

-UNITE!-
**University Network for Innovation,
Technology and Engineering**

AALTO UNIVERSITY
School of Science



unite!

Human-Centered Machine Learning Course

SUPPORT VECTOR MACHINES

Professor:
Cecilio Angulo

Authors:
Pablo De Ramon
Giang Le
Vasco Pearson
Amir Ingher
Sergi Garcia
Iván Quirante

Academic Course 2022-2023

1. Questions

Q1

For a given dataset, the SVM learning algorithm is applied leading to the optimal decision plane $H(x) = 3x_1 + 4x_2 + x_3 + 2$.

- Which is the value of the margin for this SVM?
- Find three vectors in the input space that could be support vectors of the plane.
- Demonstrate that the distance between parallel hyperplanes in SVM effectively is $\frac{2}{\|\mathbf{w}\|}$.

Solution

The margin of a given SVM is the distance between two parallel hyperplanes

$$H_{\pm}(x) = \{\mathbf{w}^T \cdot \mathbf{x} + \gamma = \pm 1\}$$

that separate a set of training data into two classes, and is the value we want to maximize in order to find the optimal decision plane - or perceptron of optimal stability - $H(x)$ that lies exactly between $H_{\pm}(x)$.

Theory tells us that the margin m is equal to $\frac{2}{\|\mathbf{w}\|}$. We also know that the weight vector \mathbf{w} is the normal vector of the decision plane. By simple inspection of $H(x)$ one can deduce the value of \mathbf{w} :

$$\mathbf{w} = (3 \ 4 \ 1)^T$$

We measure its Euclidean norm as follows:

$$\|\mathbf{w}\| := \sqrt{\langle \mathbf{w}, \mathbf{w} \rangle} = \sqrt{3^2 + 4^2 + 1^2} = \sqrt{26} \approx 5.099$$

In consequence, the value of the margin is:

$$m \equiv \frac{2}{\|\mathbf{w}\|} = \frac{2}{\sqrt{26}} \approx 0.392$$

A SVM learning algorithm aims to find the solution to the following optimisation problem:

$$\begin{aligned} \max \quad & \frac{2}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + \gamma) \geq 1 \text{ for } i = 1, \dots, n \end{aligned}$$

The \mathbf{w} and γ that solve this problem determine the classifier: $\mathbf{x} \mapsto \text{sign}(\mathbf{w}^T \mathbf{x} + \gamma)$, where $\text{sign}(\cdot)$ is the sign function. An important consequence of this geometric description is that those \mathbf{x}_i that lie in the outer hyperplanes $H_{\pm}(x)$ - also known as support vectors - completely determine the max-margin hyperplane since they are extreme points of the domain polyhedron. Thus, any vector inside these planes could be a solution for the problem.

We search for three possible support vectors by finding data points that fulfill one of the following conditions:

$$\begin{aligned} H_+(x) &= \{3x_1 + 4x_2 + x_3 + 2 = 1\} \rightarrow 3x_1 + 4x_2 + x_3 = -1 \\ H_-(x) &= \{3x_1 + 4x_2 + x_3 + 2 = -1\} \rightarrow 3x_1 + 4x_2 + x_3 = -3 \end{aligned}$$

Three possibilities are, for example:

$$\boxed{\mathbf{x}_a = (0 \ 0 \ -1)^T, \ \mathbf{x}_b = (0 \ 0 \ -3)^T, \ \mathbf{x}_c = (-1 \ 0 \ 0)^T},$$

where $\mathbf{x}_a \in H_+(x)$ and $\{\mathbf{x}_b; \mathbf{x}_c\} \in H_-(x)$.

Proof. Now we're going to demonstrate that $m = \frac{2}{\|\mathbf{w}\|}$. First, let $H_{\pm}(x) = \{\mathbf{w}^T \mathbf{x} + \gamma = \pm 1\}$ be the two outermost parallel hyperplanes of the SVM optimisation problem that maximize the margin between data points of different classes. Next, let's define $a = 1 - \gamma$ and $b = -1 - \gamma$ so that one can rewrite the previous subspaces as $H_+(x) = \{\mathbf{w}^T \mathbf{x} = a\}$ and $H_-(x) = \{\mathbf{w}^T \mathbf{x} = b\}$. Considering two data points $\mathbf{x}_+ \in H_+ \leftrightarrow \mathbf{w}^T \mathbf{x}_+ = a$ and $\mathbf{x}_- \in H_- \leftrightarrow \mathbf{w}^T \mathbf{x}_- = b$ whose distance is exactly m - that is, that both lie on the same line if it is drawn perpendicular to the decision plane - the relation between \mathbf{x}_+ , \mathbf{x}_- and \mathbf{w} can be written as follows:

$$\mathbf{x}_+ = \mathbf{x}_- + \alpha \cdot \mathbf{w}, \quad \alpha \in \mathbb{R}^+$$

Computing the dot product with $\langle \mathbf{w}^T, \cdot \rangle$ on both sides:

$$\mathbf{w}^T \mathbf{x}_+ = \mathbf{w}^T \mathbf{x}_- + \alpha \cdot \mathbf{w}^T \mathbf{w}$$

$$a = b + \alpha \cdot \|\mathbf{w}\|^2$$

Since $a = 1 - \gamma$ and $b = -1 - \gamma$, one can easily obtain the value of α :

$$\alpha = \frac{a - b}{\|\mathbf{w}\|^2} = \frac{1 - \gamma - (-1 - \gamma)}{\|\mathbf{w}\|^2} = \frac{2}{\|\mathbf{w}\|^2}$$

Finally, the value the margin m can be obtained as follows:

$$m = \|\mathbf{x}_+ - \mathbf{x}_-\| = \|\alpha \cdot \mathbf{w}\| = |\alpha| \cdot \|\mathbf{w}\| = \frac{2}{\|\mathbf{w}\|^2} \cdot \|\mathbf{w}\| = \frac{2}{\|\mathbf{w}\|}$$

□

Q2

For the 1-dimensional training set $\mathcal{T} = \{(3; -1); (2; -1); (4; -1); (6; 1); (7; 1)\}$, find the kernel machine with polynomial kernel of lower degree able to separate the data set. Draw the resulting separation function.

Solution

It often happens that the sets to discriminate are not linearly separable in the feature space, as it was assumed in Q1. In order to solve this problem, a widely-known solution is to apply the **kernel trick**, which consists of finding a function $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$ that maps a feature \mathbf{x} into a high-dimensional space \mathcal{Y} in which one can define a dot product as $K(x, y) = \langle \varphi(x), \varphi(y) \rangle$. This latter bilinear mapping is known as **kernel**. In particular, the expression of the **polynomial kernel** can be written as follows:

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + r)^d \text{¹},$$

where $d \in \mathbb{R}^+$ is the degree and $r \in \mathbb{R}$ the kernel's coefficient. As can be seen, there exist infinite combinations of d and r values, that is why the team has opted to implement a parameter-search in Python (1) in order to find the combination that gets the best accuracy.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import svm, datasets
4 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
5 from sklearn.pipeline import make_pipeline
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.metrics import accuracy_score
8
9 #Replicate datapoints to avoid errors in the fold split
10 Xp=np.array([3,3,3,2,2,4,4,4,6,6,6,6,7,7])
11 yp=np.array([-1,-1,-1,-1,-1,-1,-1,-1,-1,1,1,1,1,1])
12
13 #Data pre-processing
14 X_train,X_val,y_train, y_val = train_test_split(Xp.reshape(-1,1),yp,test_size=0.5,
15         random_state=0)
16
17 #Grid search
18 param_grid = {'kernel':['poly'], 'degree':[0,1,2,3,4], 'C':[1,10,100,1000], 'coef0'
19         :[0,1,10,100,1000]}
20 gridSearchCV =GridSearchCV(estimator = svm.SVC(),param_grid = param_grid, cv=2,scoring="
21         accuracy",refit=True,return_train_score=True)
22 gridSearchCV.fit(X_train,y_train)
23
24 print(gridSearchCV.best_estimator_)
25 print(gridSearchCV.best_score_)
26
27 #Output:
28 #SVC(C=1, coef0=0, degree=1, kernel='poly')
29 #1.0
30
31 #Plot initial and transformed data + separator hyperplane
32 def y(x):
33     return -x+10
34
35 xplot=np.linspace(2,7,100)
36 yplot=y(xplot)
37 plt.plot([2,3,4],[0,0,0], 'bo')
38 plt.plot([2,3,4],[2,3,4], 'b*')
39 plt.plot([6,7],[6,7], 'r*')
40 plt.plot([6,7],[0,0], 'ro')
41 plt.plot(xplot,yplot, 'g')
```

Code 1: Python code

¹Here \mathbf{y} refers to a data point and not to a label.

As can be seen, the best polynomial kernel obtained is:

$$K(\mathbf{x}, \mathbf{y}) = \mathbf{x} \cdot \mathbf{y} = \varphi(\mathbf{x}) \cdot \varphi(\mathbf{y}),$$

which is a linear and homogeneous kernel. Finally, lines 27-38 of [1](#) plot the current (o) and transformed (*) data as well as the classification hyperplane (in green) that best separates the two classes (blue for +1 and red for -1).

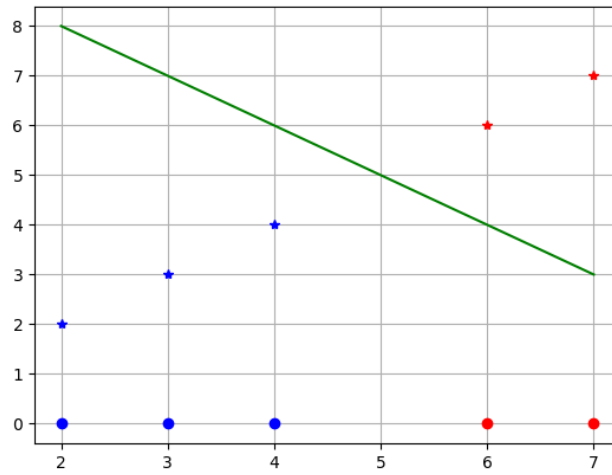


Fig. 1: Plot of (non-)transformed data and the linear classifier