

## 1 SVMs: Step by Step

A *decision rule* (or *classifier*) is a function  $r : \mathbb{R}^d \rightarrow \pm 1$  that maps a feature vector (test point) to +1 or -1. The decision rule for SVMs is

$$r(x) = \begin{cases} +1 & \text{if } w \cdot x + \alpha \geq 0, \\ -1 & \text{otherwise,} \end{cases} \quad (1)$$

where  $w \in \mathbb{R}^n$  and  $\alpha \in \mathbb{R}$  are the weights (parameters) of the SVM.

- (a) Draw a figure depicting the line  $\ell = \{u \mid u \cdot w + \alpha = 0\}$  with  $w = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$  and  $\alpha = -25$ . Include in your figure the vector  $w$ , drawn relative to  $\ell$ . Indicate in your figure the region in which data points  $x \in \mathbb{R}^2$  would be classified as 1 vs -1.

We train SVMs by maximizing the distance of the decision boundary from both positive (1) and negative ( $-1$ ) examples. The gap between the decision boundary and the closest positive and negative examples is called the margin. We can express the margin requirement by imposing the constraints

$$y_i(X_i \cdot w + \alpha) \geq c, \quad \forall i \in \{1, \dots, m\}, \quad (2)$$

where  $c$  is taken to be the maximum margin.

(b) What role does  $y_i$  play in Equation 2?

(c) The margin  $c > 0$  can be rescaled to 1 without affecting the decision rule:

$$y_i(X_i \cdot w + \alpha) \geq 1, \quad \forall i \in \{1, \dots, m\}. \quad (3)$$

Why can we rescale the margin to 1?

(d) For which examples  $i$  is  $y_i(X_i \cdot w + \alpha) = 1$ ? What is the geometric interpretation and significance of these examples?

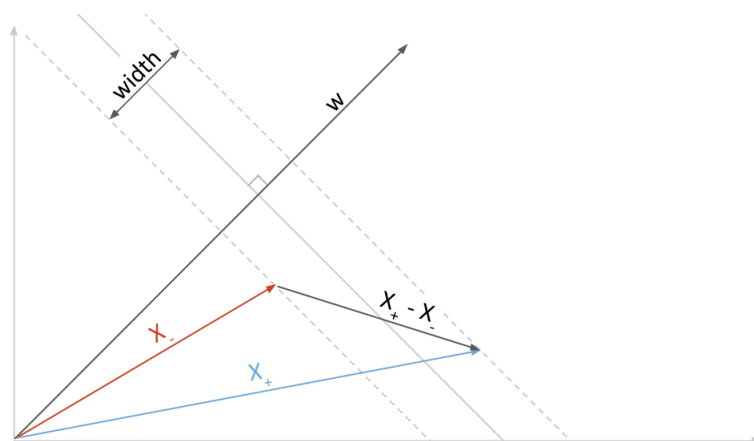


Figure 1: Diagram depicting  $X_+$ ,  $X_-$ ,  $w$ , and the width of the margins.

The constraints we obtained in the previous problem restrict the possible decision boundaries to those which separate the data with some margin that depends on  $w$  and  $b$ . We want the maximum possible margin. We'll need an objective we can optimize to obtain a maximum margin in terms of  $w$  and  $b$ . To obtain this objective, we rewrite Equation 3 as

$$y_i X_i \cdot w \geq 1 - y_i \alpha, \quad i = 1, \dots, m. \quad (4)$$

Let  $X_-$  and  $X_+$  be negative and positive examples **on the margins**, as depicted in Figure 1. The **width** is the distance from the negative margin to the decision boundary plus the distance from the decision boundary to the positive margin, as shown in Figure 1. We can compute the width in terms of  $w$  as follows.

- (e) Write down Equation 4 for  $X_-$ . Divide through by  $|w|$  to obtain a scalar projection of  $X_-$  onto  $\frac{w}{|w|}$ . Do the same for  $X_+$ .
- (f) You now have the distance between the decision boundary to each of the margins. What is the width?
- (g) We want to maximize the width. Using the answer from the previous part, show how this can be written as  $\min_{w,b} \frac{1}{2} |w|^2$ .

## 2 SVM with custom margins

In the lecture, we covered the soft-margin SVM. The objective to be optimized over the training set  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$  is

$$\min_{\mathbf{w}, b, \xi_i} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (5)$$

$$s.t. \quad y_i(\mathbf{w}^\top \mathbf{x}_i - b) \geq 1 - \xi_i \quad \forall i \quad (6)$$

$$\xi_i \geq 0 \quad \forall i \quad (7)$$

In this problem, we are interested in a modified version of the soft-margin SVM where we have a custom margin for each of the  $n$  data points. In the standard soft-margin SVM, we pay a penalty of  $\xi_i$  for each of the data point. In practice, we might not want to treat each training point equally – for example, we might know that some data points are more important than the others.

We formally define the following optimization problem:

$$\min_{\mathbf{w}, b, \{\xi_i\}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \phi_i \xi_i \quad (8)$$

$$s.t. \quad y_i(\mathbf{w}^\top \mathbf{x}_i - b) \geq 1 - \xi_i \quad \forall i \quad (9)$$

$$\xi_i \geq 0 \quad \forall i \quad (10)$$

Note that the only difference is that we have a custom weighting factor  $\phi_i > 0$  for each of the slack variables  $\xi_i$  in the objective function. These  $\phi_i$  are some constants given by the prior knowledge, and thus they can be treated as known constants in the optimization problem. Intuitively, this formulation weights each of the violations ( $\xi_i$ ) differently according to the prior knowledge ( $\phi_i$ ).

- (a) For the standard soft-margin SVM, we have shown that the constrained optimization problem is equal to the following unconstrained optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i - b)) \quad (11)$$

What's the corresponding unconstrained optimization problem for the SVM with custom margins?

(b) As seen in lecture, the dual form of the standard soft-margin SVM is:

$$\max_{\alpha} \quad \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{Q} \alpha \quad (12)$$

$$s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad (13)$$

$$0 \leq \alpha_i \leq C \quad i = 1, \dots, n \quad (14)$$

where  $\mathbf{Q} = (\text{diag } \mathbf{y}) \mathbf{X} \mathbf{X}^T (\text{diag } \mathbf{y})$ .

What's the dual form of the SVM with custom margins? To start, we provide you the Lagrangian, which is given by

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\mathbf{w}^T \mathbf{x}_i - b) + \sum_{i=1}^n (C \phi_i - \alpha_i - \beta_i) \xi_i \quad (15)$$