

Week 13: Support Vector Machines

GO GREEN. AVOID PRINTING, OR PRINT 2-SIDED MULTI-PAGE.

13.1 Introduction

Support vector machines (SVMs) are considered one of the best “out of the box” classifiers. Similar to logistic regression and random forests, their main purpose is prediction/classification, and are used for similar applications (e.g. distinguish between healthy and Alzheimers). However, SVMs are a more geometric approach. The main idea is to think of each sample as a point in a high-dimensional space, and classify it according to the region where it is located with respect to a *boundary* (see Figure 13.1). The challenge is to find such separating/classifying boundary.

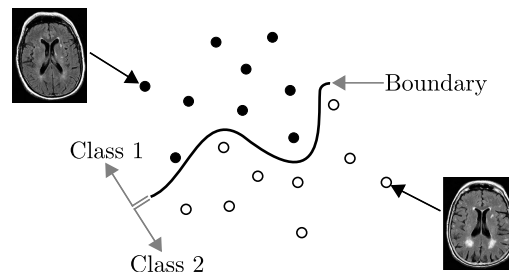


Figure 13.1: An SVM classifies each point according to the region where it is located with respect to a boundary.

13.2 Hyperplanes

Hyperplanes lie at the heart of SVMs. Intuitively, a hyperplane is the generalization of a line and a plane (in 2 and 3 dimensions) to higher dimensions. In words, hyperplanes are subspaces almost as big as the whole space. More formally, a hyperplane $\mathbb{H} \subset \mathbb{R}^d$ is a linear subspace of dimension $d - 1$ (see Figure 13.2 to build some intuition).

Recall from linear algebra that the dimensions of a subspace $\mathbb{U} \subset \mathbb{R}^d$ and its orthogonal complement $\mathbb{U}^\perp \subset \mathbb{R}^d$ add to the ambient dimension, i.e.,

$$d = \dim(\mathbb{U}) + \dim(\mathbb{U}^\perp),$$

which means that $\dim(\mathbb{H}^\perp) = 1$, i.e., \mathbb{H}^\perp is a line. As such, it is characterized by a single vector. Let $\beta \in \mathbb{R}^d$ be the vector spanning \mathbb{H}^\perp . Also recall that \mathbb{U} is the collection of all points orthogonal to \mathbb{U}^\perp . Consequently,

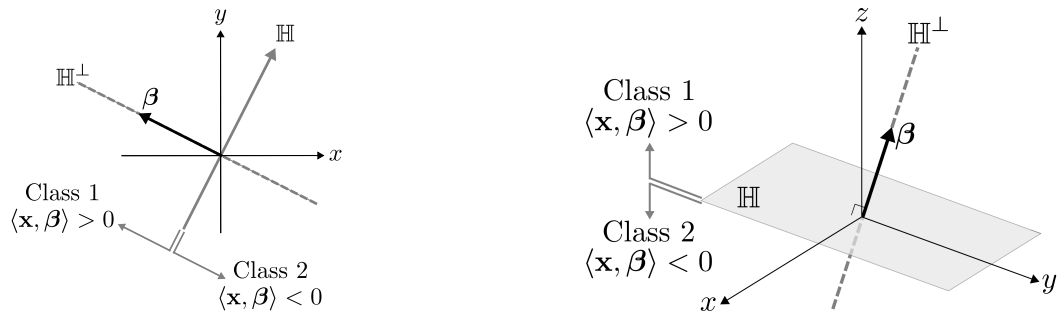


Figure 13.2: Hyperplanes are subspaces of dimension equal to the whole space minus 1. **Left:** A line in \mathbb{R}^2 is a hyperplane. **Right:** A plane in \mathbb{R}^3 is a hyperplane. Hyperplanes divide the whole space in two halves, providing a natural classifying border. **Question:** Is a line in \mathbb{R}^3 a hyperplane? Each hyperplane \mathbb{H} is characterized by the vector β spanning its orthogonal complement \mathbb{H}^\perp .

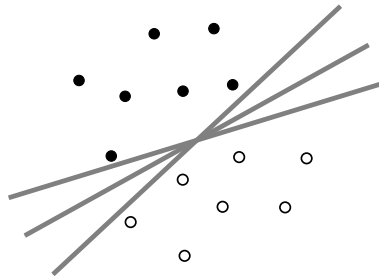
we can characterize \mathbb{H} as the collection of all points in \mathbb{R}^d that are orthogonal to β , i.e.,

$$\{ \mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \beta \rangle = 0 \},$$

where $\langle \mathbf{x}, \beta \rangle = \beta^\top \mathbf{x}$ denotes the inner product (see Figure 13.2 to build some intuition). Hyperplanes are powerful classification tools because they divide the space in two halves: the points \mathbf{x} *above* the hyperplane, meaning $\langle \mathbf{x}, \beta \rangle > 0$, and the points \mathbf{x} *below* the hyperplane, meaning $\langle \mathbf{x}, \beta \rangle < 0$ (see Figure 13.2).

13.3 Maximal Margin Classifier

SVMs happen to be a generalization of the maximal margin classifier (MMC), which uses a *hyperplane* as the separating/classifying border. Notice, however, that given some data, it is possible that there are infinitely many hyperplanes that could separate it, for example:



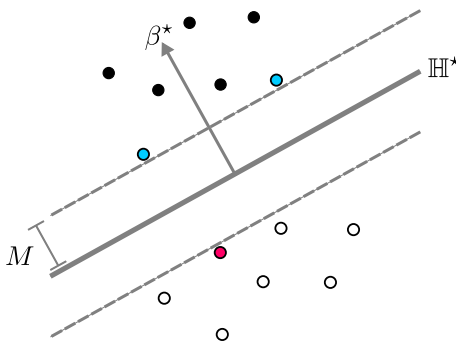
As the name suggests, the MMC finds the hyperplane that separates data with the highest possible margin M . More precisely, given data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, with labels/classes $y_1, \dots, y_n \in \{-1, 1\}$, the MMC aims to find:

$$\beta^* := \arg \max_{\substack{\beta \in \mathbb{R}^d \\ \|\beta\|=1}} M \quad \text{subject to} \quad y_i \langle \mathbf{x}_i, \beta \rangle \geq M \quad \forall i = 1, \dots, n. \quad (13.1)$$

The intuition behind (13.1) is quite simple. First, M denotes the margin that we want to maximize. The first constraint $\|\beta\| = 1$ simply guarantees that we do not cheat by making β too large. The second constraint $y_i \langle \mathbf{x}_i, \beta \rangle \geq M$ ensures that each training point is on the right side of the hyperplane (because it requires that

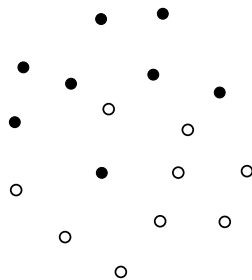
the sign of the *true* label y_i and the *classification* $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle$ are equal) and also ensures that the perpendicular distance between \mathbf{x}_i and the hyperplane (given by $y_i \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle$) is larger than the margin M .

All vectors that are exactly within M distance of the hyperplane are called *support vectors*:



13.4 Support Vector Classifier

It is pretty clear that a hyperplane may not classify two classes perfectly:

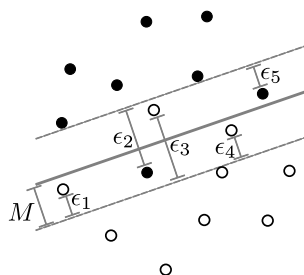


The support vector classifier (SVC), also known as the soft margin classifier, arises as an idea to address this problem. The key idea is to allow some wiggle room by letting some points to be misclassified, using the following optimization instead of (13.1)

$$\boldsymbol{\beta}^* := \underset{\substack{\boldsymbol{\beta} \in \mathbb{R}^d, \boldsymbol{\epsilon} \in \mathbb{R}^d; \\ \|\boldsymbol{\beta}\|=1, \|\boldsymbol{\epsilon}\| \leq E}}{\arg \max} M \quad \text{subject to} \quad \begin{cases} y_i \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle \geq M(1 - \epsilon_i) \\ \epsilon_i \geq 0 \end{cases} \quad \forall i = 1, \dots, n. \quad (13.2)$$

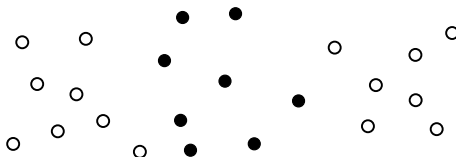
There are a few differences between (13.1) and (13.2). First, we are now also maximizing over $\boldsymbol{\epsilon}$, which models the wiggle room that will be allowed in our sample. The constraint $\|\boldsymbol{\epsilon}\| \leq E$ guarantees that the overall wiggle room does not exceed the overall allowed error E , which is a tuning parameter. The constraint $y_i \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle \geq M(1 - \epsilon_i)$ ensures that the orthogonal distance from \mathbf{x}_i and the hyperplane is larger than M , or only a *little* smaller; how *little* is determined by ϵ_i , which should be greater than zero (do you understand why?).

This type of approach is often called a *relaxation*, and will allow some points to be inside the margin, and even on the wrong side of the hyperplane, as in the following figure:



13.5 Nonlinear Boundaries

It is pretty clear that not all datasets can be divided by a linear boundary:



Hence, instead of classifying each point using the linear function $\langle \mathbf{x}, \boldsymbol{\beta} \rangle = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_d x_d$, we can use other non-linear functions, for example quadric (degree-2) polynomials:

$$f_2(\mathbf{x}) := \beta_{11}x_1^2 + \beta_{12}x_1x_2 + \beta_{13}x_1x_3 + \dots + \beta_{1d}x_1x_d \tag{13.3}$$

$$+ \beta_{22}x_2^2 + \beta_{23}x_2x_3 + \dots + \beta_{2d}x_2x_d \tag{13.4}$$

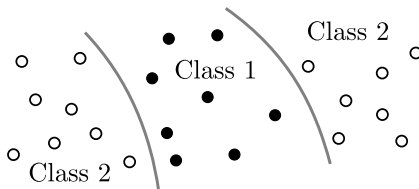
$$+ \vdots \tag{13.5}$$

$$+ \beta_{dd}x_d^2. \tag{13.6}$$

This suggests extending the SVC formulation in (13.2) to:

$$\boldsymbol{\beta}^* := \underset{\substack{\boldsymbol{\beta} \in \mathbb{R}^d, \boldsymbol{\epsilon} \in \mathbb{R}^d \\ \|\boldsymbol{\beta}\|=1, \|\boldsymbol{\epsilon}\| \leq E}}{\arg \max} M \quad \text{subject to} \quad \begin{cases} y_i f_2(\mathbf{x}_i) \geq M(1 - \epsilon_i) \\ \epsilon_i \geq 0 \end{cases} \quad \forall i = 1, \dots, n. \tag{13.7}$$

The only difference with (13.2) is that now we use $f_2(\mathbf{x}_i)$ instead of $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle$. Intuitively, this is guaranteeing that the true label y_i and the prediction $f_2(\mathbf{x}_i)$ agree by more than $(1 - \epsilon_i)$ of the margin M . This will produce a non-linear boundary:



Letting

$$\boldsymbol{\beta} \otimes \boldsymbol{\beta} := \begin{bmatrix} \beta_1 \boldsymbol{\beta} \\ \beta_2 \boldsymbol{\beta} \\ \vdots \\ \beta_d \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \beta_1 \beta_1 \\ \beta_1 \beta_2 \\ \vdots \\ \beta_1 \beta_d \\ \beta_2 \beta_1 \\ \beta_2 \beta_2 \\ \beta_2 \beta_3 \\ \vdots \\ \beta_2 \beta_d \\ \vdots \\ \beta_d \beta_1 \\ \beta_d \beta_2 \\ \vdots \\ \beta_d \beta_d \end{bmatrix},$$

and similarly for $\mathbf{x} \otimes \mathbf{x}$, we can rewrite (13.7) as

$$\boldsymbol{\beta}^* := \underset{\substack{\boldsymbol{\beta} \in \mathbb{R}^d, \boldsymbol{\epsilon} \in \mathbb{R}^d: \\ \|\boldsymbol{\beta}\|=1, \|\boldsymbol{\epsilon}\| \leq \epsilon}}{\arg \max}} M \quad \text{subject to} \quad \begin{cases} y_i \langle \mathbf{x}_i \otimes \mathbf{x}_i, \boldsymbol{\beta} \otimes \boldsymbol{\beta} \rangle \geq M(1 - \epsilon_i) \\ \epsilon_i \geq 0 \end{cases} \quad \forall i = 1, \dots, n, \quad (13.8)$$

which is linear in the *lifted* vectors $\mathbf{x} \otimes \mathbf{x}$ and $\boldsymbol{\beta} \otimes \boldsymbol{\beta}$. Consequently, by simply multiplying our data by itself, we can construct non-linear boundaries with the same linear techniques as before. The downside is that the number of parameters increases polynomially. For instance, instead of the d parameters (corresponding to the entries in $\boldsymbol{\beta}$) we now have $\binom{d}{2} = d(d - 1)/2$ parameters, corresponding to the entries $\boldsymbol{\beta} \otimes \boldsymbol{\beta}$ (Notice that $\boldsymbol{\beta} \otimes \boldsymbol{\beta} \in \mathbb{R}^{d^2}$, but some entries are duplicated; in total there are $\binom{d}{2}$ distinct entries).

13.6 Kernels and Support Vector Machines

Kernels are functions that quantify the similarity between two vectors. For example, if you remember from your linear algebra 101 class, the inner product is a kernel that measures angle as follows:

$$\angle(\mathbf{x}, \boldsymbol{\beta}) = \frac{\langle \mathbf{x}, \boldsymbol{\beta} \rangle}{\|\mathbf{x}\| \|\boldsymbol{\beta}\|}$$

There are many other kernels (functions that measure similarity between vectors), for example:

$$K(\mathbf{x}, \boldsymbol{\beta}) = \langle \mathbf{x} \otimes \mathbf{x}, \boldsymbol{\beta} \otimes \boldsymbol{\beta} \rangle,$$

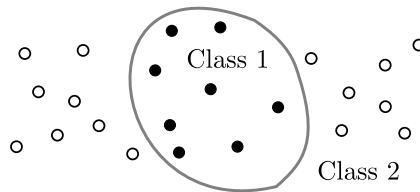
which is precisely the function that we used in (13.8), corresponding to the quadric polynomial in (13.3). The more similar \mathbf{x} and $\boldsymbol{\beta}$ are, the larger $K(\mathbf{x}, \boldsymbol{\beta})$ will be. We can extend this idea to obtain a polynomial kernel of degree r :

$$K(\mathbf{x}, \boldsymbol{\beta}) = (1 + \langle \mathbf{x}, \boldsymbol{\beta} \rangle)^r = \left(1 + \sum_{i=1}^d x_i \beta_i \right)^r.$$

Kernels are not necessarily only polynomial functions. For example, another popular choice is the *radial* kernel:

$$K(\mathbf{x}, \boldsymbol{\beta}) = e^{-\gamma \|\mathbf{x} - \boldsymbol{\beta}\|^2} = \exp\left(-\gamma \sum_{i=1}^d (x_i - \beta_i)^2\right),$$

which will produce a boundary like the following:



Depending on the dataset at hand, one kernel may be better than another. The support vector machine (SVM) is simply a generalization of the SVC that uses a kernel function:

$$\boldsymbol{\beta}^* := \underset{\substack{\boldsymbol{\beta} \in \mathbb{R}^d, \boldsymbol{\epsilon} \in \mathbb{R}^d, \\ \|\boldsymbol{\beta}\|=1, \|\boldsymbol{\epsilon}\| \leq B}}{\arg \max}} M \quad \text{subject to} \quad \begin{cases} y_i K(\mathbf{x}_i, \boldsymbol{\beta}) \geq M(1 - \epsilon_i) \\ \epsilon_i \geq 0 \end{cases} \quad \forall i = 1, \dots, n. \quad (13.9)$$