole trovare l'iperpiano a massimo margine. Questa scelta è ottima se i due insiemi risultano linearmente separabili. In particolare, il problema primale è il seguente:

- Se si vuole minimizzare l'errore in base alla **L1-Loss (modulo dell'errore)**

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$y^i(w^T x^i b) \geq 1 - \xi_i \quad i = 1, \dots, k(\lambda_i)$$

$$\xi_i \geq 0$$

- Se si vuole minimizzare l'errore in base alla **L2-Loss (modulo dell'errore al quadrato)**

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i^2$$

$$y^i(w^T x^i b) \geq 1 - \xi_i \quad i = 1, \dots, k(\lambda_i)$$

$$\xi_i \geq 0$$

Il suo duale sia se si tratta di un problema di L1-LOSS sia che si tratta di un problema con L2-Loss, è il seguente:

$$\min \frac{1}{2} \lambda^T Q \lambda - e^T \lambda_i$$

$$\sum_{i=1}^l \lambda_i y^i = 0$$

$$0 \leq \lambda_i \leq C$$

In cui:

$$Q_{ij} = \begin{cases} y^i y^j (x^i)^T x^j & \text{Kernel Lineare} \rightarrow \text{Iperpiano ottimo nello spazio di input} \\ y^i y^j K(x^i, x^j) & \text{kernel Opportuno} \rightarrow \text{SVM non lineare} \rightarrow \text{Iperpiano spazio delle features} \end{cases}$$

Si possono utilizzare pesi differenti in base alla classe $\longrightarrow \frac{1}{\# \text{Istante TS della classe } i}$

La funzione obiettivo diventa:

$$C \sum_{i=1}^n \xi_i \longrightarrow C \sum_{i \in A} w_a \xi_i + C \sum_{j \in B} w_b \xi_j$$

Importante 107. Il parametro C viene scelto in cross validation, ma di solito viene impostato a $C \approx \frac{1}{10}$. Nel caso in cui diventasse troppo grande per evitare overfitting, il problema si complica.

18 ν -Classification

Nella ν - Classification si risolve il problema:

$$\begin{cases} \min \frac{1}{2} \|w\|^2 - \nu \rho + \frac{1}{2} \sum \xi_i \\ y^i(w^T x^i + b) \geq \rho - \xi_i \quad \xi_i = 1, \dots, l \\ \xi_i \geq 0 \quad \rho \geq 0 \quad i = 1, \dots, l \end{cases}$$

Se si aumenta C diminuisce ν e viceversa.

Se $\xi_i = 0 \rightarrow y^i(w^T x^i + b) \geq \rho$. La quantità ρ rappresenta il **margin** e vale:

$$\frac{2\rho}{\|w\|}$$

Nel caso in cui $\rho > 0$, all'ottimo:

1. $\nu \geq$ frazione di errore nel margine (punti che si trovano sul margine)
2. $\nu \leq$ frazione dei vettori di supporto del TS
3. Con probabilità 1 al limite il punto 1 e 2 sono uguali.

Se $\rho > 0$ all'ottimo è equivalente alla C - Classification con $C = \frac{1}{2\rho}$.

Importante 108.

In particolare $\nu \in [\nu_{\min}, \nu_{\max}]$ o il duale diventa inammissibile:

$$\begin{cases} \min \frac{1}{2} \alpha^T Q \alpha \\ y^T \alpha > 0 \quad e^T \alpha \geq \nu \\ 0 \leq \alpha_i \leq \frac{1}{l} \end{cases}$$

$$\nu_{\max} = \frac{2 \min(|A|, |B|)}{l}$$

Con:

- $A = \{i = 1, \dots, l \mid y_i = 1\}$
- $B = \{i = 1, \dots, l \mid y_i = -1\}$

Se si utilizzasse LIBSVM l'algoritmo di tipo SMO sceglie $i \in I(\alpha)$ e j guardando le condizioni del secondo ordine. Per velocizzare l'algoritmo si usano le tecniche di **Shrinking**;

1. Se $\lambda^i = 0$ ($\alpha_i = 0$) per tante iterazioni, allora la fissa a 0 riducendo conseguentemente il numero di variabili;
2. Se $\lambda^i = C$ per tante iterazioni allora la fissa a C ;

Tuttavia si deve comunque verificare che le condizioni di ottimo siano soddisfatte. In aggiunta il costo generico di una iterazione:

$$O(|W|l)$$

Utilizzando un Kernel Lineare (es. LIBLINEAR) il costo dell'iterazione si abbatta.

Ricordando un problema con Kernel Lineare e valutando il suo duale, se N è molto grande (molte features) alcuni software fissano $b=0 \rightarrow$ i vincoli lineari scompaiono. In alternativa si può scegliere di ottimizzare rispetto ad una nuova variabile:

$$\tilde{w} = \begin{pmatrix} w \\ b \end{pmatrix} \quad \tilde{x} = \begin{pmatrix} x \\ 1 \end{pmatrix}$$

$$w^T x + b \longleftrightarrow \tilde{w}^T \tilde{x}$$

La teoria di Vapnick cade in entrambi i casi, ma il duale diventa:

$$\begin{cases} \min \frac{1}{2} \lambda^T Q \lambda - e^T \lambda \\ 0 \leq \lambda_i \leq C \\ i = 1, \dots, l \end{cases}$$

che, considerando una variabile per volta e introducendo il **metodo delle coordinate duali** $|W_k| = 1$, la soluzione del sottoproblema diventa analitica ed è possibile effettuare lo shrinking.

Nel caso di Kernel Lineare si può usare esplicitamente:

$$w = \sum_{i=1}^l \lambda_i^* y^i x^i$$

come legge di aggiornamento:

$$w^k = \sum_{i=1}^l \lambda_i^k y^i x^i$$

che al limite converge a w^* .

Il calcolo del gradiente migliora diventando $O(n) \ll O(l)$ del costo di una singola iterazione.

19 Problemi di Regressione con SVM

Le SVM possono essere utilizzati per **problemi di regressione**:

$$TS = \{(x^i, y^i), x^i \in \mathbb{R}^n, y^i \in \mathbb{R} \quad i = 1, \dots, l\}$$

Si impone che l'errore:

$$|y^i - w^T x^i - b| \leq \varepsilon$$

sia minore di una certa tolleranza ε .

Il vincolo si trasforma in una coppia di vincoli lineari del tipo:

$$y^i - w^T x^i - b \leq \varepsilon - \xi_i$$

$$y^i - w^T x^i - b \geq -\varepsilon + \hat{\xi}_i$$

$$\xi_i \geq 0 \quad \hat{\xi}_i \geq 0$$

19.1 Regressione Logistica

La **regressione logistica** è uno strumento di classificazione lineare che utilizza il concetto di distanza dell'iperpiano. In particolare, si fornisce una misura dell'importanza delle features:

$$TS = \{(x^i, y^i), x^i \in \mathbb{R}^n, y^i \in \mathbb{R} \quad i = 1, \dots, l\}$$

Si modella direttamente le probabilità di appartenenza ad una classe.

Si definisce:

$$\Theta \in \mathbb{R}^{n+1} \quad \theta^o := \text{offset} \quad \theta_i = \text{peso della feature } i$$

La probabiità viene modellata tramite la dunzione logistica:

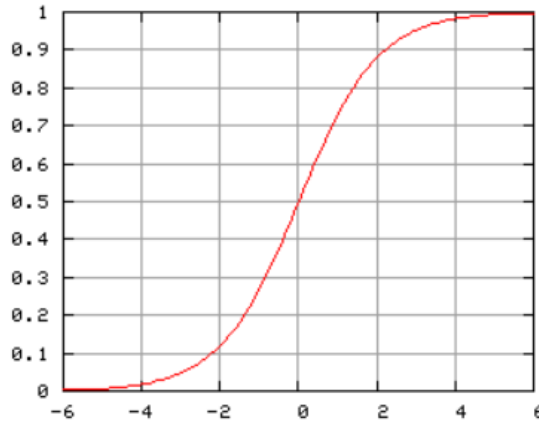


Figura 7.

che approssima la funzione gradino. Quindi:

$$P(C = 1 | x_1 = \bar{x}_1, \dots, x_n = \bar{x}_n) = \frac{1}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} = \frac{e^{(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}}{1 + e^{(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}}$$

In cui $(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)$ è un iperpiano e la distanza di un punto x è proporzionale al suo modulo.

Nota 109.

$$P(C = 1 | x_1 = \bar{x}_1, \dots, x_n = \bar{x}_n) + P(C = 0 | x_1 = \bar{x}_1, \dots, x_n = \bar{x}_n) = \frac{1}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} + \frac{e^{(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}}{1 + e^{(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} = 1$$

Corrisponde a considerare la distanze dei punti del Training Set dell'iperpiano. Maggiore probabilità fornisce una maggiore confidenza della classificazione,

19.2 Come si addestra la regressione logistica?

Per addestrare una funzione di errore di masimizzazione si prende una **funzione di massima verosimiglianza**:

$$L(\Theta) = \prod_{x^k \in \text{TS}^+} \frac{1}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} \prod_{x^k \in \text{TS}^-} \frac{1}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}}$$

Data la compelssità delle produttorie si massimizza $L(\Theta)$ a cui si applica il logaritmo:

$$\max_{\Theta} LL(\Theta) = - \sum_{x^k \in \text{TS}^+} \log(1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}) - \sum_{x^k \in \text{TS}^-} \log(1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)})$$

Questo è un problema di **massimizzazione non vincolata** e $LL(\Theta)$ è concava in Θ, θ . Allora tutti e soli i punti che annullano il gradiente sono ottimi globali.

Importante 110. $\nabla LL(\Theta)^T d \geq 0 \forall d$ ammissibile $\iff \nabla LL(\Theta) = 0$ per qualsiasi d ammissibile

Nota 111. Nel caso in cui si abbia $\min_{x \in \mathbb{R}^n} f(x)$ quello che vogliamo è una direzione d di discesa se e solo se $\nabla f(x)^T d < 0 \longrightarrow d = -\nabla f(x)$ se ho $\max f(x)$.

In particolare:

$$\begin{aligned} \frac{\partial \text{LL}}{\partial \theta_i} &= \sum_{x^k \in \text{TS}^+} \frac{e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)} x_i^k}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} - \sum_{x^k \in \text{TS}^-} \frac{e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)} x_i^k}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} = \\ &= \sum_{x^k \in \text{TS}^+} \frac{x_i^k}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} - \sum_{x^k \in \text{TS}^-} \frac{x_i^k}{1 + e^{-(\theta_o + \sum_{i=1}^n \theta_i \bar{x}_i)}} \end{aligned}$$

In pratica si stanno utilizzando le probabilità di errore per effettuare l'addestramento:

$$= \sum_{x^k \in \text{TS}^+} x_i^k \Pr(x^k \in \text{TS}^-) - \sum_{x^k \in \text{TS}^-} x_i^k \Pr(x^k \in \text{TS}^+)$$

Inoltre si pone:

$$\theta_i^k \rightarrow \theta_i^k + \alpha^k \left(\sum_{x^k \in \text{TS}^+} x_i^k \Pr(x^k \in \text{TS}^-) - \sum_{x^k \in \text{TS}^-} x_i^k \Pr(x^k \in \text{TS}^+) \right)$$

Importante 112.

Il parametro α_k deve essere opportunamente scelto.

$|\vartheta_i|$ è una misura dell'importanza della feature i – esima. Per ridurre l'overfitting si aggiunge un termine di regolarizzazione scelto in cross validation e strettamente positivo:

$$\max \text{LL}(\Theta) - \alpha \frac{1}{2} \|\Theta\|^2$$

19.3 Valutazione di un classificatore in probabilità

Per valutare un classificatore in probabilità, si passa da un assegnamento soft ad uno hard scegliendo una soglia t :

$$P(C = c | x) > t \longrightarrow C = c$$

Per ogni valore della soglia t posso definire:

$$S(t) := \text{insieme ottenuto con la soglia } t \text{ di punti nella classi } +1$$

Un'altra tecnica è quella del **G grand truth** (veri positivi) in cui si calcola al variare di t la precision:

$$\text{Precision}(t) = 100 \frac{|S(t) \cap G|}{|S(t)|} := \text{percentuale di positivi dichiarati veramente positivi}$$

Importante 113. Non è detto che all'aumentare di t la precision sia monotona

Inoltre, posso definire la recall:

$$\text{Recall}(t) = 100 \frac{|S(t) \cap G|}{|G|} := \text{percentuale di positivi effettivamente giusti}$$

A questo punto posso calcolare la **F1-Score**, è una media armonica che prende in considerazione la recall e la precision:

$$F1(t) = \frac{2 \text{Precision}(t) \text{Recall}(t)}{\text{Precision}(t) + \text{Recall}(t)}$$

19.4 ROC CURVE (receiver operating characteristic)

Per confrontare due classificatori posso utilizzare la ROC CURVE

$$\text{TPR}(t) = \text{Recall}(t) = 100 \frac{|S(t) \cap G|}{|G|} \quad \text{ordinata}$$

$$\text{FPR}(t) = \frac{|S(t) - G|}{|TS - G|} \quad \text{ascissa}$$

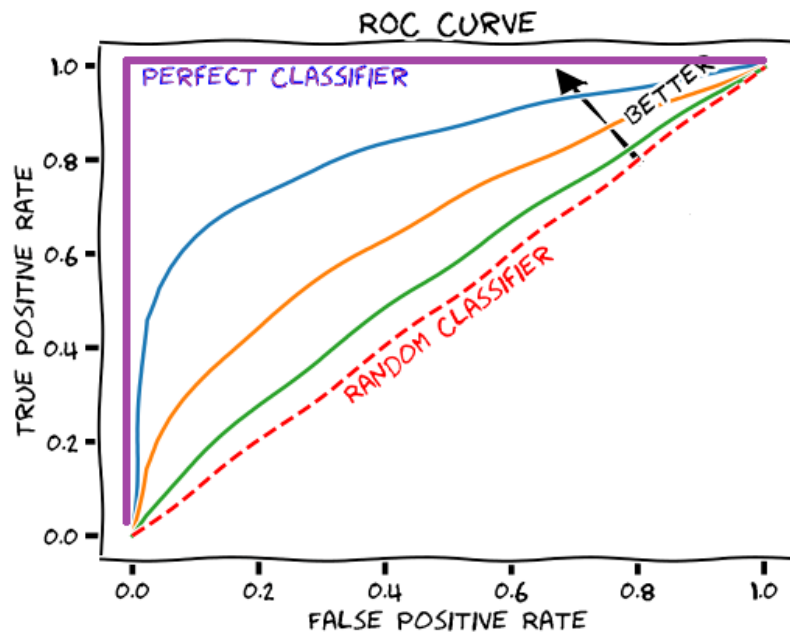


Figura 8.

Definiamo:

- **AUC**: area sotto la curva e tanto più è alta tanto è più buono il classificatore. Questo è un parametro indipendente da t .

20 Clustering

Il clustering è un tipo di classificazione **non supervisionata**, cioè ad ogni punto del Training Set non sono note né le effettive uscite né quante uscite ci si può aspettare.

Dato un Training Set $TS = \{x^i \in \mathbb{R}^n \mid i = 1, \dots, N\}$ vogliamo raggruppare i punti del TS molto simili in **clust** e in clust diversi punti non simili. In particolare, si vuole partizionare l'insieme di M cluster C_1, \dots, C_M in modo tale che:

$$C_i \cap C_j = \emptyset \quad e \quad \sum_{i=1}^M C_i = TS$$

A tale scopo, si necessita una definizione di dissimilarità tra punti e per arrivare a questa è necessario avere una nozione di distanza tra attributi:

1. **Attributi Continui** ($x_{ij} \in \mathbb{R}$): si può usare $l(|x_{ij} - x_{kj}|)$ $i, k = 1, \dots, N$
 l monotona crescente (es. norma);
2. **Attributi Ordinali**: si assumono k valori tra loro ordinati e ci si riconduce a variabili quantitative traducendo k valore numeri che rispettano l'ordinamento originale.
3. **Attributi Categorici Non Ordinati**: si deve definire la distanza ad-hoc: se ho dei valori r e r' : $L_{rr'} = 0, L_{rr'} \geq 0$

Una volta definita la distanza tra due attributi:

$$l_j(x^i, x^k) := \text{distanza da attributi } j$$

$$d(x^i, x^k) = \sum_{j=1}^n w_j l_j(x^i, x^k) \quad w_j = 1 \quad j = 1, \dots, n$$

Importante 114. I dati non vengono scalati poiché essendo non supervisionato non si hanno uscite

Definizione 115. *Variabili di Assegnamento*

$$\delta_{ij} = \begin{cases} 1 & \text{se } i \text{ è assegnato al cluster } j \\ 0 & \text{altrimenti} \end{cases} \quad \begin{matrix} i = 1, \dots, N \\ j = 1, \dots, M \end{matrix}$$

Definizione 116. *Variabile del centroide*

$z_j \in \mathbb{R}^n$ **centroide del cluster j** (è un punto significativo del cluster) $j = 1, \dots, M$

Per misurare la bontà di un cluster si utilizza la **distanza intraclassa**:

$$s_j = \sum_{i=1}^N \delta_{ij} d(x^i, z_j)$$

Il problema di clustering risulterà, quindi, del tipo:

$$\begin{cases} \min_{\delta, z} \sum_{j=1}^M \sum_{i=1}^N \delta_{ij} d(x^i, z_j) \\ \sum_{j=1}^M \delta_{ij} = 1 \quad i = 1, \dots, N \text{ (assegnamento del punto ad un solo cluster)} \\ \delta_{ij} \in \{0, 1\} \quad z_j \in \mathbb{R}^n \quad i = 1, \dots, N, j = 1, \dots, M \end{cases}$$

Prendendo in considerazione come distanza, la distanza euclidea:

$$d(x^i, z_j) = \|x^i - z_j\|^2$$

Posiamo scrivere il problema primale:

$$\begin{cases} \min_{\delta, z} \sum_{j=1}^M \sum_{i=1}^N \delta_{ij} \|x^i - z_j\|^2 \\ \sum_{j=1}^M \delta_{ij} = 1 \quad i = 1, \dots, N \text{ (assegnamento del punto ad un solo cluster)} \\ \delta_{ij} \in \{0, 1\} \quad z_j \in \mathbb{R}^n \quad i = 1, \dots, N, j = 1, \dots, M \end{cases}$$

Che può essere scritto in una forma equivalente:

$$\begin{cases} \min_{\delta, z} \sum_{j=1}^M \sum_{i=1}^N \delta_{ij} \|x^i - z_j\|^2 \\ \sum_{j=1}^M \delta_{ij} = 1 \quad i = 1, \dots, N \text{ (assegnamento del punto ad un solo cluster)} \\ \delta_{ij} \geq 0 \quad z_j \in \mathbb{R}^n \quad i = 1, \dots, N, j = 1, \dots, M \end{cases}$$

Per mostrare che i due problemi sono equivalenti deve esistere una trasformazione tra i punti ammissibili di un problema all'altro e viceversa.

Teorema 117. *Teorema Equivalenza P e Q*

Dati:

$$\begin{cases} \min_{x \in S} f(x) & (P) \\ \min_{y \in D} g(y) & (Q) \end{cases}$$

Se:

$$1. \quad \forall x \in S, \exists y \in D: g(y) \leq f(x) \longrightarrow \forall x \in S, \exists \rho(x) \in D: g(\rho(x)) \leq f(x)$$

$$2. \quad \forall y \in D, \exists x \in S: f(x) \leq g(y) \longrightarrow \forall y \in D, \exists \sigma(y) \in S: f(\sigma(y)) \leq g(y)$$

Allora il problema (P) ammette minimo globale se e solo se (Q) ammette minimo globale.

Dimostrazione.

Tenendo in considerazione il problema (Q), provando a fissare una delle due variabili:

- $\delta = \bar{\delta}$:

$$\min_z \sum_{j=1}^M \sum_{i=1}^N \bar{\delta}_{ij} \|x^i - z_j\|^2$$

Importante 118. I vincoli si levano

Lo si può decomporre in M problemi indipendenti:

$$\min_{z_j} \sum_{i=1}^N \bar{\delta}_{ij} \|x^i - z_j\|^2$$

Questo problema è coercivo e convesso in assenza di vincoli, quindi è possibile trovare il minimo globale annullando il gradiente rispetto a z :

$$\nabla_{z_j} f(\bar{\delta}, z) = 2 \sum_{i=1}^N \bar{\delta}_{ij} (x^i - z_j) = 0$$

$$\sum_{i=1}^N \bar{\delta}_{ij} x^i - \sum_{i=1}^N \bar{\delta}_{ij} z_j = 0$$

$$z_j = \frac{\sum_{i=1}^N \bar{\delta}_{ij} x^i}{\sum_{i=1}^N \bar{\delta}_{ij}}$$

In cui, il denominatore corrisponde a $|C_j|$ e la variabile z_j rappresenta il **baricentro** dei punti nel cluster j ed è anche l'ottimo globale. Inoltre, data l'assegnazione $\bar{\delta}_{ij}$ binaria non si può fare di meglio.

- $z = \bar{z}$:

$$\min_{\delta} \sum_{j=1}^M \sum_{i=1}^N \delta_{ij} \|x^i - \bar{z}_j\|^2 = \sum_{i=1}^N \delta_j^T d_j \quad \textbf{Lineare}$$

$$\delta_j = \begin{pmatrix} \delta_{1,j} \\ \vdots \\ \delta_{1,M} \end{pmatrix} \quad d_j = \begin{pmatrix} \|x^1 - \bar{z}_j\|^2 \\ \vdots \\ \|x^N - \bar{z}_j\|^2 \end{pmatrix}$$

$$s.t. \begin{cases} \sum_{j=1}^N \delta_{ij} = 1 \\ \delta_{ij} \geq 0 \end{cases}$$

è un problema di **simplexso** unitario.

Nota 119. Simplexso Unitario

Il problema è nella forma:

$$\begin{cases} \min f(x) \\ \sum_i x^i = 1 \\ x^i \geq 0 \end{cases}$$

Scriviamo le condizioni di KKT necessarie e sufficienti di ottimo globale introducendo i moltiplicatori rispettivamente μ, λ_i . Notiamo preliminarmente che l'insieme è compatto, quindi la soluzione esiste.

La lagrangiana è:

$$L(x, \lambda, \mu) = f(x) + \mu \left(\sum_i x^i - 1 \right) - \sum_{i=1}^N \lambda_i x^i$$

Annullamento della Lagrangiana:

$$\frac{\partial f}{\partial x_i}(x) + \mu - \lambda_i = 0 \quad i = 1, \dots, n$$

Complementarietà:

$$\lambda_i x_i = 0 \quad i = 1, \dots, n$$

$$\lambda_i \geq 0$$

Ammissibilità:

$$\sum_{i=1}^N x_i = 1$$

$$x_i \geq 0$$

Supponiamo che $x_i > 0 \longrightarrow \lambda_i = 0$:

$$\frac{\partial f}{\partial x_i} + \mu = 0 \longrightarrow \mu = -\frac{\partial f}{\partial x_i}$$

Sostituita nelle altre condizioni, $\forall j \neq i$:

$$\frac{\partial f}{\partial x_j} - \frac{\partial f}{\partial x_i} - \lambda_j = 0 \quad j = 1, \dots, n \quad j \neq i$$

Allora:

$$\lambda_i = \frac{\partial f}{\partial x_j} - \frac{\partial f}{\partial x_i} \geq 0$$

Quindi, se $x_i > 0$:

$$\frac{\partial f}{\partial x_j} \leq \frac{\partial f}{\partial x_i} \quad \forall j = 1, \dots, n$$

Che rappresenta la **condizioni di ottimo per problemi di semplice**.

Ricapitolando nel nostro problema:

$$\frac{\partial f}{\partial \delta_{ij}} = \|x^i - \bar{z}_j\|^2$$

Se il punto viene assegnato al cluster:

$$\delta_{ij} > 0 \longrightarrow \frac{\partial f}{\partial \delta_{ij}} \leq \frac{\partial f}{\partial \delta_{ij}} \longrightarrow \|x^i - \bar{z}_j\|^2 \leq \|x^i - \bar{z}_{\bar{j}}\|^2 \quad \forall j \neq \bar{j}$$

Si ottiene l'ottimo globale assegnando il punti i al centroide più vicino:

$$\delta_{ij}^* = \begin{cases} 1 & \text{se } j = \bar{j} \\ 0 & \text{altrimenti} \end{cases} \longrightarrow \text{Equivalenza soddisfatta} \quad \square$$

Riassumendo, (Q) è un **rilassamento** di (P) e per andare da (P) a (Q) non necessitiamo di trasformazioni poiché la funzione obiettivo è la stessa. Al contrario, se dobbiamo andare da (Q) a (P):

$$\begin{pmatrix} \bar{\delta} \\ \bar{z} \end{pmatrix} \in Q \longrightarrow \begin{pmatrix} \hat{\delta} \\ \bar{z} \end{pmatrix} \text{ dove } \hat{\delta}_{ij} = \begin{cases} 1 & \text{se } d(x^i, \bar{z}_j) \leq d(x^i, \bar{z}_k) \quad \forall k \neq j \\ 0 & \text{altrimenti} \end{cases}$$

Con la proprietà che:

$$\begin{pmatrix} \hat{\delta} \\ \bar{z} \end{pmatrix} \in (P) \quad e \quad f(\hat{\delta}, \bar{z}) \leq f(\bar{\delta}, \bar{z})$$

Dato che $\hat{\delta}$ è l'ottimo globale, allora i due problemi sono equivalenti.

20.1 Algoritmo di K-Means

Algoritmo

Dati $z^0 \in \mathbb{R}^n, k = 0$

While criterio di arresto soddisfatto

1. Assegnamento dei punti:

For $i = 1, \dots, N$ poni;

$$\delta_{ij}^{k+1} = \begin{cases} 1 & \text{se } j = j^* \text{ ovvero } \arg \min \|x^i - z_j\|^2 \\ 0 & \text{altrimenti} \end{cases}$$

2. Ricalcolo dei centroidi:

For $j = 1, \dots, M$ poni:

$$z_j^{k+1} = \frac{\sum \delta_{ij}^{k+1} x^i}{\sum_{i=1}^N \delta_{ij}^{k+1}}$$

End

Criterio di Arresto

$$\frac{f(\delta^{k+1}, z^{k+1}) - f(\delta^k, z^k)}{1 - f(\delta^k, z^k)} \leq \varepsilon$$

Teorema 120. Teorema Convergenza K-Means

L'algoritmo di K-Means con $\varepsilon = 0$ converge in un numero finito di passi in δ^*, z^* tale che:

$$f(\delta^*, z^*) \leq f(\delta, z^*) \quad \forall \delta \text{ ammissibile}$$

$$f(\delta^*, z) \leq f(\delta^*, z^*) \quad \forall z \in \mathbb{R}^m \text{ ammissibile}$$

Dimostrazione.

Supponiamo per assurdo che l'algoritmo non termini, quindi:

$$f(\delta^{k+1}, z^{k+1}) < f(\delta^k, z^k) \quad \forall k$$

Ad ogni passo:

$$f(\delta^k, z^{k+1}) \leq f(\delta^k, z) \quad \forall z \in \mathbb{R}^m$$

e:

$$f(\delta^{k+1}, z^k) \leq f(\delta, z^k) \quad \forall \delta \text{ ammissibile}$$

Allora \exists sottosuccessione k tale che:

$$\{\delta^k, z^k\}_K \quad \delta^k = \bar{\delta}$$

$$\forall k \in K \exists \text{indice } l(k) \quad k + l(k) \in K$$

$$f(\delta^k, z^k) = f(\bar{\delta}, z^k)$$

$$f(\delta^{k+l(k)}, z^{k+l(k)}) = f(\bar{\delta}, z^{k+l(k)})$$

Inoltre se:

$$f(\delta^k, z^k) = f(\bar{\delta}, z^k) \leq f(\bar{\delta}, z) \quad \forall z$$

Poiché $l(k) + k > k$ devo avere che:

$$f(\delta^{k+l(k)}, z^{k+l(k)}) < f(\delta^k, z^k) = f(\bar{\delta}, z^k) \leq f(\bar{\delta}, z = z^{k+l(k)})$$

Assurdo poiché abbiamo supposto che:

$$f(\bar{\delta}, z^{k+l(k)}) \leq f(\bar{\delta}, z = z^{k+l(k)})$$

□

20.2 Scelta del numero di Cluster

Dato che il problema di clusterizzazione è un problema di ottimizzazione non supervisionata, si ha difficoltà a scegliere il numero di cluster da utilizzare poiché non si conoscono il numero di uscite. Quindi ci sono delle euristiche che permettono di scegliere il numero di Cluster:

- **Elbow Method:**

Si definisce **elbow** il punto in cui f non decresce di molto e si hanno le seguenti proprietà:

$$f_k^* \ll f_k \text{ quando } k < k^*$$

$$f_{k+1}^* \approx f_k^* \text{ quando } k > k^*$$

- **Silhouette:**

Si fissa k e si risolve per quel valore. In particolare:

- $\forall i = 1, \dots, N \quad i \in C(i)$, calcolo:

- $a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i} d(x_i, x_j)$

- $b(i) = \min_{k \neq C(i)} \frac{1}{|C_k|} \sum_{j \in C_k} d(x_i, x_j)$

- $S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad \text{se } |C_i| > 0 \quad \text{e } S(i) = 0 \quad \text{se } |C(i)| = 1$

- Al variare di k calcolo:

$$S_k = \frac{1}{N} \sum_{i=1}^N S(i)$$

Dato che voglio il valore massimo, scelgo il k più alto.

PACIFICI

1 Alberi decisionali per la classificazione

1 Introduzione

Apprendimento statistico

L'**apprendimento statistico** sono delle tipologie di apprendimento che può essere svolta sia in maniera automatica sia tramite ML (Machine Learning).

In particolare si compone:

$$x^{(i)} \in \mathbb{R}^p \rightarrow \text{osservazione}(\text{punto})$$

$$y^{(i)} \in \mathbb{R} \rightarrow \text{risposta} \quad \text{con} \quad i = 1, \dots, n$$

Definiamo il piano/spazio ottenuto dalle osservazioni, **spazio delle feature**.

Il nostro obiettivo è quello di ottenere la funzione di uscita in funzione delle osservazioni date. Questo processo viene detto **inferire**. Formalmente:

$$y = f(x) \quad \text{con } x \in \mathbb{R}^p$$

Riassumendo il nostro problema, noti:

- $(x^{(i)}, y^{(i)})$ con $i = 1, 2, \dots, n$ osservazioni disponibili \rightarrow **Training set**;

- $x^{(i)} \in \mathbb{R}^p :=$ variabili di input indipendenti o anche dette **feature**;

Ottenere:

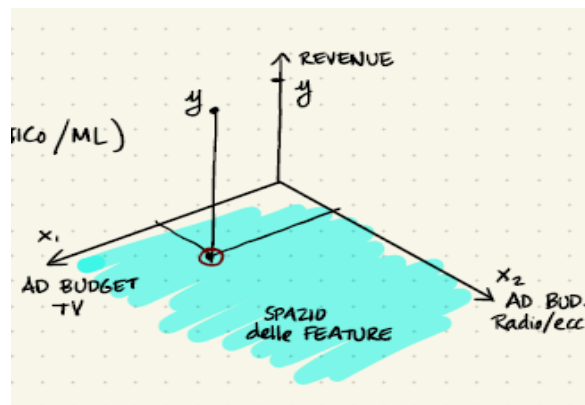
- $y \in \mathbb{R} :=$ variabile di output / **Class** / **etichetta** / **label**.

In altre parole, dato $x \neq x^{(i)}$ qual è la sua uscita? Per rispondere a questa domanda è necessario che la $f(x)$ resistisca esattamente $y^i = f(x^{(i)})$, cioè che ci sia un **fitting** con i dati del **TS (Training Set)**.

Un modello di questo tipo è un **modello di inferenza/predittivi** poiché permettono di intuire la prossima uscita e producono una y data una certa x .

Apprendimento supervisionato

L'**apprendimento supervisionato** è un tipo di apprendimento per cui l'uscita y è associata un punto dello spazio delle feature. Graficamente:



Apprendimento non supervisionato

L'**apprendimento non supervisionato** è caratterizzato dal fatto che non esiste in corrispondenza di una osservazione nello spazio delle feature la relativa uscita.

Response

Il **response** è l'uscita y che si ottiene dal fitting con il TS.

La response può essere di due tipologie:

- qualitativa: la **response** è discreta e finita e il suo risultato è una **classificazione** di punti in classi/label. Si cerca quindi di assegnare una classe ad una nuova osservazione conoscendo i dati a disposizione.

Sono metodi basati su **alberi di decisione**.

- quantitativa: tipicamente è una uscita continua ed i metodi più utilizzati sono la Regressione, Super Vector Machine, Reti Neurali etc..

Quindi il nostre Data Set sarà composto dalle righe che corrispondono alle osservazioni, dalle colonne che sono gli attributi delle varie osservazioni ed a ogni riga, a cui corrispondere un'osservazione con i suoi attributi, è associata una etichetta.

2 Alberi Decisionali

Gli alberi decisionali applicano ricorsivamente metodi euristici per creare delle divisioni nello spazio delle feature per mettere in evidenza un certo aspetto del dataset.

Per ottenere un albero decisionale dato uno spazio delle feature contenente le osservazioni, occorre:

- **Partizionare/stratificare/segmentare** lo spazio delle feature in modo tale che le partizioni siano il più omogenee possibile in base alla prevalenza o meno di etichette uguali;
- Si associa a ciascuna regione un'etichetta che corrispondendo all'etichetta predominante nella regione anche in presenza di regioni miste;
- Ipotizzando che ci arrivi una nuova osservazione, essa cadrà nello spazio delle feature all'interno di una regione R_i (ottenuta dalle precedenti partizioni) ed assumerà l'etichetta della regione corrispondente;
- Si delimitano i confini delle regioni R_i delle osservazioni e vi si assegnano dei valori;
- L'albero che si formerà sarà costituito da **nodi foglia** e **nodi interni**:
 - i **nodi interni**, sono i **test** che ci permettono di classificare le osservazioni tramite la verifica degli attributi. In particolare se si esamina una sola feature il test viene detto **univariato** ($|$); altrimenti **multivariato** ($/$);
 - i **nodi foglia** sono i risultati della nostra classificazione a cui viene associata l'etichetta della regione R_i a cui appartengono.

N.B.:

Confronto test univariato vs test multivariato:

- Nel caso di test univariato si potrebbero fare errori di classificazione, ma l'albero risultante è più semplice;
- Nel caso di test multivariato non vi sono errori di classificazione per etichette già note, ma è più soggetto ad errori rispetto al caso univariato (**errore di generalizzazione**), ma l'albero presenta una complessità maggiore data la presenza di numerosi nodi;

Gli attributi di un albero decisionale sono di due tipi:

- **categorici**: cioè qualitativi;
- **continui**: cioè quantitativi;

Dal **training data** cioè l'insieme degli attributi e dell'uscita, si ottiene l'**albero decisionale**.

Ovviamente l'albero risultante varia in base all'ordine in cui i test vengono considerati.

Ogni volta in cui si fa un test, si ha una divisione dei punti del TS e il nostro obiettivo è quello di produrre degli insiemi più omogenei possibili su cui applicare il test set.

Pro vs Contro alberi di decisione

- Facilità da spiegare/interpretare e si mima il comportamento umani;
- Facile da manovrare e le variabili sono di tipo qualitativo;

- Strumenti meno accurati;
- "Nervosismo": piccole variazioni di dati portando ad alberi differenti.

Per ovviare a questa ultima problematica si ricorre a dei **metodi di aggregazione** in cui si costruiscono diversi alberi per poi essere uniti tra loro.

TDITD Family

La **TDITD family** compende una serie di Top-Down Induction Decision Trees come il CART e ID3. Questi alberi utilizzando l'induzione, cioè la ricorsione, per creare degli split nello spazio delle feature. Tuttavia, questi potrebbero non rappresentare opportunamente le caratteristiche del dataset.

In particolare, l'albero è usato per classificare i punti del test set secondo gli split e le labels. Questi alberi vengono detti **top-down** poiché, partendo dal nodo radice, determinano uno split risolvere un problema di ottimizzazione, in genere di minimo rispetto ad una misura di impurità, per poi applicare la stessa regola alle foglie appena generate.

Questo tipo di classificazione, inoltre, viene detta **greedy**, poiché ogni divisione viene determinata in isolamento senza considerare i possibili impatti che si potrebbero avere nel futuro.

In conclusione, questa famiglia di alberi soffre problemi di complessità.

Algoritmo di Hunt

Algoritmo

1. INIT:
 $list := \{\text{training set}\}$
 $node_list := \{\text{root}\}$, $D_{root} := list$: insieme delle proprietà che soddisfano i test
 $test_pool := \{\dots\}$
2. Pick $t \in node_list$
3. Definire $D_t \subseteq list$ /*punti del TS che raggiungono il nodo t */
4. IF $D_t = \emptyset$ THEN t è LEAF label t as "*" (null)
5. IF D_t ha elementi tutti di classe K THEN t è LEAF label t as k
6. IF D_t ha elementi di classi diverse THEN
 - a. Choose $test \in test_pool$
 - b. split di D_t in sottoinsiemi $D_t^1, D_t^2, \dots, D_t^p$
 - c. update $node_list := node_list \cup \{t_1\} \cup \{t_2\} \cup \dots \cup \{t_p\}$
7. $node_list := node_list \setminus \{t\}$
8. IF $node_list = \emptyset$ THEN STOP ELSE goto 1

Algoritmo ID3

Algoritmo

ID3 (idea)

1. Selezione una "window" (sottoinsieme del TS);
2. Costruisci un DT usando i punti della window;
3. IF \exists elementi $\in TS \setminus \{\text{window}\}$ che non sono classificati correttamente THEN
 $window := window \cup \text{questi elementi}$
 ELSE Return DT and STOP

3 Errori negli alberi di classificazione

- **Errore di classificazione:** = sDato un DT, si ha errore di classificazione se il numero di elementi del TS che vengono etichettato in modo corretto. Una sua minimizzazione troppo marcata non porta ad un buon albero poiché soffre dell'**overfitting**. In aggiunta, nel caso di **overfitting**, si eliminano alcuni rami per avere un albero più “semplice”, **Tree Pruning**;
- **Errore di generalizzazione:** = è un errore atteso che il DT compie sui punti \notin TS. Viene stimato con il test set.

Indici di (in)purity

Nella TDIDT “family” la costruzioni dei DT, di norma si sceglie la **strategia greedy**: in ogni passo si sceglie il test che ottimizza un dato criterio disinteressandosi completamente di quello che può accadere in seguito.

Tuttavia, vogliamo un certo livello di **omogeneità**. A tale scopo si introducono i **purity indices** che indica quanto un sottoinsieme è omogeneo, cioè composto dalle stesse classi.

Errore di classificazione

Sia $p_i \in [0, 1]$ frequenza con cui l’etichetta i compare nell’insieme. Allora:

$$e = 1 - \max_i \{p_i\} \longrightarrow \arg_i \max \{p_i\} := \text{moda}$$

$$0 \leq e \leq 1 - \frac{1}{k} \quad \text{con } k := \# \text{ di etichette}$$

Gini Index

$$g = 1 - \sum_{i=1}^k p_i^2 \quad 0 \leq g \leq 1 - \frac{1}{k}$$

Entropia

$$E = - \sum_{i=1}^k p_i \log_k(p_i) \quad 0 \leq E \leq 1$$

Dimostrazione.

Si ha omogeneità massima $p_i = 1, p_j = 0 \quad \forall j \neq i$

Per $p = \frac{1}{k} \quad \forall i$, si ha che:

$$\begin{aligned} E &= - \sum_{i=1}^k \left(\frac{1}{k} \right) \log_k \left(\frac{1}{k} \right) = \\ &= \sum_{i=1}^k \left(\frac{1}{k} \right) (-1) \log_k \left(\frac{1}{k} \right) = \\ &= \sum_{i=1}^k \left(\frac{1}{k} \right) \log_k \left(\frac{1}{k} \right)^{-1} = \sum_{i=1}^k \left(\frac{1}{k} \right) 1 = 1 \end{aligned}$$

□

Importante 121. Noi vogliamo usare gli indici per capire il tipo di test da utilizzare per generare foglie con un grado di omogeneità maggiore. A tale scopo, si sommano gli indici di impurità dei figli:

Test

Si possono effettuare vari tipi di test, i più comuni sono i seguenti:

- **ERR_TEST** = $\sum_{j \text{ figlio}} \text{ERR_CLASS}(j)$; si sceglie il test con indice minore
- **GINI_TEST** = $\sum_{j \text{ figlio}} \text{GINI}(j) \frac{n_j}{n}$ si sceglie il test con indice minore

In cui n_j sono i punti nel nodo figlio e n i punti nel nodo padre

- **INFORMATION GAIN_TEST** = misura quando si guadagna in omogeneità negli insiemi con un certo tipo di test:

$$\text{INFORMATION_GAIN_TEST} = \text{ENTROPY}(\text{padre}) - \sum_{j \text{ figlio}} \left(\frac{n_j}{n} \right) \text{ENTROPY}(j)$$

Si sceglie il test con indice maggiore;

- **GAIN_RATIO_TEST** = equivale all'information gain test, ma con un andamento uniforme:

$$\text{GAIN_RATIO_TEST} = \frac{\text{INFORMATION_GAIN_TEST}}{\left[-\sum_{j \text{ figlio}} \left(\frac{n_j}{n} \right) \log \left(\frac{n_j}{n} \right) \right]}$$

4 Albero Ottimo di Classificazione

Introduzione

Il problema di albero decisionale ottimale tenta di risolvere il problema di classificazione di un dataset creando un albero decisione per ottenere l'ottimalità globale.

Assumiamo:

- Sia dato un Training Data (X,Y) contenente n osservazioni $(x_i, y_i), i = 1, \dots, n$ ognuna delle quali con p feature;
- Modello univariato, cioè i test sono effettuati su un attributo alla volta;
- $y_i \in \{1, \dots, k\}$ K classi/etichette;
- Attributi/feature continue (quantitative) $x_i \in [0, 1]^p$;
- Criterio: $\min(\text{errore di classificazione} + \text{complexity})$.

I metodi degli alberi decisionali tentano di effettuare una partizione ricorsiva $[0, 1]^p$ per ottenere delle regioni disgiunte che rappresentano l'albero decisionale. L'albero finale sarà formato da nodi foglia e nodi di diramazione:

- I **nodi di diramazione** si dividono con parametri **a** e **b**. Per un dato punto i , $a^T x_i < b_i$ il punto seguirà la diramazione sinistra dal nodo, altrimenti la destra;

- I **nodi foglia** sono assegnati ad una classe che determinerà la predizione di tutti i punti che cadono all'interno della foglia. La classe viene assegnata in base alla moda della foglia.

Decisioni del problema

Nella creazione dell'albero, in ogni iterazione, dobbiamo effettuare un numero fissato di decisioni:

- In ogni nodo dobbiamo decidere se dividere o fermarci;
- Dopo che abbiamo scelto di fermarci in un nodo, dobbiamo scegliere l'etichetta da assegnare al nuovo nodo foglia;
- Dopo che abbiamo scelto di dividere, dobbiamo scegliere su quali delle variabili effettuare il test;
- Quando classifichiamo i punti di training secondo la costruzione dell'albero, dobbiamo decidere a quali dei nodi foglia assegnare un punto tale che la struttura dell'albero viene rispettata.

Struttura del problema

Consideriamo il problema di costruire un albero ottimale di decisione con la massima profondità di D . Dato questo parametro:

- Albero binario completo (sotto albero di);
- Profondità fissa $D \Rightarrow \# \text{ nodi } 2^{D+1} - 1 = \mathbb{T}$

Notazione 122.

- $t \in \{T_B \cup T_L\} \setminus \{1\}$;
- $p(t) := \text{nodo padre/genitore}$;
- $A(t)$ è l'insieme degli **antenati** del nodo t ;
- $A_L(t), A_R(t)$ sono rispettivamente gli **insiemi antenati di sinistra e di destra**, presi, partendo dalla radice fino a t , scegliendo rispettivamente la direzione destra o sinistra.;
- $A(t) = A_L(t) \cup A_R(t)$

Quindi dividiamo i nodi in:

- **Branch Nodes:** sono tutti i nodi $t \in T_B = \left\{1, 2, \dots, \left\lfloor \frac{T}{2} \right\rfloor\right\}$ che applicano la divisione nella forma $a^T x < b$ se lo soddisfano seguirà il ramo sinistro, altrimenti il destro.

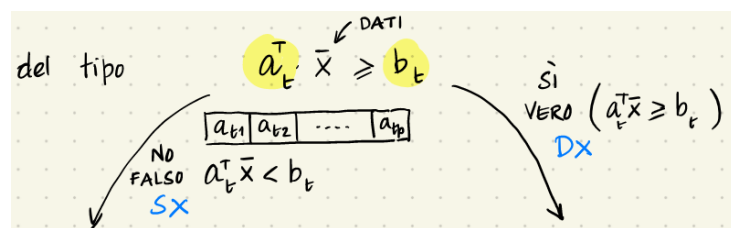


Figura 9.

- **Leaf Nodes:** Sono tutti i nodi $t \in T_L = \left\{ \left\lceil \frac{T}{2} \right\rceil, \dots, T \right\}$ formano una classe di predizione per ogni punti che cade nel nodo foglia.

Variabili

Supponiamo $t \in T_B$

- $a_t \in \mathbb{R}^p$ coefficiente del test $\longrightarrow a_t \in \{0, 1\}^p \quad \bar{x} \in [0, 1]^p \quad \sum_{i=1}^p a_t = 1$
 $(a_t \cdot \bar{x}) \in [0, 1]$
- $b_t \in \mathbb{R} \longrightarrow b_t \in [0, 1]$ è uno scalare
- $d_t \in \{0, 1\}$ variabile che indica se si effettua o no un test $d_t \in \{0, 1\}$

Nota 123.

L'idea è quella di non permettere la divisione in un branch node. A tale scopo utilizziamo la variabile d_t per sapere su quale branch node si applica la divisione. Se un nodo di diramazione non applica lo split, allora viene modellato con $a_t = 0$ e $b_t = 0$.

Di conseguenza, si forzano tutti a seguire il ramo destro del nodo t . Questo permette all'albero non smettere di crescere.

Vincoli

1. $\sum_{i=1}^p a_{t,i} = d_t \quad \forall t \in T_B$ se effettuo o no il test
2. $0 \leq b_t \leq d_t \quad \forall t \in T_B$
3. $d_t \leq d_{p(t)} \quad \forall t \in T_B$ rappresenta la coerenza con d_t se non si effettua lo SPLIT no split.

Allocazione

Il problema di **allocazione** riguarda l'assegnazione dei punti del training set alle foglie dell'albero decisionale e associarli agli errori che sono indotti da questa struttura.

Variabili

- Indichiamo con $z_{it} = \{0, 1\} \quad i = 1, \dots, n, t \in T_L$ se il punto x_i appartiene o meglio alla foglia t ;
- Per verificare se una foglia possiede dei punti. Utilizziamo:

$$l_t \in \{0, 1\} \quad t \in T_L \quad l_t = \begin{cases} 1 & \text{se la foglia } t \text{ ha assegnati dei punti} \\ 0 & \text{altrimenti} \end{cases}$$

- Il numero minimo di punti in ogni foglia N_{\min} ,

Vincoli

$$\begin{aligned} z_{it} &\in \{0, 1\} \quad i = 1, \dots, n \quad t \in T_L \\ \sum_{t \in T_L} z_{it} &= 1 \quad i = 1, \dots, n \quad \text{deve } \exists \text{ una foglia per il punto } i \\ \sum_{i=1}^n z_{it} &\geq N_{\min} l_t \quad t \in T_L \quad \text{Una foglia deve essere sufficientemente popolata} \end{aligned}$$

A questo punto dobbiamo aggiungere dei vincoli di coerenza che impongono le divisioni richieste fino al raggiungimento della foglia:

- Per i test **non soddisfatti**, deve valere: (9)

$$a_m^T x_i < b_i + M_1(1 - z_{it}) \quad i = 1, \dots, n \quad \forall t \in T_B, \forall m \in A_L(t)$$

- Per i test **soddisfatti**, deve valere:

$$a_m^T x_i \geq b_i - M_2(1 - z_{it}) \quad i = 1, \dots, n \quad \forall t \in T_B, \forall m \in A_R(t)$$

Nota 124.

Il vincolo •(9) non è supportato dai risolutori. Quindi occorre convertirlo in una forma che non utilizza la disuguaglianza stretta. A tale scopo, aggiungiamo un fattore piccolo ε e il vincolo diventa:

$$a_m^T x_i + \varepsilon \leq b_i + M_1(1 - z_{it}) \quad i = 1, \dots, n \quad \forall t \in T_B, \forall m \in A_L(t)$$

Tuttavia, se ε è troppo piccolo, potrebbe casare **instabilità** numeriche, quindi occorre renderlo il più grande possibile senza influenzare la soluzione in qualsiasi problema valido.

Occorre definire un e_j per ogni feature j ; il più grande valore valido è la distanza diversa da 0 più piccola tra i valori adiacenti della feature j :

$$\varepsilon_j = \min \{ |x_j^{(i+1)} - x_j^i| \mid x_j^{(i+1)} \neq x_j^i \quad i = 1, \dots, n-1 \}$$

In cui x_j^i è il più grande valore della feature j -esima.

Dalla nota 124 possiamo utilizzare ε nel vincolo, dove ε_j corrisponde alla feature sul quale effettuiamo lo split:

$$a_m^T(x_i + \varepsilon) \leq b_m + M_1(1 - z_{it}) \quad i = 1, \dots, n \quad \forall t \in T_B, \forall m \in A_L(t)$$

Scelta M_1 e M_2

M_1 e M_2 sono dette costanti big-M. Valutando i vincoli, sappiamo che:

- Il massimo valore di $a_m^T(x_i + \varepsilon) - b_m$ è $1 + \varepsilon_{\max}$ in cui $\varepsilon_{\max} = \max_j \{\varepsilon_j\}$. Quindi vogliamo:

$$M_1 \geq 1 + \varepsilon_{\max}$$

- Il massimo valore di $b_t - a_t^T x_i$ è 1 e quindi possiamo scegliere:

$$M_2 \geq 1$$

Ottenendo, in conclusione gli ultimi due vincoli:

$$a_m^T x_i + \varepsilon \leq b_i + (1 + \varepsilon_{\max})(1 - z_{it}) \quad i = 1, \dots, n \quad \forall t \in T_B, \forall m \in A_L(t)$$

$$a_m^T x_i \geq b_i - (1 - z_{it})$$

Funzione Obiettivo

La funzione obiettivo che vogliamo determinare si basa su due criteri di minimizzazione riguardo:

- la **complessità dell'albero**: corrisponde alla dimensione dell'albero decisionale ed è fissata a priori da D ;
- **Errore di classificazione**: assegnamo ad un'etichetta non corretta il costo di 1 e ad una corretta il costo 0;

Dati

Sia:

$$y_{ik} = \begin{cases} +1 & \text{se } y_i = k \text{ il punto } i \text{ ha etichetta } k \\ -1 & \text{altrimenti} \end{cases} \quad i = 1, \dots, n; k = 1, \dots, K$$

Inoltre, $N_{kt} :=$ il numero di punti per l'etichetta k nel nodo t e $N_t :=$ il numero di punti nel nodo t :

$$N_t = \sum_{i=1}^n z_{it} \quad t \in T_L$$

$$N_{kt} = \frac{1}{2} \sum_{i=1}^n (1 + y_{ik}) z_{ik} \quad k = 1, \dots, K$$

Task - obiettivo

Vogliamo assegnare alle foglie $t \in T_L$ con $l_t=1$ l'etichetta a lui associata. Essa è la moda delle etichette assegnate nel nodo t .

Inoltre dobbiamo assegnare un'etichetta ad ogni nodo t nell'albero, definita come:

$$c_t \in \{1, \dots, K\}$$

Ovviamente l'etichetta ottima è quella che corrisponde alla moda di in quel nodo, cioè alla frequenza in cui quella determinata etichetta compare nel nodo:

$$c_t = \arg \max_{k=1, \dots, K} \{N_{kt}\}$$

Al fine di tenere traccia della predizione di ogni nodo utilizziamo la variabile binaria c_{kt} :

$$c_{kt} = \begin{cases} 1 & \text{se l'etichetta della foglia } t \text{ è } k \\ 0 & \text{altrimenti} \end{cases}$$

A questo punto ci dobbiamo assicurare che in ogni nodo ci sia una sola predizione. A tale scopo introduciamo il vincolo:

$$\sum_{k=1}^K c_{kt} = l_t \quad \forall t \in T_L$$

Dato che sappiamo come scegliere l'etichetta in maniera ottimale in ogni nodo, definiamo L_t come l'**errore di classificazione** in ogni nodo che risulta uguale al numero di punti nel nodo minore del numero di punti dell'etichetta più comune:

$$L_t = N_t - \max_{k=1, \dots, K} \{N_{kt}\} = \min_{k=1, \dots, K} \{N_t - N_{kt}\} \quad t \in T_L$$

Che può essere linearizzato nel seguente modo:

$$\begin{aligned} L_t &\geq N_t - N_{kt} - M(1 - c_{kt}) & k=1, \dots, K & \quad \forall t \in T_L \\ L_t &\leq N_t - N_{kt} + M c_{kt} & k=1, \dots, K & \quad \forall t \in T_L \\ L_t &\geq 0 & & \quad \forall t \in T_L \end{aligned}$$

Dove M è una costante sufficientemente grande che rende il vincolo intattivo dipendendo dal valore di c_{kt} . Un possibile valore che possiamo prendere è $M = n$.

Il costo totale dell'errore di classificazione è dato da $\sum_{t \in T} L_t$ e la complessità è il numero di split nell'albero, dato da $\sum_{t \in T_b} d_t$. A questo punto possiamo normalizzare l'errore di classificazione con una **soglia di accuratezza** \hat{L} , ottenuta dalla semplice predizione dell'etichette più popolari dell'intero dataset rendendolo di conseguenza, indipendente dalla quantità α .

Finalmente possiamo dunque scrivere la funzione obiettivo:

$$\min \frac{1}{\hat{L}} \sum_{t \in T_L} L_t + \alpha \sum_{t \in T_B} d_t$$

Riepilogo

Mettendo insieme tutte queste formulazioni delle sezioni precedenti possiamo scrivere il nostro modello OCT:

$$\begin{aligned} \min \frac{1}{\hat{L}} \sum_{t \in T_L} L_t + \alpha \sum_{t \in T_B} d_t \\ \\ \begin{aligned} L_t &\geq N_t - N_{kt} - M(1 - c_{kt}) & k=1, \dots, K & \quad \forall t \in T_L \\ L_t &\leq N_t - N_{kt} + M c_{kt} & k=1, \dots, K & \quad \forall t \in T_L \\ L_t &\geq 0 & & \quad \forall t \in T_L \\ N_{kt} &= \frac{1}{2} \sum_{i=1}^n (1 + y_{ik}) z_{ik} & k=1, \dots, K \\ N_t &= \sum_{i=1}^n z_{it} & \forall t \in T_L \\ \sum_{k=1}^K c_{kt} &= l_t & \forall t \in T_L \\ a_m^T x_i &\geq b_i - M_2(1 - z_{it}) & i=1, \dots, n & \quad \forall t \in T_B, \forall m \in A_R(t) \\ a_m^T x_i + \varepsilon &\leq b_i + (1 + \varepsilon_{\max})(1 - z_{it}) & i=1, \dots, n & \quad \forall t \in T_B, \forall m \in A_L(t) \\ \sum_{t \in T_L} z_{it} &= 1 & i=1, \dots, n \\ \sum_{i=1}^n z_{it} &\leq l_t & t \in T_L \\ \sum_{i=1}^n z_{it} &\geq N_{\min} l_t & t \in T_L \\ \sum_{i=1}^p a_{t,i} &= d_t & \forall t \in T_B \\ 0 &\leq b_t \leq d_t & \forall t \in T_B \\ d_t &\leq d_{p(t)} & \forall t \in T_B \\ z_{it}, l_t &\in \{0, 1\} & i=1, \dots, n & \quad \forall t \in T_L \\ a_{jt}, d_t &\in \{0, 1\} & j=1, \dots, p & \quad \forall t \in T_B \end{aligned} \end{aligned}$$

Importante 125.

La difficoltà del modello è determinata dal numero di variabili z_{it} in quale sono $n \cdot 2^D$.

Inoltre, occorre specificare a priori tre parametri: la massima profondità dell'albero D , la dimensione minima della foglia N_{\min} e il parametro di complessità. Questa scelta viene detta **tuning** (sincronizzazione).

Il tempo di calcolo di un OCT è nell'ordine dei minuti.

Warm Starts

I risolutori traggono grande beneficio quando gli viene fornito una soluzione intera ammissibile come un **warm start** per il processo di soluzione.

Inserendo una forte soluzione di **warm start** il risolutore incrementa notevolmente la velocità con cui è in grado di generare soluzioni fortemente ammissibili. Quindi fornisce un upper-bound iniziale sulla soluzione ottima che permette di effettuare l'azione di **pruning** e inoltre fornisce un punto di partenza per la ricerca locale euristica.

Se noi avessimo una soluzione generata per la profondità D , questa soluzione è un valido **warm start** per la profondità $D + 1$. Questo è importante poiché la difficoltà del problema aumenta con il crescere della profondità e quindi potrebbe risultare svantaggioso eseguire un problema MIO (Mixed-Integer-Optimization) con una piccola profondità per generare un forte warm start.

Quindi, data una soluzione iniziale ammissibile è possibile confrontare la velocità del solutore nel determinare una soluzione ottima.

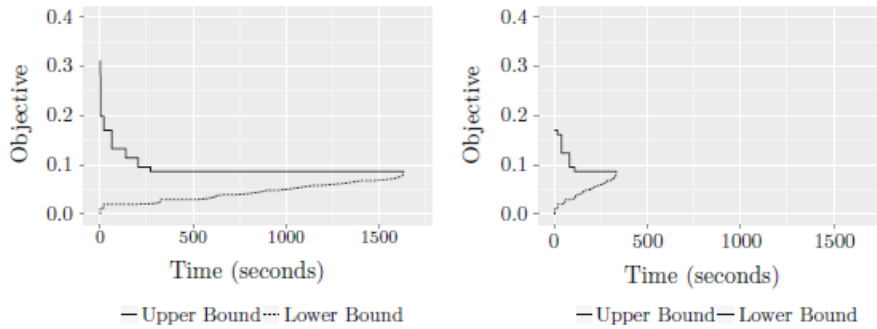


Figura 10. A sinistra senza Warm star; a destra con

Comi si può evincere dall'immagine, si ha convergenza quando upperbound e lower bound coincidono.

Tipicamente il warm start viene fornito da un algoritmo diverso.

5 Addestramento di un OCT

Addestrare una OCT consiste nello scegliere i parametri adatti al modello che vogliamo classificare. Questa azione viene detta **tuning** dei parametri e coinvolge la scelta di N_{\min} , D , α per minimizzare l'errore di training e la complessità dell'albero

$$\min \frac{1}{L} \sum_{t \in T_L} L_t + \alpha \sum_{t \in T_B} d_t$$

Importante 126.

- $\alpha :=$ è il fattore di importanze se $\alpha = 0 \rightarrow$ overfitting

Nota 127.

Un tipico approccio per la scelta di α è quello di discretizzare lo spazio di ricerca e i test di ogni valore.

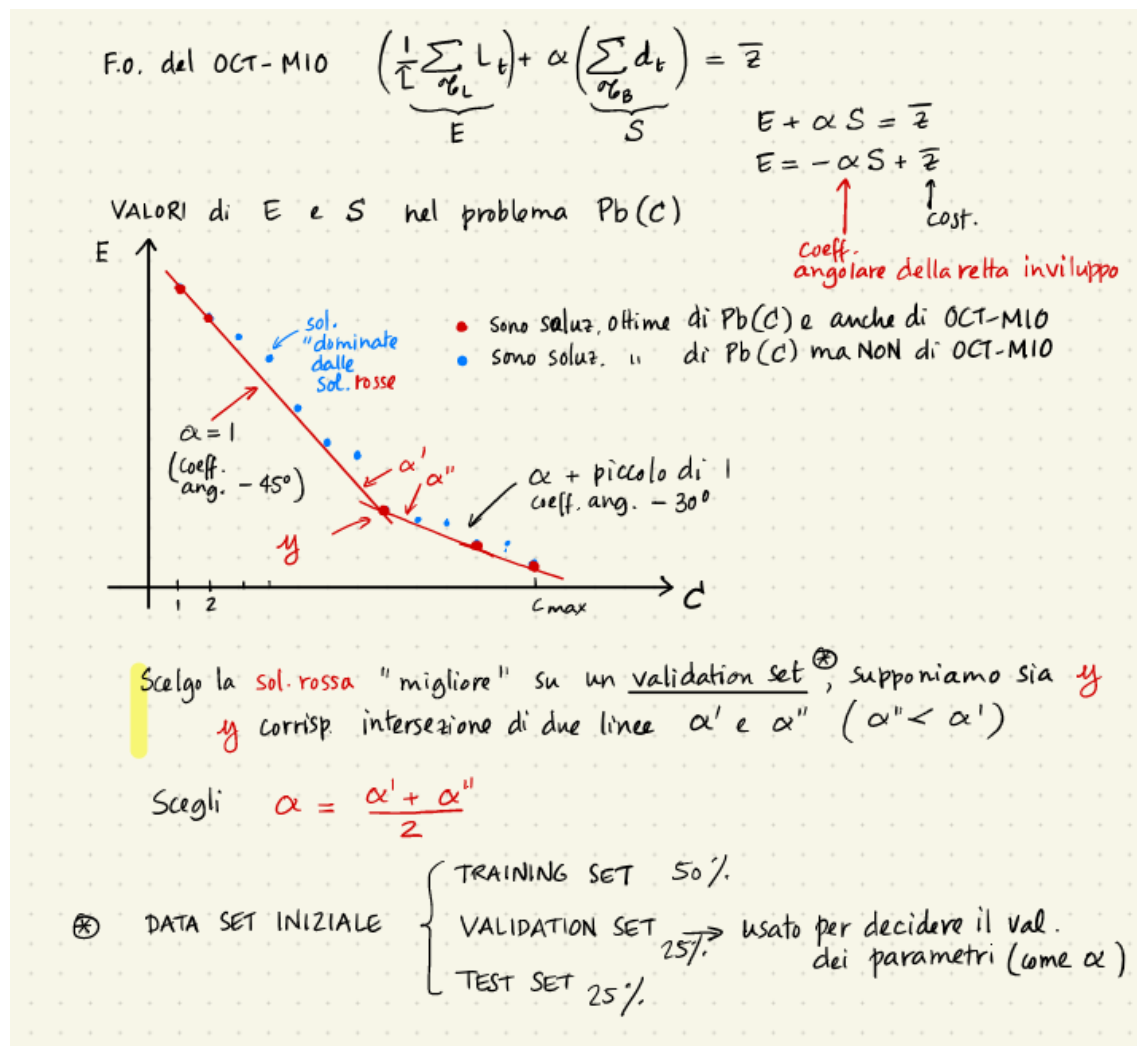
Nell'addestramento di un OCT occorre scegliere la **prodondità** dell'alber D_{\max} da cui si genereranno da $D = 1$ (3 nodi) fino a D_{\max} .

Possiamo riscrivere il problema di minimizzazione nel seguente modo. Posto $E :=$ errore di training
 $S :=$ complessità dell'albero:

$$\min E$$

$$S \leq C$$

con C è una costante che corrisponde al massimo numero di split $(2^D - 1)$.



Algoritmo 2

1. Scegli D_{\max} e N_{\min} ;
2. For $D=1, \dots, D_{\max}$ do:
 - For $C = 1, \dots, 2^D - 1$ do:
 - Run TDITD con $\alpha = 0$ e N_{\min} . Effettua prune a profondità D , max numero di split = C . Inserisci la soluzione nel pool di warm start;
 - Scegli il candidato più accurato sul validation set nel pool di warm start;
 - Risolvi $Pb(C)$ con profondità D e C split usando il warm start selezionato. Inserisci la soluzione nel pool di warm start;
3. Post-process: rimuovi le soluzioni non ottime per OCT-MIO per alcun valore di α . Determina quindi le soluzioni non dominate e rimuovi le altre.
4. Selezione le soluzioni migliori sul validation set e determina il range di α .

6 OCT: Modelli Multivariati

6.1 Introduzione

Fino ad ora abbiamo considerato alberi di decisione che utilizzano una singola variabile nei loro split ad ogni nodo. Questi alberi vengono detti **alberi decisionali univariati**. Ora cerchiamo di estendere questo insieme di alberi decisionali al caso **multivariato**, cioè nel caso in cui negli split di ogni nodo si prendono in considerazione più variabili. Questi alberi decisionali multivariati, vengono detti **OCT-H**.

6.2 Formulazione di OCT-H

Negli alberi decisionali multivariati, non siamo più vincolari a scegliere una singola variabile negli split, ma possiamo scegliere un generico iperpiano negli split di ogni nodo. Andiamo, ora, a riscrivere i vincoli del nostro problema...

Le variabili a_t sono usate per modellare lo split in ogni nodo ed in particolare $\in [-1, 1]^p$:

$$a_{tj} \in [-1, 1]$$

$$b_t \in [-1, 1]$$

Come nel caso univariato, le variabili d_t rappresentano la presenza o meno di uno split.

$$\sum_{j=1}^p |a_{tj}| \leq d_t$$

La problematica di questo vincolo è che è diventato non più lineare. A tale scopo linearizzo tramite delle variabili ausiliarie $\hat{a}_{jt} \geq 0 \quad \forall t \in T_B, j = 1, \dots, p$ ottenendo:

$$\hat{a}_{jt} \geq a_{jt}; \hat{a}_{jt} \geq -a_{jt} \quad j = 1, \dots, p \quad \forall t \in T_B$$

Quindi:

$$\sum_{j=1}^p \hat{a}_{jt} \leq d_t$$

Inoltre, abbiamo la variabile b_t che ci indica la presenza di punti dopo la diramazione:

$$-d_t \leq b_t \leq d_t \quad \forall t \in T_B$$

I **vincoli di consistenza**, mi garantiscono l'assegnamento del punto nel nodo e diventano:

$$\begin{aligned} a_m^T x_i &\geq b_m - M(1 - z_{it}) & \forall i = 1, \dots, n & \quad \forall t \in T_B & \quad \forall m \in A_L(t) \\ a_m^T x_i + \mu &\leq b_m + M(1 - z_{it}) & \forall i = 1, \dots, n & \quad \forall t \in T_B & \quad \forall m \in A_R(t) \end{aligned}$$

$\mu = 5 \cdot 10^{-3}$ quantità piccola e di norma si utilizza $M = 2 + \mu$;

$$\begin{aligned} M \leq \max(a_m^T x_i - b_m) &= 2 \\ a_m^T x_i + \mu &\leq b_m + (2 + \mu)(1 - z_{it}) & \forall i = 1, \dots, n & \quad \forall t \in T_L & \quad \forall m \in A_L(t) \end{aligned}$$

Varia anche la funzione obiettivo poiché la complessità dell'albero varia: vogliamo tenere in considerazione l'uso di test con più attributi.

Introduco $s_{jt} \in \{0, 1\} \quad \forall t \in T_B, j = 1, \dots, p$:

$$s_{jt} = \begin{cases} 1 & \text{se al nodo } t \in T_B \text{ utilizzo } j \text{ per la definizione dello split } t \\ 0 & \text{altrimenti} \end{cases}$$

$$s_{jt} \geq |a_{jt}|$$

Nota 128.

Nel caso in cui $a_{jt} \neq 0 \longrightarrow s_{jt} = 1$

Nel caso in cui $a_{jt} = 0 \longrightarrow s_{jt} = 0$

Linearizzando il vincolo non lineare:

$$-s_{jt} \leq a_{jt} \leq s_{jt} \quad \forall t \in T_B, j = 1, \dots, p$$

Al fine di ottenere una maggiore efficienza si aggiungono altri due vincoli per rendere s_{jt} compatibile con d_{jt} :

$$s_{jt} \leq d_t$$

$$\sum_{j=1}^p s_{jt} \geq d_t$$

Quindi la funzione obiettivo nel caso OCT-H diventa:

$$\min \frac{1}{\hat{L}} \sum_{t \in T_L} L_t + \alpha \sum_{t \in T_B} \sum_{j=1}^p s_{jt}$$

6.3 Warm Start OCT-H

Negli OCT-H come Warm start si utilizza un'euristica greedy appartenente alla TDIDT family. Per determinare il migliore split, invece di utilizzare gli impurity index, si utilizza, sui punti "sopravvissuti" al nodo $t \in T_B$ un OCT-H a profondità $D = 1$.

6.4 Tuning degli Iperparametri

Nella scelta degli iperparametri si utilizza la stessa procedura di OCT con la differenza che:

$$\text{OCT} \longrightarrow C_{\max} = 2^D - 1$$

$$\text{OCT} - H \longrightarrow C_{\max} = p(2^D - 1)$$

7 Alberi di Decisione Pro e Contro

La caratteristica degli alberi di decisione è che sono semplici da spiegare/giustificare ai non esperti. Rappresentano in modo più verosimile il processo decisionale umano rispetto ad altri metodi e modellano in maniera immediata variabili qualitative. Tuttavia, il livello di accuratezza è inferiore e sono poco robusti poiché piccole variazioni nei dati producono cambiamenti significativi nell'albero risultante,

8 Bagging, Random Forests, Boosting

Gli alberi decisionali DT soffrono di **varianza elevata**:

- Estrai in modo casuale due o più dataset da una popolazione;
- Ricava alberi di decisione sulla base dei nuovi dataset;
- alberi di decisione molto dissimili;

Per ridurre la varianza si può ricorrere al **bagging** (Bootstrap Aggregation). Questo è un metodo general-purpose utile a ridurre la variazione aumentando l'accuratezza della predizione tramite l'aggregazione di un insieme di osservazioni.

Idoneamente si ricava un certo numero B di training set della popolazione osservata e se ne considera la media come predizione del modello aggregato.

8.1 Bagging

Si applica l'idea della sezione precedente agli alberi di classificazione. In particolare si costruisce un albero di classificazione a partire da ciascuno dei training set bootstrap $1, \dots, B$.

A questo punto, si classifica sulla base del voto di maggioranza cioè dalla classe di predizione più frequente tra i risultati ottenuti dai B alberi.

In aggiunta, con questo tipo di tecnica non si effettua pruning degli alberi generati a partire dai training set bootstrap. La varianza viene ridotta prendendo il risultato a maggioranza sulle B classificazioni.

A questo punto occorre definire una grandezza per valutare l'errore che stiamo commettendo tramite bagging. Questa grandezza di errore viene detta **stima out-of-bag**.

Un DT nel bagging usa circa $\frac{2}{3}$ dei dati originali. Quindi $\frac{1}{3}$ delle osservazioni non vengono utilizzate per generare un singolo DT e proprio questi punti/osservazioni rappresentano **Out-Of-Bag**.

L'idea è quella di utilizzare i risultati prodotti dai DT che hanno x_i OOB per ottenere una predizione su x_i . Si ripete questa procedura per tutti gli n punti x_i $i = 1, \dots, n$ e si calcola l'errore di classificazione complessivo.

Ricapitolando il bagging aumenta l'accuratezza della predizione al costo di una minore interpretabilità e utilizza l'indice di Gini come misura dell'importanza di una variabile indipendente j :

Algoritmo 3

Per tutti gli alberi T_b dal training set di bootstrap $b = 1, \dots, B$:

Calcola la diminuzione $d_{b,j}$ dell'indice di Gini associata allo split sulla variabile j

Output = $\frac{1}{B} \sum_b d_{b,j}$ è la misura dell'importanza relavita della variabile j

8.2 Random Forests

I random forests sono un miglioramento del bagging basato sulla scelta di alberi meno correlati tra loro. In particolare, si costruiscono B alberi di decisione e ad ogni split:

- si considerano soltanto un **campione casuale** di $m < p$ variabili candidate;
- La scelta di p viene effettuata in cross validation, ma sperimentalmente si è visto che si pone $p \approx \sqrt{m}$

La conseguenza della riduzione delle alternative ad ogni split è che si riducono i rischi dovuti alla forte correlazione tra alberi e si aumenta la riduzione della varianza promuovendo la costruzione di alberi poco correlati tra loro.

In particola, ponendo $m = p \rightarrow$ bagging

Con un numero alto di attributi fortemente correlati, scegliere m piccolo rispetto a p può migliorare le performance di queste tecniche rispetto al bagging.

8.3 Boosting

Le tecniche di **boosting** possono essere applicabili a metodi diversi di apprendimento statistico di regressione e classificazione; ogni albero viene creato a partire dalle informazioni sugli alberi creati in precedenza e si combina un numero sufficientemente elevato di B alberi di decisione a partire da B training set.

In questo caso, i training set non sono bootstrap, ma sono creati modificando in maniera opportuna i dati originali.

I campioni successivi sono determinati pesando i dati dei campioni precedenti.

Il peso di una osservazione x_i al passo b è maggiore se il DT non ha classificato correttamente x_i al passo $(b-1)$. La classificazione viene effettuata pesando opportunamente i risultati dei vari DT.

8.4 Adaboost ST

Assumo che le label $\in \{-1, 1\}$ ed ho M weak learners k_j a stamp:

$$C(x_i) = \sum_{j=1}^M \alpha_j k_j(x_i) \quad x_i \in \text{TS} \quad x_i \in \mathbb{R}^p$$

La risposta:

$$y = \frac{C(x_i)}{|C(x_i)|}$$

Il suo segno determina la risposta del sistema e il denominatore è proporzionale alla confidenza, cioè maggiore è il suo valore, maggiore sarà la confidenza della risposta per il punto del TS.

Inoltre:

$$C_m(x_i) = \sum_{j=1}^M \alpha_j k_j(x_i) \longrightarrow C(x_i) = C_M(x_i); C_m(x_i) = C_{m-1}(x_i) + \alpha_m k_m(x_i)$$

A questo punto si necessita di una funzione per determinare l'errore di C_m : **exponential loss**.

8.5 Errore di classificazione al passo m-esimo

$$E = \sum_{i=1}^n e^{y_i C_m(x_i)} \equiv \sum_{i=1}^m e^{-y_i C_{m-1}(x_i)} e^{\alpha_i k_m(x_i) y_i} = \sum_{i=1}^m W_i^{(m)} e^{\alpha_i k_m(x_i) y_i}$$

8.6 Scelta del Weak Learner al passo m

Il mio obiettivo è quello di avere E minimo:

$$\begin{aligned} E &= \sum_{i: K_m(x_i) \neq y_i}^m W_i^{(m)} e^{-\alpha_m k_m(x_i) y_i} + \sum_{i: K_m(x_i) = y_i}^m W_i^{(m)} e^{\alpha_m k_m(x_i) y_i} = \\ &= \sum_{I^=} W_i^{(m)} e^{-\alpha_m} + \sum_{I^{\neq}} W_i^{(m)} (e^{\alpha_m} - e^{-\alpha_m}) \end{aligned}$$

In cui il primo termine non dipenda da $k_m(\cdot)$, mentre il secondo non dipende né da i né da $k_m(\cdot)$. Quindi, per minimizzare viene scelta maniera tale che:

$$\sum_{I^{\neq}} W_i^{(m)}$$

sia minimo.

8.7 Scelta del parametro α_m

Voglio minimizzare E come forma di α_m . Quindi per trovare il minimo derivo $\frac{\partial E}{\partial \alpha_m} = 0$:

$$\begin{aligned} E &= \sum_{i \in I^{\neq}}^m W_i^{(m)} e^{-\alpha_m k_m(x_i) y_i} + \sum_{i \in I^=}^m W_i^{(m)} e^{\alpha_m k_m(x_i) y_i} = \sum_{i \in I^{\neq}}^m W_i^{(m)} e^{-\alpha_m} + \sum_{i \in I^=}^m W_i^{(m)} e^{\alpha_m} = \\ &= e^{-\alpha_m} W^E - e^{\alpha_m} W^C \end{aligned}$$

Quindi:

$$\frac{\partial E}{\partial \alpha_m} = 0 \iff e^{\alpha_m} W^E = e^{-\alpha_m} W^C$$

$$\alpha_m + \ln W^E = -\alpha_m + \ln W^C$$

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W^C}{W^E} \right) = \frac{1}{2} \ln \left(\frac{1 - \varepsilon}{\varepsilon} \right)$$

Rappresenta un errore complessivo pesato che può essere scritto come:

$$\varepsilon = \frac{W^E}{W^E + W^C}$$

2 Classificazione Robusta

1 Introduzione

Nel caso in cui ci fosse incertezza sui dati, si ha un peggioramento delle prestazioni del classificatore. A tale scopo vogliamo delle procedure che rendano il classificatore più robusto rispetto a delle piccole incertezze. L'approccio standard è la **regolarizzazione** che si basa sulle funzioni di penalità per diminuire l'overfitting. UN altro approccio è quello di utilizzare l'**ottimizzazione robusta**.

2 Robust Classification

Il problema di ottimizzazione è del tipo:

$$\max c(x, y)$$

$$g(x, u) \leq 0^m$$

$$x \in X$$

In cui:

- x := variabili di decisione;
- u = parametri;
- $g(x, u) \leq 0^m$:= vincoli di disuguaglianza.

I **parametri** u possono essere:

- fissati → problema di ottimizzazione incerti;
- incerti → $u \in U$ insieme di incertezze

.Cerchiamo la soluzione ottima nel caso peggiore.

3 Approccio MaxMin

$$\max_{x \in X} \min_{u \in U} \{c(x, u) : g(x, u) \leq 0^m\}$$

Se $|U| = +\infty$ allora si hanno infiniti vincoli che possono essere ridotti ad un numero finito.

Definizione 129. *Problema Norma Duale*

La norma duale è un numero reale associato ad una funzione lineare definita sullo spazio vettoriale X .

Se $X = \mathbb{R}^n$: una funzione lineare genera $f = a^T x$

$$z_q = \sup \{ |f(x)| : \|x\|_q < 1 \}$$

$$z_q = \max \{ a^T x : \|x\|_q \leq 1, x \in \mathbb{R}^n \}$$

Si può dimostrare che la soluzione di questo problema è:

$$z_q = \|a\|_{q^*} \quad \text{con } q^* = \frac{q}{q-1} \quad \text{se } q \neq 1$$

$$z_q = \|a\|_\infty \quad \text{Se } q = 1$$

Si può estendere questo problema nella norma di $f(x)$ ristretta da un qualsiasi numero $\rho > 0$. Ottenendo:

$$\begin{aligned} z_q &= \max \{ a^T x : \|x\|_q \leq \rho, x \in \mathbb{R}^n \} = \\ &= \max \{ a^T \rho y : \|y\|_q \leq 1, y \in \mathbb{R}^n \} = \\ &= \rho \|a\|_{q^*} \end{aligned}$$

4 Uncertainty Set

L'**insieme di incertezza** può essere utilizzato per modellare l'**incertezza delle feature**. Poiché l'incertezza nelle feature è dovuto a mancanza di dati, errori di manipolazioni e errori di misura vogliamo costruire un modello di ottimizzazione per il classificatore che sia una **controparte robusta** del modello OCT.

A tale scopo, supponiamo che:

$$x_i \in \mathbb{R}^p \text{ è il valore esatto}$$

$$(x_i + \Delta x_i) \quad \forall i = 1, \dots, n \quad \text{Training data effettivo}$$

con $\Delta x_i :=$ perturbazione dell' i -esimo dato.

A questo punto, possiamo definire un insieme delle perturbazioni:

$$\Delta X = \{ \Delta x_1, \dots, \Delta x_n \} \in \mathbb{R}^{p \times n}$$

Ottenendo:

$$U_x = \{ \Delta x \in \mathbb{R}^{p \times n} : \|\Delta X_i\|_q \leq \rho, i = 1, \dots, n \}$$

Il termine ρ viene detto **parametro di magnitudo(incertezza)** e viene scelto tramite cross validation.

5 Modello MIP (Programmazione Intera Mista)

Se consideriamo un albero binario non completo, esso sarà formato da $\left\lceil \frac{\#D}{2} \right\rceil$. Fissata la struttura, il numero di nodi T è dispari ed abbiamo:

$$\# \text{Foglie} = \left\lceil \frac{T}{2} \right\rceil$$

I nodi sono numerati in modo che $k = 1, 2, \dots, \left\lfloor \frac{T}{2} \right\rfloor, \left\lceil \frac{T}{2} \right\rceil, \dots, T$.

Al posto delle variabili $d_t = \begin{cases} 1 & \text{split nel nodo } t \\ 0 & \text{altrimenti} \end{cases}$, si utilizza:

$$\delta_k = \begin{cases} 1 & \text{se non si ha lo split al nodo } k \\ 0 & \text{altrimenti} \end{cases}$$

5.1 Vincoli Strutturali

1. $\delta_k = 1 \quad k = \left\lceil \frac{T}{2} \right\rceil, \dots, T;$
2. $\delta_k \geq \delta_n \quad k = 1, \dots, T; u \in A(k);$
3. $\delta_k + \sum_{i=1}^P a_{k,i} = 1 \quad k = 1, 2, \dots, T$

5.2 Allocazione Punti nelle Foglie

$$z_{i,k} \in \{0, 1\} \quad i = 1, \dots, n \quad k = 1, \dots, T$$

Importante 130.

$$z_{ik} = 1 \iff \text{assegno punto } i \text{ al nodo } k$$

Vincoli

- $\sum_{k=1}^T z_{ik} = 1 \quad i = 1, \dots, n$
- $z_{ik} \leq \delta_k \quad i = 1, \dots, n \quad k = 1, \dots, T$
- $\sum_{i=1}^n z_{ik} \geq N_{\min} l_k \quad l_k \in \{0, 1\} \quad k = 1, \dots, T \quad l_k = 1 \text{ se la foglia popolata da punti}$
- $z_{ik} + \delta_u \leq 1 \quad i = 1, \dots, n; k = 1, \dots, T; u \in A(k)$
- $l_k + \sum_{u \in A(k)} \delta_u \geq \delta_k \quad k = 1, \dots, T.$

5.3 Funzione Obiettivo

Assumiamo 2 classi $y_i \in \{-1; +1\} \quad \forall i = 1, \dots, n$

Variabili

$g_k, h_k \in \mathbb{Z}_{\geq 0} (\mathbb{R}_{\geq 0})$ contano il numero di punti assegnati a k che assumo ± 1

$$(\text{numero di } -1 \text{ in } k) \quad g_k = \frac{1}{2} \sum_{i=1}^n z_{ik}(1 - y_i) \quad \forall k \in T$$

$$(\text{numero di } 1 \text{ in } k) \quad h_k = \frac{1}{2} \sum_{i=1}^n z_{ik}(1 + y_i)$$

Vincoli:

$$(\alpha) f_k \leq g_k + M[w_k + (1 - l_k)] \quad (\gamma) f_k \geq -M[1 - w_k + (1 - l_k)]$$

$$(\beta) f_k \leq h_k + M[1 - w_k + (1 - l_k)] \quad (\varepsilon) f_k \geq h_k - M[w_k + (1 - l_k)]$$

Se $l_k = 0 \longrightarrow$ vincoli soddisfatti

Se $l_k = 1$:

$$w_k \in \{0, 1\} \quad k = 1, \dots, T$$

Se $w_k = 0$ $\beta e \gamma$ sono soddisfatti banalmente, ma i vincoli attivi α, ε impongono:

$$h_k \leq f_k \leq g_k$$

Se $w_k = 1$ α, ε vengono soddisfatti e rimangono attivi gli altri due e impongono:

$$g_k \leq f_k \leq h_k$$

Quindi l'insieme f_k contiene tutti i punti malclassificati e quindi può rappresentare la funzione da minimizzare:

$$\min \left[\sum_{k=1}^T f_k + \alpha \sum_{k=1}^T (1 - \delta_k) \right]$$

Coerenza dei Test

$$a_u^T x_i \leq b_u + M(1 - z_{ik}) \quad u \in A_l(k) \quad \forall i = 1, \dots, l$$

$$a_u^T x_i \geq b_u - M(1 - z_{ik}) \quad u \in A_R(k) \quad \forall i = 1, \dots, l$$

6 Robustezza vs. Incertezza delle feature

$$U_x = \{ \Delta x \in \mathbb{R}^{p \times n} : \|\Delta x_i\|_q \leq \rho, i = 1, \dots, n \}$$

Nei vincoli 5.3 della coerenza dei testi, **controparte robusta**:

$$a_u^T(x_i + \Delta x_i) + \varepsilon \leq b_u + M(1 - z_{ik}) \quad \forall i, k, u \in A_L(k) \quad \forall \Delta x \in U_x$$

Possiamo riscriverlo come:

$$a_u^T \Delta x_i \leq -a_u^T x_i - \varepsilon + b_u + M(1 - z_{ik})$$

Quindi:

$$\max_{\Delta x \in U_x} \{a_u^T \Delta x_i\} \leq -a_u^T x_i - \varepsilon + b_u + M(1 - z_{ik})$$

$$\Updownarrow$$

$$\max_{\|\Delta x = \rho\|} \{a_u^T \Delta x : \|\Delta x_i\| \leq \rho\} = \rho \|a_u\|_q = \rho$$

La versione robustificata diventa:

$$a_u^T x_i + \varepsilon + \rho \leq b_u + M(1 - z_{ik})$$

7 Robustezza vs. Incertezza nelle label

7.1 Dualità PL

Il problema primale di minimizzazione in un contesto di classificazione è il seguente:

$$z = \min c^T x$$

$$A x \leq b$$

$$x \in \mathbb{R}_{\geq 0}^n$$

Il suo duale:

$$w = \max u^T b$$

$$u^T A \leq c^T$$

$$u \in \mathbb{R}_{\leq 0}^m$$

Se $\bar{x} \in \mathbb{R}_{\geq 0}^n$ tale che $A\bar{x} \leq b$, cioè ammissibile per il primale e $\bar{u} \in \mathbb{R}_{\leq 0}^m$ tale che $\bar{u}^T A \leq c^T$, ammissibile per il problema duale allora:

$$c^T \bar{x} \geq z = w \geq \bar{u}^T b$$

7.2 Robustezza nelle label

Se vi è incertezza nelle label, esse sono affette da errore e quindi occorre robustificare il modello e definire l'insieme di incertezza. Sostanzialmente esse sono dei **flip**: la label da 1 diventa -1 e viceversa. Quindi definiamo:

$$\Delta y_i = \begin{cases} 0 & \text{Se non vi è flip} \\ 1 & \text{altrimenti} \end{cases}$$

Posso definire l'insieme di incertezza come:

$$U_y = \left\{ \Delta y_i \in \{0, 1\}^n : \sum_{i=1}^n \Delta y_i \leq \Gamma \right\} \quad \Gamma \in \mathbb{Z}_+$$

La funzione obiettivo, tenendo conto dell'incertezza sulle label diventa:

$$g_k = \frac{1}{2} \sum_{k=1}^n z_{ik} (1 - y_i (1 - 2 \Delta y_i))$$

$$h_k = \frac{1}{2} \sum_{i=1}^n z_{ik} (1 + y_i (1 - 2 \Delta y_i))$$

I vincoli:

$$(\alpha) f_k \leq g_k + M[w_k + (1 - l_k)] \quad (\gamma) f_k \geq -M[1 - w_k + (1 - l_k)]$$

$$(\beta) f_k \leq h_k + M[1 - w_k + (1 - l_k)] \quad (\varepsilon) f_k \geq h_k - M[w_k + (1 - l_k)]$$

Considerando $\alpha + \beta + \gamma + \varepsilon \longrightarrow$ Se $l_k = \begin{cases} h_k \leq f_k \leq g_k & \text{se } w_k = 0 \\ g_k \leq f_k \leq h_k & \text{se } w_k = 1 \end{cases}$

Prendiamo in considerazione solo il primo vincolo, le restanti si svolgono allo stesso modo $\forall k = 1, \dots, T, \Delta y_i \in U_y$

$$f_k \leq \frac{1}{2} \sum_{i=1}^n z_{ik} (1 - 2 \Delta y_i) + M[w_k + (1 - l_k)]$$

$$\frac{1}{2} \sum_{i=1}^n (1 - y_i) z_{ik} + \sum_{i=1}^n z_{ik} y_i \Delta y_i \geq f_k - M[w_k + (1 - l_k)]$$

$$\frac{1}{2} \sum_{i=1}^n (1 - y_i) z_{ik} + \min_{\Delta y \in U} \{z_{ik} y_i \Delta y_i\} \geq f_k - M[w_k + (1 - l_k)]$$

La funzione di minimo è un problema di programmazione lineare PL, che può essere rilassato:

$$\xi_k = \min_{\Delta y \in U} \{z_{ik} y_i \Delta y_i\}$$

$$\sum_{i=1}^n \Delta y_i \leq \Gamma$$

$$\Delta y_i \in \{0, 1\} \quad i = 1, \dots, n$$

Che rilassando il vincolo binario:

$$\xi_k = \min_{\Delta y \in U} \{z_{ik} y_i \Delta y_i\}$$

$$\sum_{i=1}^n \Delta y_i \leq \Gamma$$

$$0 \leq \Delta y_i \leq 1 \quad i = 1, \dots, n$$

Quindi, risolto questo problema, otteniamo ξ_k :

$$\frac{1}{2} \sum_{i=1}^n (1 - y_i) z_{ik} + \xi_k \geq f_k - M[w_k + (1 - l_k)]$$

7.3 Proprietà

$$\xi_k^R \leq \xi_k \quad \forall k = 1, \dots, T$$

Allora posso considerare il duale del rilassamento:

$$\begin{aligned} \xi_k &= \min_{\Delta y \in U} \{z_{ik} y_i \Delta y_i\} & \tau_k &= \max \left(\Gamma \nu_k + \sum_{i=1}^n \mu_{ik} \right) \\ \sum_{i=1}^n \Delta y_i &\leq \Gamma & \iff & \mu_{ik} + \nu_k \leq z_{ik} y_i \\ 0 \leq \Delta y_i &\leq 1 \quad i = 1, \dots, n & \mu_{ik} &\leq 0; \nu_k \leq 0 \end{aligned}$$

Se sono ammissibili per il duale allora:

$$\left(\Gamma \bar{\nu}_k + \sum_{i=1}^n \bar{\mu}_{ik} \right) \leq \tau_k \leq \xi_k^R \leq \xi_k$$

Allora 7.2 diventa:

$$\frac{1}{2} \sum_{i=1}^n (1 - y_i) z_{ik} + \left(\Gamma \bar{\nu}_k + \sum_{i=1}^n \bar{\mu}_{ik} \right) \geq f_k - M[w_k + (1 - l_k)]$$

$$\mu_{ik} + \nu_k \leq z_{ik} y_i \quad \forall i = 1, \dots, n$$

$$\mu_{ik} \leq 0 \quad \forall k = 1, \dots, T$$

$$\nu_k \leq 0$$

Avendo così, in conclusione, $n \cdot T$ vincoli invece di T .