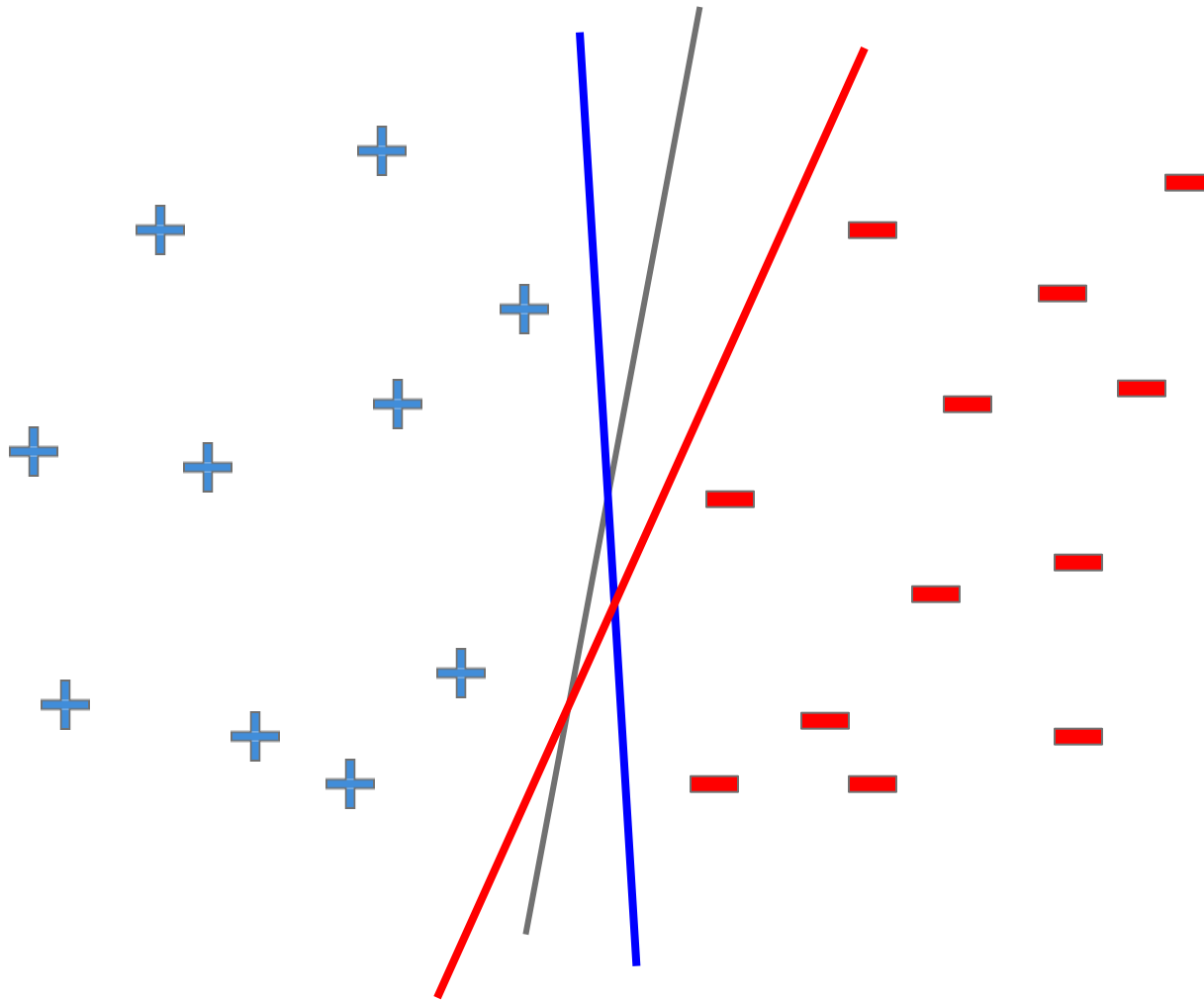


# 机器学习

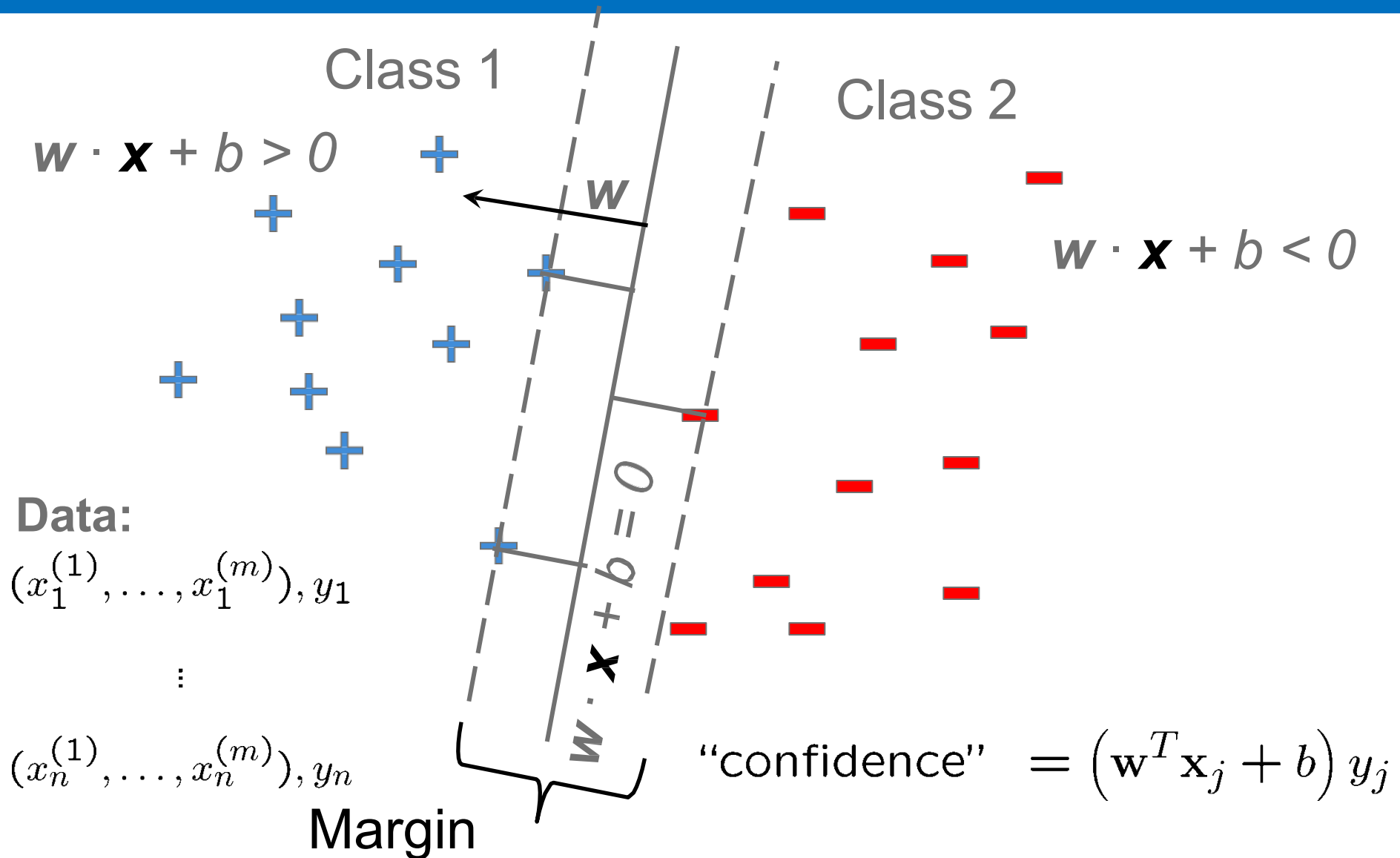
## 实验课题1-Support Vector Machines

# Linear classifiers which line is better?

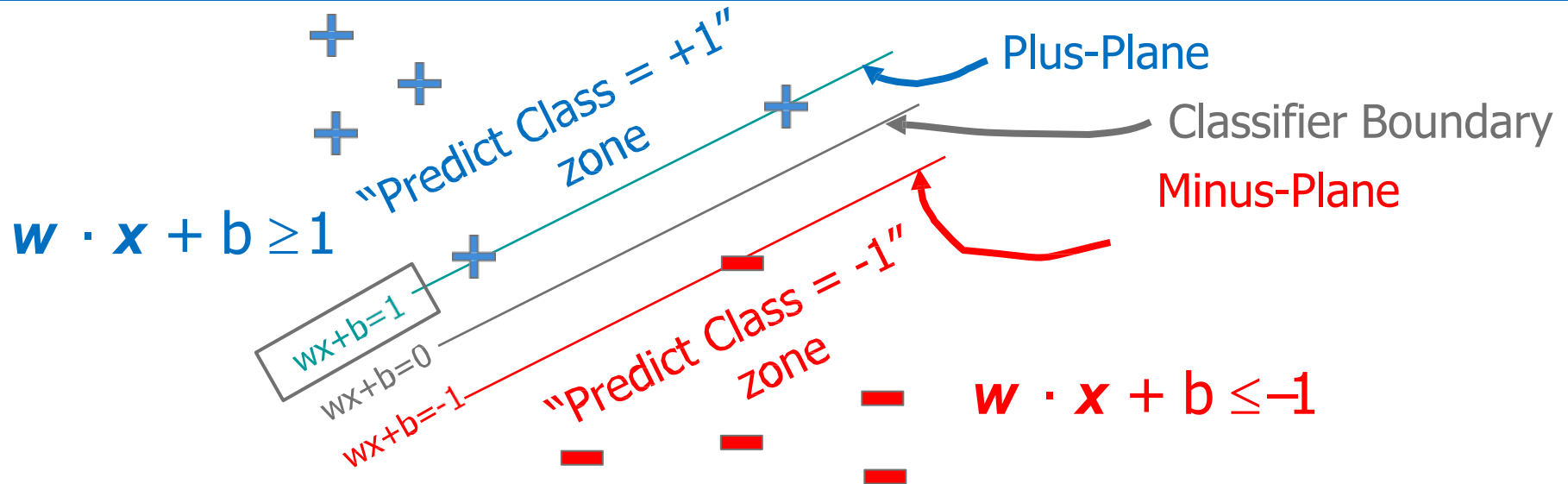


**Which decision boundary is better?**

# Pick the one with the largest margin!



# Scaling



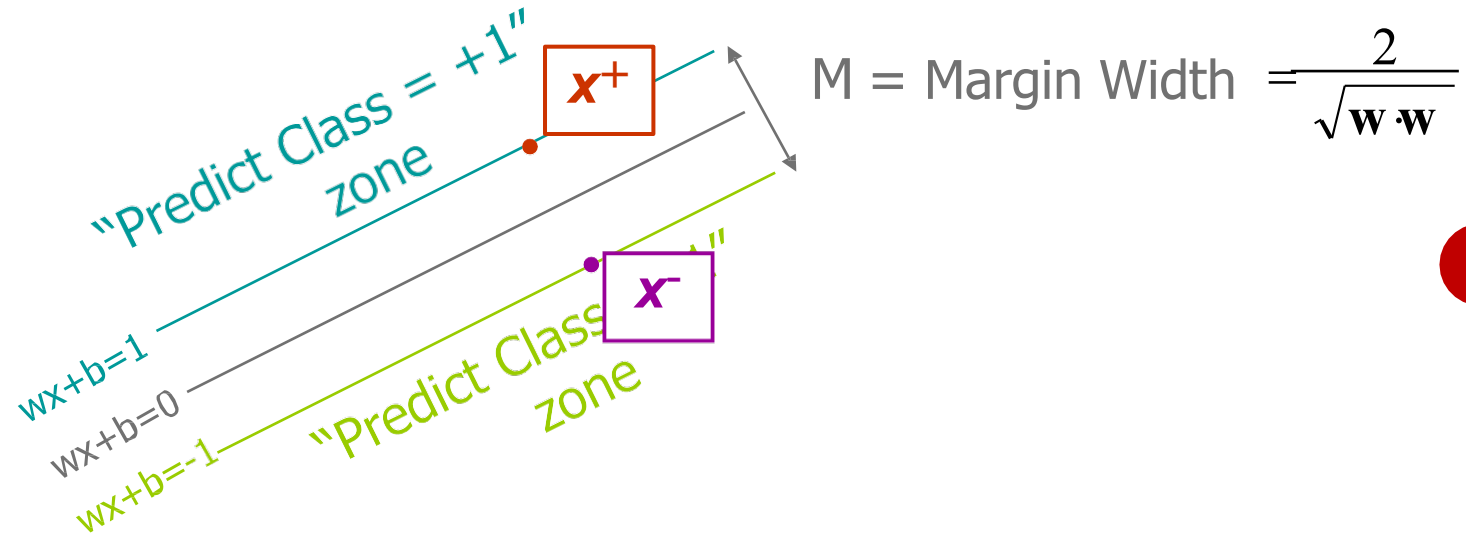
Classification rule:

|               |                   |    |                          |
|---------------|-------------------|----|--------------------------|
| Classify as.. | <b>+1</b>         | if | $w \cdot x + b \geq 1$   |
|               | <b>-1</b>         | if | $w \cdot x + b \leq -1$  |
|               | Universe explodes | if | $-1 < w \cdot x + b < 1$ |

How large is the margin of this classifier?

**Goal:** Find the maximum margin classifier

# Computing the margin width



Let  $\mathbf{x}^+$  and  $\mathbf{x}^-$  be such that

- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$

Maximize  $M \equiv \text{minimize } \mathbf{w} \cdot \mathbf{w} !$

# Observations

We can assume  $b=0$

Classify as..  $+1$  if  $\mathbf{w} \cdot \mathbf{x} + b \geq 1$

$-1$  if  $\mathbf{w} \cdot \mathbf{x} + b \leq -1$

Universe  
explodes if  $-1 < \mathbf{w} \cdot \mathbf{x} + b < 1$

This is the same as  $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n$

# The Primal SVM

- Given  $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$  training data set.
- Assume that  $D$  is **linearly separable**.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n$$

**Prediction:**  $f_{\hat{\mathbf{w}}}(\mathbf{x}) = \text{sign}(\langle \hat{\mathbf{w}}, \mathbf{x} \rangle)$

**This is a QP problem (m-dimensional)**  
**(Quadratic cost function, linear constraints)**

# The Primal SVM

- Given  $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$  training data set.
- Assume that  $D$  is **linearly separable**.

$$\text{argmin: } \frac{1}{2} |\mathbf{w}|^2 + C \sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

The cost function is often referred to as **the primal optimization problem**, which seeks to minimize the sum of hinge losses over all training examples, subject to a constraint that the magnitude of the weight vector is bounded.



# The Primal SVM

- Given  $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$  training data set.
- Assume that  $D$  is **linearly separable**.

When we assume  $b \neq 0$

$$\operatorname{argmin}: \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

The cost function is often referred to as **the primal optimization problem**, which seeks to minimize the sum of hinge losses over all training samples, subject to a constraint that the magnitude of the weight vector is bounded.

# Implementation--Training

1. Any programming language can be used.
2. Write the function of the gradient computation with respect to  $\mathbf{w}$  and  $b$  (if you assume  $b \neq 0$ ). (**manually derive the gradients on papers**).
3. Write the function of computing the loss of the objective function over all training samples.
4. Construct a loop structure to compute gradients, update weights ( $\mathbf{w}$ ,  $b$ ) according to the gradients, and update the loss in each iteration.

When we assume  $b \neq 0$

$$\operatorname{argmin}: \frac{1}{2} |\mathbf{w}|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T x_i + b))$$



## Algorithm: gradient descent

Initialize  $\mathbf{w} = [0, \dots, 0]$

For  $t = 1, \dots, T$ :

$$\mathbf{w} \leftarrow \mathbf{w} - \underbrace{\eta}_{\text{step size}} \underbrace{\nabla_{\mathbf{w}} \text{TrainLoss}(\mathbf{w})}_{\text{gradient}}$$

# Implementation--Testing

1. Load the trained model (load the parameters of  $\mathbf{w}$  and  $b$  learned).
2. Load the testing data
3. Input each testing sample into the below prediction function to make a prediction.
4. Report the accuracy by comparing the prediction outputs and the true labels.

$$f_{\hat{\mathbf{w}}}(x) = \text{sign}(\langle \hat{\mathbf{W}}, \mathbf{x} \rangle)$$

**Prediction:**

$$f_{\hat{\mathbf{w}}}(x) = \text{sign}(\langle \hat{\mathbf{W}}, \mathbf{x} \rangle + b)$$