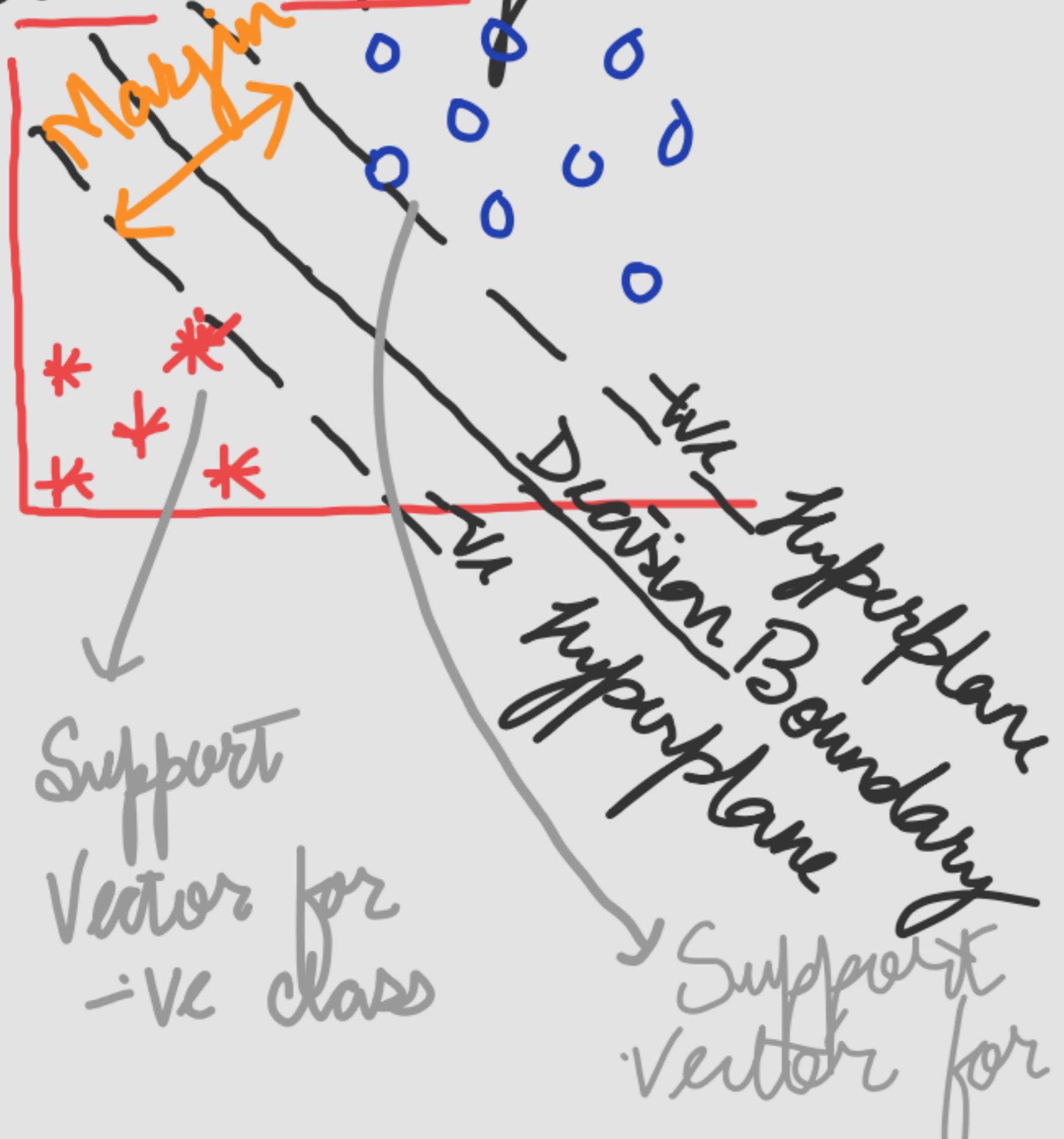


# Hard Margin SVM

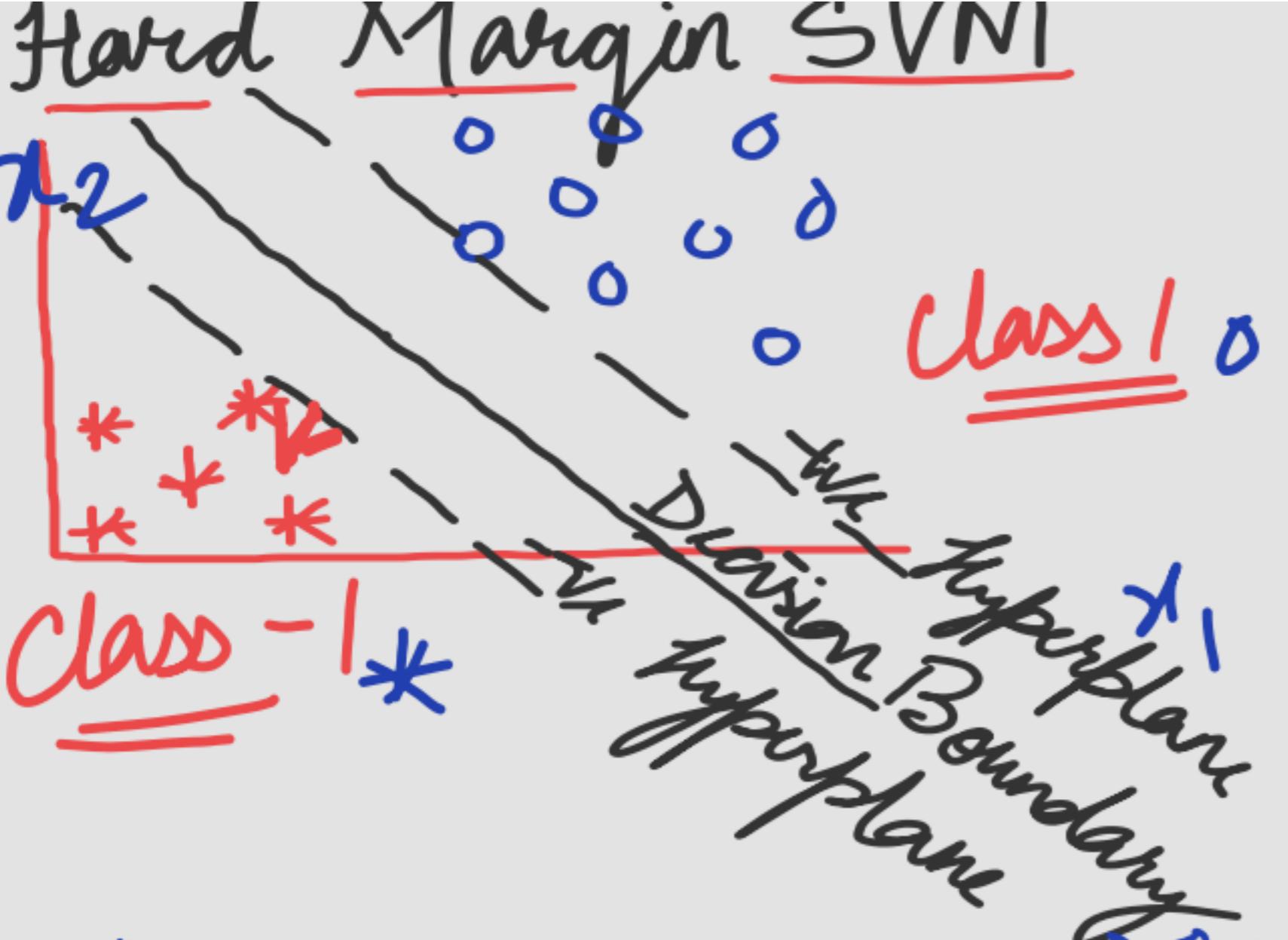
\*NOTE : Refer class notes



AIM

To come up with the division boundary that can act as max. margin classifier i.e. distance of support vectors from decision boundary is max.

Margin : Distance of Support Vectors (closest pt. from division boundary) from the decision boundary



Eq'n of Decision Boundary

$$w_1 x_1 + w_2 x_2 + b = 0$$

or

$$\vec{w} \cdot \vec{x} + b = 0$$

where

$$w = [w_1 \ w_2]$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Eq'n of positive hyperplane :  $\vec{w} \cdot \vec{x} + b = 1$

Eq'n of negative hyperplane :  $\vec{w} \cdot \vec{x} + b = -1$

Thus, the Constraints that training data should satisfy :

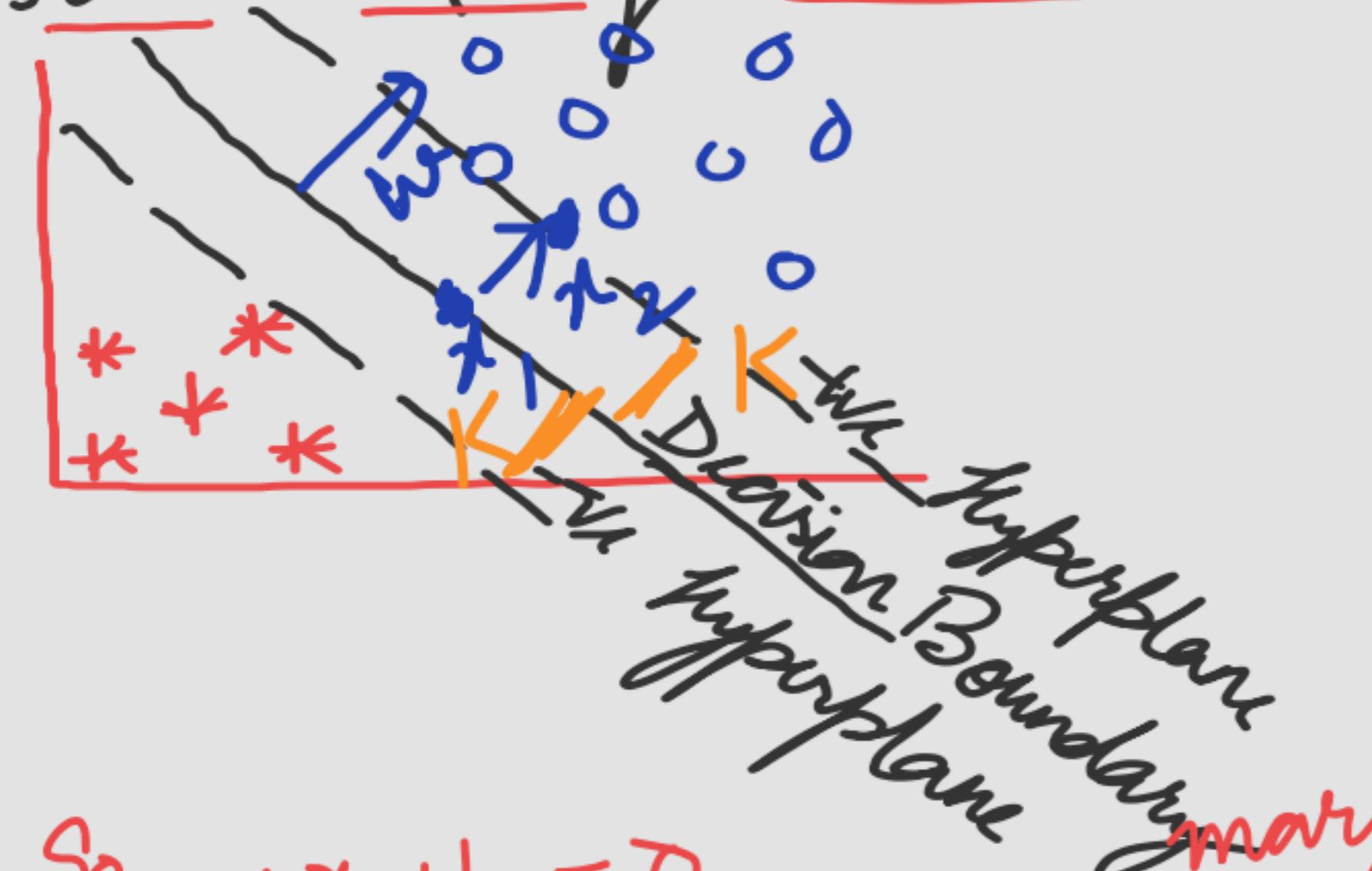
If  $y=1$ , then  $\vec{w} \cdot \vec{x} + b \geq 1$

4 If  $y=-1$ , then  $\vec{w} \cdot \vec{x} + b \leq -1$

(OR)

$$(y) \cdot (\vec{w} \cdot \vec{x} + b) \geq 1$$

## Hard Margin SVM



$w$ : vector which is normal  
Or perpendicular to the decision boundary

Consider pt  $x_1$  on decision boundary &  $x_2$  on the hyperplane which is  $K$  units away

$$\text{So, } w \cdot x_1 + b = 0$$

To reach pt  $x_2$  on the  $K$ . Can be found as

$$w \left( x_1 + K \frac{w}{\|w\|} \right) + b = 1 \quad \left( \because x_2 = x_1 + K \frac{w}{\|w\|} \right)$$

$K$  units in direction of  $w$

Refr  
class  
Notes |

$$\Rightarrow \vec{w}\vec{x}_i + \frac{K \vec{w} \cdot \vec{w}}{\|\vec{w}\|} + b = 1 \quad (\text{so } \vec{w}\vec{x} + b = 0)$$

$$\Rightarrow K = \frac{\|\vec{w}\|}{\vec{w} \cdot \vec{w}} \quad (\because \vec{w} \cdot \vec{w} = \|\vec{w}\|^2)$$

$$= \frac{\|\vec{w}\|}{\|\vec{w}\|^2} = \frac{1}{\|\vec{w}\|}$$

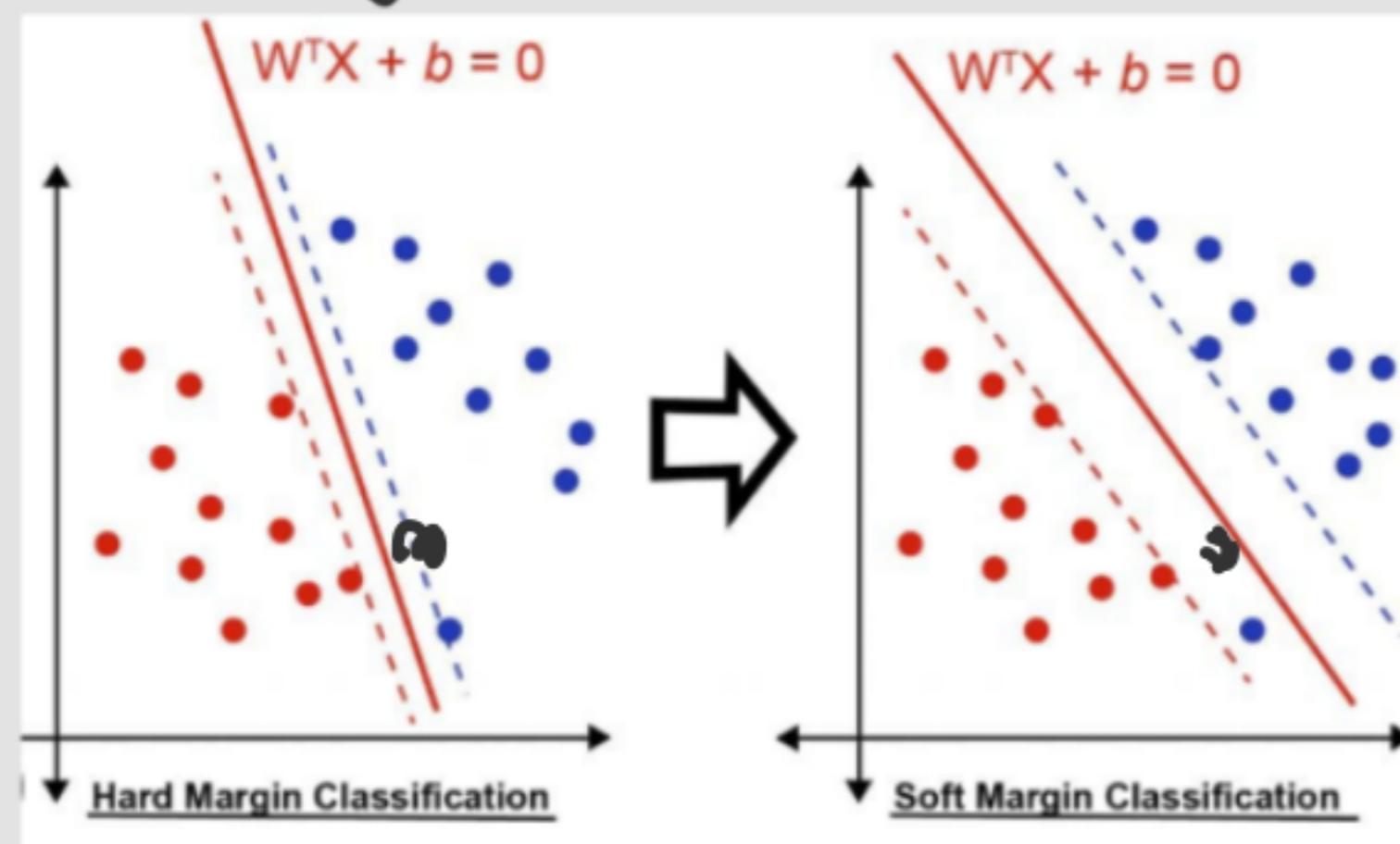
Total Margin :  $2K = \frac{2}{\|\vec{w}\|}$

Optimization Problem of Hard margin SVM:

$$\text{Max } \frac{2}{\|\vec{w}\|} \quad \text{or Min. } \frac{1}{2} \|\vec{w}\|^2$$

subject to  $y_i(\vec{w}\vec{x}_i + b) \geq 1 \quad \forall i=1,2,3,\dots,N$

# Hard Margin vs Soft Margin

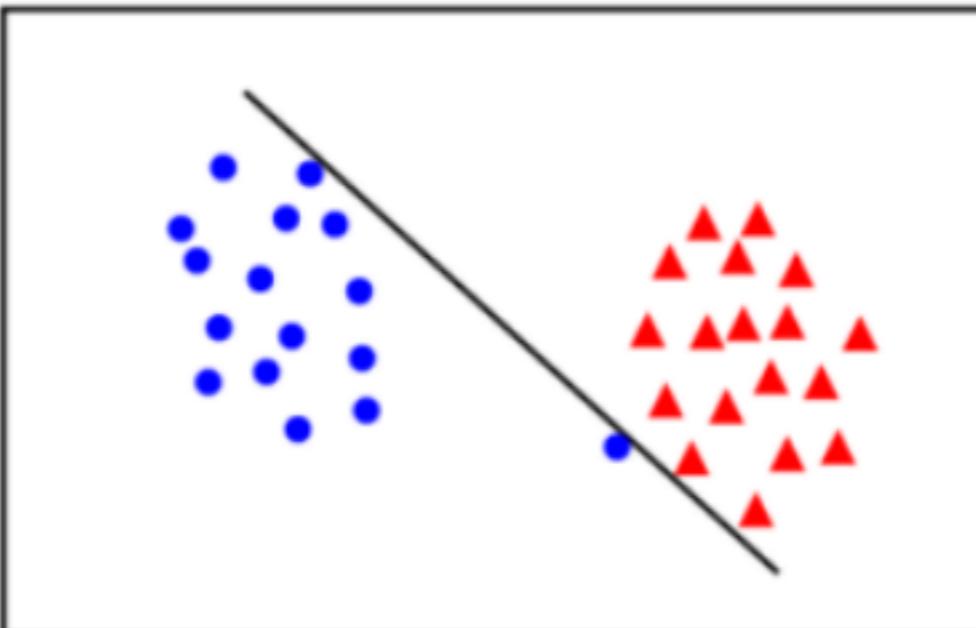


Soft Margin: try to find a line to separate, but tolerate one or few misclassified dots (e.g. the dots circled in red)

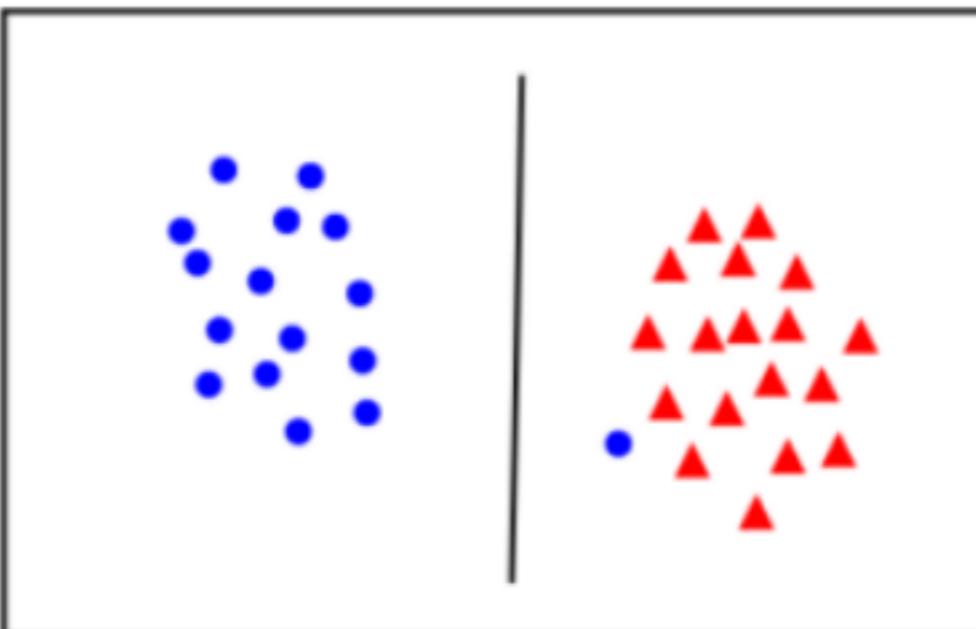
which case denotes overfitting?

→ Hard Margin  
→ Soft Margin

\* Generally data is not linearly separable



- the points can be linearly separated but there is a very narrow margin



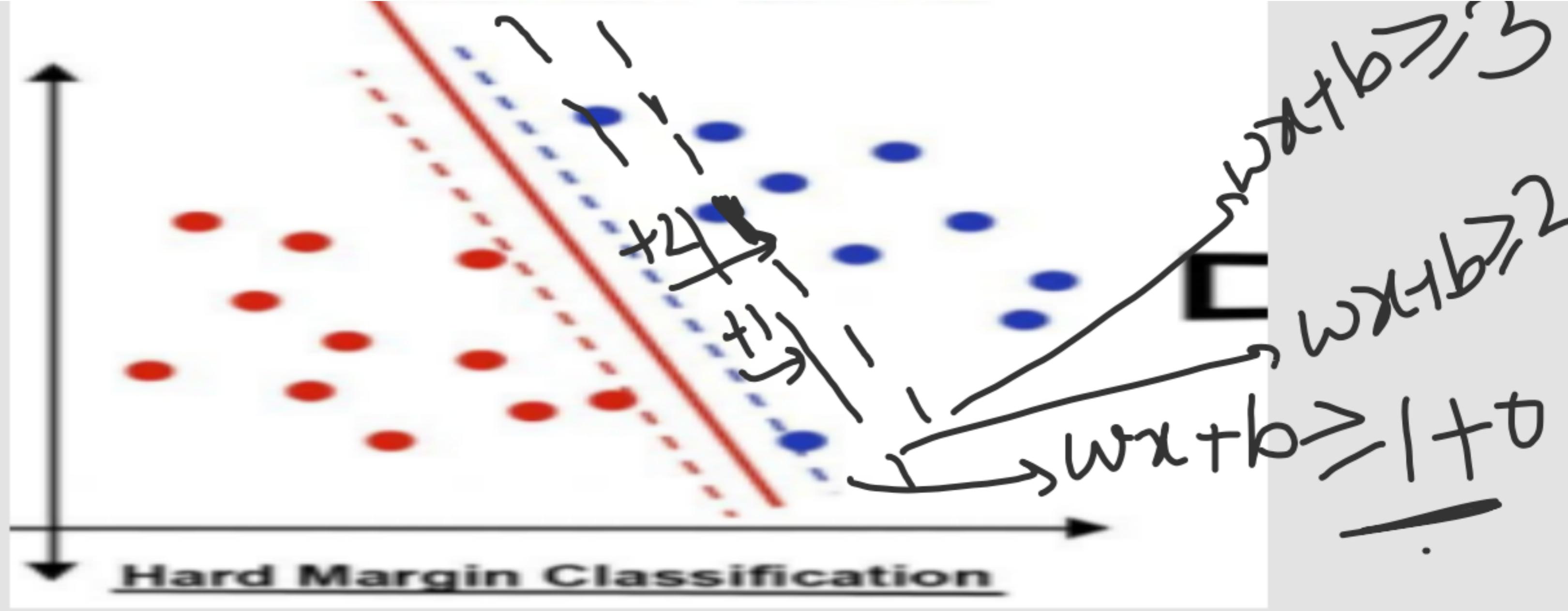
- but possibly the large margin solution is better, even though one constraint is violated

In general there is a trade off between the margin and the number of mistakes on the training data

Idea: Give certain penalties for mistakes based on how big they are

→ Avoid overfitting of our model by introducing slack var  $\xi_i$  one for each example, which allow some of the points to be inside margin

Aim: Max. margin 4 min. misclassifications margin error

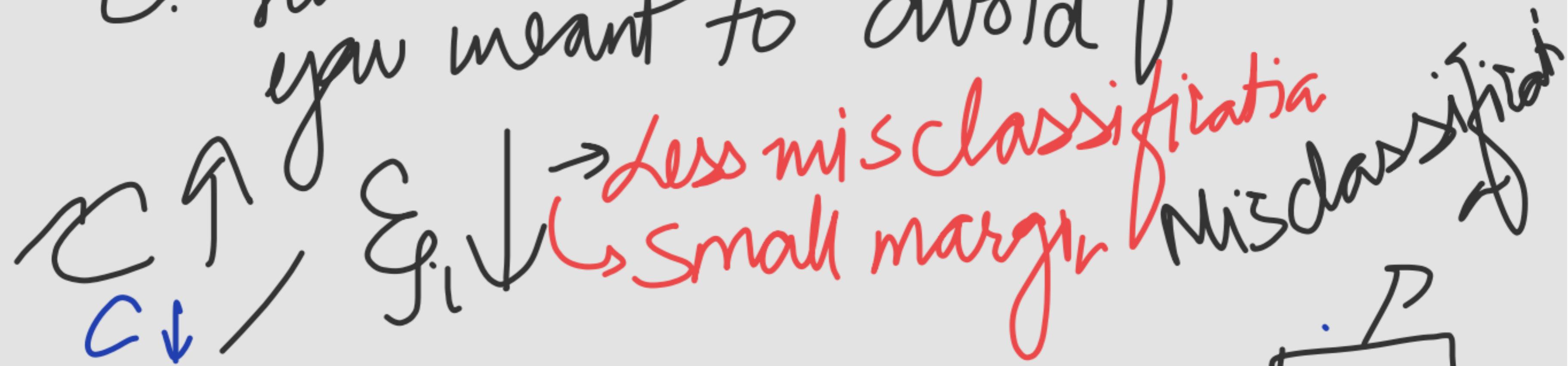


Optimization Problem of Soft margin SVM:

Max or Min.  $\frac{1}{2} \|w\|^2 + C \sum_{i=0}^N \xi_i$

subject to  $y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 + \xi_i \quad \forall i = 1, 2, 3, \dots, N$ .

C: How much misclassifications  
you meant to avoid



Optimization Problem of Soft margin SVM

or Min.  $\frac{1}{2} \|w\|^2 + C \sum_{i=0}^P \epsilon_i$

subject to  $y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 + \epsilon_i \quad \forall i = 1, 2, 3, \dots, N$ .

$C$ : large  $\rightarrow E_g \downarrow$

- ↳ Lesser misclassification
- Optimizer will choose smaller margin hyperplane

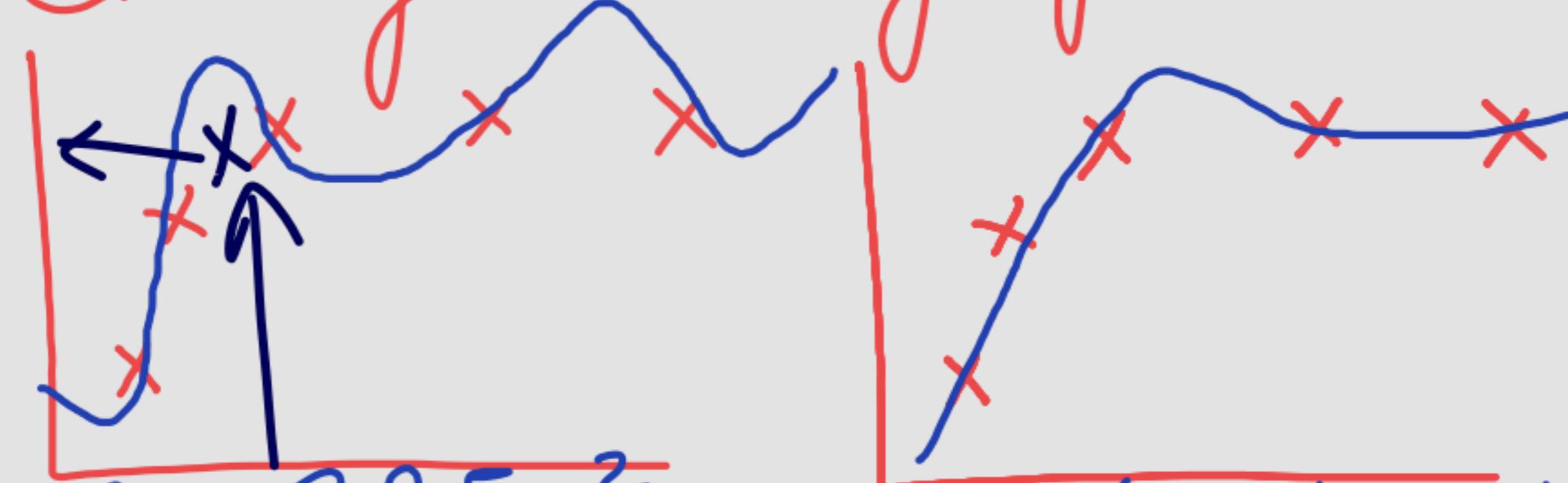
$C$ : small

- ↳ Entire focus on first term
- Optimizer will look for large margin hyperplane.

# Regularization

< Refer Class Notes >

→ Technique to avoid overfitting by shrinking the value of parameters.



$$f(x) = 29x - 385x^2 + 39x^3 - 174x^4 + 10$$

$$f(x) = 13x - 0.23x^2 + 0.00014x^3 - 0.0001x^4 + 10$$

Idea:  
Update cost func. by penalizing  
weight parameters

## L1 Regularization (Lasso Regression)

$$J(w, b) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}))^2 + \frac{\lambda}{2N} \sum_{j=1}^d |w_j|$$

## L2 Regularization (Ridge Regression)

$$J(w, b) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}))^2 + \frac{\lambda}{2N} \sum_{j=1}^d w_j^2$$

Orange: Optional

Ques: Why  $\lambda$  should not be too small or too large?  
If  $\lambda = 0$   
→) No regularization  
→) Fit overly wiggly, overly complex curve.

↳ Effect of overfitting data  
↳ less Training Error, large Generalization / Prediction error

If  $\lambda \rightarrow$  Very Very Large  
⇒ Placing heavy weight on regularization term  
⇒ Learning algo will choose  $w_1, w_2, \dots, w_n$  extremely close to 0.  
⇒  $F(x')$  is basically equal to  $b$ .

Large Train Error & Large Generalization Prediction Error  
↳ Learning algo will fit a straight line if it underfits

Choose  $\lambda$  that balances first & second terms of the order of  
↳ Then make the parameters small.

## The differences between L1 and L2 regularization:

L1 regularization penalizes the sum of absolute values of the weights

- ↳ Lasso shrinks weights of less imp. features to zero.
- ↳ Has the effect of feature Selection

L2 regularization penalizes the sum of squares of the weights.

$$J = \frac{1}{2N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2N} \sum_{j=1}^n w_j^2$$

Gradient Descent for

Linear Regression

repeat until convergence

$$b = b - \alpha \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)})$$

$$w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

$$w_1 = w_1 - \frac{\partial J}{\partial w_1}$$

Regularized



SAME

$$w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

$$\frac{\partial J}{\partial w_1} = \alpha \left( \frac{1}{N} \sum_{i=1}^N (w_1^2 + w_2^2) \right)$$

$$\beta = \frac{1}{2N} \sum_{j=1}^N (\omega_j)^2$$

$$\frac{\partial B}{\partial w_1} = \frac{\partial}{\partial w_1} \left( \frac{1}{2N} (\omega_1^2 + \omega_2^2 + \dots + \omega_N^2) \right)$$

$$= \frac{1}{2N} \times 2w_1 = \frac{1}{N} \omega_1$$

$$J = -\frac{1}{N} \sum_{i=1}^N \left( y^{(i)} \log(f(x^{(i)})) + (1-y^{(i)}) \log(1-f(x^{(i)})) \right) + \frac{\lambda}{2N} \sum_{j=1}^d w_j^2$$

Gradient Descent for  $f(x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + \dots + w_dx_d + b)}}$

## Logistic Regression

repeat until Convergence

$$b = b - \alpha \frac{1}{N} \sum (f(x^{(i)}) - y^{(i)})$$

$$w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^N (f(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

if

## Regularized Logistic Reg

SAME

$$w_j = w_j - \alpha \frac{1}{N} \sum_{i=1}^N (f(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{N} w_j$$