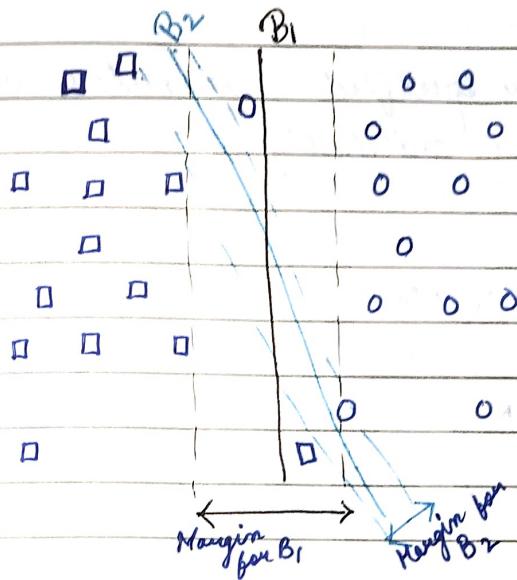


Date: \_\_\_\_\_

## DECISION BOUNDARY OF SVM FOR NONSEPARABLE CASE



In above figure we see that the classifier  $B_2$  successfully classifies all the sample but the decision boundary is small. But the other classifier  $B_1$  makes a slight error while classification but it has a greater decision margin.

So we add a slack variable to tolerate this small error. Our inequality constraint becomes:

$$w \cdot x_i + b \geq 1 - \xi_i \quad \text{if } y_i = 1$$

$$w \cdot x_i + b \leq -1 + \xi_i \quad \text{if } y_i = -1$$

where  $\xi_i > 0$  and  $\xi_i$  is the slack variable

Objective func:

$$f(w) = \frac{\|w\|^2}{2} + C \left( \sum_{i=1}^N \xi_i \right)$$

where  $C$  ~~and~~ <sup>is an</sup> user specified parameter which represents the penalty for misclassification

Lagrangian dual problem formulation

$$L_p = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \lambda_i \{ y_i (w \cdot x_i + b) - 1 + \xi_i \} - \sum_{i=1}^N \mu_i \xi_i \quad (1)$$

where first two terms are the objective function to be minimized, the third term represents the inequality constraints associated with the slack variable, and the last term is the result of the non-negativity requirement on the value of  $\xi_i$ 's.

KKT condition

$$\xi_i \geq 0, \lambda_i \geq 0, \mu_i \geq 0$$

$$\lambda_i \{ y_i (w \cdot x_i + b) - 1 + \xi_i \} = 0$$

$$\mu_i \xi_i = 0$$

First order derivative of  $L$  w.r.t  $w, b, \xi_i$  to zero.

$$\frac{\partial L}{\partial w_j} = w_j - \sum_{i=1}^N \lambda_i y_i x_{ij} = 0$$

$$\Rightarrow w_j = \sum_{i=1}^N \lambda_i y_i x_{ij}$$
(2)

$$\frac{\partial L}{\partial \xi_i} = (-\lambda_i - \mu_i) = 0 \Rightarrow \lambda_i + \mu_i = C \quad (3)$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^N \lambda_i y_i = 0 \Rightarrow \sum_{i=1}^N \lambda_i y_i = 0 \quad (4)$$

Putting (5), (4), (3) in (1)

$$L_D = \frac{1}{2} \sum_{ij} \lambda_i \lambda_j y_i y_j x_i \cdot x_j + C \sum_i \xi_i$$

$$- \sum_i \lambda_i \left\{ y_i \left( \sum_j \lambda_j y_j x_i \cdot x_j + b \right) - 1 + \xi_i \right\}$$

$$- \sum_i (C - \lambda_i) \xi_i$$

$$= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j y_i y_j x_i \cdot x_j$$

Pointing out the margins in primal and dual problem:

In the derivation here we have taken a **soft margin** since we are tolerating a small error to maximise the margins as oppose to a **hard margin** where we tolerate zero error.

*In primal problem:*

The width of the margins is represented by  $2/\|W\|$  which we are trying to maximise in the model.

The objective function  $\frac{1}{2}W^2 + C \sum \zeta_i$ . Here C is penalty we are incurring where we make kind of a trade off between maximising the margins or reducing the misclassified points.

If we set the value C to be very high, we would then tolerate no misclassification, and on setting the value of C to be very low, we get a wider margin.

*In dual problem:*

We need to maximise the margin  $(\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i x_j)$  with respect to  $\lambda$  which

has been derived in the previous section, subject to the constraints,

$$\sum_{i=1}^N \lambda_i y_i = 0, \lambda_i \geq 0.$$

---

The Benefits for **maximising the margins** are:

- I. A larger margin will help in generalising the model when we use sample outside the training data.
  - II. It prevents overfitting of the model as discussed in the bias-variance tradeoff. Thus instead of hard margin where training error is zero we prefer soft margin which generalises more. It is a method of regularisation.
- 

**Characterizing support vectors:**

$(\lambda_i = 0 \text{ and } \xi_i = 0)$  referring that the datapoint  $x_i$  has been correctly classified

$(0 < \lambda_i < C, \xi_i = 0, y_i(w^T x_i + b) = 1)$  which concludes that  $x_i$  is a support vector. All the support vectors which satisfy the condition  $0 < \lambda_i < C$  are **free or unbounded support vectors**.

$(\lambda_i = C)$ :  $y_i(w^T x_i + b) = 1 - \xi_i$ ,  $\xi_i \geq 0$ , implies  $x_i$  is a support vector. When  $\lambda_i = C$  the support vectors are **bounded support vectors**; that is, they lie inside the margin. Furthermore, for  $0 \leq \xi_i < 1$ ,  $x_i$  is correctly classified, but if  $\xi_i \geq 1$ ,  $x_i$  is misclassified.

---

### Benefits of Solving Dual rather than Primal Problem:

- According to the KKT conditions described in the derivation one of them states  $\lambda_i(y_i(w^T x_i + b) - 1 + \xi_i) = 0$ . Thus for all points not lying on support vectors must have  $\lambda$  to be zero. This reduces parameters in dual problems.
- Kernels map points to high dimensional plane. As we know the dual problem relies on the product  $x_i^T x_j$ . This saves a lot of computation when kernels are applied as they need not compute mapping of each point but just their dot product. Thus dual helps in classifying non linearly separable points with less computational effort.
- Microsoft Research also introduced Sequential Minimal Optimisation for solving dual Lagrangian SVM problem.

### References:

[https://www.stat.berkeley.edu/~arturof/Teaching/EE127/Notes/support\\_vector\\_machines.pdf](https://www.stat.berkeley.edu/~arturof/Teaching/EE127/Notes/support_vector_machines.pdf)

### Introduction to Data Mining:

Pang-Ning Tan, Michael Steinbach, and Vipin Kumar. 2005. Introduction to Data Mining, (First Edition). Addison-Wesley Longman Publishing Co., Inc., USA.

<https://homepage.cs.uri.edu/faculty/hamel/courses/2014/spring2014/csc581/lecture-notes/12-dual-maximum.pdf>

3. Formulate the primal problem and derive its dual if there are multiple classes.  
 Grammer and Singer proposed an approach for multi-class problems.  
 Assuming there are  $K$ -classes and  $\ell$ -training examples

Solving the primal problem

$$\min_{w_m \xi_i} \frac{1}{2} \sum_{m=1}^K w_m^T w_m + C \sum_{i=1}^{\ell} \xi_i$$

$$w_i^T \phi(x_i) - w_m^T \phi(x_i) \geq e_i^m - \xi_i, \quad i = 1, \dots, \ell$$

where  $e_i^m \equiv 1 - \delta_{y_i, m}$  and

$$\delta_{y_i, m} \equiv \begin{cases} 1 & \text{if } y_i = m \\ 0 & \text{if } y_i \neq m \end{cases}$$

∴ Therefore the decision function is

$$\arg \max_{m=1, \dots, K} w_m^T \phi(x)$$

The dual problem is

$$\min_{\alpha} f(\alpha) = \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} K_{i,j} \alpha_i^{-T} \bar{x}_j + \sum_{i=1}^{\ell} \bar{\alpha}_i^T \bar{e}_i$$

$$\sum_{m=1}^K \alpha_i^m = 0, \quad i = 1, \dots, \ell,$$

$$\alpha_i^m \leq 0 \quad \text{if } y_i \neq m$$

$$\alpha_i^m \leq C \quad \text{if } y_i = m$$

$$i = 1, \dots, \ell, \quad m = 1, \dots, k$$

$$\text{where } K_{i,j} \equiv \phi(x_i)^T \phi(x_j)$$

$$\bar{x}_i \equiv [\bar{x}_i^1, \dots, \bar{x}_i^K]^T \text{ and } \bar{e}_i \equiv [e_i^1, \dots, e_i^K]^T$$

Then

$$w_m = \sum_{i=1}^l \alpha_i m \phi(x_i)$$

If we write  $\alpha = [d_1^1, d_2^1, \dots, d_K^1, \dots, d_1^l, \dots, d_l^l]^T$

and  $e = [e_1^1, \dots, e_1^K, \dots, e_2^1, \dots, e_2^K]^T$

then the dual objective function can be written as

$$\frac{1}{2} \alpha^T (K \otimes I) \alpha + e^T \alpha$$

where  $I$  is an  $K$  by  $K$  identity matrix and  $\otimes$  is  
Kronecker product. Since  $K$  is positive semi-definite,  
 $K \otimes I$ , is the Hessian of the dual objective function  
is also positive semi definite

Thus the decision function is

$$\text{argmax}_{m=1, \dots, K} \sum_{i=1}^l \alpha_i m K(x_i, n)$$

## References

- A comparison of Methods for Multiclass Support Vector Machines  
Chih-Wei Hsu and Chih-Jen Lin
- On the algorithmic implementation of multiclass kernel based  
vector machines  
Katy Crammer, Yoram Singer