



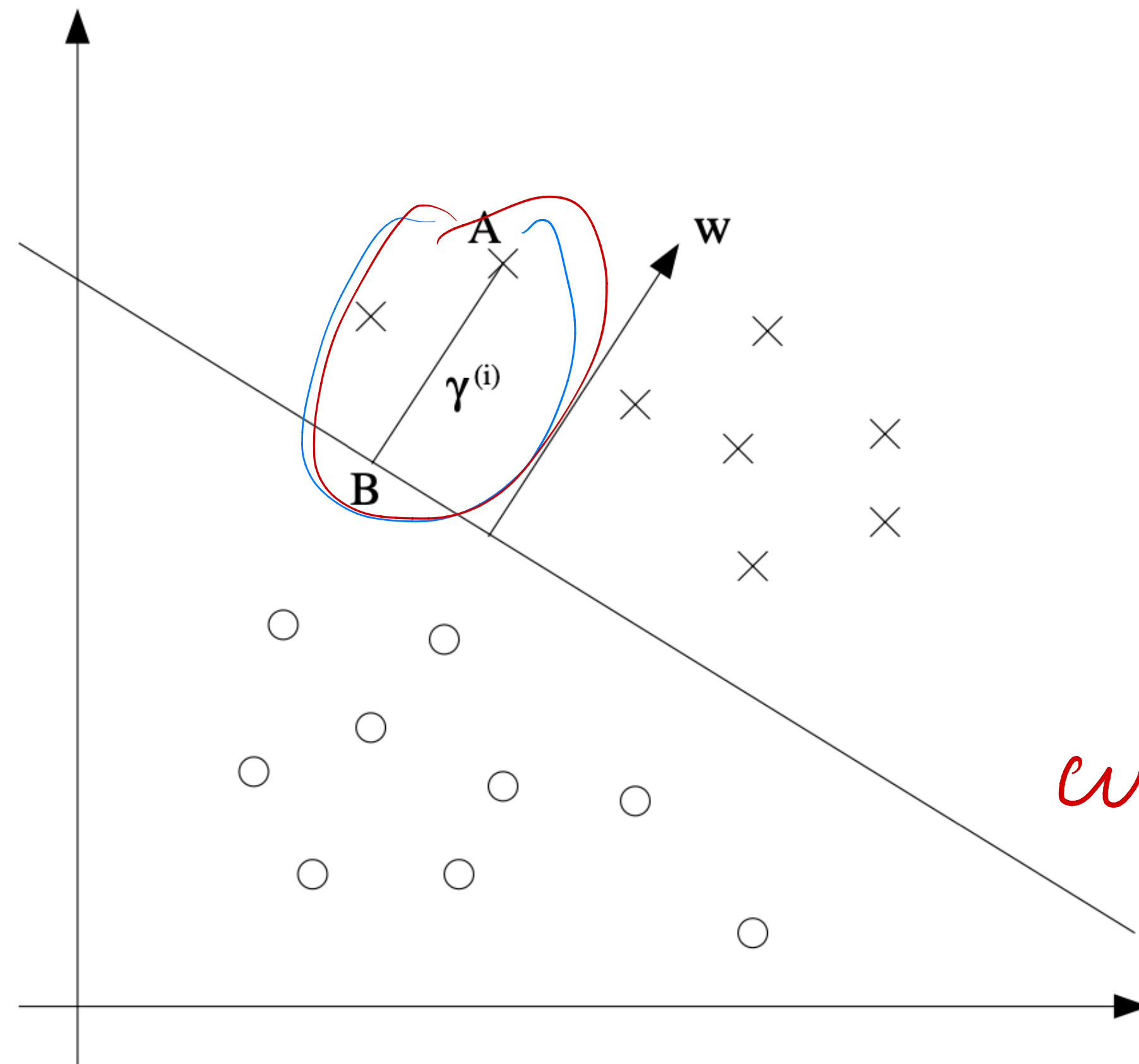
香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212  
Machine Learning  
Lecture 7

# Support Vector Machines

Junxian He  
Sep 26, 2024

# Recap: Support Vector Machines



$$w^T x + b = 0 \quad w, b$$

$$w \rightarrow Kw$$

$$b \rightarrow Kb$$

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

# Recap: The Optimization Problem

$$\max_{w,b} \min_{i=1,\dots,n} \gamma^{(i)}$$

Rewrite

$$\begin{aligned} & \max_{\gamma,w,b} \gamma \\ & \text{s.t. } y^{(i)} \left( \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right) \geq \gamma, \quad i = 1, \dots, n \end{aligned}$$

Linear constraint

$$\begin{aligned} & \max_{\hat{\gamma},w,b} \frac{\hat{\gamma}}{\|w\|} \\ & \text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq \hat{\gamma}, \quad i = 1, \dots, n \end{aligned}$$

$\frac{\hat{\gamma}}{\|w\|} = \gamma$

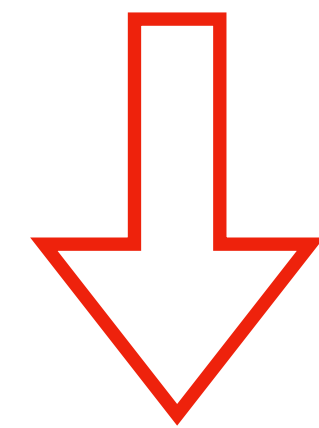
Infinite solutions, as  $\hat{\gamma}$  can be at any scale without changing the classifier

$\|w\|$  is not easy to deal with, non-convex objective

# Recap: The Optimization Problem

$$\begin{aligned} \max_{\hat{\gamma}, w, b} \quad & \frac{\hat{\gamma}}{\|w\|} \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, \quad i = 1, \dots, n \end{aligned}$$

$\alpha^*$  →  $w^*$



Add constraint  $\hat{\gamma} = 1$

This is a standard quadratic problem that can be directly solved with quadratic problem solvers

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

not kernel

Assumption: the training dataset is linearly separable

$y^{(i)}(w^T x^{(i)} + b) \geq 1$  only if prediction correct

# Recap: The Dual Problem

$$\mathcal{L}(w, \alpha, \beta) = \frac{1}{2} \|w\|^2 + \alpha \cdot g(w) + \beta h(w)$$

$\theta_D(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta)$

$\beta h(w)$   
equality

The dual optimization problem

$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta)$$

$$\min_w \mathcal{L}(w, \alpha, \beta)$$

The primal optimization problem

$$\min_w \theta_P(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta)$$

$\alpha, \beta$

$f(w)$  if constraint are satisfied  
 $\infty$  otherwise

What is the relation of the two problems?

# Recap: The Dual Problem

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1]$$

remove w, b.

The dual optimization problem

$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta)$$

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

$$\theta(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}$$

$$L(w, \alpha, \beta) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i (y^{(i)} w^T x^{(i)} + b) - 1$$

$$\min_w L(w, \alpha, \beta)$$

$$\max_{\alpha \geq 0} \min_w L(w, \alpha, \beta)$$

$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$\left( \sum_i \alpha_i y^{(i)} = 0 \right)$$

$$\min_w L(w, \alpha) = \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right)^T \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right) - \sum_i \alpha_i (y^{(i)} \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right)^T x^{(i)} + b) - 1$$

$$-\sum_i \alpha_i y^{(i)} b = 0$$

# Recap: The Dual Problem

$$\theta(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad i = 1, \dots, n$$

$\langle x^{(i)}, x^{(j)} \rangle$   
 $\langle \langle x^{(i)}, x^{(j)} \rangle \rangle$   
 $\Rightarrow \alpha^*$

$\langle \langle x^{(i)}, x^{(j)} \rangle \rangle$

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

What is the relation between solving this dual problem and solving the original problem

$w^* = \sum_{i=1}^n \alpha_i^* y^{(i)} x^{(i)}$



# The Dual Problem

*dual*

$$d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*$$

*primal*

$$\max_x \min_y f(x, y) \leq \min_y \max_x f(x, y)$$

Under certain conditions:  $d^* = p^*$  Zero-duality Gap (Strong Duality)

What are the conditions?

$$\max_x \min_y f(x, y) \leq \min_y \max_x f(x, y)$$

$\Downarrow$   
 $(x_0, y_0)$

$\Downarrow$   
 $(x_1, y_1)$

$$f(x, y_0) \leq f(x, y) \text{ for any } x, y$$

$$f(x_1, y) \geq \underline{f(x, y)} \text{ for any } x, y$$

$$f(x_0, y_0) \leq \underline{f(x_0, y_1)} \leq f(x_1, y_1)$$

# Slater's Condition

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k \\ & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

$$\underline{g_i(w) \leq 0}$$

- $f(w)$  and  $g(w)$  are convex
- $h_i(w)$  is affine (i.e. linear)
- $g_i(w)$  are strictly feasible for all  $i$ , which means there exists some  $w$  so that  $g_i(w) < 0$  for all  $i$

$\frac{1}{2} \|w\|^2$   $g(w)$  linear  
no  $h_i(w)$  for SVM

$$1 - y^{(i)} (w^T x^{(i)} + b) \leq 0$$

If Slater's condition holds, then  $d^* = p^*$

The primal optimization problem of SVM satisfies the Slater's condition

Slater's condition  $\Rightarrow$  zero duality gap  $\iff$  KKT

# KKT Conditions

Denote the solution to the primal problem as  $w^*$ , the solution to the dual problem as  $\alpha^*, \beta^*$ , then zero duality gap is sufficient and necessary (i.e. equivalent) to satisfy KKT Conditions:

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, d$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

Normal Lagrange multiplier equations

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

The original constraints

$g_i(w) \leq 0$

---

$\alpha \geq 0$

# KKT Conditions

Denote the solution to the primal problem as  $w^*$ , the solution to the dual problem as  $\alpha^*, \beta^*$ , then zero duality gap is sufficient and necessary (i.e. equivalent) to satisfy KKT Conditions:

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

$w = [a, b]$

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, d$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

If  $\alpha_i^* > 0$ , then

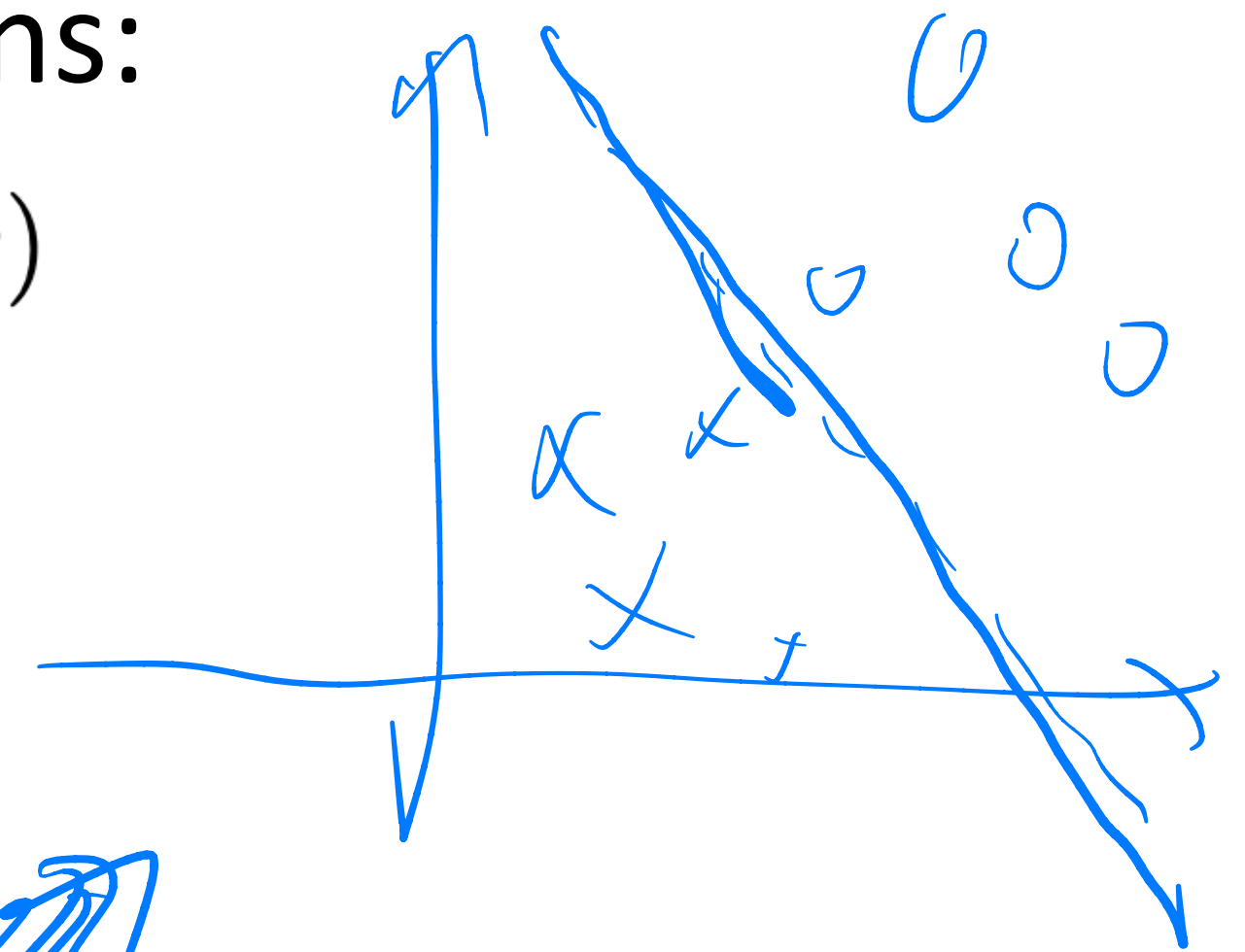
$g_i(w^*) = 0$ , the inequality is actually equality

$g_i(w^*) < 0, \alpha_i^* = 0$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha_i^* \geq 0, \quad i = 1, \dots, k$$



inequality

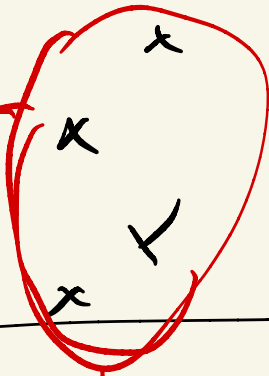
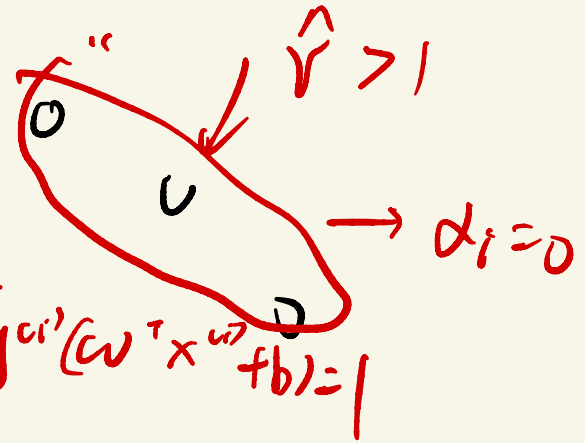
$$g_i(w) = 0$$

$$= 1 - g^{cov}(w^T x + b) \leq 0$$

$$d_i g_i(\omega^*) = 0$$

$d_i = 0$

$d_i \neq 0$



$\hat{\gamma} > 1$

$g_i(\omega) < 0$

$g_i(\omega) < 0$

$$1 - y^{(i)}(\omega^T x^{(i)} + b) < 0$$

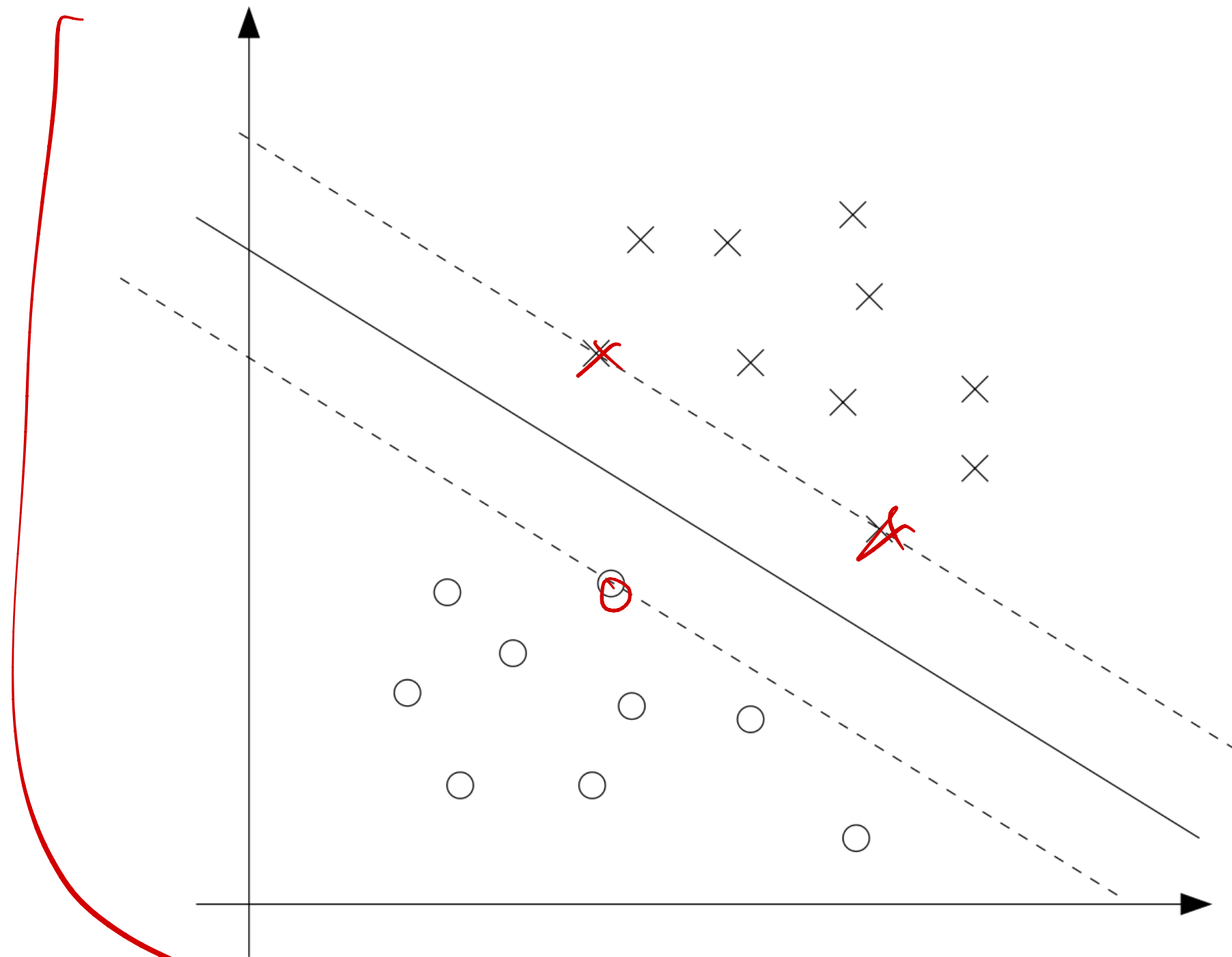
# Supporting Vectors

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$\alpha_i = 0 \quad \alpha_i \neq 0$$

imply?

Only the 3 points have non-zero  $\alpha_i$ , and they are called supporting vectors



# Lagrangian for SVM

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1]$$

The dual optimization problem

$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta)$$

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

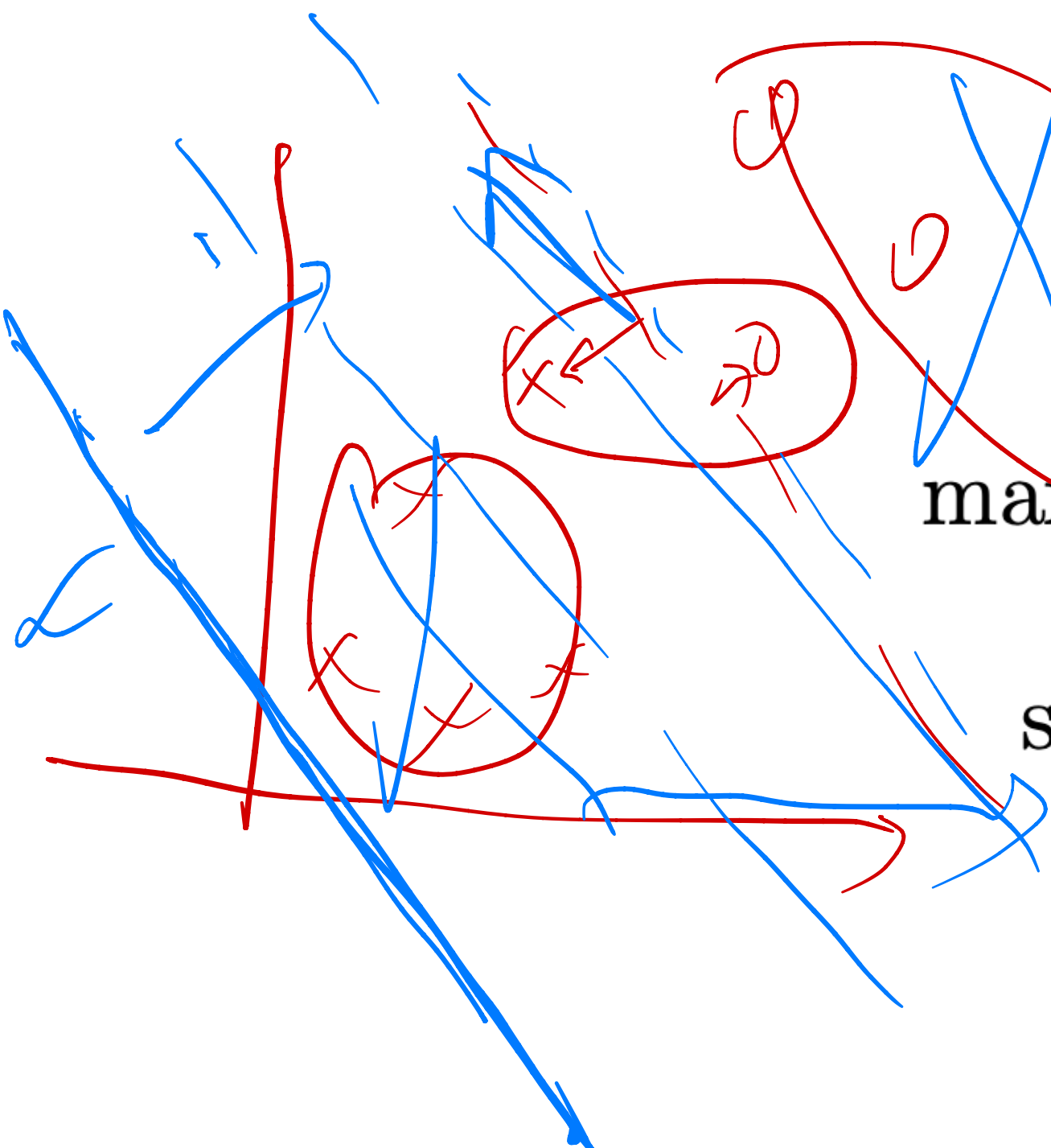
$$\theta(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}$$

$\sum_{i=1}^n \alpha_i y^{(i)} = 0$

$\alpha_i = 0$   
 $x^{(i)}$  is irrelevant



# The Dual Problem of SVM



$$\max_{\alpha} W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } \alpha_i \geq 0, \quad i = 1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0,$$

Kernel is all we need!

$\checkmark \quad \langle C x^{(i)}, x^{(j)} \rangle$

After solving  $\alpha$  (coordinate ascent with clipping, 6.8.2 of the CS229 Notes)

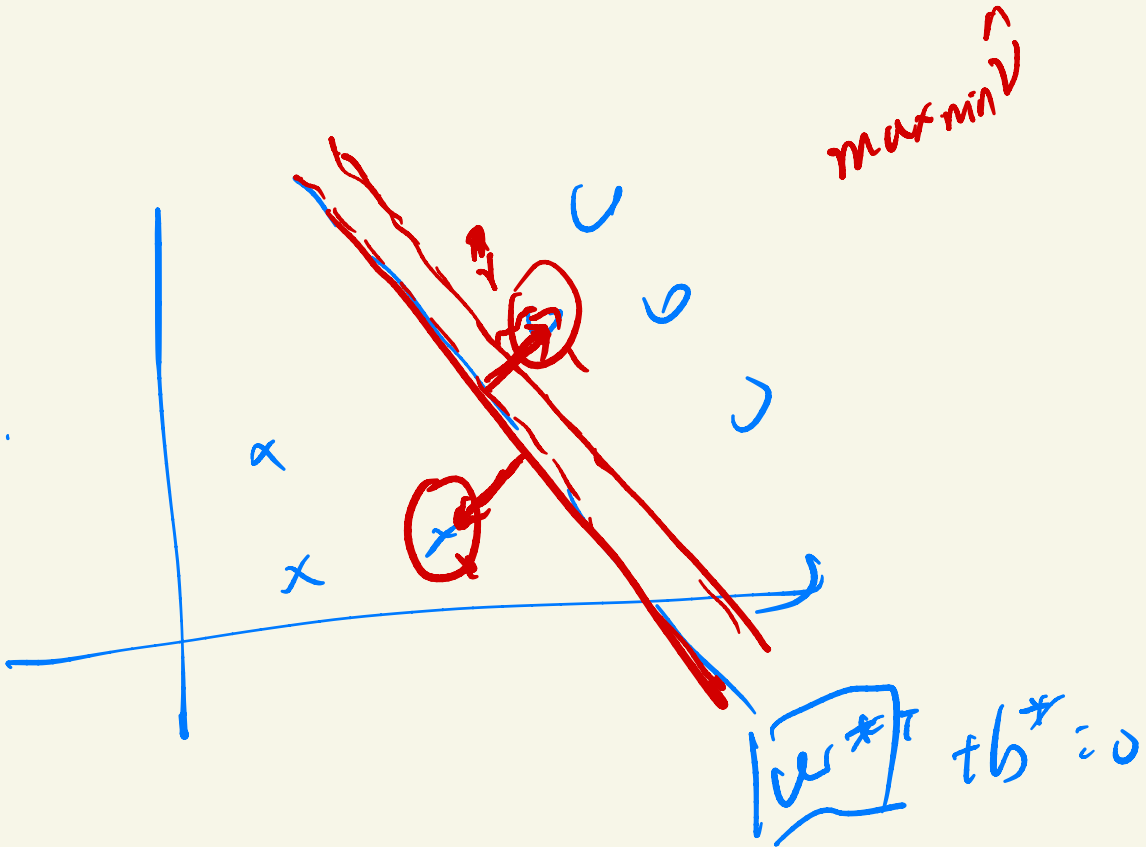
$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

From KKT Conditions

$$b^* = \frac{\max_{i:y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i:y^{(i)}=1} w^{*T} x^{(i)}}{2}$$

From the original constraints

$$y^{(i)} (w^T x^{(i)} + b) = 1$$



# Inference

$$w^T x + b = \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right)^T x + b$$
$$= \sum_{i=1}^n \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b.$$

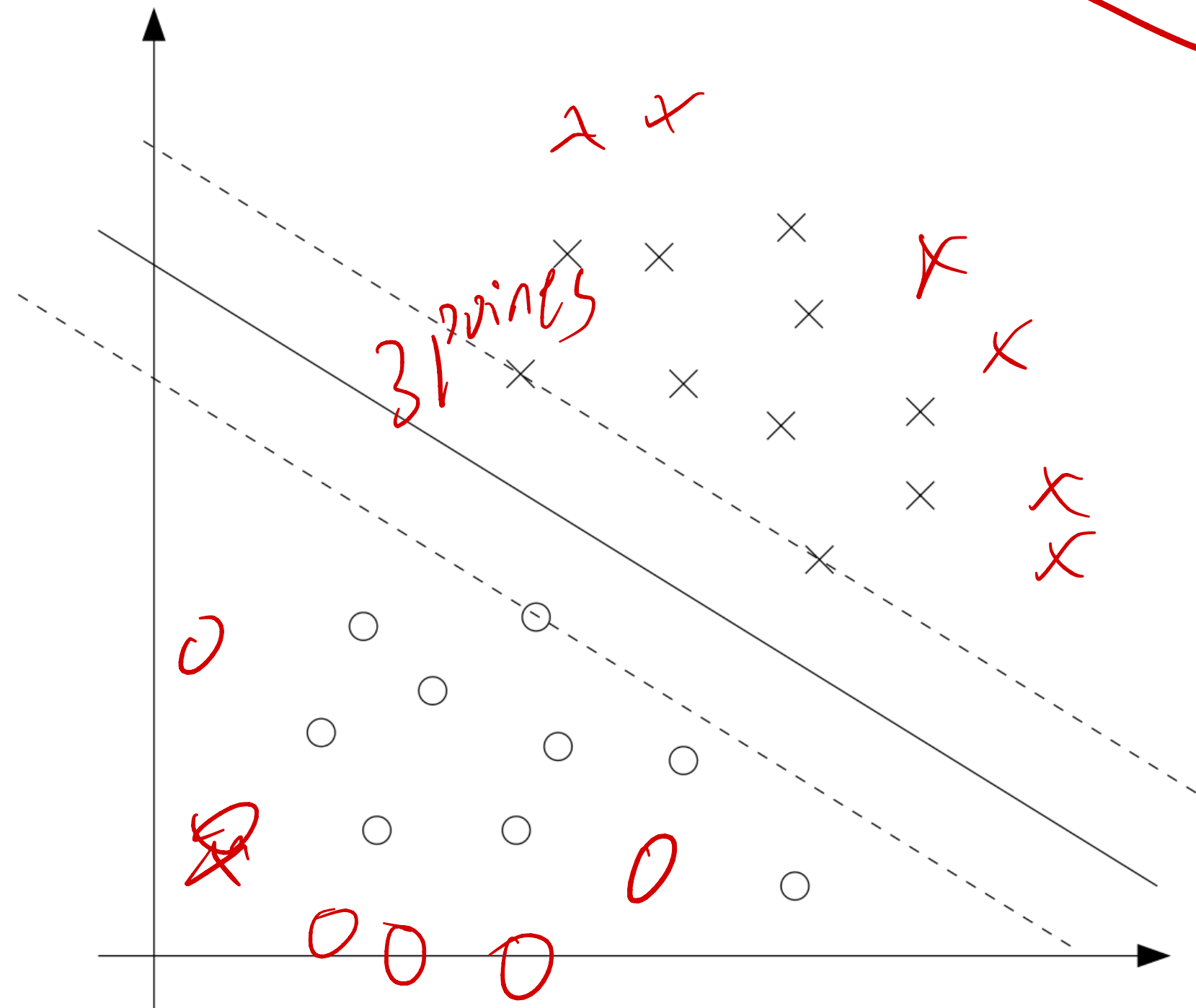
Support

3 terms

We never need to really compute  $w$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

Most  $\alpha_i$  are 0, only the supporting examples will influence the final prediction

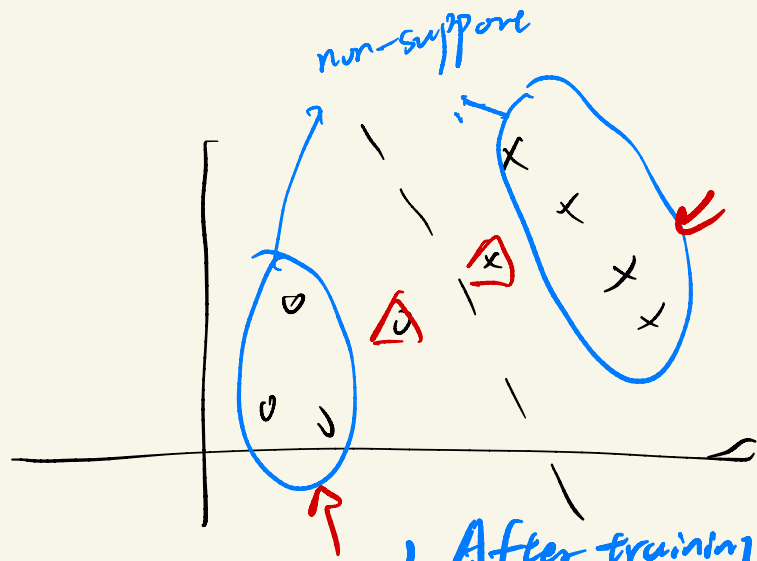




$$\hat{y} = y (\omega^T x + b)$$

$$\hat{y} = 1$$

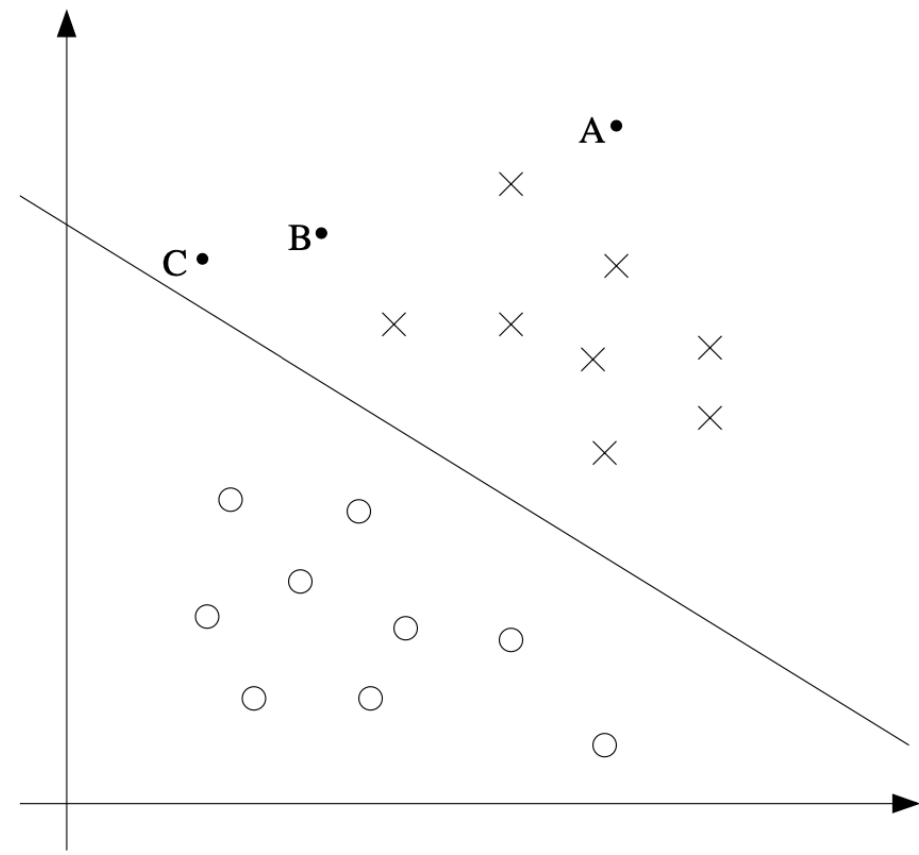
$$r = \frac{1}{\|\omega\|}$$



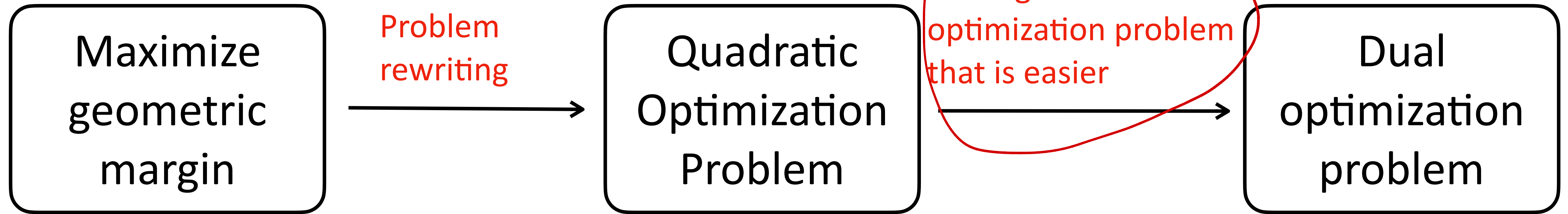
1. After training, if we delete all the non-support data samples. prediction change? **No**
2. If we know the non-support samples in advance, and we delete in the beginning, **Yes** decision boundary change? **No**
3. Same to 2, whether training process change? **Yes** dynamics?



# Review of the High-Level Logic



$$h_{w,b}(x) = g(w^T x + b).$$



$$\gamma^{(i)} = y^{(i)} \left( \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right)$$

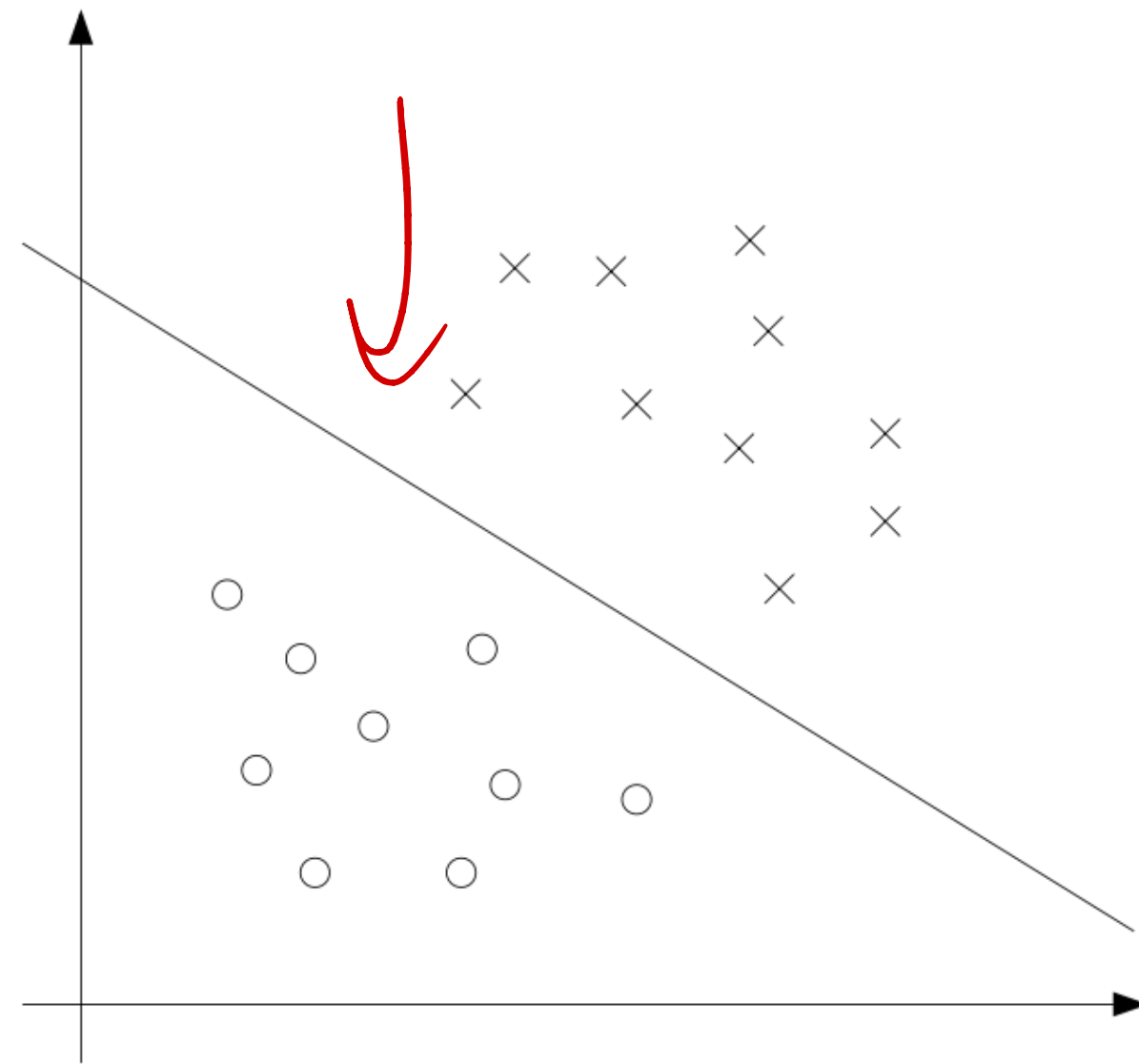
$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)} (w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \end{aligned}$$

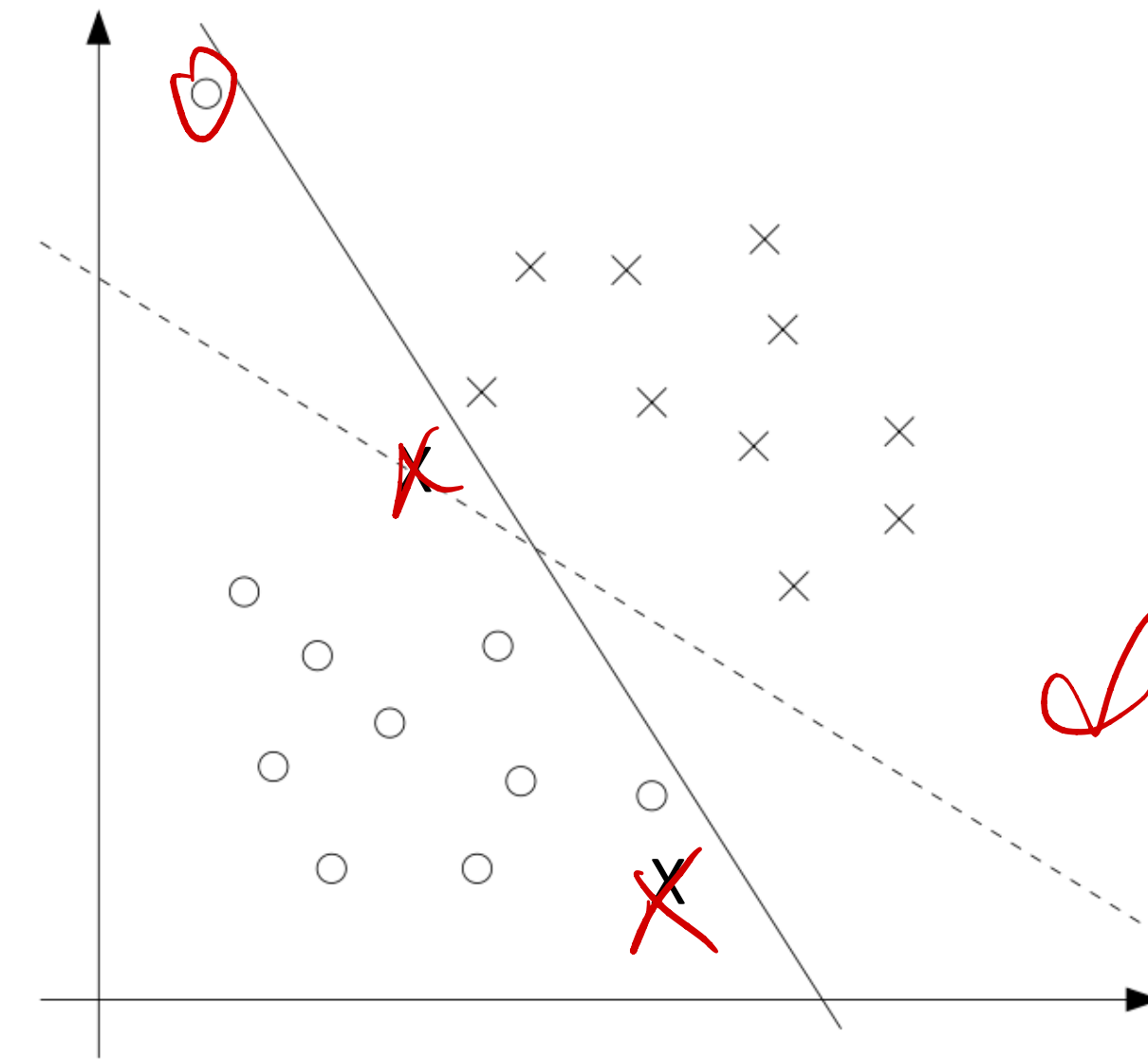
Not suitable for non-linear cases (high-dim feature map)

Kernel makes it very flexible in non-linear cases!

# The Non-Separable Case



Linearly Separable



Linearly Non-Separable



# The Non-Separable Case

Primal opt problem:

$$\min_{\gamma, w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$y^{(i)} (w^T x^{(i)} + b) < 0$

$\xi_i$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n.$$

Dual opt problem

You will prove this in your hw

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0,$$

$0 \leq \alpha_i \leq C$

**Thank You!**  
**Q & A**