

# 1. Maximum Margin Classifiers.

Our tasks & assumptions

① Our raw data  $X$  will be mapped to another space by kernel function  $\phi(x)$ .

② We know before hand that there will be two classes.

③ We will use a hyperplane to separate them

④ Our dataset is linearly separable

We have a vector of true labels

$$\{t_i\} \in \{-1, 1\}$$

We have a set of input data

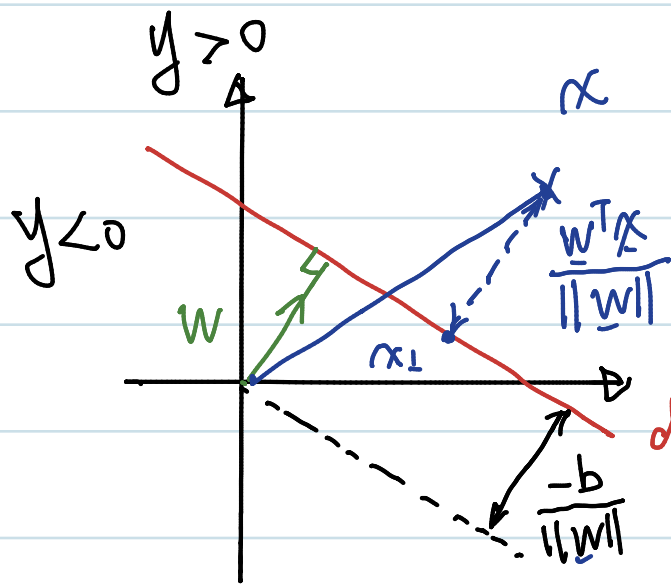
$$\{x_i\}$$

Our linear discriminant model is

$$y(x) = w^T \phi(x) + b$$

and  $w, b$  will represent the decision hyperplane that will linearly separate the two classes.

Let's remind ourselves about some geometries.



If  $x$  is on decision plane,

$$y = w^T x + b = 0$$

$$\Rightarrow \frac{w^T}{\|w\|} x = - \frac{b}{\|w\|}$$

For a point  $x$ , its distance to the decision plane with kernel  $\phi(\cdot)$  is

$$\left| \frac{y(x_n)}{\|w\|} \right|$$

now, w.l.o.g, let's assume that  $w$ 's sign is carefully chosen, so  $\text{tn } y(x_n)$  is always non-negative. So the distance for point  $x$  to the decision plane is

$$d_n = \frac{\text{tn } y(x_n)}{\|w\|} = \frac{\text{tn}(w^T \phi(x_n) + b)}{\|w\|}$$

Now, margin is defined to be the distance of the closest point to decision plane, mathematically it is.

$$\min_n \left\{ \frac{1}{\|w\|} \cdot \text{tn}(w^T \phi(x_n) + b) \right\}$$

Our final goal is to maximize this gap, or say a "maxmin" problem

$$\arg \max_{w, b} \left\{ \frac{1}{\|w\|} \min_n [\text{tn}(w^T \phi(x_n) + b)] \right\}$$

However, solving this optimization problem directly is difficult. We want to convert it to one equivalent simple problem.

Because we don't care the magnitude of  $w$  and  $b$ , say multiply them both by a non-zero vector  $k$ , we won't change our final results for the "min" problem.

Based on this freedom, we can choose a factor  $k$  s.t.

$$t^* (W^T \phi(x^*) + b) = 1$$

where  $x^*$  is the point that is closest to the decision plane.

So, we have

$$t_n (W^T \phi(x_n) + b_n) \geq 1 \quad \forall n \quad (*)$$

Actually, (\*) can be viewed as  $n$  inequality constraints. For point which makes the equality holds can be regarded as an "active" constraint. There will be at least 2 active constraints.

The new optimization problem is

$$\min_{W, b} \frac{\|W\|^2}{2}$$

$$\text{s.t. } t_n (W^T \phi(x_n) + b_n) \geq 1 \quad \forall n$$

So far, we simplify our original target. New thing is to solve it.

We will use the Lagrangian.

$$\begin{aligned} & \mathcal{L}(W, b, \lambda) \\ &= \frac{1}{2} \|W\|^2 - \sum_{n=1}^N \lambda_n \{ \text{tr}(W^T \phi(x_n) + b) - 1 \} \quad \lambda \geq 0 \end{aligned}$$

Let's compute dual function  $g(\lambda)$

$$\begin{aligned} \nabla_W \mathcal{L}(W, b, \lambda) &= W - \sum_{n=1}^N \lambda_n \text{tr} \phi(x_n) = 0 \\ \Rightarrow W^* &= \sum_{n=1}^N \lambda_n \text{tr} \phi(x_n) \end{aligned}$$

$$\nabla_b \mathcal{L}(W, b, \lambda) = \sum_{n=1}^N \lambda_n \text{tr} = 0$$

Replace  $W$  with  $W^*$ , we have

$$\begin{aligned} g(\lambda) &= \frac{1}{2} \left( \sum_{n=1}^N \lambda_n \text{tr} \phi(x_n) \right)^T \left( \sum_{m=1}^N \lambda_m \text{tr} \phi(x_m) \right) \\ &\quad - \left( \sum_{n=1}^N \lambda_n \text{tr} \phi(x_n) \right)^T \left( \sum_{m=1}^N \lambda_m \text{tr} \phi(x_m) \right) + \sum_{n=1}^N \lambda_n \\ &= \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m \text{tr} \text{tr} \underbrace{\phi(x_m)^T \phi(x_n)}_{K(x_m, x_n)} \end{aligned}$$

$$= \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m \text{tr} \text{tr} K(x_m, x_n)$$

and  $\lambda_n \geq 0 \quad \forall n$

$$\sum \lambda_n \text{tr} = 0$$

Then, we simply assume that "according to some optimization technique", the dual problem is solved, and we have the optimal  $\lambda^*$ .

From the formula of  $w^*$  when we minimize  $L(w, b, \lambda)$  w.r.t  $w$  &  $b$ , we have.

$$w^* = \sum_{n=1}^N \lambda_n t_n \phi(x_n)$$

For the optimal value for  $b$ , the step is not that difficult.

$$t_n (w^T \phi(x) + b) \geq 1 \quad t_n \in \{-1, 1\}, \forall n$$

$$\Rightarrow t_n b \geq 1 - t_n w^T \phi(x)$$

$$\Rightarrow t_n b = \min_n (1 - t_n w^T \phi(x))$$

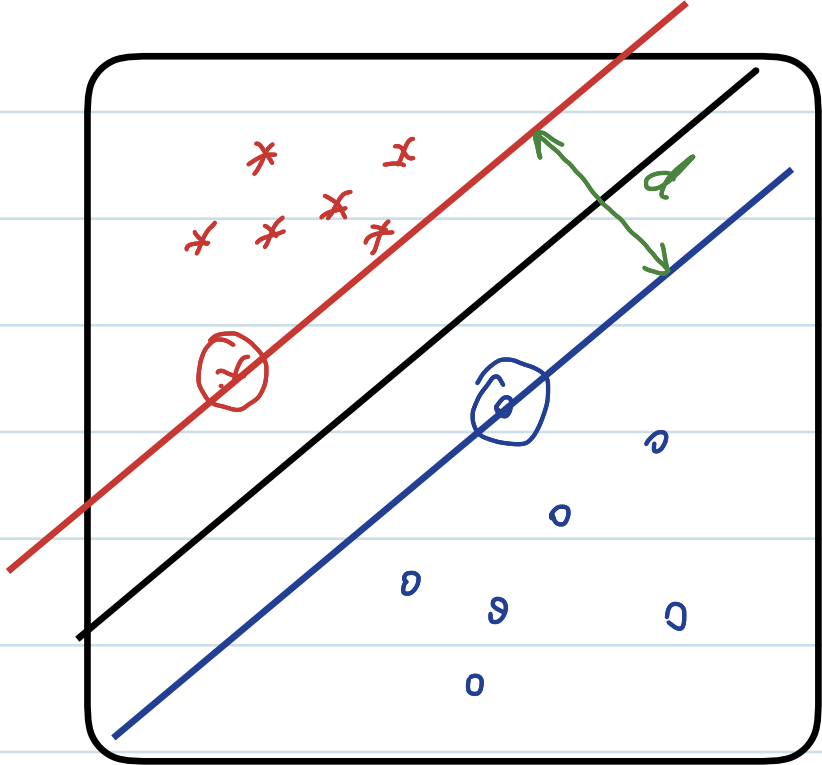
$$\Rightarrow b^* = \min_n (1 - t_n w^T \phi(x))$$

Later, we can use this formula &  $\lambda^*$ , we have

$$y(x) = \sum_{n=1}^N \lambda_n^* t_n k(x, x_n) + b^*$$

We can simply determine the class of sample  $x$  by checking the sign of  $f(x)$ .

Now, let's talk about the intuition for maximum margin classifier.



First, we only care the so called "Support Vector".

In left case, the circled points are regarded as support vectors, because they support the boundaries of red and blue samples.

The middle black line is regarded

as the decision hyperplane. The direction of planes is chosen by maximizing the distance "d" between two boundary planes.

So we know that, only the support vectors that really care. So, this algorithm won't be affected by weird outliers.

E.g, points extremely far away from common class central won't affect the decision plane.