

## Recitation 6

Lecturer: Regev Schweiger

Scribe: Regev Schweiger

## 6.1 SVM - Realizable case

In the lecture we saw the following optimization problem, for a maximum margin classifier.

$$\begin{aligned} \min_{\mathbf{w}, b} & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} & y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n = 1, \dots, N \end{aligned}$$

where  $\mathbf{w} \in \mathbb{R}^d$  is the weight vector,  $b \in \mathbb{R}$  is the bias, and  $(\mathbf{x}_n, y_n)$  are the examples and  $\mathbf{x}_n \in \mathbb{R}^d$  and  $y_n \in \{+1, -1\}$ .

The first step is to write the Lagrangian. In general, for a program

$$\begin{aligned} \min_{\mathbf{z}} & f(\mathbf{z}) \\ \text{s.t.} & g_i(\mathbf{z}) \leq 0 \quad \forall i = 1, \dots, N \end{aligned}$$

the Lagrangian is

$$L(\mathbf{z}, \boldsymbol{\alpha}) = f(\mathbf{z}) + \sum_{i=1}^N \alpha_i g_i(\mathbf{z})$$

where  $\boldsymbol{\alpha}$  are called the *Lagrangian multipliers*. The KKT conditions (on which we do not elaborate here) tell us that that if  $\boldsymbol{\alpha}^*$  is a solution for the dual program:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \min_{\mathbf{z}} & L(\mathbf{z}, \boldsymbol{\alpha}) \\ \text{s.t.} & \alpha_i \geq 0 \quad \forall i = 1, \dots, N \end{aligned}$$

then  $\mathbf{z}^* = \arg \min_{\mathbf{z}} L(\mathbf{z}, \boldsymbol{\alpha}^*)$  is a solution to the original program (for all the cases which we will consider).

For our SVM program we get

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n(\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

The first step is to assumed  $\boldsymbol{\alpha}$  is fixed, and minimize over  $\mathbf{w}$  and  $b$ . We now take the derivative of  $L$  and equate it with zero to minimize over  $\mathbf{w}$  and  $b$ .

$$\nabla_{\mathbf{w}}L = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0 \implies \mathbf{w}^*(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

this gives us a way to compute the  $\mathbf{w}$  that achieves the minimal point, given  $\boldsymbol{\alpha}$ . We call this the  $\mathbf{w}$ -constraint. For  $b$  we have

$$\frac{d}{db}L = - \sum_{n=1}^N \alpha_n y_n = 0 \implies \sum_{n=1}^N \alpha_n y_n = 0$$

We call this the  $b$ -constraint. This effectively tells us that there are two classes of  $\boldsymbol{\alpha}$ , and that the behavior of the Lagrangian's minimal point differs between them. If  $\sum_{n=1}^N \alpha_n y_n \neq 0$ , then there is no minimal point; we can take arbitrarily large (or small, depending on the sign of  $\sum_{n=1}^N \alpha_n y_n$ ) values of  $b$ . Therefore, the minimum value (technically, infimum), is  $-\infty$ . However, when  $\sum_{n=1}^N \alpha_n y_n = 0$ , then the value of  $b$  doesn't matter, so there is an finite minimum. This stems from the fact that the Lagrangian is a linear function of  $b$ . Since we are interested, at the next step, at the maximum over all of these values, we are not interested in the case  $L(\mathbf{w}^*, b^*, \boldsymbol{\alpha}) = -\infty$ , so we limit ourselves only to the case  $\sum_{n=1}^N \alpha_n y_n = 0$ .

Plugging the constraints back in  $L$  we have

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\alpha}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \underbrace{\mathbf{w}^T \left( \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right)}_{\mathbf{w}} - \underbrace{b \left( \sum_{n=1}^N \alpha_n y_n \right)}_0 + \left( \sum_{n=1}^N \alpha_n \right) \\ &= -\frac{1}{2} \mathbf{w}^T \mathbf{w} + \left( \sum_{n=1}^N \alpha_n \right) \\ &= -\frac{1}{2} \left( \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^T \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right) + \sum_{n=1}^N \alpha_n \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{n=1}^N \alpha_n \end{aligned}$$

where we have the constraints  $\sum_{n=1}^N \alpha_n y_n = 0$  and  $\forall n$  we have  $\alpha_n \geq 0$ .

Formally, the dual problem is

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha}) = \\ \max_{\boldsymbol{\alpha}} L(\mathbf{w}^*(\boldsymbol{\alpha}), b^*(\boldsymbol{\alpha}), \boldsymbol{\alpha}) = \min_{\boldsymbol{\alpha}} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{n=1}^N \alpha_n \\ \text{s.t. } \sum_{n=1}^N \alpha_n y_n = 0 \\ \forall n \quad \alpha_n \geq 0 \end{aligned}$$

This is an instance of quadratic programming, for which there are efficient algorithms. Suppose we solved this, and got a solution  $\boldsymbol{\alpha}$ . How do we get the solution for the original problems?

For  $\mathbf{w}^*$ , recall we have the  $w$ -constraint  $\mathbf{w}^*(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ , which gives us an explicit formula of  $\mathbf{w}$  as a function of the Lagrange multipliers.

When  $\alpha_n^* > 0$  (for a specific  $n$ ), this means the constraint  $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$  must be satisfied (otherwise,  $\alpha_n = 0$  would give a smaller solution when trying to minimize the Lagrangian). Therefore, the support vectors are those with  $\alpha_n^* > 0$ . This allows us also to get the solution for  $b^*$  - choose an  $n$  for which  $\alpha_n^* > 0$ ; then  $b^* = y_n - (\mathbf{w}^*)^T \mathbf{x}_n$ .

## 6.2 Unrealizable Case

We add slack variables  $\xi_n$  to ensure feasibility. We have,

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ \text{s.t. } y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \quad \forall n = 1, \dots, N \\ \xi_n \geq 0 \end{aligned}$$

We can now write the Lagrangian

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \mathbf{r}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (y_n(\mathbf{w}^T \mathbf{x}_n + b) - 1 + \xi_n) - \sum_{n=1}^N r_n \xi_n$$

Following the same derivation, we now take the derivatives

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0 \quad \implies \quad \mathbf{w}^* = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

identically as before. For  $b$  we have

$$\frac{d}{db}L = -\sum_{n=1}^N \alpha_n y_n = 0 \implies \sum_{n=1}^N \alpha_n y_n = 0$$

also as before. For  $\xi_n$  we have

$$\frac{d}{d\xi_n}L = C - \alpha_n - r_n = 0 \implies \alpha_n = C - r_n$$

Substituting the constraints in  $L$ , we get

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\alpha}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \mathbf{w}^T \underbrace{\left( \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right)}_{\mathbf{w}} - b \underbrace{\left( \sum_{n=1}^N \alpha_n y_n \right)}_0 + \left( \sum_{n=1}^N \alpha_n \right) + \sum_{n=1}^N \xi_n \underbrace{(C - \alpha_n - r_n)}_0 \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{n=1}^N \alpha_n \end{aligned}$$

identically to before. The only difference is that now we have two additional constraints,  $r_n \geq 0$  and  $\alpha_n = C - r_n$ . Since  $r_n$  does not appear in the optimization, we can drop it, and join the two constraints to  $\alpha_n \leq C$ . (For any solution of  $\alpha_n$  we can set  $r_n = C - \alpha_n$ .)

Note that when we have an error in classification or in the margin, then  $\xi_n > 0$  and therefore  $r_n = 0$ , which implies that  $\alpha_n = C$ .

If  $C > \alpha_n > 0$ , this means as before that  $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$  and  $\xi_n = 0$ , and thus  $\mathbf{x}_n$  is a support vector.