



MACHINE LEARNING

- Support Vector Machine -

Notazione

| | | |
|---|---|--------------------------------------|
| $x_j, j = 1, \dots, d$ | → | feature |
| $\mathbf{x} = [x_1, \dots, x_d]$ | → | feature vector |
| y | → | variabile target |
| $T = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$ | → | dataset di training |
| $[w_0, w_1, \dots, w_d] = [w_0, \mathbf{w}]$ | → | vettore dei parametri |
| $[v_0, v_1, \dots, v_d] = [v_0, \mathbf{v}]$ | → | vettore dei parametri |
| \mathbf{w}^T | → | vettore \mathbf{w} trasposto |
| $(\mathbf{x}^{(i)})^T$ | → | vettore $\mathbf{x}^{(i)}$ trasposto |
| $X = \mathbb{R}^d$ | → | feature space |
| $W = \mathbb{R}^{d+1}$ | → | parameter space |

Oggi parleremo di **Support Vector Machine (SVM)**, che sono dei metodi discriminativi parametrici molto robusti rispetto al problema dell'overfitting (interpretabili anche come non-parametrici)

Le SVM sono dei classificatori ma possono essere estese per trattare task di regressione

La loro caratteristica è una definizione basata sugli aspetti geometrici dei loro bordi di decisione

Prima di trattare le SVM, però, è necessario approfondire l'interpretazione geometrica dei classificatori parametrici e dei loro decision boundary

Il classificatore (decision rule inclusa) della logistic regression può essere espresso come:

$$C(\mathbf{x}) = \text{If } g(\mathbf{w}^T \mathbf{x} + w_0) > 0.5 \text{ then } y_1 \\ \text{else } y_2$$

dove $[w_0, \mathbf{w}]$ sono i parametri, $g()$ è la logistic function e

$$f_{[w_0, \mathbf{w}]}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

la chiameremo *score function* (in analogia con le score function della Softmax regression)

Dato \mathbf{x} , il suo «score» $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ è, appunto, un «punteggio» che $f()$ associa ad \mathbf{x} (più alto è il punteggio, più «probabile» è l'appartenenza alla classe y_1 secondo $f()$). Lo score viene mappato in un valore di probabilità usando $g()$

La parte parametrica, ovvero la score function $f_{[w_0, \mathbf{w}]}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$, è lineare

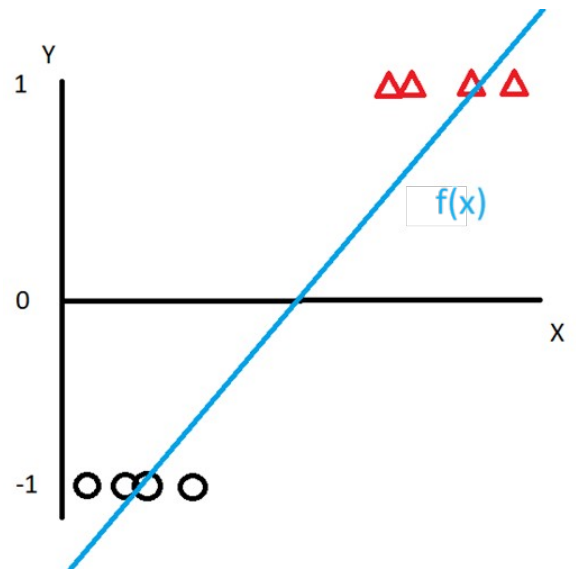
Esaminiamone le proprietà geometriche. Per semplicità, faremo a meno della funzione logistica $g()$ e ci concentreremo solo sulla parte parametrica lineare

Consideriamo un task di classificazione binaria e tralasciamo il requisito di voler ottenere una stima di probabilità in output

Non siamo quindi più costretti a rappresentare la ground truth (y) con 0 o 1, perciò useremo etichette positive e negative per la variabile target: $y \in \{-1, 1\}$ (cosa che semplifica la trattazione successiva)

Un *classificatore lineare* è basato su una decision rule e una *score function lineare* e può essere espresso come:

$$C(\mathbf{x}) = \text{If } \mathbf{w}^T \mathbf{x} + w_0 > 0 \text{ then } y_1 \\ \text{else } y_2$$

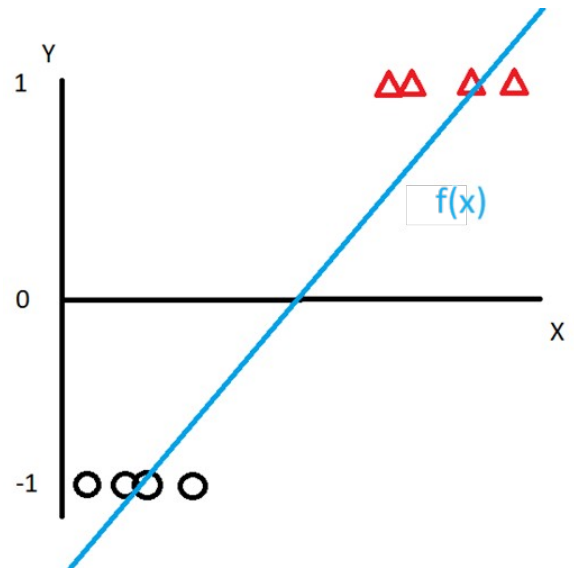


La parte parametrica di $C()$, ovvero:

$$f_{[w_0, w]}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

(senza decision rule) è lineare, ovvero è un polinomio di primo grado e, quindi, geometricamente, corrisponde ad una funzione di predizione rappresentata da un (iper-)piano (lo dimostreremo formalmente tra un po'), esattamente come nel caso della linear regression: si tratta della stessa classe di funzioni!

Ma che forma ha il suo decision boundary?

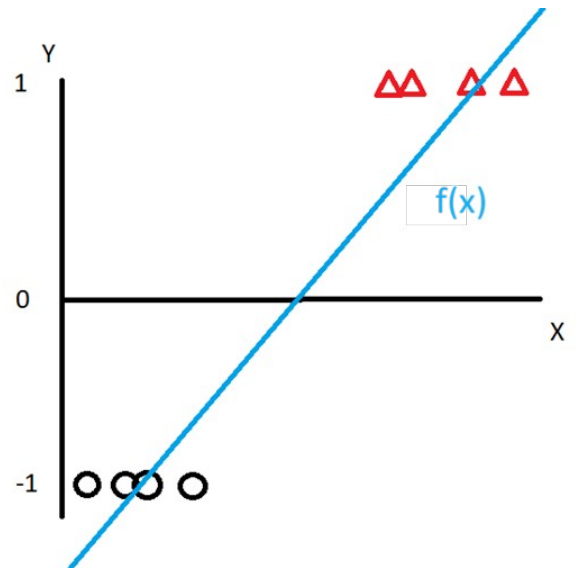


Il decision boundary (H) corrisponde ai punti in cui la score function assume valore 0, perché questi sono i punti del feature space in cui $C()$ passa dal predire una classe al predire l'altra (sono i bordi delle decision region) :

$$H = \{\mathbf{x} \mid f_{[w_0, \mathbf{w}]}(\mathbf{x}) = 0\}$$

Se il feature space è composto da una sola feature, come nel grafico accanto, H è un punto

Nel caso generale di d feature, H è una (iper-)superficie in R^d



Classificatori Lineari: iperpiano di decisione

$$H = \{\mathbf{x} \mid f_{[w_0, \mathbf{w}]}(\mathbf{x}) = 0\}$$

Per la precisione, H è un iperpiano. Dimostrazione:

Se $\mathbf{x}_1, \mathbf{x}_2 \in H$, allora:

$$\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0 = 0$$

quindi:

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

Ricordandoci che $\cos \theta = \mathbf{x}^T \mathbf{y} / \|\mathbf{x}\| \|\mathbf{y}\|$, ne deduciamo che l'angolo tra \mathbf{w} e il vettore $(\mathbf{x}_1 - \mathbf{x}_2)$ è 90°

Ne segue che \mathbf{w} è ortogonale *ad ogni* vettore $(\mathbf{x}_1 - \mathbf{x}_2)$ giacente su H , quindi H è una superficie piana

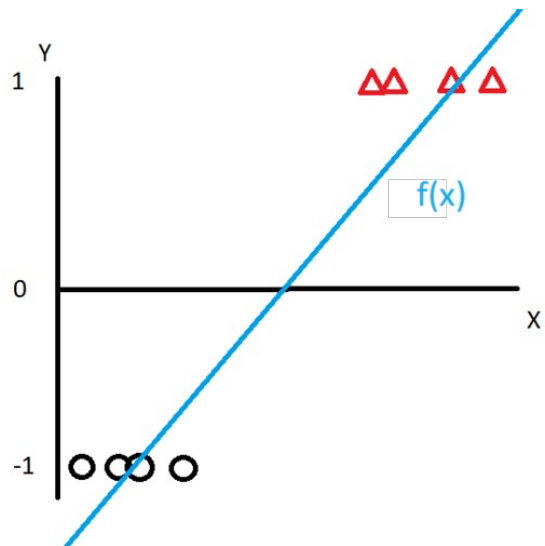
Classificatori Lineari: iperpiani corrispondenti

Un ragionamento analogo si può applicare all'equazione $f_{[w_0, w]}(\mathbf{x}) = y$, dimostrando che la superficie definita da $F = \{(\mathbf{x}, y) \in R^{d+1} \mid f_{[w_0, w]}(\mathbf{x}) = y\} = \{(\mathbf{x}, y) \in R^{d+1} \mid f_{[w_0, w]}(\mathbf{x}) - y = 0\}$ è un'iper-piano in R^{d+1} , dove R^{d+1} è lo spazio comprendente sia le feature (\mathbf{x}) che la variabile target (y)

Quando $f(\mathbf{x})$ interseca lo spazio delle feature alle coordinate $y = 0$, si ottiene H , che è anch'esso un iper-piano (in R^d)

Esempi:

- Se $d=1$, $f(\mathbf{x})$ è una retta in R^2 e H è un punto in R
- Se $d=2$, $f(\mathbf{x})$ è un piano in R^3 e H è una retta in R^2
- Se $d=3$, $f(\mathbf{x})$ è un iper-piano in R^4 e H è un piano in R^3

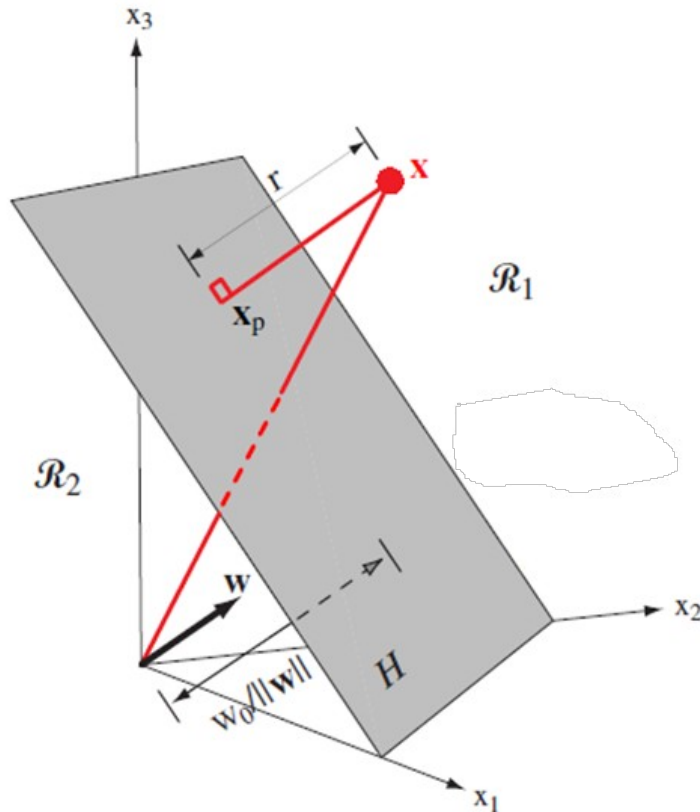


Classificatori Lineari: distanza dall'iperpiano di decisione

Si può dimostrare inoltre che la distanza r di un generico punto \mathbf{x} da H è data da:

$$r = |f(\mathbf{x})| / \|\mathbf{w}\| = |\mathbf{w}^T \mathbf{x} + w_0| / \|\mathbf{w}\|$$

dove $\|\mathbf{w}\|$ è la norma L_2 di \mathbf{w}



Classi linearmente separabili

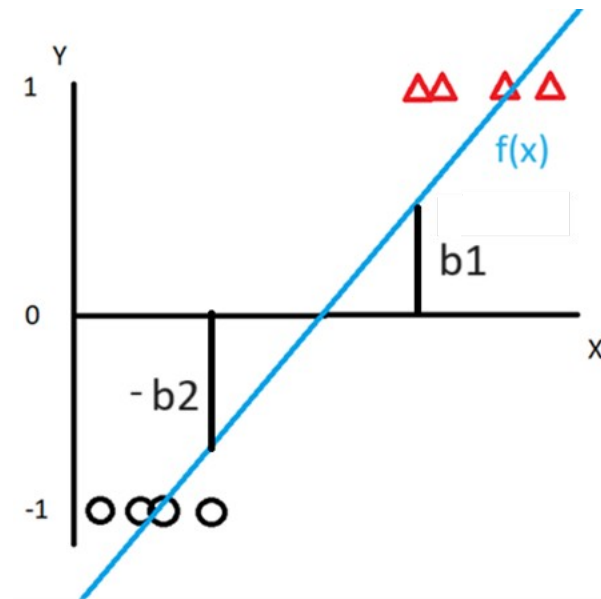
Due classi di punti si dicono «linearmente separabili» quando esiste un iper-piano H che li separi o, equivalentemente, quando esiste un classificatore lineare che separi completamente le due classi

Matematicamente, ciò è dato dall'esistenza di una score function $f_{[w_0, w]}(\mathbf{x})$ e di due valori reali positivi b_1 e b_2 tali che:

Per ogni $(\mathbf{x}^{(i)}, y^{(i)}) \in T$:

- $f_{[w_0, w]}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}^{(i)} + w_0 \geq b_1$ iff $y^{(i)} = +1$
- $f_{[w_0, w]}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}^{(i)} + w_0 \leq -b_2$ iff $y^{(i)} = -1$

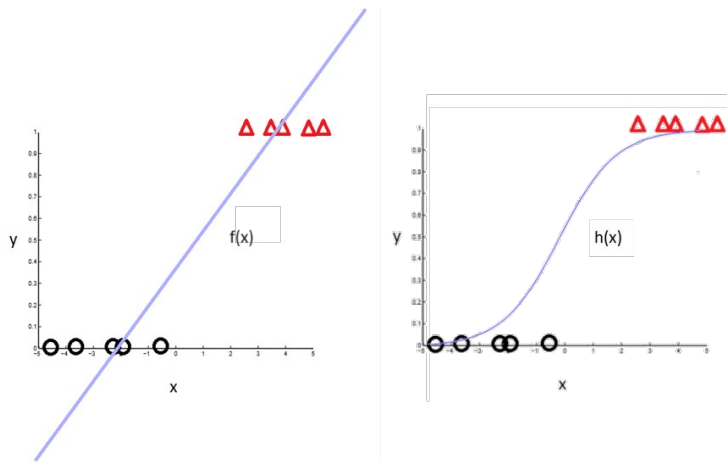
Se il bordo di decisione di $f_{[w_0, w]}(\mathbf{x})$ è a metà strada tra i punti delle due classi ad esso più vicini, allora $b_1 = b_2 = b$



Decision boundary della Logistic Regression

Nel caso della logistic regression la score function $f()$ viene combinata con la logistic function $g()$, e la funzione ipotesi finale $h(\mathbf{x}) = g(f(\mathbf{x}))$ non è più lineare (né lo è la superficie corrispondente a $\{(\mathbf{x}, y) \in \mathbb{R}^{d+1} \mid h(\mathbf{x}) = y\}$)

Tuttavia, è possibile dimostrare che il decision boundary della logistic regression continua ad essere un iper-piano



Classificatori Lineari: ricerca dei parametri

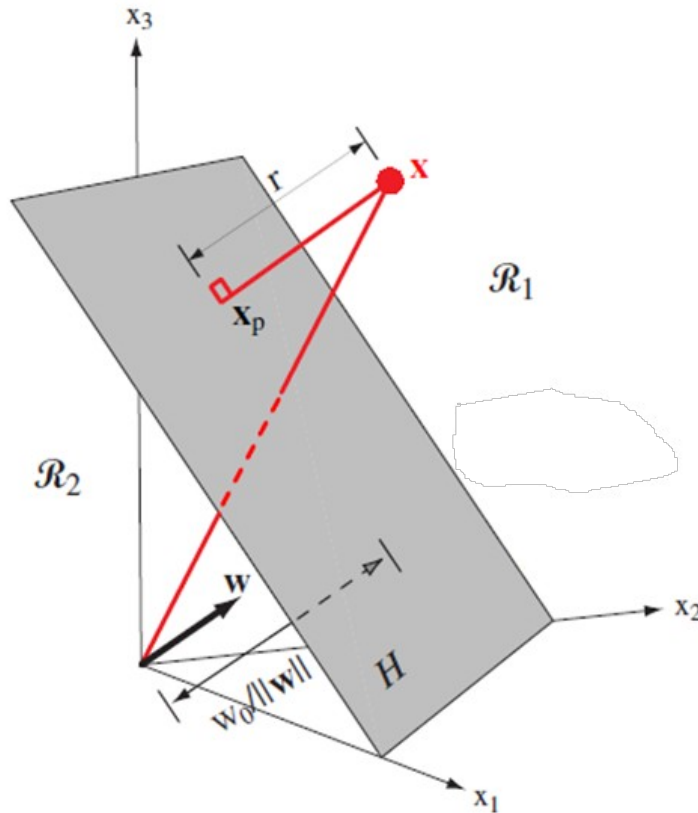
Come trovare $[w_0, \mathbf{w}]$?

Abbiamo visto che, nel caso della logistic regression, $g()$ ha la funzione di convertire gli score in valori di probabilità

Ciò permette di definire una (log-)likelihood function, massimizzando la quale troviamo $[w_0, \mathbf{w}]$

Nel caso delle SVM, invece, la ricerca del «miglior» valore per $[w_0, \mathbf{w}]$ è effettuata diversamente e si basa sulla distanza r dei punti del training set rispetto ad un generico H , dove H è il decision boundary della funzione ipotesi (dettagli tra un po'...)

Ora possiamo introdurre le SVM

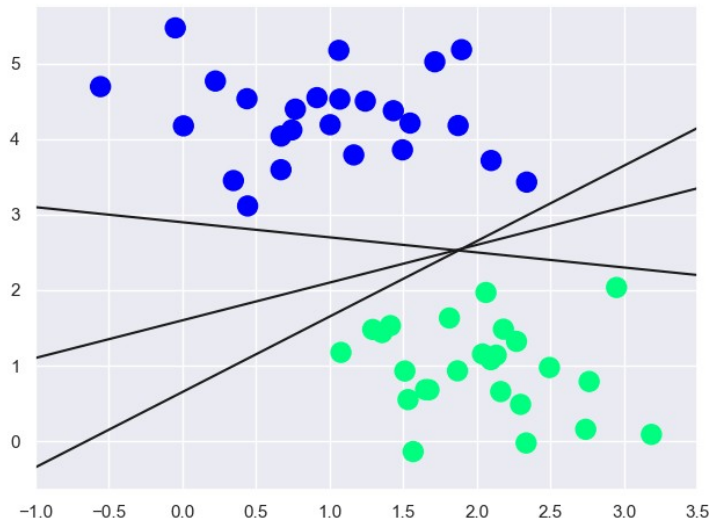


Support Vector Machine

Le **Support Vector Machine (SVM)** sono dei classificatori *binari lineari* che cercano l'iperpiano separatore tra due classi con distanza massima tra gli esempi più vicini delle due classi

Intuitivamente, l'iperpiano con distanza massima dagli esempi di training delle due classi è quello che probabilmente generalizza meglio a inference time

Le SVM possono essere estese per problemi non binari e trattare classi non linearmente separabili (e anche task di regressione)



SVM: trovo i parametri impostando un problema geometrico

Nel caso delle SVM la ricerca di $[w_0, \mathbf{w}]$ viene impostata come un problema geometrico

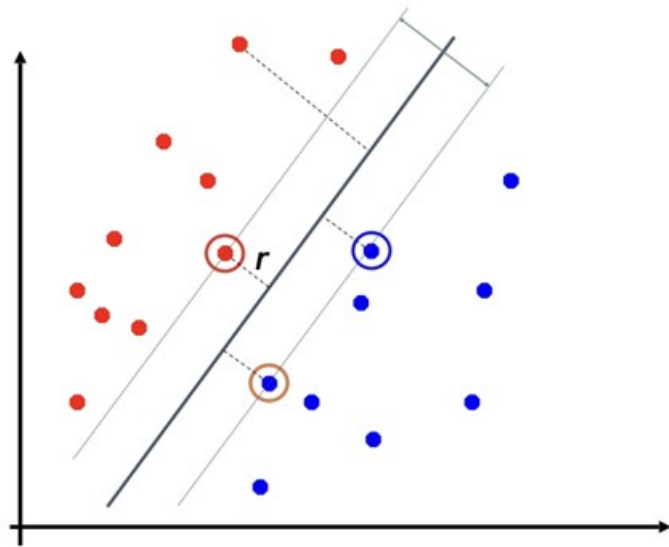
L'obiettivo è trovare l'iperpiano H che meglio separi le due classi

Il «miglior» H possibile è definito come quello che massimizza il «margine», ovvero la distanza dagli esempi \mathbf{x}_s ($s = 1, \dots, k$) più vicini ad H

Tali esempi «più vicini ad H » sono chiamati *vettori di supporto* (cerchiati nella figura a lato)

N.B.: Di solito $k < n$ ($n = |T|$)

Nelle SVM assumo che H sia situato a metà strada tra i support vector delle due classi



SVM: problema geometrico

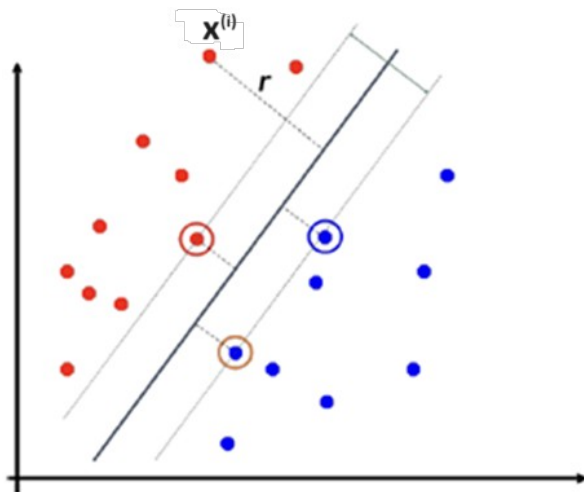
Dato che le SVM non modellano delle distribuzioni di probabilità, posso usare valori diversi da 0 e 1 per la ground truth di y . Assumeremo che $y \in \{-1, 1\}$

Se H esiste, ovvero le due classi sono *linaramente separabili*, allora posso assumere che esista un classificatore lineare $C(\mathbf{x})$ basato su una score function che assegni valori positivi/negativi ad ogni $(\mathbf{x}^{(i)}, y^{(i)}) \in T$ a seconda del valore di $y^{(i)}$:

Per ogni $(\mathbf{x}^{(i)}, y^{(i)}) \in T$:

- $\mathbf{v}^T \mathbf{x}^{(i)} + v_0 \geq b$ iff $y^{(i)} = +1$
- $\mathbf{v}^T \mathbf{x}^{(i)} + v_0 \leq -b$ iff $y^{(i)} = -1$

Dove b è un valore positivo qualsiasi e $[v_0, \mathbf{v}]$ sono i parametri che devo cercare



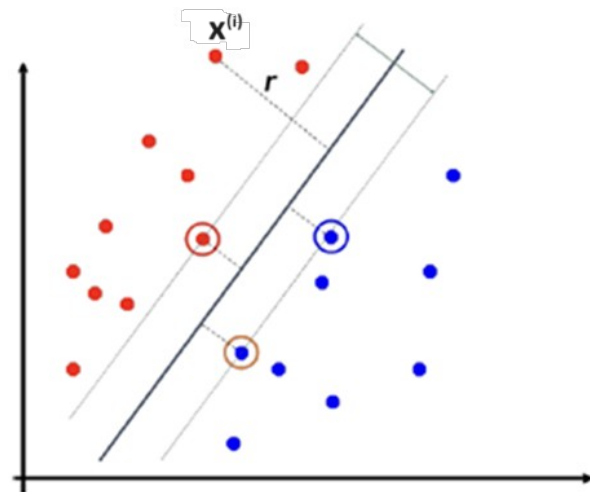
SVM: problema geometrico

Tuttavia, se uso: $\mathbf{w}^T = \mathbf{v}^T / b$ e $w_0 = v_0 / b$ posso semplificare i calcoli ed assumere:

- $\mathbf{w}^T \mathbf{x}^{(i)} + w_0 \geq 1$ iff $y^{(i)} = +1$
- $\mathbf{w}^T \mathbf{x}^{(i)} + w_0 \leq -1$ iff $y^{(i)} = -1$

Ovvero, per ogni $(\mathbf{x}^{(i)}, y^{(i)}) \in T$:

$$y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1$$

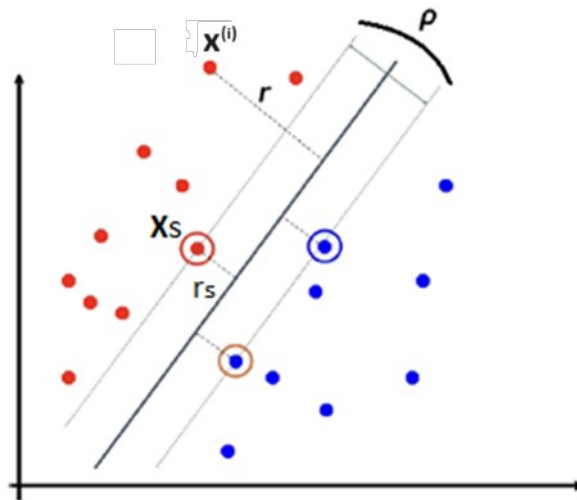


SVM: margine

Per i support vector \mathbf{x}_s la disuguaglianza diventa un'eguaglianza (ricordiamoci che i support vector di ambo le classi sono a distanza uguale da H per definizione e quindi hanno lo stesso score):

$$y_s (\mathbf{w}^T \mathbf{x}_s + w_0) = 1$$

Per cui abbiamo che $f(\mathbf{x}_s) = \mathbf{w}^T \mathbf{x}_s + w_0 = 1/y_s = y_s$
($y_s = 1$ o $y_s = -1$)



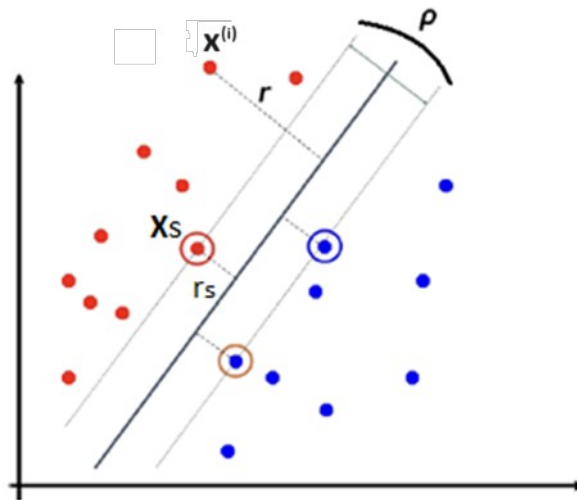
SVM: margine

Ne segue che la distanza (r_s) di \mathbf{x}_s da H sarà:

$$r_s = \frac{|\mathbf{w}^T \mathbf{x}_s + w_0|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

Il *margine* è definito come due volte tale distanza (v. figura) :

$$\rho = 2r_s = \frac{2}{\|\mathbf{w}\|}$$

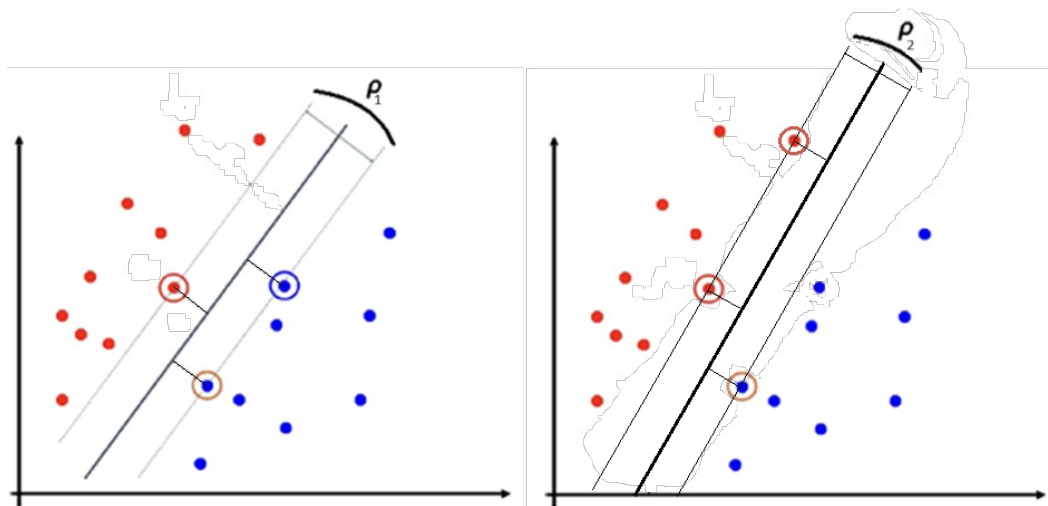


SVM: margine

$$\rho = 2r_s = \frac{2}{\|w\|}$$

L'obiettivo è trovare la score function definita da $[w_0, w]$ il cui corrispondente margine sia massimo

Attenzione: cambiando $[w_0, w]$ cambia l'iper-piano e, quindi, anche il sottoinsieme di punti di T che costituiscono i support vector, come mostra la figura qui sotto



Usando la definizione di margine è possibile impostare la seguente funzione obiettivo:

$$[w_0^*, \mathbf{w}^*] = \arg \max_{[w_0, \mathbf{w}] \in W} \frac{2}{\|\mathbf{w}\|} \quad s.t.$$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

← Vincoli

o, equivalentemente:

$$[w_0^*, \mathbf{w}^*] = \arg \min_{[w_0, \mathbf{w}] \in W} \frac{1}{2} \|\mathbf{w}\|^2 \quad s.t.$$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

← Vincoli

SVM: primal solution

$$[w_0^*, \mathbf{w}^*] = \arg \min_{[w_0, \mathbf{w}] \in W} \frac{1}{2} \|\mathbf{w}\|^2 \quad s.t.$$
$$y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

← Vincoli

Il problema di ottimizzazione così impostato è detto *soluzione primale*

La funzione obiettivo, quadratica rispetto a $[w_0, \mathbf{w}]$, è convessa, quindi ammette un unico minimo globale (potrebbe essere risolta in closed form solution)

Tuttavia, la presenza di *vincoli* ($\{y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1\}_{i=1, \dots, n}$) rende l'ottimizzazione più difficile

Tipicamente la soluzione viene trovata utilizzando tecniche di ottimizzazione quadratica (detta anche *quadratic programming*)

Una formulazione alternativa è data introducendo i cosiddetti “moltiplicatori di Lagrange” ($\alpha_i \geq 0$), uno per ogni esempio in T , i quali permettono di includere i vincoli direttamente nella funzione obiettivo, ottenendo la *soluzione duale*:

$$\boldsymbol{\alpha}^* = [\alpha_1^*, \dots, \alpha_n^*] = \arg \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} \quad s.t.$$
$$\alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

Nella soluzione duale ho ancora dei vincoli, ma sono più semplici da risolvere

Inoltre, la soluzione duale ci permette di individuare i vettori di supporto, che sono gli $\mathbf{x}^{(i)}$ associati ad $\alpha_i > 0$

SVM: confronto tra soluzioni primale e duale

Primale: $[w_0^*, \mathbf{w}^*] = \arg \min_{[w_0, \mathbf{w}] \in W} \frac{1}{2} \|\mathbf{w}\|^2 \quad s.t.$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

Duale: $[\alpha_1^*, \dots, \alpha_n^*] = \arg \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} \quad s.t.$

$$\alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

- Entambe le soluzioni possono essere ottimizzate usando tecniche di ottimizzazione quadratica
- I metodi più recenti, tuttavia, utilizzano delle varianti del Gradient Descent/Ascent
- La soluzione primale consiste nel risolvere un problema di ottimizzazione in $d+1$ variabili
- La soluzione duale consiste nel risolvere un problema di ottimizzazione in n variabili
- N.B.: nella soluzione duale i parametri della funzione ipotesi (score function) sono scomparsi...

SVM: modello predittivo (inference time)

A inference time, il classificatore finale sarà:

$$C(\mathbf{x}) = \text{If } f(\mathbf{x}) > 0 \text{ then } y_1 \\ \text{else } y_2$$

dove, usando la soluzione primale, avremo: $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

che è un modello parametrico lineare

Usando invece la soluzione duale avremo : $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x} + w_0$

dove è possibile calcolare la sommatoria limitandosi ai soli $\alpha_i > 0$, ovvero usando solo i support vector

w_0 , nel caso duale, può essere ricavato utilizzando un qualsiasi support vector \mathbf{x}_s e sfruttando il fatto che $f(\mathbf{x}_s) = y_s$:

$$y_s = f(\mathbf{x}_s) = \sum_{i=1}^n \alpha_i y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x}_s + w_0$$

$$w_0 = y_s - \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^{(i)})^T \mathbf{x}_s$$

SVM: modello predittivo (inference time)

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

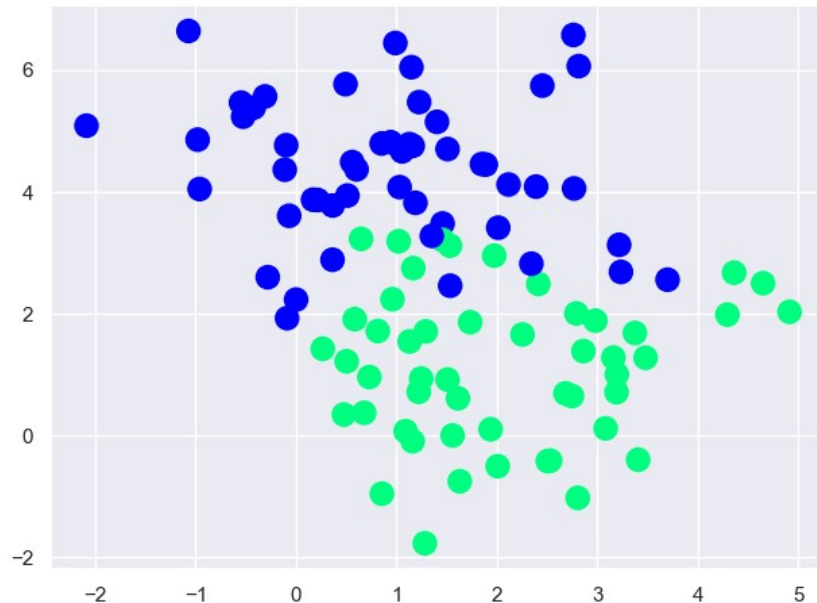
$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x} + w_0 = \sum_{i=1}^n \alpha_i y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x} + y_s - \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^{(i)})^T \mathbf{x}_s$$

Le due soluzioni, primale e duale, sono equivalenti (e sono entrambi modelli parametrici lineari), per cui il vettore dei parametri \mathbf{w} è tale che:

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

Classi parzialmente sovrapposte

Quando le due classi sono parzialmente sovrapposte, come nella figura a lato, un iper-piano separatore non esiste



Classi parzialmente sovrapposte

Lasciamo allora che la SVM commetta un certo numero di «errori», violando i vincoli per alcuni samples (ma facendole pagare una penalità per ogni violazione):

$$w_0^*, \mathbf{w}^*, \xi_1^*, \dots, \xi_n^* = \arg \min_{[w_0, \mathbf{w}] \in W, \xi_i \geq 0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad s.t.$$

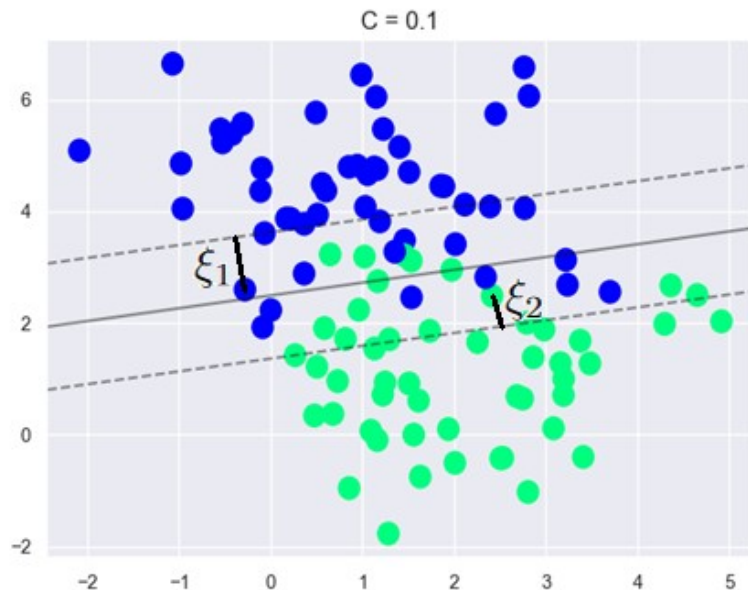
$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 - \xi_i \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

Confronto con la soluzione primale senza penalità:

$$[w_0^*, \mathbf{w}^*] = \arg \min_{[w_0, \mathbf{w}] \in W} \frac{1}{2} \|\mathbf{w}\|^2 \quad s.t.$$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

Le n variabili ξ_1, \dots, ξ_n sono chiamate «slack variables»



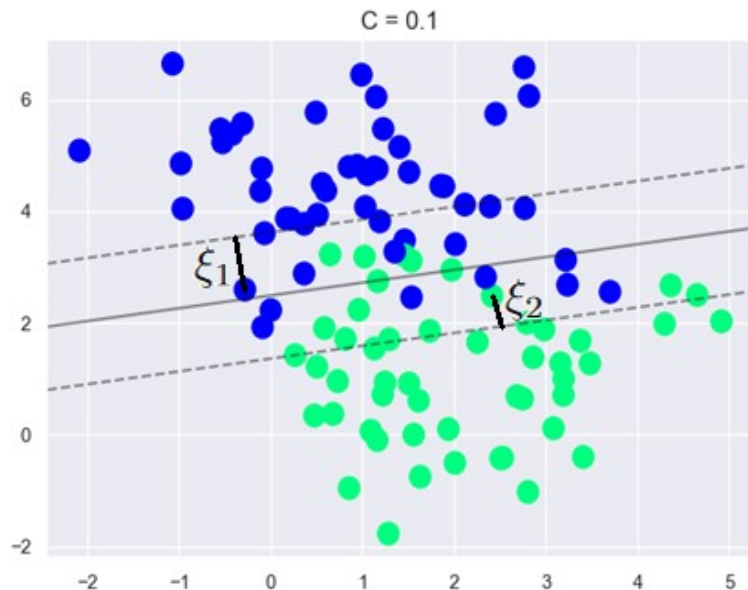
Classi parzialmente sovrapposte

$$w_0^*, \mathbf{w}^*, \xi_1^*, \dots, \xi_n^* = \arg \min_{[w_0, \mathbf{w}] \in W, \xi_i \geq 0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad s.t.$$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 - \xi_i \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in T$$

Anche questa è una forma di *regolarizzazione* che serve ad alleviare l'overfitting, diminuendo la dipendenza dell'iper-piano dall'esatta collocazione spaziale dei sample di frontiera tra le due classi

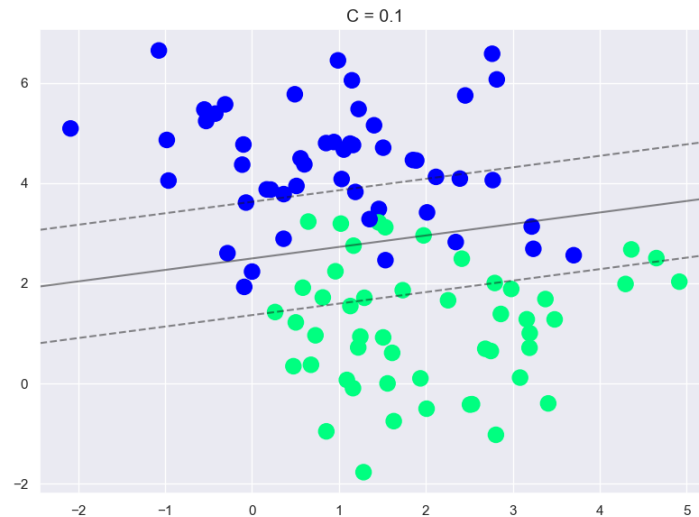
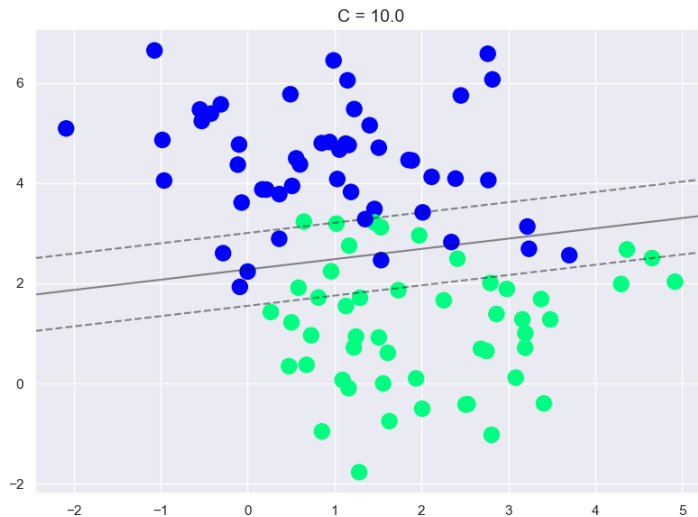
Usando la soluzione duale si può seguire un ragionamento analogo



Classi parzialmente sovrapposte

L'iper-parametro C determina il trade-off tra la massimizzazione del margine e la minimizzazione del numero di sample che violano i vincoli:

- Con C grande, ci si concentra di più sull'evitare errori nella separazione lineare delle due classi, al costo di avere un margine più piccolo (v. figura sotto a sinistra)
- Con C piccolo, si dà meno importanza agli errori di separazione e ci si concentra di più sul massimizzare il margine rispetto ai samples linearmente separabili (cioè quelli che non corrispondono ad «errori»)

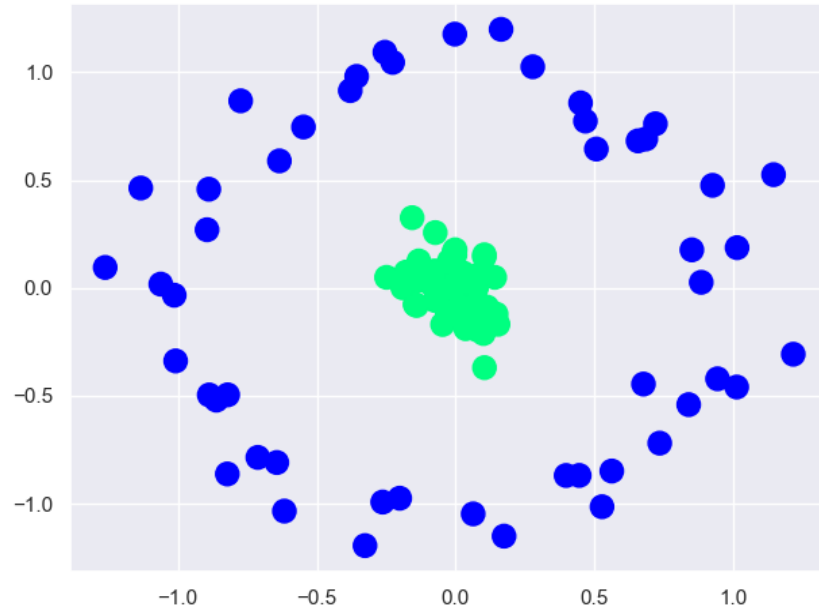


Non-linear Support Vector Machine

Se le due classi, anzichè essere solo parzialmente sovrapposte, sono distribuite in modo tale da non poter essere separate linearmente (e.g., v. figura accanto), allora il problema non può essere risolto con una semplice regolarizzazione

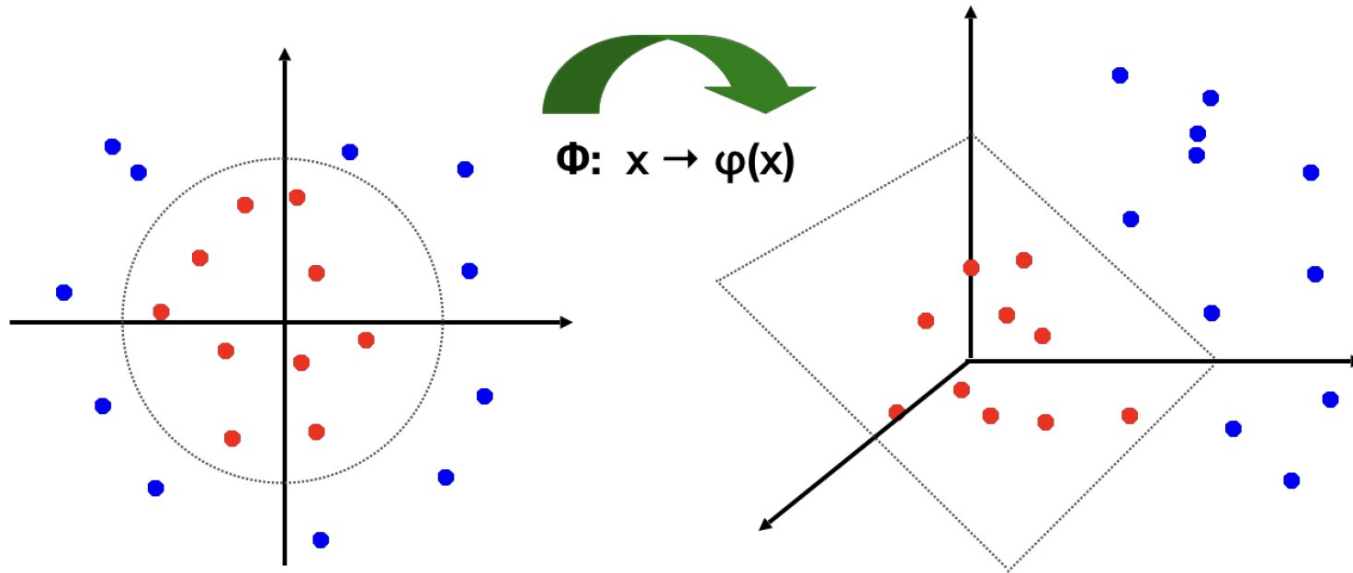
Come faccio a capire se le classi sono linearmente separabili, solo parzialmente sovrapposte o non sono linearmente separabili?

Uso il validation set (lo vedremo meglio nella prossima lezione)



Non-linear Support Vector Machine

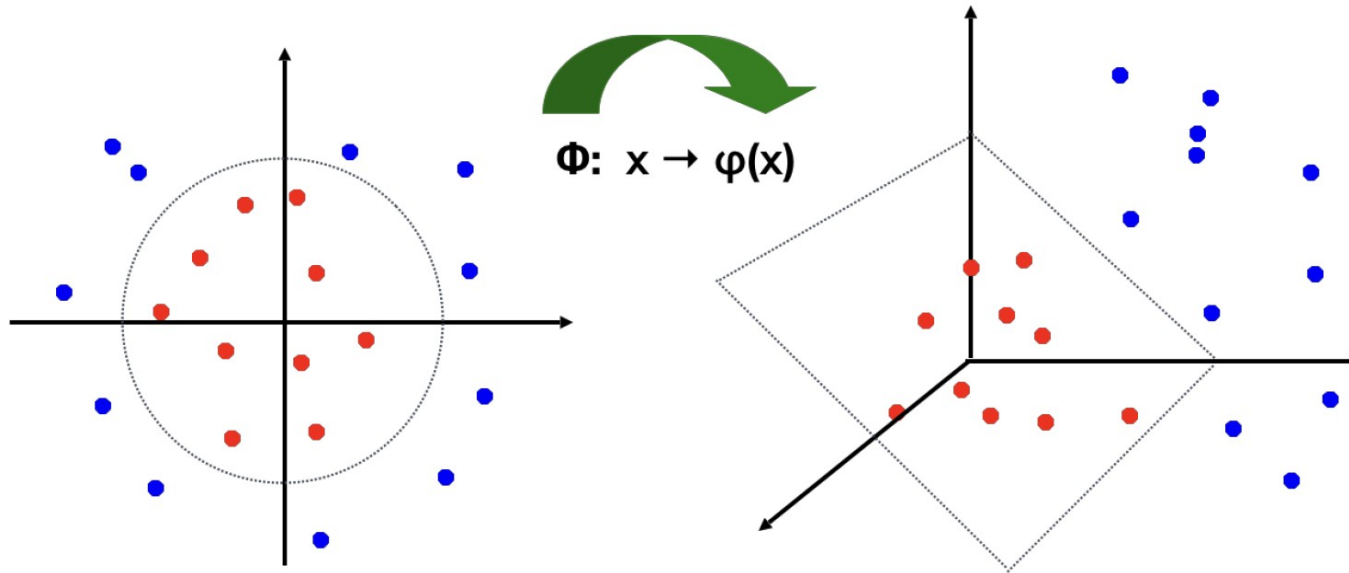
Nel caso le classi non siano linearmente separabili possiamo applicare una trasformazione non lineare $\phi(\mathbf{x})$ per mappare i feature vector $\mathbf{x} \in X = R^d$ in uno spazio con una dimensione più grande ($\phi(\mathbf{x}) = \mathbf{x}' \in X' = R^D, D > d$), dove i sample saranno (sperabilmente) linearmente separabili



Non-linear Support Vector Machine

Benchè, in teoria, una tale trasformazione non lineare $\phi()$ esiste sempre, in pratica trovarla (automaticamente) non è semplice

È però possibile aggirare questo problema ed evitare di dover calcolare $\phi()$ esplicitamente usando il cosiddetto «kernel trick» (ora vedremo come...)



Ricordiamoci che, usando la soluzione duale, si ha che,

- in fase di training:

$$[\alpha_1^*, \dots, \alpha_n^*] = \arg \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} \quad s.t.$$

$$\alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

-
- e a inference time:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x} + w_0$$

Notiamo che i vettori delle feature $(\mathbf{x}, \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \in R^d)$ compaiono solo in prodotti scalari

Questo ci permette di evitare di trasformare direttamente \mathbf{x} , ovvero di dover trovare un nuovo feature space X' , e di limitarci a definire il solo prodotto scalare in uno spazio delle feature «ipotetico» X'

In altre parole, invece di definire esplicitamente $\phi()$, ci basta definire una funzione che calcola il prodotto scalare di due feature vector in uno spazio a dimensioni maggiori di d , senza che questo feature space nuovo sia effettivamente definito

Tale funzione è detta «funzione kernel». Dati due feature vector $\mathbf{x}, \mathbf{z} \in X$, il kernel restituisce un valore scalare che rappresenta il loro prodotto scalare in un feature space a dimensioni maggiori di X :

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z}) \in \mathbb{R}$$

In pratica, $K()$ stima una qualche misura di «similarità» (non lineare) tra \mathbf{x} e \mathbf{z} (tra un po' vedremo degli esempi di kernel)

Non-linear SVM con kernel trick

Con la soluzione duale si ha, in fase di training:

$$[\alpha_1^*, \dots, \alpha_n^*] = \arg \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(j)})$$

che, usando il kernel, diventa:

$$[\alpha_1^*, \dots, \alpha_n^*] = \arg \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

Non-linear SVM con kernel trick

Mentre, ad inference time, ho:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i) \phi(\mathbf{x}) + w_0$$

che, usando il kernel, diventa:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0$$

dove, come nel caso senza kernel, se \mathbf{x}_s è un support vector, sfruttando il fatto che $f(\mathbf{x}_s) = y_s$, posso calcolare:

$$w_0 = y_s - \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}^{(i)}, \mathbf{x}_s)$$

Quindi non ho bisogno di definire $\phi(\mathbf{x})$: sia in fase di training che di inferenza, mi basta definire $K(\mathbf{x}, \mathbf{z})$. Questo modo di procedere è chiamato “kernel trick”

Non-linear Support Vector Machine

I kernel più usati sono:

- **Kernel lineare.** Prodotto scalare standard:
 $K(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$ (i.e., nessun kernel e nessuna deformazione di R^d)
- **Kernel polinomiale.** Permette di «curvare» R^d :

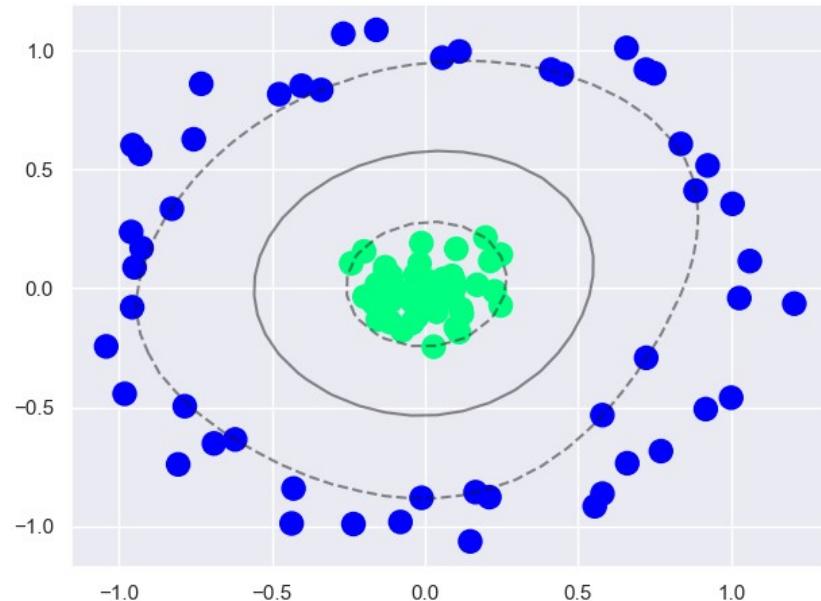
$$K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + a)^m$$

dove a ed m sono iper-parametri

- **Radial Basis Function (RBF) o Gaussian kernel.** Crea deformazioni complesse di R^d :

$$K(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right)$$

E' il kernel più usato e σ è il suo iper-parametro



Non-linear SVM: modello non parametrico

Se uso i kernel, però, non posso usare la soluzione primale, in quanto, non avendo esplicitamente definito lo spazio delle feature $X' = R^D$, non ho neanche lo spazio dei parametri ($W' = R^{D+1}$) e il vettore dei parametri w non è definito

Similmente, a inference time, non posso calcolare direttamente il vettore dei parametri: non avendo definito $\phi()$, non posso *sommare* i vettori trasformati e quindi non posso calcolare:

$$w = \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i)$$

Non-linear SVM: modello non parametrico

Se uso i kernel, quindi, devo necessariamente memorizzare i support vector, perchè, a inference time, dovrò calcolare:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0 \quad w_0 = y_s - \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}^{(i)}, \mathbf{x}_s)$$

Questa cosa ha due implicazioni

La prima è di natura pragmatica: finito il training, devo memorizzarmi tutti i support vector e, a inference time, applicare il kernel ad ognuno di essi

Considerando che, di solito, il numero dei support vector non è piccolissimo, ciò comporta un certo svantaggio computazionale in spazio e tempo

Non-linear SVM: modello non parametrico

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0$$

La seconda implicazione è di natura, per così dire, più “filosofica”, nel senso che le SVM con kernel non lineari non possono più essere considerate come metodi parametrici (senza kernel invece sono parametrici a tutti gli effetti)

Questo perchè, per implementare la formula di sopra, che esprime la funzione di score finale, non posso fare a meno di basarmi su un algoritmo che iteri per tutti i support vector, e questo loop non può essere espresso con una combinazione di funzioni analitiche note, quindi non rientra nella definizione di “metodo parametrico”

In altri termini, $f()$ ora dipende dagli specifici support vector estratti dallo specifico T e corrispondente task che stiamo esaminando, e non c'è modo di esprimerli con un numero finito di parametri (il parameter vector) di una funzione analitica

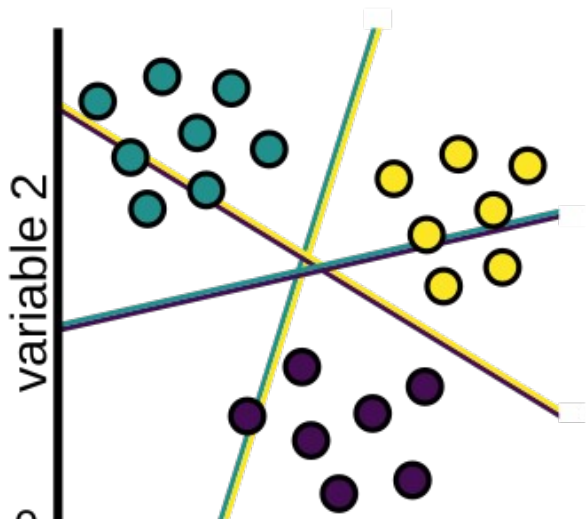
$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0$$

In un certo senso, la soluzione duale con kernel non lineari può essere vista come una specie di k-Nearest Neighbors (k-NN), in cui:

- a inference time, il vettore di testing \mathbf{x} viene confrontato solo con i k vettori di supporto (indipendentemente dal valore di \mathbf{x} la scelta dei k samples per il confronto è sempre la stessa e sono sempre i k vettori di supporto),
- il kernel è usato come misura di *similarità* tra due vettori (maggiore è il valore $K(\mathbf{x}^{(i)}, \mathbf{x})$, maggiore è la loro similarità) al posto (dell'inverso) della distanza Euclidea,
- α_i è usato come peso fisso per determinare l'importanza relativa del vettore di supporto $\mathbf{x}^{(i)}$,
- i valori di ground truth y_i vengono pesati da $\alpha_i K(\mathbf{x}^{(i)}, \mathbf{x})$ e la media pesata viene restituita in output come nel caso del k-NN usato per regressione

SVM: estensione a problemi multiclasse ($y \in \{1, \dots, k\}$)

One versus one



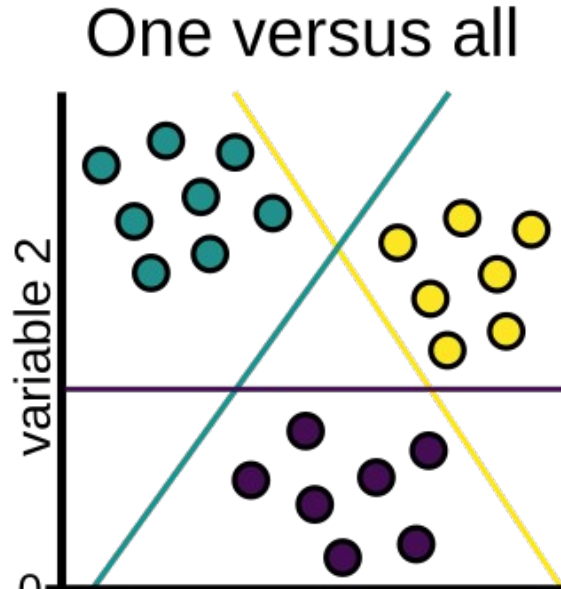
In **One-vs-One** (OvO) si addestrano M classificatori binari per discriminare ogni possibile paio di classi, per un totale di

$$M = \frac{k(k-1)}{2}$$

classificatori binari

A inference time si applica uno schema di voto (eventualmente pesato usando lo score): la classe che riceve più voti «vince»

SVM: estensione a problemi multiclasse ($y \in \{1, \dots, k\}$)



In **One-vs-All** (OvA), o **One-vs-Rest** (OvR) si addestrano M classificatori binari per discriminare ciascuna delle k classi da tutte le altre $k-1$ classi, per un totale di

$$M = k$$

A inference time si sceglie la classe corrispondente al classificatore più “confident”, ovvero quello che, per la “sua” classe positiva ha stimato lo score col valore maggiore rispetto agli altri classificatori

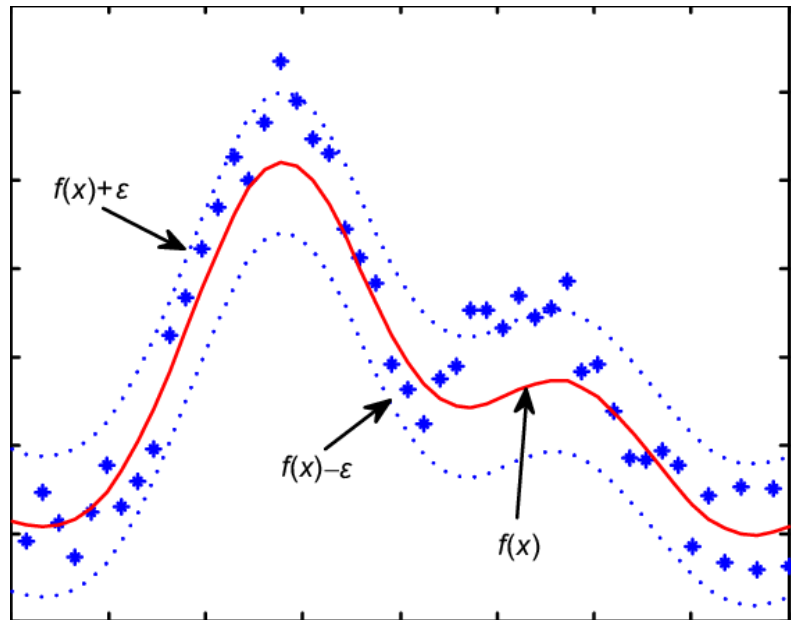
Support Vector Regression (SVR)

Come sappiamo, in un task di regressione vogliamo trovare una funzione che approssimi la relazione tra le feature e la variabile target, dove quest'ultima assume valori numerici continui

Le SVM possono essere usate come metodo di regressione, mantenendo tutte le loro caratteristiche principali. In tal caso, si chiamano Support Vector Regression (SVR)

Nel caso della regressione i vincoli delle SVM vengono riformulati imponendo che, per ogni $(\mathbf{x}_i, y_i) \in T$, la differenza tra il valore di ground truth y_i e la predizione $f(\mathbf{x}_i)$ sia inferiore o uguale ad un margine massimo di errore definito dall'iper-parametro ε

Usando il kernel trick la funzione di regressione risultante può essere altamente non lineare (e.g., v. figura a fianco)



Anche con le SVM/SVR è utile scalare i dati prima di iniziare il training perché:

- L'ottimizzazione, sia per la soluzione primale che duale, avviene con algoritmi iterativi come il Gradient Descent
- Sono metodi definiti rispetto ad un «margine», che è una distanza (Euclidea), il cui calcolo è influenzato dalle unità di misura delle varie feature
- Anche il calcolo del kernel è influenzato dalle unità di misura delle varie feature

Vantaggi:

- Strumento molto potente, con una notevole resistenza all'overfitting (spesso maggiore ad altri metodi)
- Possono trattare problemi non lineari (se uso i kernel)

Limitazioni:

- Non c'è nessuna garanzia che si riesca a trovare (facendo delle prove col validation set...) il kernel e i relativi iper-parametri giusti che permettano di separare due classi non linearmente separabili

Svantaggi:

- Nel caso duale con kernel non-lineare bisogna memorizzare tutti i support vector, che potrebbero essere tanti con un dataset di training enorme
- Non esiste un'interpretazione probabilistica diretta per l'appartenenza di un sample a una classe (ma questa potrebbe essere ricavata come post-processing a partire dallo score)

Riferimenti

- <https://see.stanford.edu/materials/aimlcs229/cs229-notes3.pdf>
- **A tu per tu col Machine Learning. L'incredibile viaggio di un developer nel favoloso mondo della Data Science**, Alessandro Cucci, The Dot Company, 2017 [cap.5]
- **Pattern Classification, second edition**, R. O. Duda, P. E. Hart, D. G. Stork, Wiley-Interscience, 2000
- **An introduction to statistical learning**, Gareth M. James, Daniela Witten, Trevor Hastie, Robert Tibshirani, New York: Springer, 2013 [cap. 4]