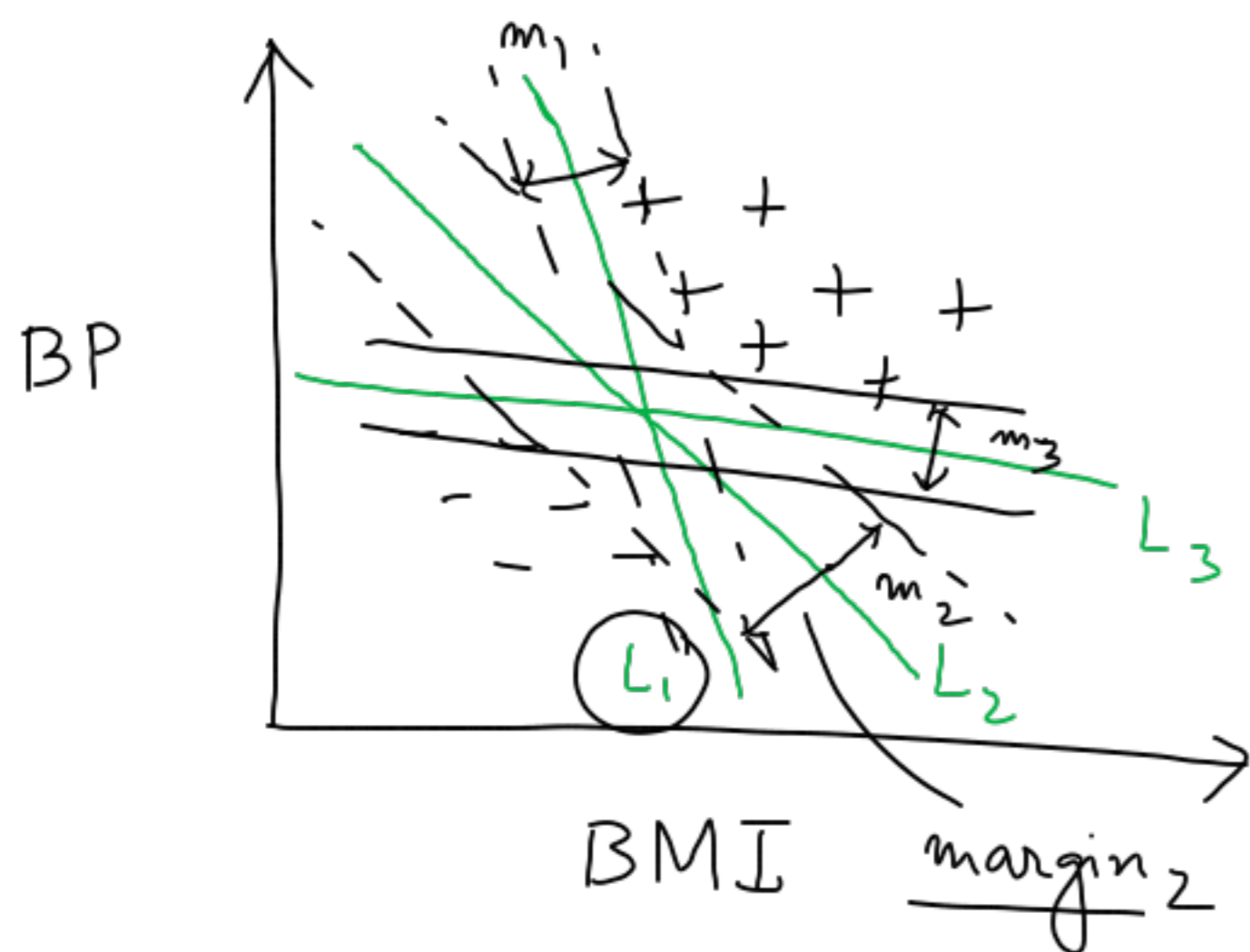


# Lec 15: Support Vector Machines (SVM)

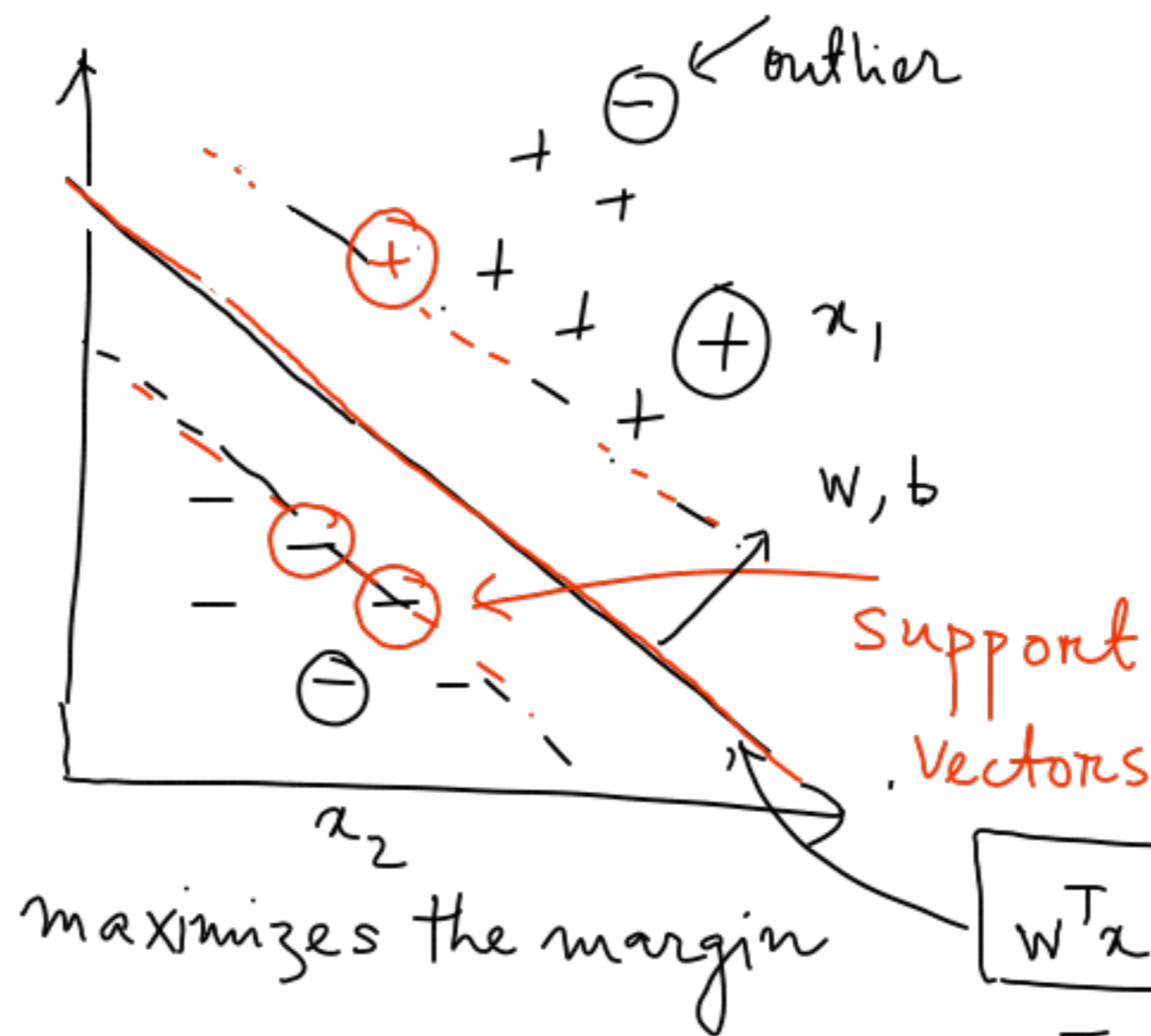


margin-based classifier.  
 $m_2 > m_1, m_3$

Two variants of SVMs:

① Hard margin SVMs: allows no misclassification.

② Soft margin SVMs: misclassifications allowed but to a "limited" extent.



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

$$w^T x_1 + b > 0$$

$$w^T x_2 + b < 0$$

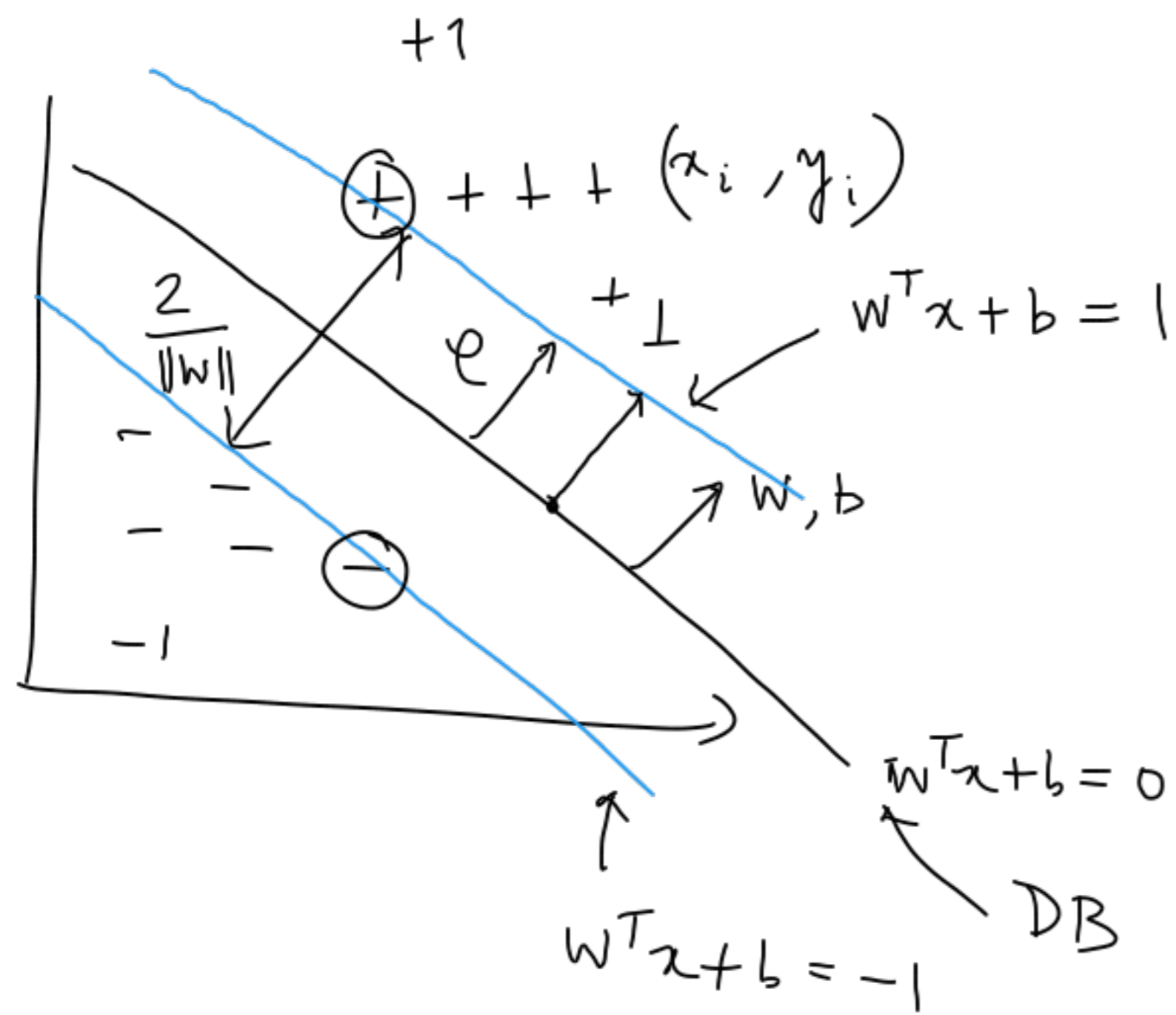
$$D = \left\{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \right\}$$

$$x_i \in \mathbb{R}^d$$

$$d \gg n$$

$$y_i \in \{+1, -1\}$$

# Formal description



$$\frac{w}{m} \rightarrow w, \frac{b}{m} \rightarrow b$$

Pick a point  $x$  on DB:  $w^T x + b = 0$

$$w^T \left( x + \rho \frac{w}{\|w\|} \right) + b = 1$$

$$\Rightarrow \underbrace{w^T x + b}_{=0} + \rho \|w\| = 1 \quad \Rightarrow \rho = \frac{1}{\|w\|}$$

# Convex optimization

$$\max \frac{2}{\|w\|} \Rightarrow \min \frac{1}{2} \|w\|^2$$

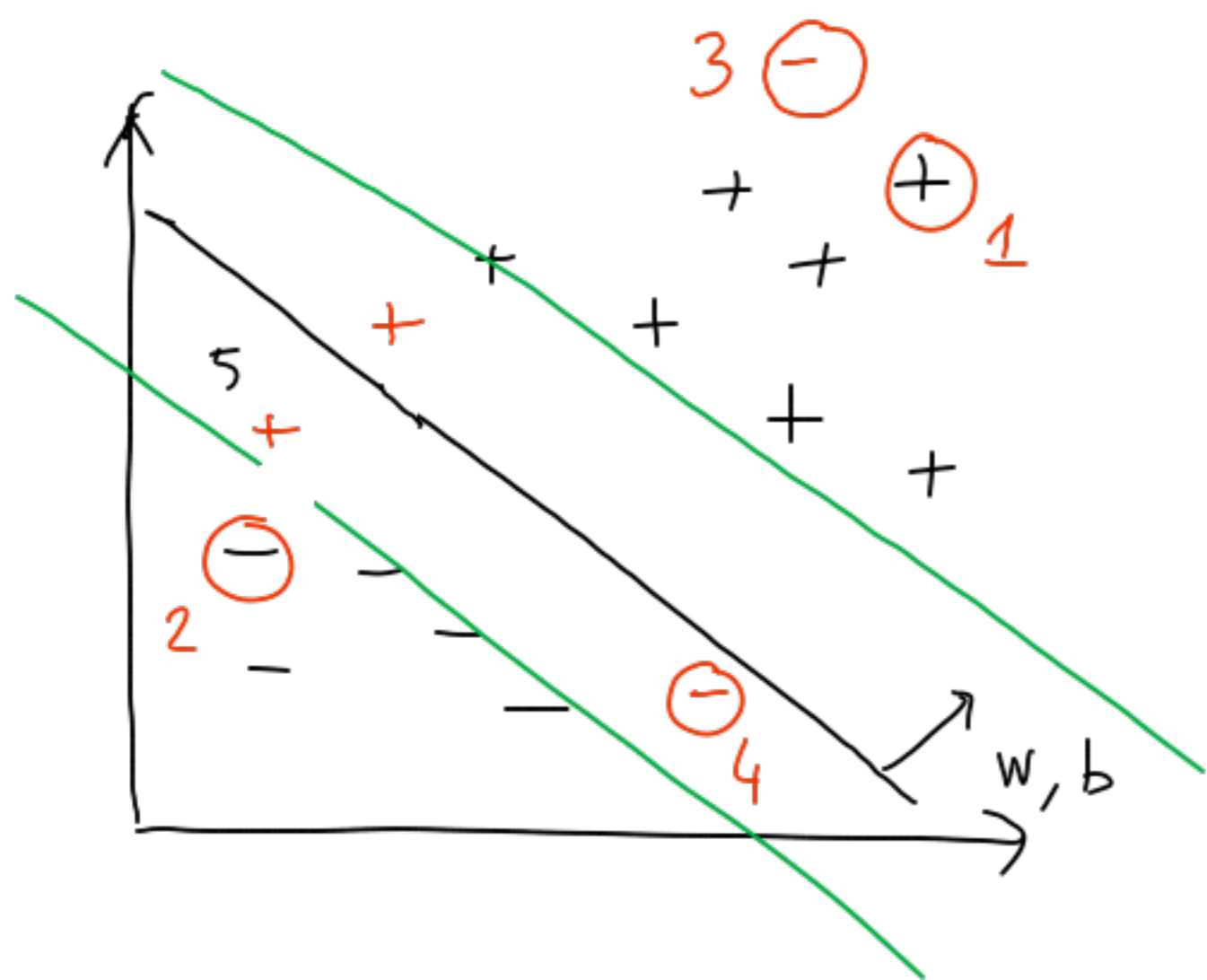
Constraints:  $s.t. \quad y_i (w^T x_i + b) \geq 1$   
 $\forall i = 1, \dots, n.$

$$\left. \begin{aligned} y_i = +1, & \Rightarrow w^T x_i + b \geq 1 \\ y_j = -1, & \Rightarrow w^T x_j + b \leq -1 \end{aligned} \right\}$$

$$y_i (w^T x_i + b) \geq 1 \quad \forall i = 1, \dots, n.$$

Remark: for  $i \in$  support vectors

$$y_i (w^T x_i + b) = 1.$$



## Soft margin SVMs

Data points are not linearly separable

Q: How to capture the degree of misclassification?

A: via a loss function.

Correct classification:  $y_i (w^T x_i + b) \geq 1$

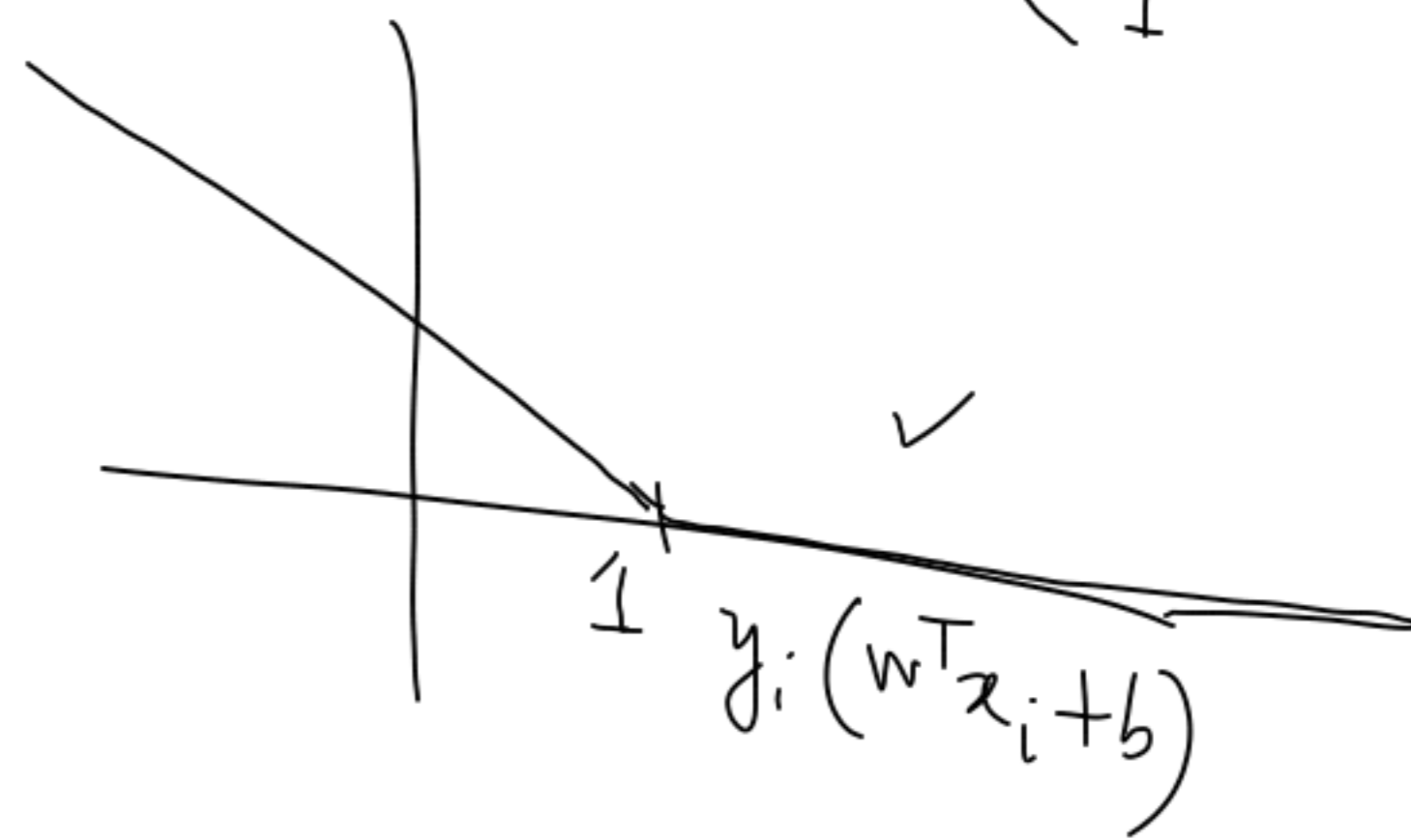
$$\max \left\{ 0, 1 - y_i (w^T x_i + b) \right\} = L_i(w, b)$$

$\leq 0$  for correct

$> 0$  for misclassification

(hinge loss)

misclassification:  $< 1$





min  
 $w, b$

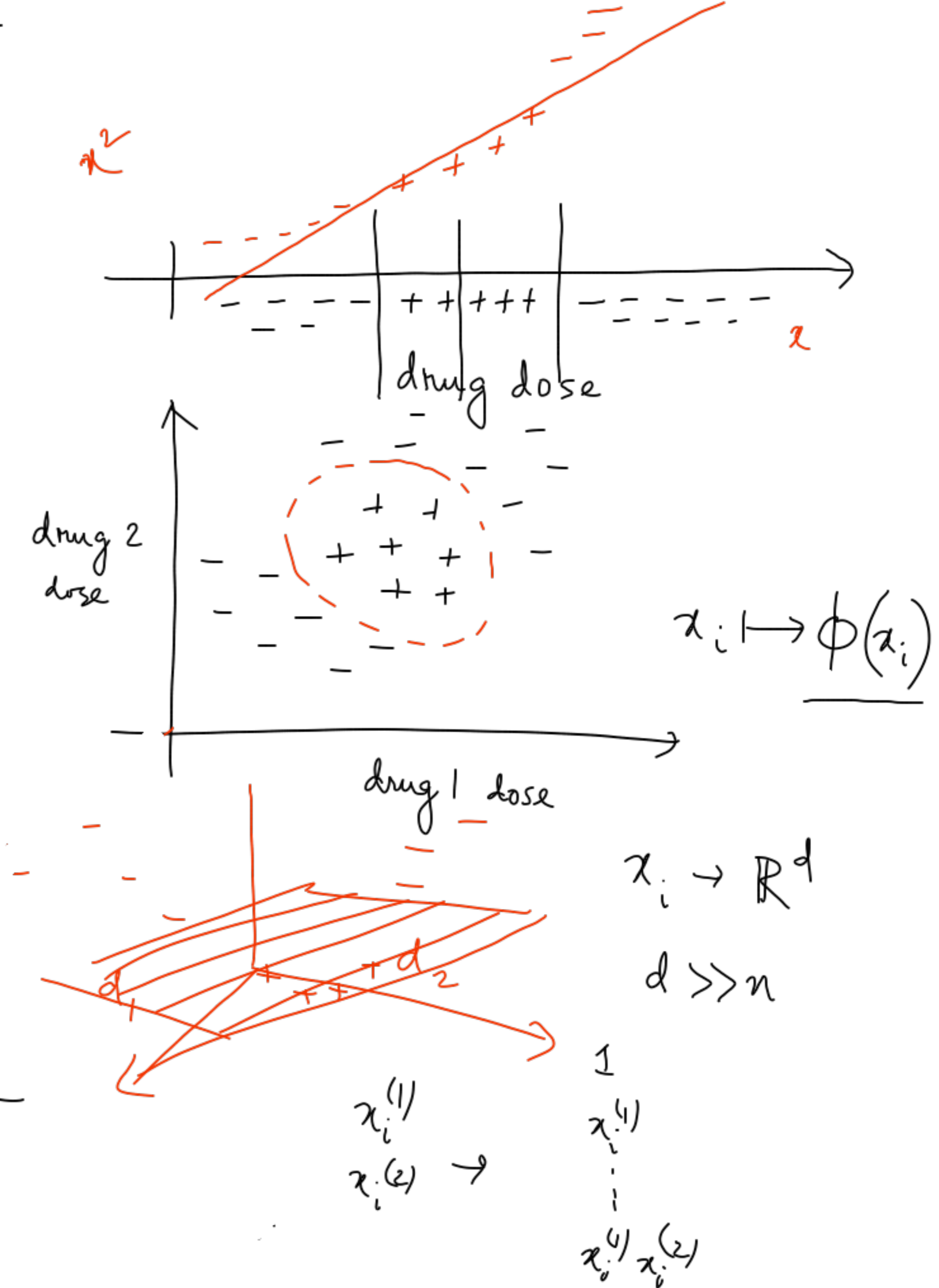
$$\frac{1}{n} \sum_{i=1}^n L_i(w, b) + \lambda \|w\|^2$$

↑ hyperparameter

Soft margin SVM.

$\lambda$  determines the tradeoff between misclassification and margin maximization.

Kernelization: a method to calculate the higher dimension trans. in a computationally efficient manner.



## Background of Kernelization

Hard margin SVM:

$$\min \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1$$

$$\forall i = 1, \dots, n.$$

$$\mathcal{L}(w, b, \lambda)$$

$$= \frac{1}{2} \|w\|^2 - \sum \lambda_i (y_i (w^T x_i + b) - 1)$$

$\lambda_i \geq 0$

In general

$$\min f_0(x)$$

PRIMAL

$$\text{s.t. } f_i(x) \leq 0, i = 1, \dots, m \leftarrow \lambda_i$$

$$h_i(x) = 0, i = 1, \dots, p \leftarrow \gamma_i$$

domain  $D$  where the feasible  $x$ 's live

$$D = \left\{ x : f_i(x) \leq 0, \forall i, h_j(x) = 0, \forall j \right\}$$

## Lagrangian

$$\mathcal{L}(x, \lambda, \gamma) = \underbrace{f_0(x)} + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^k \gamma_j \underline{h_j(x)}$$

$\lambda_i \geq 0$

Consider the problem

$$\max_{\lambda \geq 0, \gamma} \left[ \min_{x \in D} \mathcal{L}(x, \lambda, \gamma) \right]$$

EQUIVALENT to PRIMAL.

## Lagrange dual

$$\underline{g(\lambda, \gamma)} = \min_{x \in D} \mathcal{L}(x, \lambda, \gamma)$$

Dual problem  $\leq p^*$

$$\max_{\lambda \geq 0, \gamma} g(\lambda, \gamma)$$

primal optimal  
Convex opt:  
dual opt =  $p^*$



$$g(\lambda) = \min_{W, b} \mathcal{L}(W, b, \lambda)$$

$$\frac{\partial \mathcal{L}}{\partial W} = 0 \Rightarrow W = \sum_{i=1}^n \lambda_i y_i x_i \quad \text{--- ①}$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \lambda_i y_i = 0 \quad \text{--- ②}$$

use ① and ② to simplify

$$\max_{\substack{\lambda_i \geq 0 \\ \forall i}} g(\lambda) = \sum \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \underbrace{x_i^T x_j}_{\substack{\text{Kernel} \\ -i-j}}$$

$n \times n$

$d \gg n$   
 $nd \gg n^2$

$$\sum_i \lambda_i y_i = 0$$

missed this constraint in the class

Why dual?

- ① comp eff  
 $nd \gg n^2$
- ② Kernel friendly.