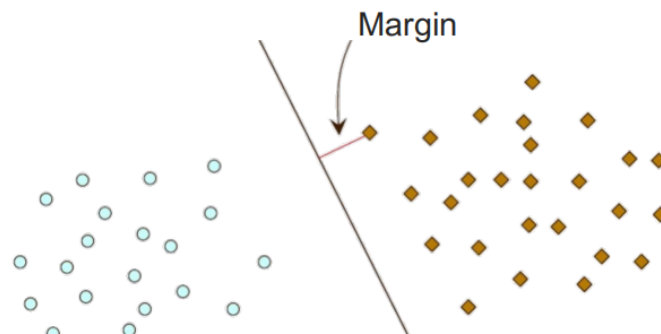




# ch7: Support Vector Machine

**The Idea:** Find a decision boundary that **maximizes the margin** between two classes



## Problem Formulation(P6-9)

### Maximizing the Margin(P7)

$$\mathbf{w}^T \mathbf{x} + w_0 = \begin{cases} 1 & \text{for the closest points on one side} \\ -1 & \text{for the closest points on the other} \end{cases}$$

- Let  $x^{(1)}$  and  $x^{(2)}$  be two closest points on each side of the hyperplane.
- Note that

$$\mathbf{w}^T x^{(1)} + w_0 = +1$$

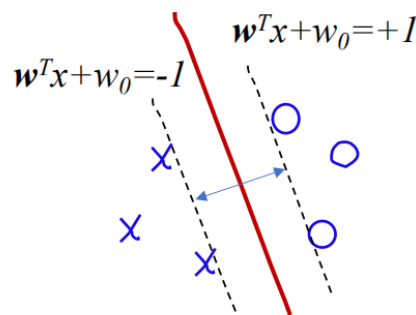
$$\mathbf{w}^T x^{(2)} + w_0 = -1$$

which imply

$$\mathbf{w}^T (x^{(1)} - x^{(2)}) = 2.$$

Hence, the margin can be given by

$$\text{margin} = \frac{\mathbf{w}^T (x^{(1)} - x^{(2)})}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$



**Maximizing** the margin is equivalent to **minimizing**  $\frac{1}{2} \|\mathbf{w}\|$

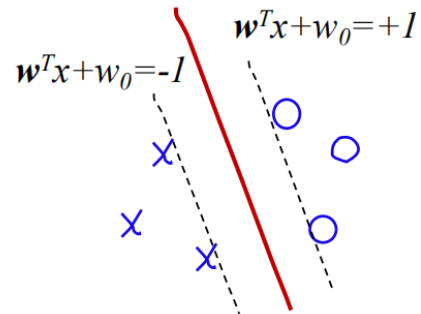
## Inequality Constraints(P8)

- Given:  $D = \{(x^{(\ell)}, y^{(\ell)}), \dots, (x^{(N)}, y^{(N)})\}$ , we want  $\mathbf{w}$  and  $w_0$  to satisfy

$$\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0 \begin{cases} \geq +1 & \text{if } y^{(\ell)} = +1 \\ \leq -1 & \text{if } y^{(\ell)} = -1 \end{cases}$$

- Or, equivalently,

$$y^{(\ell)} (\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1$$



## SVM = Solving an Optimization Problem(P9)

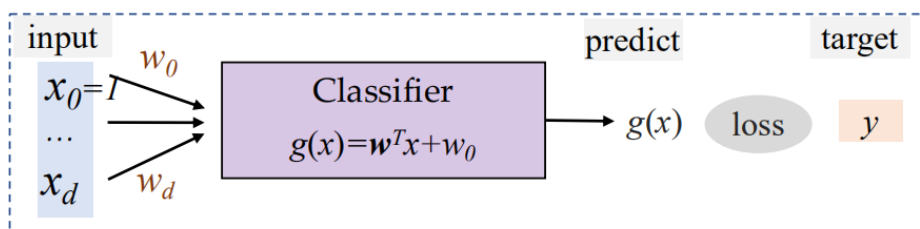
In summary, SVM aims to solve a constrained optimization problem:

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ &\text{subject to} \quad y^{(\ell)} (\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1, \ell = 1, \dots, N \end{aligned}$$

- This is a **quadratic programming** (QP) problem, which is one type of convex optimization problem.
- The complexity depends on the dimensionality  $d$  of inputs

## Model Architecture(P10)

- The same architecture as Perceptron.



- Train:**
  - optimize the parameters  $\mathbf{w}$  and  $w_0$  using data
- Test:**
  - calculate  $g(x) = \mathbf{w}^T \mathbf{x} + w_0$  and choose  $C_1$  if  $g(x) > 1$  or choose  $C_2$  if  $g(x) < -1$ .

## Loss Function(P11)

- For a given input  $x$ , the model outputs a score  $g(x) = \mathbf{w}^T x + w_0$ . Let  $y \in \{-1, +1\}$  be the label of the real class ( $y = +1: x \in C_1, y = -1: x \in C_2$ ):
  - if  $y(\mathbf{w}^T x + w_0) < 1$ : we aim to maximize  $y(\mathbf{w}^T x + w_0)$  until reaching 1, cost is  $1 - y(\mathbf{w}^T x + w_0)$
  - if  $y(\mathbf{w}^T x + w_0) > 1$ : outlier points, no need to optimize, cost is 0
- Given:  $D = \{(x^{(1)}, r^{(1)}), \dots, (x^{(N)}, r^{(N)})\}$ , the loss over the dataset is defined as:

$$L(\mathbf{w}, w_0 | D) = \frac{1}{N} \sum_{\ell=1}^N \max(0, 1 - y^{(\ell)}(\mathbf{w}^T x^{(\ell)} + w_0)) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

## Optimization – Gradient Descend(P12)

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \eta \frac{\partial L}{\partial \mathbf{w}}$$

What is  $\frac{\partial L}{\partial \mathbf{w}}$ ?

$$L(\mathbf{w}, w_0 | D) = \begin{cases} \frac{\lambda}{2} \|\mathbf{w}\|^2 & \text{if } y^{(\ell)}(\mathbf{w}^T x^{(\ell)} + w_0) \geq 1 \\ \frac{1}{N} \sum_{\ell=1}^N 1 - y^{(\ell)}(\mathbf{w}^T x^{(\ell)} + w_0) + \frac{\lambda}{2} \|\mathbf{w}\|^2 & \text{otherwise} \end{cases}$$

$$\text{For each } \mathbf{w}_j \ (j=0, \dots, d): \quad \frac{\partial L}{\partial w_j} = \begin{cases} \lambda w_j & \text{if } y^{(\ell)}(\mathbf{w}^T x^{(\ell)} + w_0) \geq 1 \\ \lambda w_j - y^{(\ell)} x^{(\ell)} & \text{o.w.} \end{cases}$$

$$\frac{\partial L}{\partial w_0} = \begin{cases} 0 & \text{if } y^{(\ell)}(\mathbf{w}^T x^{(\ell)} + w_0) \geq 1 \\ -y^{(\ell)} & \text{o.w.} \end{cases}$$

## Algorithm(P13)

## Gradient Descend for SVM

```
Input:  $D = \{(x^{(l)}, y^{(l)})\} \ (l=1:N)$   
for  $j = 0, \dots, d$   
     $w_j \leftarrow \text{rand}(-0.01, 0.01)$   
repeat  
    for  $j = 0, \dots, d$   
         $\Delta w_j \leftarrow 0$   
    for  $l = 1, \dots, N$   
         $a \leftarrow 0$   
        for  $j = 0, \dots, d$   
             $a \leftarrow a + w_j x_j^{(l)}$   
        for  $j = 0, \dots, d$   
             $\Delta w_j \leftarrow \Delta w_j + \lambda w_j$   
             $\Delta w_0 \leftarrow 0$   
            if  $y^{(l)}a < 1$ :  $\Delta w_j \leftarrow \Delta w_j - y^{(l)}x^{(l)}$   
         $\Delta w_j = \Delta w_j / N$   
    for  $j = 0, \dots, d$   
         $w_j \leftarrow w_j + \eta \Delta w_j$   
until convergence
```

## Dual Problem

### Lagrangian(P15)

- 定义

$$\begin{aligned} &\text{minimize} && f_0(x) \\ &\text{subject to} && f_i(x) \leq 0, \quad \text{for } i = 1, \dots, m. \\ &&& h_i(x) = 0, \quad \text{for } i = 1, \dots, k. \end{aligned}$$

**Idea:** relax/soften the constraints by replacing each with a linear “penalty term” or “cost” in the objective.

### The Lagrangian Function

$$L(x, \lambda, v) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^k v_i h_i(x)$$

- $\lambda_i$  is the **lagrange multiplier** for the  $i$ -th inequality constraint
  - required to be nonnegative
- $v_i$  is the **lagrange multiplier** for the  $i$ -th equality constraint
  - allowed to be arbitrary sign
- **Lagrange Dual Problem**

This is the problem of finding the **best lower bound** on  $\text{OPT}(\text{primal})$  obtained from the Lagrange dual function

$$\begin{aligned} &\underset{\lambda, \nu}{\text{maximize}} && g(\lambda, \nu) \\ &\text{subject to} && \lambda \geq 0. \end{aligned}$$

- $(\lambda^*, v^*)$  solving the above are referred to as the **dual optimal solution**.
- **Note:** this is a **convex** optimization problem, regardless of whether primal problem was convex

$$\begin{aligned} &\text{minimize} && f_0(x) \\ &\text{subject to} && f_i(x) \leq 0, \quad \text{for } i = 1, \dots, m. \\ &&& h_i(x) = 0, \quad \text{for } i = 1, \dots, k. \end{aligned}$$

**Idea:** creates the **lower bound** of the primal optimum subject to the softened constraints.

### The Lagrange Dual Function

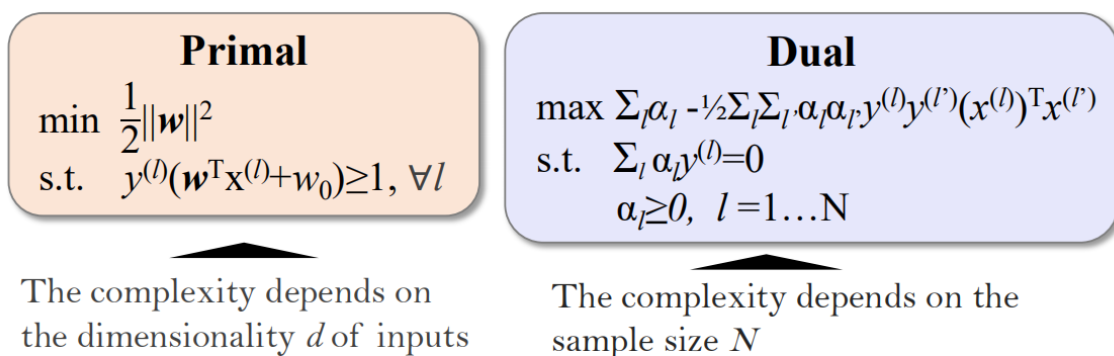
$$g(\lambda, v) = \inf_{x \in \mathcal{D}} L(x, \lambda, v) = \inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^k v_i h_i(x))$$

## The Dual Problem for SVM(P16-17)

### Dual optimization problem:

$$\begin{aligned} &\text{maximize} && \sum_{\ell=1}^N \alpha_{\ell} - \frac{1}{2} \sum_{\ell=1}^N \sum_{\ell'=1}^N \alpha_{\ell} \alpha_{\ell'} y^{(\ell)} y^{(\ell')} (x^{(\ell)})^T x^{(\ell')} \\ &\text{subject to} && \sum_{\ell=1}^N \alpha_{\ell} y^{(\ell)} = 0 \\ &&& \alpha_{\ell} \geq 0, \quad \ell = 1 \dots N \end{aligned}$$

- This is also a QP problem, but its complexity depends on the sample size  $N$  (rather than the input dimensionality  $d$ )



- **It turns out to be more convenient to solve the dual problem** than solving the **primal problem** ( $N < d$ )
- We can firstly solve **Dual** to obtain  $\{\alpha_{\ell}\}$ , and then obtain the  $\mathbf{W}$  in **Primal**

## Support Vectors(P19)

Suppose the optimal  $\{\alpha_\ell\}$  have been obtained

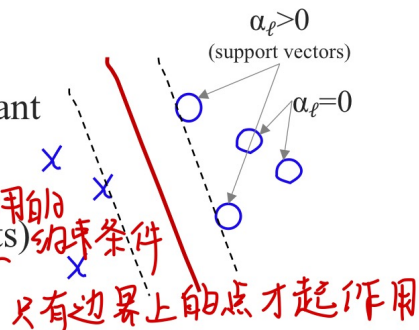
- Patterns for which  $y^{(\ell)}(wx^{(\ell)} + w_0) > 1$

$\alpha_\ell = 0$  (inactive constraints)  $\Rightarrow x^{(\ell)}$  irrelevant

recall complimentary slackness:  $\lambda_i^* g_i(x^*) = 0$

- Patterns that have  $\alpha_\ell > 0$  (active constraints)

$y^{(\ell)}(wx^{(\ell)} + w_0) = 1 \Rightarrow x^{(\ell)}$  lies on margin



- Most of the dual variables vanish with  $\alpha_\ell = 0$ . They are points lying beyond the margin with **no effect** on the hyperplane.

- Solution is only determined by the examples on the margin (support vectors), i.e.,  $x^{(\ell)}$  with  $\alpha_\ell > 0$ , hence the name **support vector machine (SVM)**.

在边界上的点支持向量