

# CS 181 Spring 2016 Section 5 Notes

## (Support Vector Machines)

### 1 Motivation

The idea for Support Vector Machines is that, for all the linear hyperplanes that exist, we want one that will create the largest space, or “margin”, with the data. This gives us more insurance about our classifications being correct. To define the margin, we consider a hyperplane of the form

$$\mathbf{w}^T \mathbf{x} + b = 0$$

What is the perpendicular distance from a data point  $\mathbf{x}_n$  to the decision boundary  $y_w(x)$ ?

For two points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  on the hyperplane, consider the projection with  $\mathbf{w}$ :

$$\mathbf{w}(\mathbf{x}_1 - \mathbf{x}_2) = \mathbf{w}\mathbf{x}_1 - \mathbf{w}\mathbf{x}_2 = -b - (-b) = 0$$

Therefore,  $\mathbf{w}$  is orthogonal to the hyperplane. So to get the distance between an arbitrary data point  $\mathbf{x}$ , we just need the length of the vector in the direction of  $\mathbf{w}$  between the point and the vector. We let  $r$  signify the distance between a point and the hyperplane. Then  $\mathbf{x}_\perp$  is defined as the vector on the hyperplane to satisfy the equation

$$\mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|_2} = \mathbf{x}_n$$

Left multiply by  $\mathbf{w}^T$ :

$$\mathbf{w}^T \mathbf{x}_\perp + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|_2} = \mathbf{w}^T \mathbf{x}_n \Rightarrow r = \frac{\mathbf{w}^T \mathbf{x}_n + b}{\|\mathbf{w}\|_2}$$

$r$  then gives the distance between a point and the hyperplane. To make sure it's positive, we multiply by the label. Therefore, the margin of the dataset is the minimum such margin over all data:

$$\min_n \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$$

What is the maximization problem for support vector machines?

We want the  $\mathbf{w}$  and  $b$  that maximize the margin:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_n t_n(\mathbf{w}^T \mathbf{x}_n + b)$$

We can observe that  $\mathbf{w}$  and  $b$  are invariant to changes of scale; therefore, without loss of generality, we can set

$$\min_n t_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$$

So our optimization becomes

$$\arg \max_{w,b} \frac{1}{\|w\|_2} \text{ s.t. } \forall n \ t_n(w^T x_n + b) \geq 1$$

What is the corresponding minimization problem for support vector machines?

Explain (at a high level) why the minimization problem for support vector machines is equivalent to the max-margin problem

We can invert  $w$  to change the max to a min:

$$\arg \min_{w,b} \frac{1}{2} \|w\|_2^2 \text{ s.t. } \forall n \ t_n(w^T x_n + b) \geq 1$$

Minimizing  $\frac{1}{2} \|w\|^2$  is equivalent to maximizing  $\frac{1}{\|w\|}$  because  $\|w\| \geq 0$ . Furthermore, we note that the normalized orthogonal distance from a point to the decision boundary is invariant under scalar multiplication. To see this, we have

$$\frac{t_n(w^T \Phi(x_n) + b)}{\|w\|^2} = \frac{\beta}{\beta} \cdot \frac{t_n(w^T \Phi(x_n) + b)}{\|w\|^2} = \frac{t_n(\beta w)^T \Phi(x_n) + (\beta b)}{\|\beta w\|^2}$$

Thus, since the data is linearly separable, so there exists a decision boundary with non-zero, positive margin for each example, we do not lose any generality in imposing the restraint  $t_n(w^T \Phi(x_n) + b) \geq 1$  (because we can just scale  $w$  until the minimal value is  $\geq 1$ ).

## 2 Dual Form

Use Lagrange multipliers to convert the new minimization problem into dual form.

Our original optimization problem, which looks for the parameters that maximize the margin, is

$$w^*, b^* = \arg \min_{w,b} \left\{ \frac{1}{2} \|w\|_2^2 \right\} \text{ s.t. } t_n \cdot (w^T \phi(x_n) + b) \geq 1, \quad \forall n \in \{1, \dots, N\}. \quad (1)$$

We introduce *Lagrange multipliers*,  $\alpha_1, \dots, \alpha_N \geq 0$ , one for each inequality in Equation 1, i.e., one per datum, to obtain the Lagrangian function (Bishop appendix E does a good job of explaining Lagrange multipliers):

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 - \sum_{n=1}^N \alpha_n (t_n \cdot (w^T \phi(x_n) + b) - 1) \quad (2)$$

Our ideal values are

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \max_{\alpha \geq 0} L(\mathbf{w}, b, \alpha)$$

Since our objective is quadratic, and we have linear constraints, we can (by Slater's condition) change the order to derive:

$$\max_{\alpha \geq 0} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)$$

Taking derivatives, we see

$$\frac{\partial L}{\partial w_j} = w_j - \sum_{n=1}^N \alpha_n t_n \phi(x_n) = 0 \Rightarrow \mathbf{w}^* = \sum_{n=1}^N \alpha_n t_n \phi(x_n)$$

$$\frac{\partial L}{\partial b} = - \sum_{n=1}^N \alpha_n t_n = 0 \Rightarrow \sum_{n=1}^N \alpha_n t_n = 0$$

Plugging these optimal values into our Lagrange multiplier, we get:

$$\begin{aligned} L(\mathbf{w}, b, \alpha) &= \frac{1}{2} \left( \sum_{i=1}^N \alpha_i t_i \phi(x_i) \right)^2 - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j t_i t_j \phi(x_i)^T \phi(x_j) - b \sum_{i=1}^N \alpha_i t_i + \sum_{i=1}^N \alpha_i \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j t_i t_j \phi(x_i)^T \phi(x_j) \end{aligned}$$

We then solve for  $\alpha$ , maximizing  $L$ , conditioned on  $\sum_{n=1}^N \alpha_n t_n = 0$ . If desired, we can then plug in the optimal  $\alpha$  to obtain  $\mathbf{w}^*$  and  $b^*$ .

### 3 Why bother with the dual form?

Solving the primal form, we obtain the optimal  $\mathbf{w}$ , but know nothing about the  $\alpha_n$ . To classify a point, we must compute  $\mathbf{w}^T \phi(x)$ , which is expensive if the  $J$  feature dimension, produced by  $\phi$ , is large.

Solving the dual problem, we obtain  $\alpha_n$  (where  $\alpha_n = 0$  for all but a few points, i.e. the support vectors). To classify a point, we compute  $\sum_{n=1}^N \alpha_n^* t_n K(x_n, x) + b^*$  which is efficient if there are few support vectors. Further, now that we have a scalar product involving only the data vectors, we can apply the **kernel trick**.

### 4 Kernel trick

Given features  $\phi(x)$  for each data point  $x$ , computing the features  $\phi$  may be expensive if the dimensionality is large. It may be more efficient to directly compute the dot product  $\phi(x)^T \phi(x')$ . This is known as the kernel trick:

$$K(x, x') = \phi(x)^\top \phi(x') \quad (3)$$

As long as the **kernel** (or **similarity measure**)  $K$  is a valid inner product, the dual SVM problem can thus be solved without actually computing  $\phi$ .

## 5 Practice Problems

### 1. (Berkeley, Fall '11)

Suppose that we train two SVM's, the first containing all of the sample data points and the second trained on a data set constructed by removing some of the support vectors from the first set. How does the size of the optimal margin change between the first and second SVM?

### 2. Weights of Hard Margin SVM (MIT 6.867, Fall '12)

Consider a hard-margin SVM classification scenario where we have linearly separable data  $\{\phi_n, t_n\}_{n=1}^N$ , where  $t_n \in \{-1, 1\}$ . As usual, for a new data-point  $\phi_{new}$ , we predict 1 if  $w^\top \phi_{new} + b \geq 0$ , and -1 otherwise. Now, assume  $N = 4$ , and that our data is  $\phi_1 = \langle 1, 1 \rangle$ ,  $\phi_2 = \langle 2, 2 \rangle$ ,  $\phi_3 = \langle -1.5, -1.5 \rangle$ ,  $\phi_4 = \langle 4, 4 \rangle$ . Show that for any labeling  $t_1, \dots, t_4$  of our 4 data-points, the  $w^*$  we learn by optimizing the hard-margin SVM criterion will have  $w_1^* = w_2^*$ .

### 3. SVM with Only Positive Training Examples (MIT 6.867, Fall '12)

Suppose we have data  $D = \{\phi_n, t_n\}_{n=1}^N$ , where  $t_n \in \{-1, 1\}$ , and that we would like to learn a linear classification boundary by optimizing an SVM-like criterion that completely ignores the negative training examples. In particular, letting  $D^+ = \{\phi_n \in D \mid t_n = 1\}$  (i.e., the set of positive training examples), we will look for

$$\arg \min_w \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad w^\top \phi_i \geq 1, \forall \phi_i \in D^+. \quad (4)$$

Note that the above optimization problem does not include an offset parameter  $b$ !

- (a) If instead our optimization problem above *did* include an offset parameter  $b$  (i.e., we minimized  $\|w\|^2$  subject to  $w^\top \phi_i + b \geq 1, \forall \phi_i \in D^+$ ), what would the minimizing weight-vector  $w^*$  be?
- (b) If we've found a  $w^*$  that minimizes Equation (4), what is the value of  $\min_{\phi_i \in D^+} w^{*\top} \phi_i$ ?
- (c) Again, assuming we've found a  $w^*$  that minimizes Equation (4), suppose that for a new data-point  $\phi_{new}$  we predict 1 if  $w^{*\top} \phi_{new} \geq (\min_{\phi_i \in D^+} w^{*\top} \phi_i) - \epsilon$  for some small  $\epsilon > 0$ , and -1 otherwise. Will this decision rule guarantee that all the training examples in  $D$  (both positive and negative) are classified correctly?

#### 4. SVM with Three Points (MIT 6.867, Fall '12)

Consider the following dataset consisting of three points on the real line

$$(x_1 = -1, t_1 = 1), (x_2 = 0, t_2 = -1), (x_3 = 1, t_3 = 1),$$

which we will attempt to separate with a linear hyperplane through feature-space *that goes through the origin*. That is, our discriminant function will be  $w^\top \phi(x) \geq 0$ , with no offset term  $b$ . The primal form of this problem is:

$$\begin{aligned} \text{Minimize} \quad & \frac{1}{2} \|w\|^2 + C \sum_{n=1}^3 \xi_n \\ \text{s.t.} \quad & t_n (w^\top \phi(x_n)) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad i = 1, 2, 3 \end{aligned}$$

The dual form is:

$$\begin{aligned} \text{Maximize} \quad & \sum_{n=1}^3 a_n - \frac{1}{2} \sum_{n=1}^3 \sum_{m=1}^3 a_n a_m t_n t_m k(x_n, x_m) \\ \text{s.t.} \quad & 0 \leq a_n \leq C, \quad i = 1, 2, 3 \end{aligned}$$

- Suppose our kernel function is  $k(x, x') = 1 + |xx'|$ , where  $|\cdot|$  is the absolute value. What feature mapping  $\phi$  results in this kernel function?
- Using the feature mapping  $\phi$  from your answer to the previous question, are the three training examples linearly separable in feature space?
- Consider any pair of kernels  $k_1(x, x')$  and  $k_2(x, x')$  such that the training points are linearly separable with  $k_1$  but not with  $k_2$ . Are the data linearly separable if we use the kernel  $k(x, x') = k_1(x, x') + k_2(x, x')$ ? Justify your answer.
- Using the kernel  $k(x, x') = 1 + |xx'|$  that we had in part (a), express the value of  $w^\top \phi(x_2)$  in terms of the  $a_n$  (i.e., in dual form).
- If in our optimization we set  $C < 1$ , will we necessarily have  $\xi_2 > 0$  (i.e., will the slack variable for the second training example exceed 0)? (Hint: you will probably want to use the expression you derived in the previous question for  $w^\top \phi(x_2)$  in answering this).

5. (MIT 6.867, Fall '12)

Suppose that  $K_1$  and  $K_2$  are valid kernel. Recall that a valid kernel has the matrix  $(K)_{ij} = K(x_i, x_j)$  is a positive semi-definite matrix for any finite set of examples  $x_1, x_2, \dots, x_n$ . Show that  $K(x_i, x_j) = K_1(x_i, x_j) + K_2(x_i, x_j)$  is a valid kernel if  $K_1, K_2$  are both valid kernels.

6. (Berkeley, Fall '11)

Recall that the equation of an ellipse in the 2-dimensional plane is  $c(x_1 - a)^2 + d(x_2 - b)^2 - 1 = 0$ . Show that an SVM using the polynomial kernel of degree 2,  $K(x, z) = (1 + x \cdot z)^2$  is equivalent to a linear SVM in the feature space  $(x_1^2, x_2^2, x_1, x_2, x_1x_2, 1)$ . This shows that SVM's with this kernel can separate any elliptic region from the rest of the plane.