

# Miniproject: Documentation for Support Vector Machines

October 24, 2013

## 1 Sequential Minimal Optimization

In this note, we describe the sequential minimal optimization (SMO) algorithm to solve the soft margin support vector machine (SVM) binary classification problem, which is to be implemented as part of the miniproject. In order to understand our arguments here, the reader is advised to study the SVM chapter of the course notes, in particular the “Solving the Support Vector Machine” section (the section “Lagrange Multipliers and Lagrangian Duality” is not required). The SMO algorithm is due to Platt [?], but we make use of algorithmic improvements proposed by Keerthi *et.al.* [?]. The interested reader is encouraged to study these references as well.

The notation we employ here is the same as in the SVM chapter of the course notes. SMO solves the dual of the soft margin SVM problem:

$$\begin{aligned} \min_{\alpha} \left\{ \Phi(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j t_i t_j K_{ij} - \sum_{i=1}^n \alpha_i \right\}, \\ \text{subj. to } \alpha_i \in [0, C], i = 1, \dots, n, \quad \sum_{i=1}^n t_i \alpha_i = 0. \end{aligned} \quad (1)$$

Here,  $\mathbf{K} = [K_{ij}] = [K(\mathbf{x}_i, \mathbf{x}_j)]$  is the kernel matrix for the training dataset  $\mathcal{D} = \{(\mathbf{x}_i, t_i) \mid i = 1, \dots, n\}$ , where  $t_i \in \{-1, +1\}$ . SMO is essentially a constrained coordinate descent algorithm. In each iteration,  $\Phi$  is minimized w.r.t.  $\alpha_i$  and  $\alpha_j$  for some  $i \neq j$ . Note that due to the equality constraint, we have to update at least two variables in order to stay in the feasible set, so SMO is “minimal” in that sense. There are two questions we have to deal with:

- How to pick  $\{i, j\}$  for the next update? When are we done?
- How to minimize  $\Phi$  w.r.t.  $\alpha_i, \alpha_j$ , keeping all other  $\alpha$  coordinates fixed?

The SVM discriminant function is given by the kernel expansion

$$y(\mathbf{x}) = \sum_{i=1}^n \alpha_i t_i K(\mathbf{x}, \mathbf{x}_i) - b,$$

where  $b$  is a bias term, which can be determined from the final  $\alpha$ . **Attention:** In the course notes, the bias parameter has a different sign: it is added to the kernel expansion to get  $y(\mathbf{x})$ , while here, we subtract  $b$ . We do this in order to keep our notation the same as in [?]. This does not matter much, since  $b$  is determined after convergence of SMO only. But take care to rewrite the formula for  $b$  given in the course notes in your code.

Recall from the course notes that a solution of the dual (1) is determined by the KKT conditions:

$$\begin{aligned}\alpha_i = 0 : \quad & t_i y(\mathbf{x}_i) - 1 \geq 0 \\ \alpha_i \in (0, C) : \quad & t_i y(\mathbf{x}_i) - 1 = 0 \\ \alpha_i = C : \quad & t_i y(\mathbf{x}_i) - 1 \leq 0\end{aligned}$$

If we define

$$f_i = \sum_{j=1}^n \alpha_j t_j K_{ij} - t_i, \quad (2)$$

then

$$(f_i - b)t_i = (y(\mathbf{x}_i) - t_i)t_i = t_i y(\mathbf{x}_i) - 1,$$

since  $t_i^2 = 1$ . This means we can rewrite the KKT conditions as follows:

$$\begin{aligned}\alpha_i = 0 : \quad & (f_i - b)t_i \geq 0 \\ \alpha_i \in (0, C) : \quad & (f_i - b)t_i = 0 \\ \alpha_i = C : \quad & (f_i - b)t_i \leq 0\end{aligned}$$

Let us partition  $\{1, \dots, n\}$  into  $I_0, I_+, I_-$ , where

$$\begin{aligned}I_0 &= \{i \mid 0 < \alpha_i < C\}, \quad I_+ = \{i \mid t_i = +1, \alpha_i = 0 \text{ or } t_i = -1, \alpha_i = C\}, \\ I_- &= \{i \mid t_i = -1, \alpha_i = 0 \text{ or } t_i = +1, \alpha_i = C\}.\end{aligned}$$

Moreover, let  $I_{\text{up}} = I_0 \cup I_+$ ,  $I_{\text{low}} = I_0 \cup I_-$ . Then, the KKT conditions can be written as

$$b \leq f_i \quad \forall i \in I_{\text{up}} \quad \text{and} \quad b \geq f_i \quad \forall i \in I_{\text{low}}.$$

If we define

$$b_{\text{up}} = \min \{f_i \mid i \in I_{\text{up}}\}, \quad b_{\text{low}} = \max \{f_i \mid i \in I_{\text{low}}\},$$

the KKT conditions are fulfilled at  $\alpha$  iff  $b_{\text{low}} \leq b_{\text{up}}$  (as we can squeeze  $b$  in between if and only if this inequality holds). For numerical reasons, and also to attain convergence in a reasonable number of iterations, we will be satisfied with approximately fulfilled conditions:

$$b_{\text{low}} \leq b_{\text{up}} + 2\tau, \quad \tau > 0. \quad (3)$$

Here,  $\tau$  is a small positive number, such as  $10^{-8}$ . This is the termination condition for the SMO algorithm. As long as it does not hold, there will always be a violated pair,

$$f_i > f_j + 2\tau, \quad i \in I_{\text{low}}, \quad j \in I_{\text{up}} \quad (4)$$

on which an update has to be done. As a heuristic to minimize the number of updates, we may select the most violated pair ( $f_i = b_{\text{low}}, f_j = b_{\text{up}}$ ). Before discussing the pair selection in detail, we derive the minimization for a chosen pair.

### 1.1 Minimization w.r.t. $\alpha_i, \alpha_j$

Suppose that  $(i, j)$  is a violated pair. The fact which drives SMO is that we can minimize  $\Phi$  analytically w.r.t.  $\alpha_i, \alpha_j$ , keeping all other coordinates of  $\alpha$  fixed. The feasible set of this bivariate minimization is given by the box constraints and the requirement to keep the equality constraint valid. We summarize results at the end of this section.

Denote the old values by  $\tilde{\alpha}_i, \tilde{\alpha}_j$ , the minimizers by  $\alpha_i, \alpha_j$ , and  $\Delta_i = \alpha_i - \tilde{\alpha}_i$ ,  $\Delta_j = \alpha_j - \tilde{\alpha}_j$ . Moreover, let  $\sigma = t_i t_j \in \{-1, +1\}$ . The feasible set is given by

$$\alpha_i, \alpha_j \in [0, C], \quad t_i \Delta_i + t_j \Delta_j = 0.$$

The last equation is obtained by subtracting the equality constraint for the new from that for the old values. Multiplying it by  $t_i$ , we have

$$\Delta_i + \sigma \Delta_j = 0 \quad \Rightarrow \quad \alpha_i = \tilde{\alpha}_i - \sigma(\alpha_j - \tilde{\alpha}_j) = w - \sigma \alpha_j, \quad w := \tilde{\alpha}_i + \sigma \tilde{\alpha}_j.$$

Therefore, the problem boils down to a *univariate* quadratic minimization for  $\alpha_j$  only, subject to  $\alpha_j \in [L, H]$ . In order to work out  $L$  and  $H$ , note that  $\alpha_j \in [0, C]$  and

$$\alpha_i = w - \sigma \alpha_j \in [0, C] \quad \Rightarrow \quad w - C \leq \sigma \alpha_j \leq w.$$

Distinguishing between  $\sigma = +1$  and  $\sigma = -1$ , we obtain

$$L = \max \{0, \sigma w - \mathbf{I}_{\{\sigma=+1\}} C\}, \quad H = \min \{C, \sigma w + \mathbf{I}_{\{\sigma=-1\}} C\}, \quad \sigma w = \tilde{\alpha}_j + \sigma \tilde{\alpha}_i. \quad (5)$$

We solve for  $\alpha_j$  in two steps. First, we determine the unconstrained minimum point, then clip it according to  $[L, H]$ . Details for the first step can be found in [?] (exercise!). If we plug  $\alpha_i = w - \sigma \alpha_j$  into  $\Phi$ , then

$$\frac{\partial \Phi}{\partial \alpha_j} = (K_{ii} + K_{jj} - 2K_{ij})\alpha_j - \sigma(K_{ii} - K_{ij})w - t_j(v_i - v_j) + \sigma - 1.$$

Here,

$$v_i = \sum_{p \neq i, j} \alpha_p t_p K_{ip} = f_i + t_i - \tilde{\alpha}_i t_i K_{ii} - \tilde{\alpha}_j t_j K_{ij}, \quad v_j = f_j + t_j - \tilde{\alpha}_i t_i K_{ij} - \tilde{\alpha}_j t_j K_{jj}.$$

Note that  $f_i$  and  $f_j$  in these equations denote the values for the old values  $\tilde{\alpha}_i, \tilde{\alpha}_j$ . Also,  $\eta := K_{ii} + K_{jj} - 2K_{ij}$  must be positive, since the kernel matrix  $\mathbf{K}$  is positive definite:

$$\eta = (\delta_i - \delta_j)^T \mathbf{K} (\delta_i - \delta_j) > 0.$$

Plugging  $w = \tilde{\alpha}_i + \sigma \tilde{\alpha}_j$  into  $(\partial \Phi)/(\partial \alpha_j) = 0$ , some algebra gives

$$\alpha_j^{\text{unc}} = \tilde{\alpha}_j + \frac{t_j(f_i - f_j)}{\eta} \quad (6)$$

for the unconstrained minimizer. Finally, if  $\alpha_j^{\text{unc}}$  does not lie in  $[L, H]$ , we clip it accordingly:

$$\alpha_j = \left\{ \begin{array}{l|l} \alpha_j^{\text{unc}} & \alpha_j^{\text{unc}} \in [L, H] \\ L & \alpha_j^{\text{unc}} < L \\ H & \alpha_j^{\text{unc}} > H \end{array} \right\}. \quad (7)$$

Finally, in rare cases it can happen that  $\eta$  is zero or very close to zero. In these cases, we evaluate  $\Phi$  (up to a constant) at both ends of the segment  $[L, H]$  and update  $\alpha_j = L$  or  $\alpha_j = H$ , depending on which value is lower. Since the relevant part of  $\Phi$  is given by

$$\Phi = \frac{1}{2} (K_{ii} \alpha_i^2 + K_{jj} \alpha_j^2) + \sigma K_{ij} \alpha_i \alpha_j + t_i \alpha_i v_i + t_j \alpha_j v_j - \alpha_i - \alpha_j + \text{const},$$

we can compute the value  $\Phi_L$  for  $\alpha_j = L$  as

$$\Phi_L = \frac{1}{2} (K_{ii}(L_i)^2 + K_{jj}L^2) + \sigma K_{ij}L_iL + t_iL_i v_i + t_jLv_j - L_i - L, \quad L_i = w - \sigma L, \quad (8)$$

where  $w$ ,  $v_i$  and  $v_j$  are defined above. The value  $\Phi_H$  for  $\alpha_j = H$  is obtained in the same way, with  $L \rightarrow H$  and  $L_i \rightarrow H_i = w - \sigma H$ .

To summarize this section, the minimization w.r.t.  $\alpha_i$ ,  $\alpha_j$  works as follows. First, we determine the feasible range  $[L, H]$  for  $\alpha_j$  (5), and  $\eta = K_{ii} + K_{jj} - 2K_{ij}$ . If  $\eta > 10^{-15}$ , we determine the minimizer  $\alpha_j$  via (6) and (7). Otherwise, we determine  $\Phi_L$  and  $\Phi_H$ , setting  $\alpha_j = L$  if  $\Phi_L < \Phi_H$ , otherwise  $\alpha_j = H$ . Finally, we set  $\alpha_i = w - \sigma\alpha_j$ . Importantly, all these computations are  $O(1)$  and require a single kernel evaluation ( $K_{ij}$ ; we assume that the diagonal of the kernel matrix is precomputed), given that we maintain the values  $[f_p]$  up to date.

## 1.2 Heuristics for choosing violated pair $(i, j)$

Based on [?], we discuss how to select the most violated pair  $(i, j)$ . As mentioned above, we will use the optimality conditions (3) and (4) to check for optimality at  $\alpha$ . For efficiency reasons we maintain a vector  $\mathbf{f}$  with entries as defined in (2). We define new indices  $i_{\text{low}}$  and  $i_{\text{up}}$  such that

$$f_{i_{\text{up}}} = b_{\text{up}} = \min \{f_i \mid i \in I_{\text{up}}\}, \quad f_{i_{\text{low}}} = b_{\text{low}} = \max \{f_i \mid i \in I_{\text{low}}\}.$$

Obviously, the pair  $(i_{\text{low}}, i_{\text{up}})$  is the most violated one. We now use (3) to check for optimality. If (3) holds, then we have converged to the optimal solution and we exit the main loop. Otherwise, we return the pair  $(i_{\text{low}}, i_{\text{up}})$ .

## 2 Pseudocode

Now we are in position to provide the complete pseudocode for the SVM optimization algorithm. The reader should be able to write down a complete program from this pseudocode.

---

**Algorithm 1** SMO Implementation

---

**Input:**

$\mathbf{t}$ : Target vector (entries +1 or -1).

$\mathbf{K}$ : Kernel matrix.

// Initialization

$\boldsymbol{\alpha} = \mathbf{0}$ ;

$\mathbf{f} = -\mathbf{t}$ ;

Compute index sets  $I_{\text{low}}$  and  $I_{\text{up}}$ ;

**while** (1) **do** // Main Loop

$(i, j) = \text{select\_pair}$ ; // Select most violated pair

**if**  $j = -1$  **then**

        break;

**end if**

$\sigma = t_i t_j$ ;

    Compute  $L, H$  from (5);

$\eta = K_{ii} + K_{jj} - 2K_{ij}$ ;

**if**  $\eta > 10^{-15}$  **then**

        Compute the minimum along the direction of the constraint from (6);

        Clip unconstrained minimum to the ends of the line segment according to (7);

**else** // the second derivative is negative

        Compute  $\Phi_H$  and  $\Phi_L$  according to (8);

**if**  $\Phi_L > \Phi_H$  **then**

$\alpha_j^{\text{new}} = H$ ;

**else**

$\alpha_j^{\text{new}} = L$ ;

**end if**

**end if**

$\alpha_i^{\text{new}} \leftarrow \alpha_i + \sigma(\alpha_j - \alpha_j^{\text{new}})$ ; // compute new  $\alpha_i$  from the new  $\alpha_j$

    Update  $\boldsymbol{\alpha}$  vector with the new values  $\alpha_i^{\text{new}}, \alpha_j^{\text{new}}$ ;

$\mathbf{f} \leftarrow \mathbf{f} + t_i(\alpha_i^{\text{new}} - \alpha_i)\mathbf{K}_{\cdot,i} + t_j(\alpha_j^{\text{new}} - \alpha_j)\mathbf{K}_{\cdot,j}$ ;

    //  $\mathbf{K}_{\cdot,i}$  denotes the  $i$ -th column of the kernel matrix  $\mathbf{K}$

    Update sets  $I_{\text{low}}, I_{\text{up}}$ ;

**end while**

---

```

procedure select_pair
// Compute  $i_{\text{up}}$  and  $i_{\text{low}}$ 
 $i_{\text{up}} = \operatorname{argmin}_i \{f_i \mid i \in I_{\text{up}}\};$ 
 $i_{\text{low}} = \operatorname{argmax}_i \{f_i \mid i \in I_{\text{low}}\};$ 
// Check for optimality
if  $f_{i_{\text{low}}} \leq f_{i_{\text{up}}} + 2 \cdot \tau$  then
     $i_{\text{low}} = -1;$ 
     $i_{\text{up}} = -1;$ 
end if
// return the most violated pair
return  $(i_{\text{low}}, i_{\text{up}});$ 
end procedure

```

## 2.1 Hints for Implementation

Apart from the main loop, you do not need any loops in order to implement SMO. In particular, make sure to understand why the update of  $(\alpha_i, \alpha_j)$  for a violated pair  $(i, j)$  requires  $O(1)$  operations only.

The update of  $\mathbf{f} = [f_p]$  requires two columns  $\mathbf{K}_{:,i}$  and  $\mathbf{K}_{:,j}$  of the kernel matrix  $\mathbf{K}$ . If your training dataset has a size of up to 5000 or so, it makes sense to compute the entire kernel matrix  $\mathbf{K}$  up front. Otherwise, you can compute required columns on demand. Advanced implementations maintain a cache of kernel columns. For example, you will find that indices  $i \in I_0$  are frequently part of the most violated pair. However, such optimizations are not required for this project.

Whatever you do, make sure not to use any loops for evaluating the kernel matrix  $\mathbf{K}$  or columns  $\mathbf{K}_{:,i}$  thereof. Here is how you do it. Whether you use the Gaussian or the polynomial kernel, you first compute a matrix of all pairwise inner products or squared norms, then apply the kernel function elementwise to this matrix (elementwise application of a scalar function is supported in `Matlab` and is very efficient). For the Gaussian kernel, suppose that  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is your data matrix ( $n$  patterns,  $d$  attributes). Let  $\mathbf{d} = \operatorname{diag}(\mathbf{X}\mathbf{X}^T) = [\|\mathbf{x}_i\|^2] \in \mathbb{R}^n$ , which in `Matlab` you compute as `d = sum(x.*x,2)`. Then,

$$\mathbf{A} = [(1/2)\|\mathbf{x}_i - \mathbf{x}_j\|^2] = (1/2)\mathbf{d}\mathbf{1}_n^T + (1/2)\mathbf{1}_n\mathbf{d}^T - \mathbf{X}\mathbf{X}^T.$$

And if  $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\tau}{2}\|\mathbf{x}_i - \mathbf{x}_j\|^2}$ , then  $\mathbf{K} = f(\mathbf{A})$ , where  $f(a) = e^{-\tau a}$ . In the same way, you avoid loops when computing  $\mathbf{K}_{:,i}$ .