

## Support vector machines (with offset)

As we had discussed before, maximum margin linear separators can be found by solving the following primal quadratic programming problem

$$\text{(primal)} \quad \min \frac{1}{2} \|\theta\|^2 \quad \text{subject to} \quad y^{(t)}(\theta \cdot x^{(t)} + \theta_0) \geq 1, \quad t = 1, \dots, n \quad (1)$$

If we wish to solve the problem in the dual, i.e., represent the parameters  $\theta$  in terms of examples as  $\theta = \sum_{t=1}^n \alpha_t y^{(t)} x^{(t)}$ , we can solve the dual

$$\text{(dual)} \quad \max \sum_{t=1}^n \alpha_t - \frac{1}{2} \sum_{t=1}^n \sum_{t'=1}^n \alpha_t \alpha_{t'} y^{(t)} y^{(t')} (x^{(t)} \cdot x^{(t')}) \quad (2)$$

$$\text{subject to} \quad \alpha_t \geq 0, \quad \sum_{t=1}^n \alpha_t y^{(t)} = 0 \quad (3)$$

where the additional constraint  $\sum_{t=1}^n \alpha_t y^{(t)} = 0$  pertains to including the offset parameter  $\theta_0$  in the primal. But  $\theta_0$  does not appear anywhere in the dual. How do we set it? After we solve the dual, we can see that  $\hat{\alpha}_t > 0$  for some examples (support vectors) and  $\hat{\alpha}_t = 0$  for others (non-support vectors). Since for support vectors, the classification constraints must be satisfied with equality, we have that when  $\hat{\alpha} > 0$ ,

$$y^{(t)}(\hat{\theta} \cdot x^{(t)} + \hat{\theta}_0) = y^{(t)} \left( \sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) = 1 \quad (4)$$

As  $y^{(t)} \in \{-1, 1\}$  or  $(y^{(t)})^2 = 1$  we can multiply both sides by  $y^{(t)}$  and get

$$\hat{\theta}_0 = y^{(t)} - \left( \sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) \right) \quad (5)$$

This should hold for all support vectors. However, solving the dual quadratic programming problem numerically does introduce errors and relying on any particular constraint may be unwise. Instead, we can simply take the median of all  $\hat{\theta}_0$  estimates from support vector constraints.

## Support vector machines (with errors)

If the labels for training examples contain errors, finding the maximum margin linear classifier may be problematic. The resulting decision boundary is potentially drastically affected by a single mislabeled point. Consider, for example, how the decision boundary changes by the single positive (perhaps mislabeled) point included in Figure 1b).

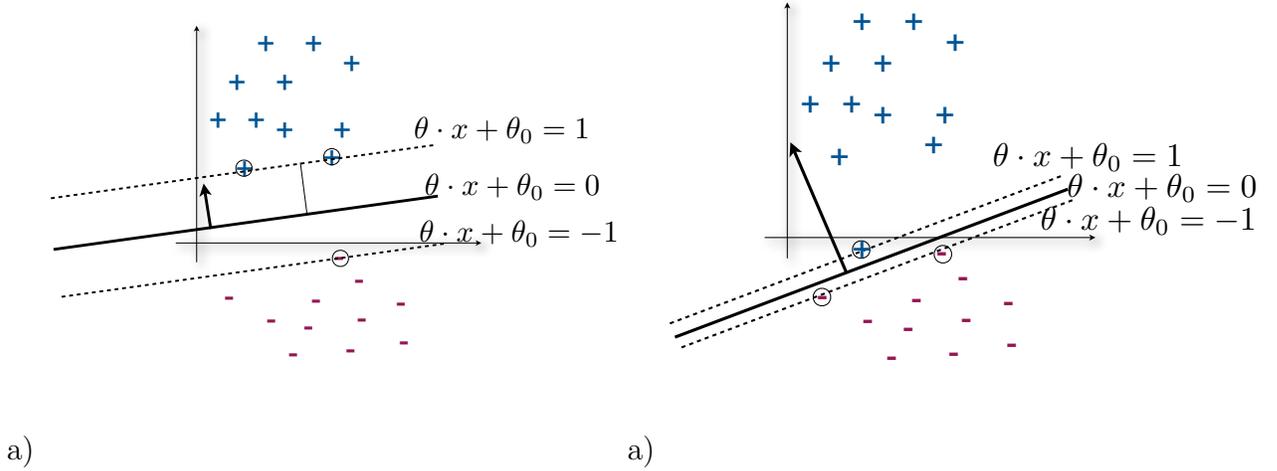


Figure 1: The maximum margin separator is strongly affected by individual points

In order to remedy the situation, we should allow for misclassified points, yet maximize the margin. How do we achieve this? The simplest way is to treat the classification constraints as “soft” constraints rather than hard constraints. In other words, we permit them to be violated but specify the cost of such discrepancies in relation to increasing the margin (decreasing the norm  $\|\theta\|$ ). Specifically, we will add “slack” variables  $\xi_t \geq 0$  to the primal optimization problem as follows

$$\text{(primal)} \quad \min \quad \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum_{t=1}^n \xi_t \quad (6)$$

$$\text{subject to } y^{(t)}(\theta \cdot x^{(t)} + \theta_0) \geq 1 - \xi_t, \quad \xi_t \geq 0, \quad t = 1, \dots, n \quad (7)$$

We now minimize with respect to the parameters  $\theta$ ,  $\theta_0$  as well as the slack variables  $\xi_t \geq 0$ . Note that we had to introduce a parameter<sup>1</sup>  $\lambda$  (regularization parameter) to balance how much we favor increasing the margin over satisfying the classification constraints. Larger values of  $\lambda$  will push the margin boundaries and potentially the decision boundary past the examples. The margin boundaries are still defined as points satisfying  $\theta \cdot x + \theta_0 = 1$  or  $\theta \cdot x + \theta_0 = -1$ . Figure 2 illustrates the choice of  $\lambda$  and the resulting maximum margin solutions when the examples are still linearly separable.

<sup>1</sup>In the literature you will often see a parameter  $C$  multiplying the sum of slack variables instead. This is the same formulation so long as  $C = 1/\lambda$ .

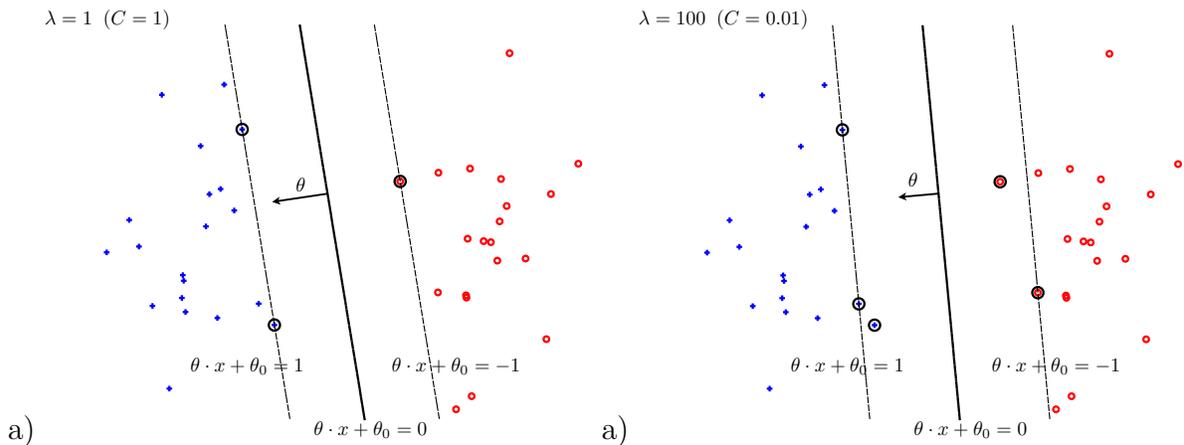


Figure 2: The effect of slack when examples are still linearly separable

The other advantage of the slack variables is that we can now solve problems that are no longer linearly separable. This is illustrated in Figure 2 with different values of the regularization parameter  $\lambda$ .

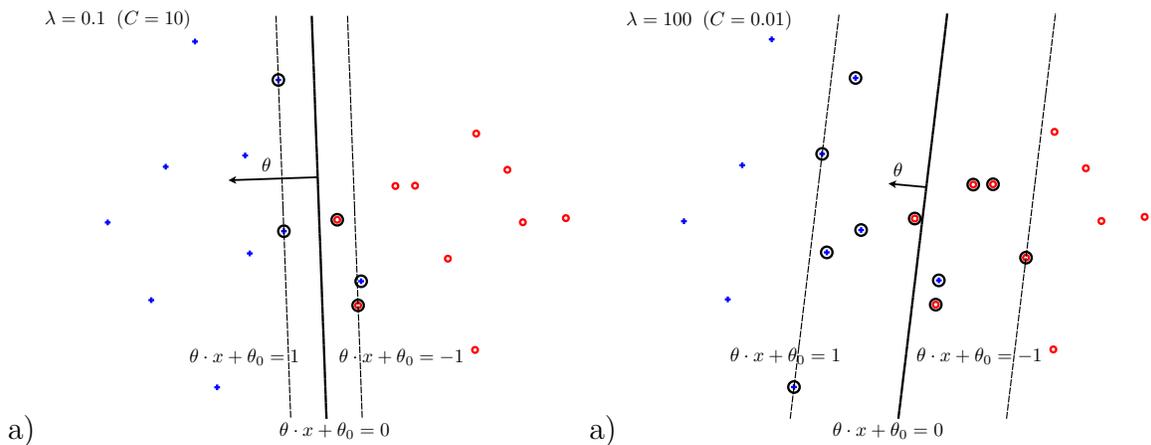


Figure 3: The effect of slack when examples are no longer linearly separable

Our current formulation of the support vector machines can be written in a familiar form. Indeed, the primal version is exactly the same as minimizing

$$\frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum_{t=1}^n \text{Loss}_h(y^{(t)}(\theta \cdot x^{(t)} + \theta_0)) \quad (8)$$

with respect to  $\theta$  and  $\theta_0$ . Here  $\text{Loss}_h(z) = \max\{1 - z, 0\}$  is the hinge loss. To see that

these are equivalent, let us define

$$\xi_t = \text{Loss}_h(y^{(t)}(\theta \cdot x^{(t)} + \theta_0)) = \max\{1 - y^{(t)}(\theta \cdot x^{(t)} + \theta_0), 0\} \quad (9)$$

Clearly,  $\xi_t \geq 0$  and

$$\xi_t \geq 1 - y^{(t)}(\theta \cdot x^{(t)} + \theta_0) \quad \text{or} \quad y^{(t)}(\theta \cdot x^{(t)} + \theta_0) \geq 1 - \xi_t \quad (10)$$

In other words, the slack variables are simply encoding the hinge loss in the primal formulation.

How will the dual formulation change in light of the slack variables in the primal? We will simply limit how large the Lagrange multipliers  $\alpha_t$  can become. In other words, the adversary “gives up” on trying to further satisfy some of the classification constraints. More formally,

$$\text{(dual)} \quad \max \sum_{t=1}^n \alpha_t - \frac{1}{2} \sum_{t=1}^n \sum_{t'=1}^n \alpha_t \alpha_{t'} y^{(t)} y^{(t')} (x^{(t)} \cdot x^{(t')}) \quad (11)$$

$$\text{subject to } 0 \leq \alpha_t \leq 1/\lambda, \quad \sum_{t=1}^n \alpha_t y^{(t)} = 0 \quad (12)$$

Note that the larger the value of  $\lambda$  is, i.e., the more we wish to increase the margin at the expense of the constraints, the smaller the resulting  $\alpha_t$  must be.

We have to be a bit more careful now in reconstructing the offset parameter  $\theta_0$ . Not all the support vectors, i.e., points with non-zero  $\alpha$ 's, will lie exactly on the margin boundaries. Indeed, there will be points that lie on the wrong side of the margin boundaries or points that are misclassified altogether. The *complementary slackness* constraints that characterize the dual solution are now given by

$$\hat{\alpha}_t = 0 \quad \Rightarrow \quad y^{(t)} \left( \sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) \geq 1 \quad (\text{non-support vectors}) \quad (13)$$

$$\hat{\alpha}_t \in (0, 1/\lambda) \quad \Rightarrow \quad y^{(t)} \left( \sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) = 1 \quad (\text{SVs, on the margin}) \quad (14)$$

$$\hat{\alpha}_t = 1/\lambda \quad \Rightarrow \quad y^{(t)} \left( \sum_{t'=1}^n \hat{\alpha}_{t'} y^{(t')} (x^{(t')} \cdot x^{(t)}) + \hat{\theta}_0 \right) \leq 1 \quad (\text{SVs, margin violations}) \quad (15)$$

We can use the points for which  $\alpha_t$  lies in the interior of the possible values, i.e.,  $\alpha_t \in (0, 1/\lambda)$ , to reconstruct  $\hat{\theta}_0$ .

## Non-linear support vector machines – kernels

Let's start by considering how we can use linear classifiers to make non-linear predictions. The easiest way is to first map all the examples  $x \in \mathcal{R}^d$  into a different feature

representation  $\phi(x) \in \mathcal{R}^p$  where typically  $p$  is much larger than  $d$ . We would then simply use a linear classifier on the new (higher dimensional) feature vectors, pretending that they were the original input vectors. There are many ways to create such feature vectors. For example, we can build  $\phi(x)$  by concatenating polynomial terms of the original coordinates. For example, in two dimensions, we could map  $x = [x_1, x_2]^T$  to a five dimensional feature vector

$$\phi(x) = [x_1, x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^T \quad (16)$$

We can then train a “linear” classifier (linear in the new  $\phi$ -coordinates)

$$y = \text{sign}(\theta \cdot \phi(x) + \theta_0) \quad (17)$$

by mapping each training example to the corresponding feature vector. In other words, our training set is now  $S_n^\phi = \{(\phi(x^{(t)}), y^{(t)}), t = 1, \dots, n\}$ . The resulting parameter estimates  $\hat{\theta}, \hat{\theta}_0$  define a linear decision boundary in the  $\phi$ -coordinates but a non-linear boundary in the original  $x$ -coordinates

$$\hat{\theta} \cdot \phi(x) + \hat{\theta}_0 = 0 \Leftrightarrow \hat{\theta}_1x_1 + \hat{\theta}_2x_2 + \hat{\theta}_3\sqrt{2}x_1x_2 + \hat{\theta}_4x_1^2 + \hat{\theta}_5x_2^2 + \hat{\theta}_0 = 0 \quad (18)$$

The non-linear boundary can represent, e.g., an ellipse in the original two dimensional space.

The main problem with the above procedure is that the feature vectors  $\phi(x)$  can become quite high dimensional. For example, if we start with  $x \in \mathcal{R}^d$ , where  $d = 1000$ , then compiling  $\phi(x)$  by concatenating polynomial terms up to the 2nd order would have dimension  $d + d(d+1)/2$  or about 500,000. Using higher order polynomial terms in such situations becomes quickly infeasible. However, it may still be possible to *implicitly* use such feature vectors. If training and prediction problems can be formulated only in terms of inner products between examples, then the relevant computation for us is  $\phi(x) \cdot \phi(x')$ . Depending on how we define  $\phi(x)$ , this computation can be carried out efficiently even if using  $\phi(x)$  explicitly is not. For example, when  $\phi(x) = [x_1, x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^T$ , we see that (check!)

$$\phi(x) \cdot \phi(x') = (x \cdot x') + (x \cdot x')^2 \quad (19)$$

So the inner product is obtained easily from the original input vectors.

One of the main advantages from considering the dual form of support vector machines is that it can be expressed entirely in terms of inner products. For example, in solving for the Lagrange multipliers  $\hat{\alpha}_t$ , the dual problem in Eq.(11) only requires us to evaluate inner products between the training examples, i.e.,  $(x^{(t)} \cdot x^{(t')})$ . Similarly, the parameter  $\hat{\theta}_0$  can be reconstructed from the margin constraints in Eq.(14) based on inner products. Finally, we can predict labels for new examples  $x$  with access only to the inner products between the training examples  $x^{(t)}$  and the new points:

$$y = \text{sign} \left( \hat{\theta} \cdot x + \hat{\theta}_0 \right) = \text{sign} \left( \sum_{t=1}^n \hat{\alpha}_t y^{(t)} (x^{(t)} \cdot x) + \hat{\theta}_0 \right) \quad (20)$$

When we map examples into feature vectors, we replace  $(x^{(t)} \cdot x^{(t')})$  with  $\phi(x^{(t)}) \cdot \phi(x^{(t')})$ . Instead of explicitly using feature vectors, we will focus on evaluating *kernel functions*  $K(x^{(t)}, x^{(t')}) = \phi(x^{(t)}) \cdot \phi(x^{(t')})$ .