

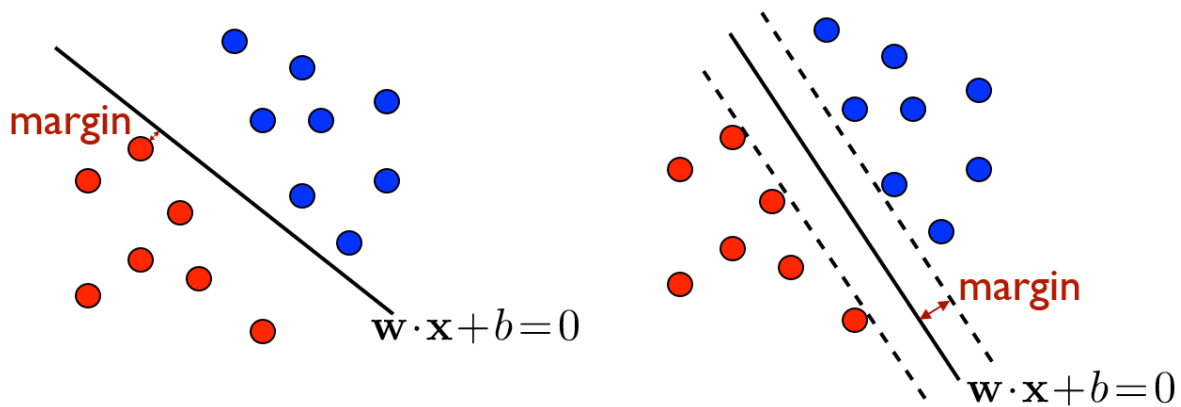
SVM Problem with SVO Solver

By Minghe Ren, Xinqiao Wei

EC500 (HPC) Project

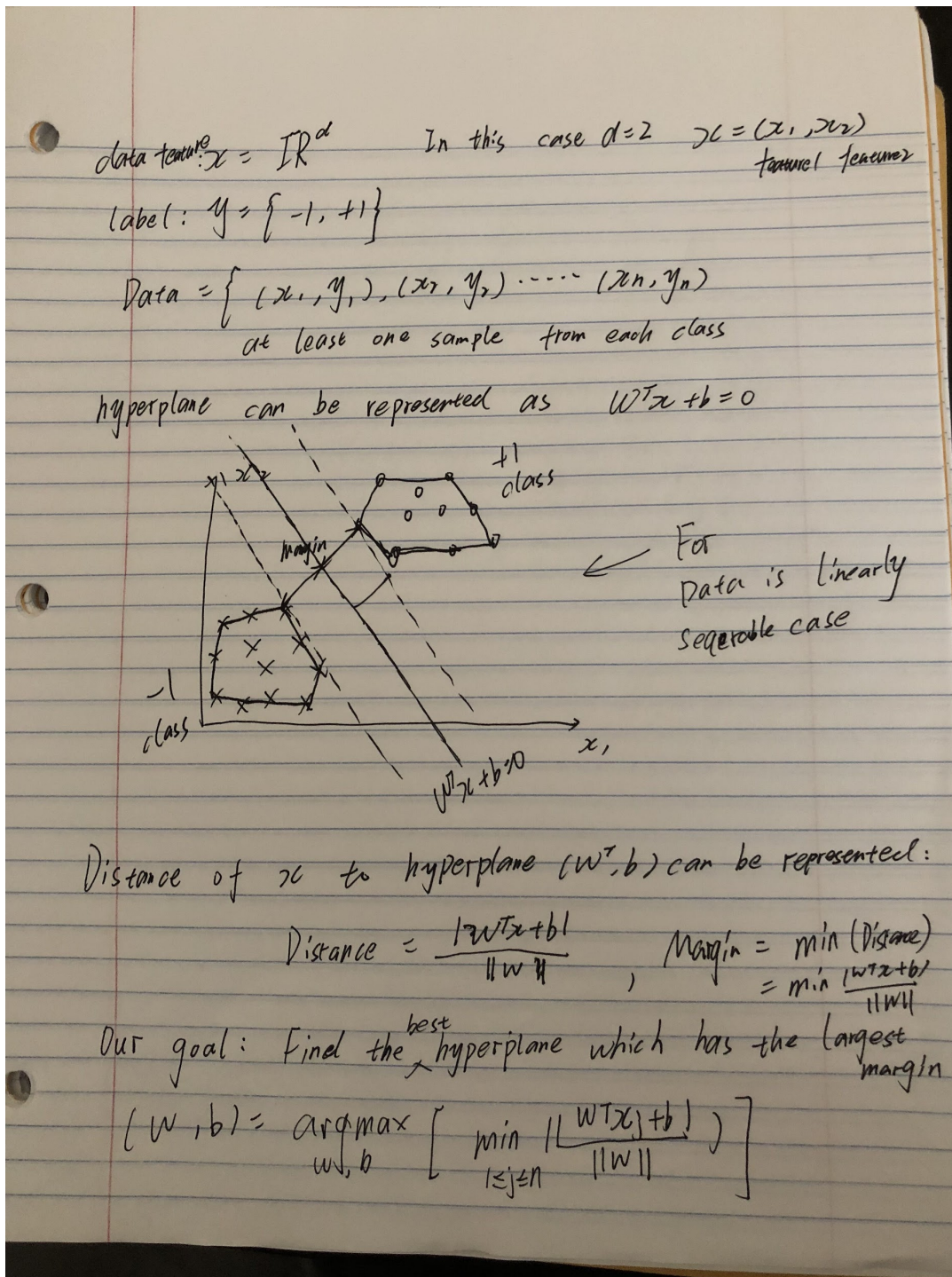
1. SVM problem Introduction

In machine learning algorithm, SVM (support vector machine) is often used in supervised classification. The idea is, giving a data set, the algorithm trying to find the best hyperplane that separating the data. For example, in two dimension case, suppose you have a dataset which has two different labels. And the data points has two features (two dimension). The general idea can be represented as figure below, $W \cdot x + b = 0$ represents the hyperplane.



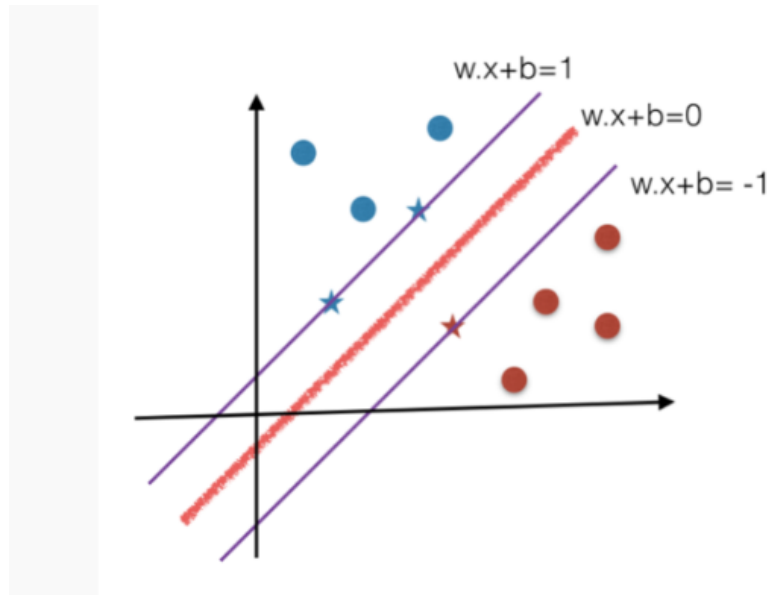
In our project, we decided to write a two dimensional SVM classifier in C++. That is, writing our SVM solver code for two dimensional dataset to find the best w and b to separate the data, which is the hyperplane.

The idea can be shown as:



The hyperplane is the classifier. Above it, it's class +1. Below it, it's class -1.

If we let the $\min |W \cdot x + b| = 1$, we can represent the margin as:



By introducing Lagrangian Multipliers and KKT conditions, we are actually solving this dual optimization problem:

■ **Lagrangian:** for all $w, b, \alpha_i \geq 0$,

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1].$$

■ **KKT conditions:**

$$\begin{aligned} \nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 &\iff \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i. \\ \nabla_b L = - \sum_{i=1}^m \alpha_i y_i = 0 &\iff \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

$$\forall i \in [1, m], \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1] = 0.$$

Substitute expression for w back into Lagrangian

We get
$$\mathcal{L} = \frac{1}{2} \left\| \sum_{j=1}^n \alpha_j y_j x_j \right\|^2 + \sum_{j=1}^n \alpha_j \left(1 - y_j \sum_{k=1}^n \alpha_k y_k x_k^T x_j \right) + b \sum_{j=1}^n \alpha_j$$

$$= \frac{1}{2} \left\| \sum_{j=1}^n \alpha_j y_j x_j \right\|^2 + \sum_{j=1}^n \alpha_j - b \sum_{j=1}^n \alpha_j y_j$$

$$- \sum_{j=1}^n \alpha_j y_j \sum_{k=1}^n \alpha_k y_k (x_j^T x_k)$$

$$- \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k y_j y_k (x_j^T x_k)$$

$$\underbrace{\quad}_{\|w\|^2 = w^T w}$$

$$\mathcal{L} = \frac{1}{2} \|w\|^2 + \sum_{j=1}^n \alpha_j - b \sum_{j=1}^n \alpha_j y_j - \|w\|^2$$

$$= \left(\sum_{j=1}^n \alpha_j \right) - \frac{1}{2} \|w\|^2$$

$$= \sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k y_j y_k (x_j^T x_k)$$

Dual optimization problem

$$\max_{\alpha} \mathcal{L}(\alpha)$$

$$\alpha_j \geq 0$$

$$\sum \alpha_j y_j = 0$$

$$\alpha^{sum} = \arg \max \left(\sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k y_j y_k (x_j^T x_k) \right)$$

$$\alpha_j \geq 0$$

$$\sum_{j=1}^n \alpha_j y_j = 0$$

$$\alpha^{SVM} = \underset{\alpha \geq 0}{\operatorname{argmax}} \left[\sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k y_j y_k (x_j^T x_k) \right]$$

may not be unique. $\sum_{j=1}^n \alpha_j y_j = 0$ $\| \sum_{j=1}^n \alpha_j y_j x_j \|^2$

This Dual Optimization problem is solved via an iterative numerical method, eg - Quadratic Program (QP) - Sequential Minimal Optimization (SMO)

Complementary Conditions

$$\alpha_j^{SVM} [1 - y_j ((W^{SVM})^T x_j + b^{SVM})] = 0$$

Support Vector (SV)

A feature x_j for which $\alpha_j^{SVM} > 0$ is called a support vector

\Rightarrow if x_j is a SV then

$$1 - y_j ((W^{SVM})^T x_j + b^{SVM}) = 0$$

$$\Rightarrow (W^{SVM})^T x_j + b^{SVM} = y_j$$

\Rightarrow all SV's are close to the SVM hyperplane.

Our problem becomes: a dual optimization problem with boundary condition:

$$\alpha^{SVM} = \underset{\alpha \geq 0}{\operatorname{argmin}} \left[\sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k y_j y_k (x_j^T x_k) \right]$$

The condition is $\sum_{j=1}^n \alpha_j y_j = 0$ and $0 \leq \alpha_j \leq C$ (C is hyper parameter). This can be solved by Sequential minimal Optimization.

2. SMO

SMO breaks this problem into a series of smallest possible sub-problems, which are then solved analytically. Because of the linear equality constraint involving the Lagrange multipliers α_j the smallest possible problem involves two such multipliers. Then, for any two multipliers α_1 and α_2 , the constraints are reduced to:

$$0 \leq \alpha_1, \alpha_2 \leq C$$

$$y_1\alpha_1 + y_2\alpha_2 = k$$

This reduced problem can be solved analytically: one needs to find a minimum of a one-dimensional quadratic function. K is the negative of the sum over the rest of terms in the equality constraint, which is fixed in each iteration.

The linear classifier $f(x) = w^T x + b$ can be express as:

$$f(x) = \sum_{i=1}^m \alpha_i y^i \langle x^i, x \rangle + b$$

The KKT condition for this problem:

$$\alpha_i = 0 \rightarrow y^i(w^T x_i + b) \geq 1$$

$$\alpha_i = C \rightarrow y^i(w^T x_i + b) \leq 1$$

$$0 < \alpha_i < C \rightarrow y^i(w^T x_i + b) = 1$$

Next step, optimizing α_i and α_j . First, find the bounds L and H such that $L \leq \alpha_j \leq H$ for the constraint $0 \leq \alpha_j \leq C$:

If $y_i \neq y_j$, $L = \max(0, \alpha_j - \alpha_i)$, $H = \min(C, C + \alpha_j - \alpha_i)$

If $y_i = y_j$, $L = \max(0, \alpha_i + \alpha_j - C)$, $H = \min(C, \alpha_i + \alpha_j)$

The optimal α_j is given by:

$$\alpha_j := \alpha_j - \frac{y_i(E_i - E_j)}{\eta}$$

where $E_k = f(x_k) - Y_k$. It is the error between the SVM output on the k th example and the true label Y_k . And $\eta = 2\langle x_i, x_j \rangle - \langle x_i, x_i \rangle - \langle x_j, x_j \rangle$

If α_j ends up lying outside the bounds L and H , we modify α_j :

$$\alpha_j := \begin{cases} H & \text{if } \alpha_j > H \\ \alpha_j & \text{if } L \leq \alpha_j \leq H \\ L & \text{if } \alpha_j < L \end{cases}$$

Then we solve for α_i :

$$\alpha_i := \alpha_i + y_i y_j (\alpha_j^{\text{old}} - \alpha_j)$$

Select the threshold b satisfying KKT for x_i and x_j . If α_i is not at the bounds then b_1 is valid:

$$b_2 = b - E_i - y_i(\alpha_i - \alpha_i^{\text{old}}) < \langle x_i, x_i \rangle - y_i(\alpha_j - \alpha_j^{\text{old}}) < \langle x_i, x_j \rangle$$

Similarly, if $0 < \alpha_j < C$, b_2 is valid:

$$b_2 = b - E_j - y_j(a_i - a_i^{old}) < x_i, x_j > -y_j(\alpha_j - \alpha_j^{old}) < x_j, x_j >$$

If both b_1 and b_2 are both valid, the value would be the same. If both α s are at the bounds then the $b := (b_1 + b_2)/2$ would meet the KKT condition.

3. Code

Here is our pseudo-code:

```

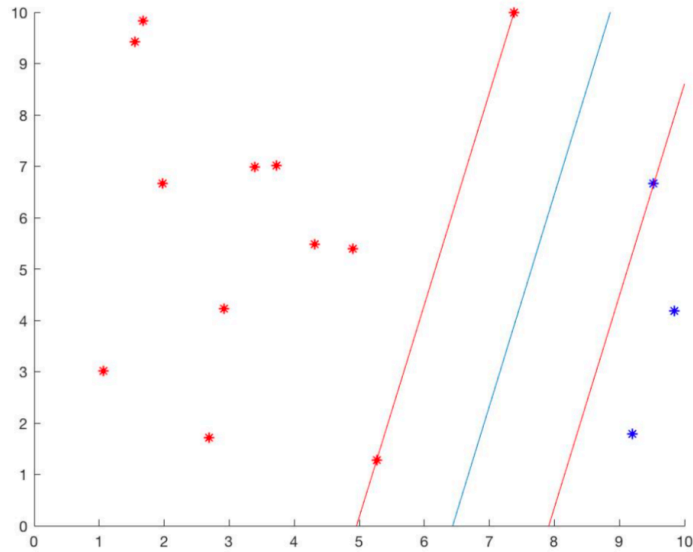
1: Initial:  $\alpha_i = 0, \forall i, b = 0, passes = 0$ 
2: while passes < max_passes do
3:   num_changed_alphas = 0;
4:   for  $i = 1, \dots, m$  do
5:     Calculate  $E_i = f(x^i) - y^i$ 
6:     if  $((y^i E_i < -tol \& \& a_i < C) \parallel (y^i E_i > tol \& \& a_i > 0))$  then
7:       Select  $j \neq i$  randomly.
8:       Calculate  $E_j = f(x^j) - y^j$ 
9:       Save old  $\alpha$ 's:  $\alpha_i^{old} = \alpha_i, \alpha_j^{old} = \alpha_j$ 
10:      Compute L and H.
11:      if  $(L == H)$  then
12:        continue to next  $i$ .
13:      Compute  $\eta$ .
14:      if  $(\eta \geq 0)$  then
15:        continue to next  $i$ .

16:      Compute and clip new value for  $\alpha_j$ 
17:      if  $(|\alpha_j - \alpha_j^{old}| < 10^{-5})$  then
18:        continue to next  $i$ .
19:      Determine value for  $\alpha_i$ 
20:      Compute  $b_1$  and  $b_2$ 
21:      Compute  $b$ 
22:      num_changed_alphas := num_changed_alphas + 1
23:    end if
24:  end for
25:  if (num_changed_alphas == 0) then
26:    passes := passes + 1
27:  else
28:    passes := 0
29: end while

```

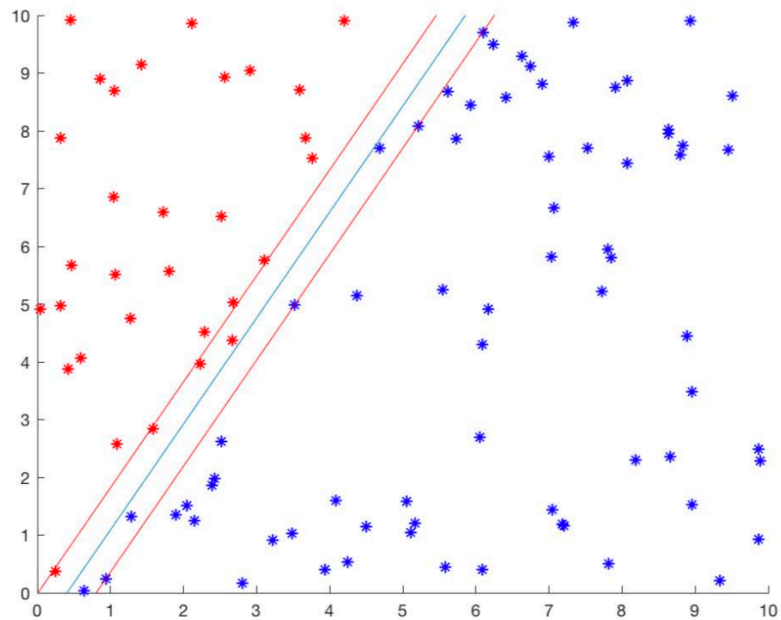

4. Result

For 15 points dataset:



Green line is the hyperplane and the red lines represents the margin. The hyperplane we calculate perfectly separates the data.

For 100 data points:



The green hyperplane also perfectly solve the problem and with some data points living the boundary.

The output of our code shown below, where most of the Laplacian multipliers are zero.

```
The calculated alphas are:  
0 0 0 0 0 0 0 0 0 0 1.6 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1.6 0 0.00401323 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 1.14711 0 0 0 0 0 1.6 0 0 1.6 0 1.6 0 0 0 0 0 0 0 0 0 0 1.6 1.6 0 0 0 0 0 0 0 0 0 0 0  
0.0158122 0 0 0.433067 0 0 0 0 0  
The calculated W is:  
2.51426 -1.37197  
  
The calculated b is:  
-1.02552  
  
The calculated values of support vector are:  
-0.923299 0.374325 1.00948 1 0.520782 -0.849337 -0.325389 -0.949446 0.184743 1.00134 1 [Finished in 1.0s]
```

For 1000 data points, our code can't converge due to some unknown reason.

5. Future work

- a.** Test on a much bigger dataset like 100,000 points and check the performance.
- b.** Different algorithms to achieve SVM like Quadratic Programming.
- c.** Test on non-linear separable dataset and non-separable dataset.
- d.** Try to parallelize our code

6. Reference:

<http://cs229.stanford.edu/materials/smo.pdf>

Our Github link:

<https://github.com/rmhsawyer/EC500-HPC-Final-Project>

