

# Support Vector Machines

In this section we briefly explain the Support Vector Machines (SVM) classification algorithm for binary classification. We consider a set of training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$  where each vector  $\mathbf{x}_i \in \mathbb{R}^N$  is a collection of  $N$  features and is labelled  $y_i \in \{-1, +1\}$ . The basic idea of SVMs is that we can find some boundary in  $\mathbb{R}^N$  which "best" separates the vectors labelled  $-1$  from the vectors labelled  $+1$ . We first consider boundaries that are hyperplanes, before we go on to more complicated non-linear boundaries.

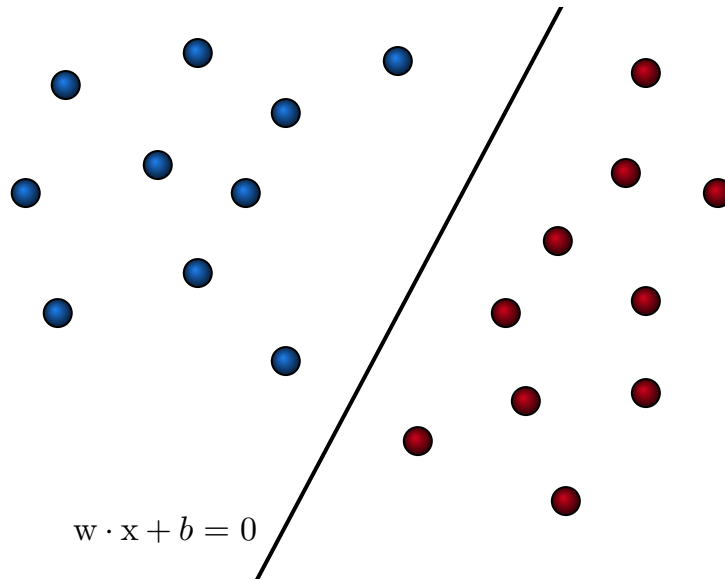
## Linearly Separable Case

To start out simple we will first consider the case where the vectors  $\mathbf{x}_i$  are *linearly seperable*. This means that we can find a hyperplane  $\mathbf{w} \cdot \mathbf{x} + b = 0$ ,  $\mathbf{w} \neq 0$  such that all the vectors labelled  $-1$  fall on one side and all the vectors labelled  $+1$  fall on the other. More specifically, we want to find  $\mathbf{w}$  and  $b$  such that all the vectors labelled  $-1$  satisfy  $\mathbf{w} \cdot \mathbf{x}_i + b < 0$  and all the vectors labelled  $+1$  to satisfy  $\mathbf{w} \cdot \mathbf{x}_i + b \geq 0$ . We can write this condition compactly as

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 0, \quad \forall i \in [m]$$

Once we've found  $\mathbf{w}$  and  $b$  that satisfies this condition for all the training data, we can label *new* data using a linear classifier

$$h(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$



However, from the figure we realize that there are infinitely many hyperplanes which satisfy the condition, so which one do we choose? It might be tempting to use *any* hyperplane that satisfies the condition, but misclassification of new data is always a possibility. And intuitively, it would make sense that in order to avoid misclassification, we must not pick a hyperplane too close to either the vectors labelled  $-1$  or  $+1$ , but square in the middle. This extra condition restricts us to a single hyperplane. A useful quantity for us in order to solve this problem is the so-called geometric margin.

**Definition of the Geometric margin:** The geometric margin  $\rho_h(x)$  of a linear classifier  $h : x \mapsto \text{sgn}(w \cdot x + b)$  is the shortest Euclidean distance between  $x$  and the hyperplane  $w \cdot x + b = 0$ .

To find an explicit expression of  $\rho_h(x)$ , consider the orthogonal projection  $\hat{x}$  of the point  $x$  onto the hyperplane. The distance between  $\hat{x}$  and  $x$  is  $\rho_h(x)$  by definition. The direction of the vector  $x - \hat{x}$  is normal to the hyperplane, so it is given by the normal vector  $w$ . Thus we have

$$\rho_h(x) \frac{w}{|w|} = x - \hat{x}$$

Since the point  $\hat{x}$  lies in the hyperplane, it satisfies the equation  $w \cdot \hat{x} + b = 0$ , or  $w \cdot \hat{x} = -b$ . Taking the dot product of  $w$  and both sides of the equation we get

$$\rho_h(\mathbf{x}) \frac{\mathbf{w} \cdot \mathbf{w}}{|\mathbf{w}|} = \mathbf{w} \cdot \mathbf{x} - \mathbf{w} \cdot \hat{\mathbf{x}}$$

$$\rho_h(\mathbf{x}) |\mathbf{w}| = \mathbf{w} \cdot \mathbf{x} + b$$

$$\rho_h(\mathbf{x}) = \frac{\mathbf{w} \cdot \mathbf{x} + b}{|\mathbf{w}|}$$

Since  $\rho_h(\mathbf{x})$  represents a distance it should be a positive function. Since  $\mathbf{w} \cdot \mathbf{x} + b$  may be negative we can take the absolute value and redefine

$$\rho_h(\mathbf{x}) \equiv \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{|\mathbf{w}|} = \frac{y_i(\mathbf{w} \cdot \mathbf{x} + b)}{|\mathbf{w}|}$$

Specifically, we'll use the smallest geometric margin over the training set

$\rho_h \equiv \min_{i \in [m]} \rho_h(\mathbf{x}_i)$ . The SVM solution is the hyperplane  $\mathbf{w} \cdot \mathbf{x} + b = 0$  which maximizes  $\rho_h$  under the constraint [EQUATION HERE]. The distance between this hyperplane and the closest training vector is given by

$$\rho = \max_{\mathbf{w}, b} \min_{i \in [m]} \frac{y_i(\mathbf{w} \cdot \mathbf{x} + b)}{|\mathbf{w}|} : y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 0, \forall i \in [m]$$

An observation that will simplify the optimization problem substantially is that the geometric margin is invariant under multiplication of  $(\mathbf{w}, b)$  by a positive scalar. This freedom allows us to choose  $(\mathbf{w}, b)$  such that we can set  $\min_{i \in [m]} y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$ . Under this constraint the optimization problem is

$$\rho = \max_{\mathbf{w}, b} \frac{1}{|\mathbf{w}|} : y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \forall i \in [m]$$

Since  $|\mathbf{w}|^{-1}$  is maximized when  $|\mathbf{w}|^2$  is minimized, which in contrast is infinitely differentiable everywhere and thus easier to deal with, we can reformulate the optimization problem as

$$\min_{\mathbf{w}, b} \frac{1}{2} |\mathbf{w}|^2 : y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \forall i \in [m]$$

We can solve the optimization problem using the Lagrange multiplier method. The Lagrangian to be minimized is given by

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2}|\mathbf{w}|^2 - \sum_i^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1]$$

where the Lagrange multipliers are  $\alpha_i \geq 0 \ \forall \ i \in [m]$ . Taking the gradient of  $\mathcal{L}$  with respect to  $\mathbf{w}$  and the derivative with respect to  $b$  and setting both to zero we obtain

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_i^m \alpha_i y_i \mathbf{x}_i = 0 \rightarrow \mathbf{w} = \sum_i^m \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_i^m \alpha_i y_i = 0 \rightarrow \sum_i^m \alpha_i y_i = 0$$

Furthermore, according to Karush-Kuhn-Tucker's theorem (see Theorem B.30 in [BOOK REFERENCE HERE]), the following equation is satisfied at the minimum of  $\mathcal{L}$ .

$$\alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1] = 0 \ \forall \ i \in [m] \rightarrow \alpha_i = 0 \vee y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 \ \forall \ i \in [m]$$

This means that all the  $\alpha_i$ 's corresponding to all *but* the vectors  $\mathbf{x}_i$  closest to the hyperplane are zero. These vectors are called support vectors and they satisfy  $\mathbf{w} \cdot \mathbf{x}_i + b = \pm 1$ .

The *dual optimization problem* can be obtained by putting [EQUATIONS HERE] back into the Lagrangian and *maximizing* with respect to  $\alpha$ . I.e

$$\max_{\alpha} \mathcal{L}(\alpha) = \max_{\alpha} \sum_i^m \alpha_i - \frac{1}{2} \sum_i^m \sum_j^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \ : \ \alpha_i \geq 0 \ \forall \ i \in [m], \ \sum_i^m \alpha_i y_i = 0$$

Once the optimal  $\alpha$  has been found, we can calculate  $\mathbf{w}$  using [EQUATION HERE]. If  $\mathbf{x}_j$  is a support vector we can calculate  $b$  using the equation

$$y_j(\mathbf{w} \cdot \mathbf{x}_j + b) = 1 \rightarrow b = y_j - \mathbf{w} \cdot \mathbf{x}_j$$

where we've used that  $(y_j)^{-1} = y_j$ . Once the optimal  $w$  and  $b$  has been found, we classify new data according to the equation

$$h(\mathbf{x}) = \text{sgn}\left(\sum_i^m \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)$$

At this point we can make two important observations: The first observation is that since  $\alpha_i = 0$  except for all the support vectors, the solution depends only on the support vectors. The second observation is that the solution [REFERENCE TO THE DUAL PROBLEM] depends only on inner products between vectors. This observation leads to a major simplification when we'll consider non-linear boundaries which, after all these mathematics for just the simplest form of boundary, would otherwise seem outright frightening.