



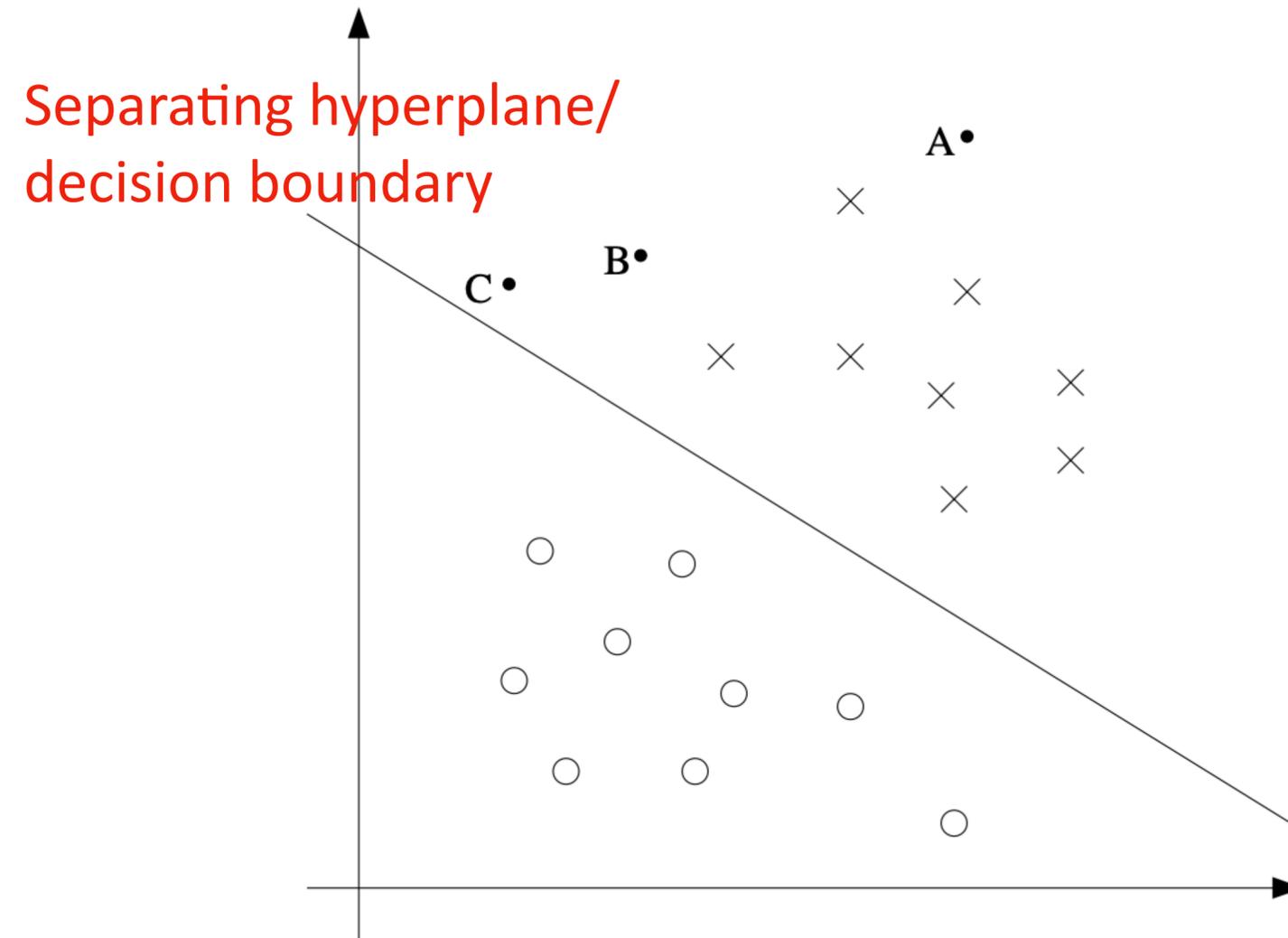
香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 5212
Machine Learning
Lecture 6

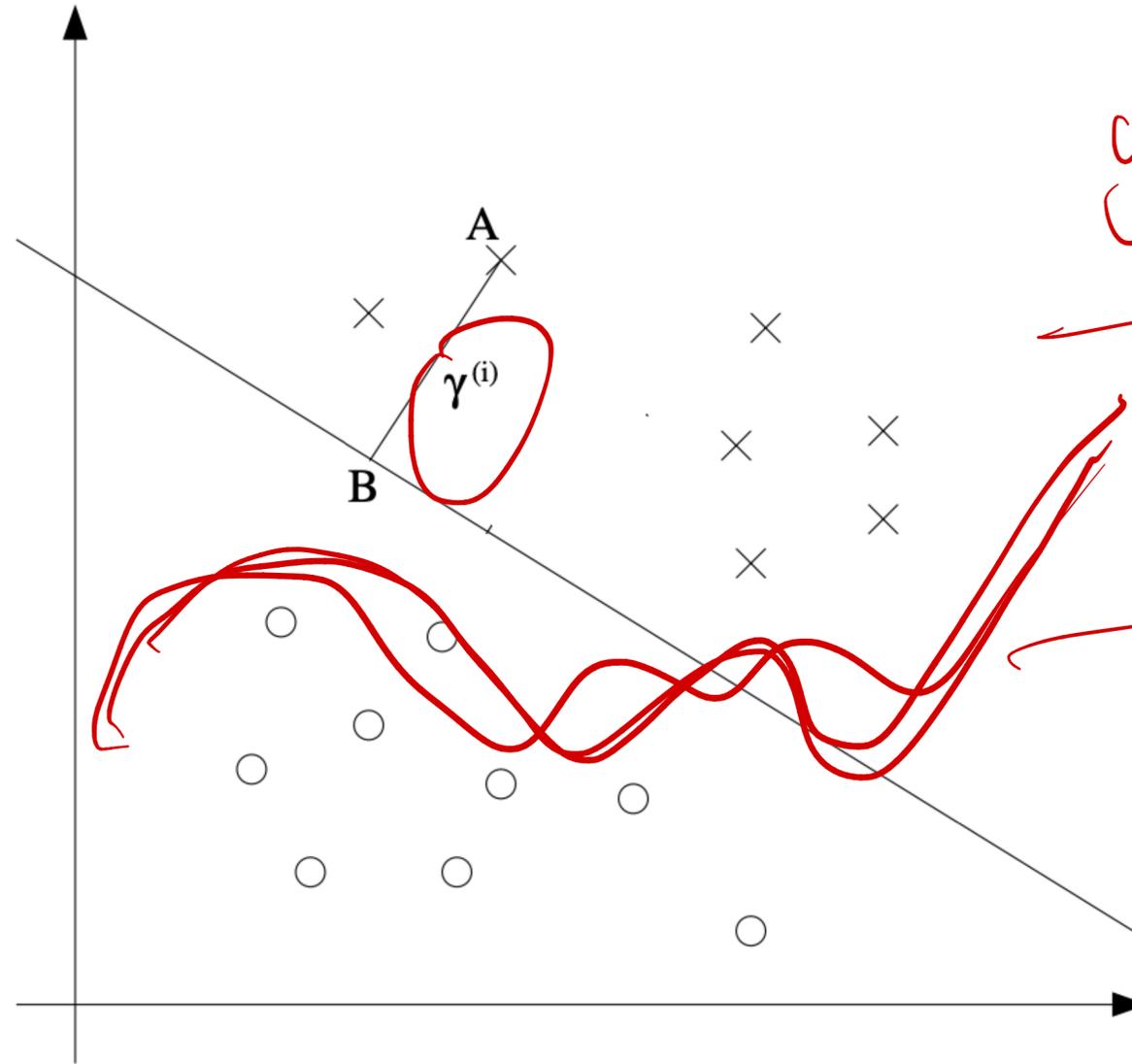
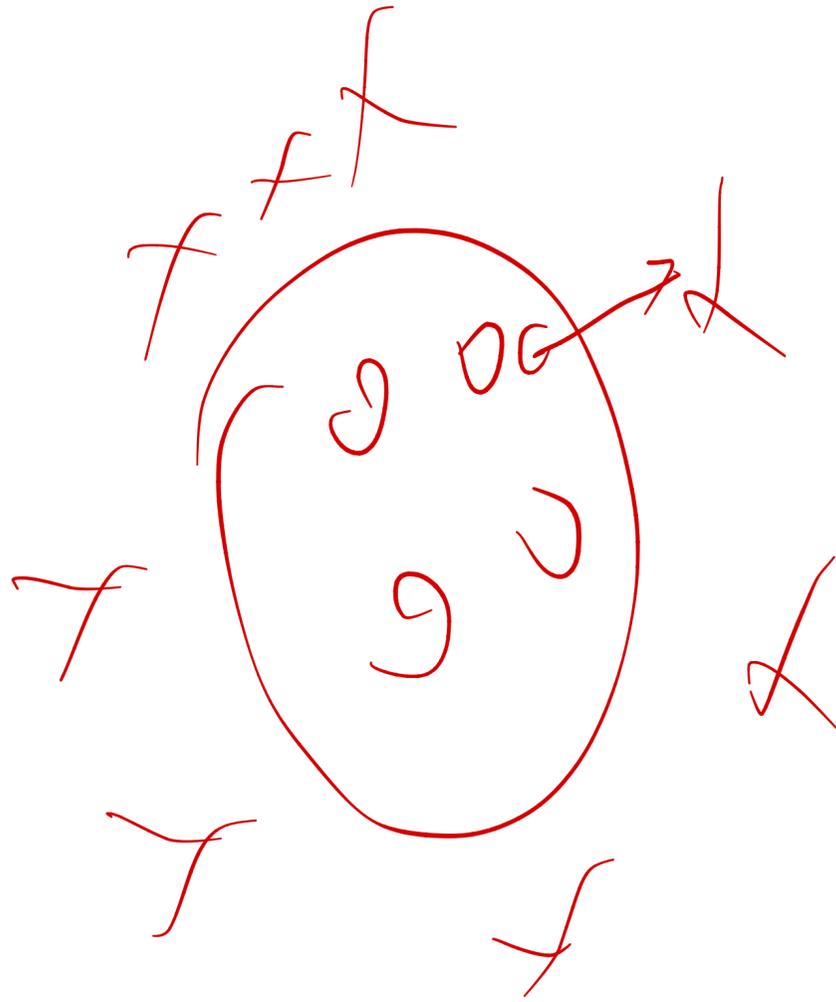
Support Vector Machines

Junxian He
Mar 3, 2025

Recap: Support Vector Machines



Recap: Geometric Margin



$$y(w \cdot x + b)$$

$$\|w\|$$

What is the geometric margin?

Recap: Functional Margin

Given a training example $(x^{(i)}, y^{(i)})$

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b).$$

Recap: Functional Margin

Given a training example $(x^{(i)}, y^{(i)})$

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b).$$

Handwritten red equation: $\frac{\hat{\gamma}}{\|w\|} = \gamma$. The γ is circled.

Given a training set $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$

$$\hat{\gamma} = \min_{i=1, \dots, n} \hat{\gamma}^{(i)}$$

Recap: Functional Margin

Given a training example $(x^{(i)}, y^{(i)})$

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b).$$

$$w \rightarrow 2w$$

$$b \rightarrow 2b$$

Given a training set $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$

$$\hat{\gamma} = \min_{i=1, \dots, n} \hat{\gamma}^{(i)}$$

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

Functional margin changes rescaling parameters, making it a bad objective, e.g. when $w \rightarrow 2w$, $b \rightarrow 2b$, the functional margin changes while the separating plane does not really change

Recap: Geometric Margin

Given a training set $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$

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

max γ
min $i \in S, \gamma^{(i)}$

Recap: The Optimization Problem

Recap: The Optimization Problem

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

Recap: The Optimization Problem

Infinite solutions, as \hat{y} can be at any scale without changing the classifier
 $||w||$ is not easy to deal with, non-convex objective

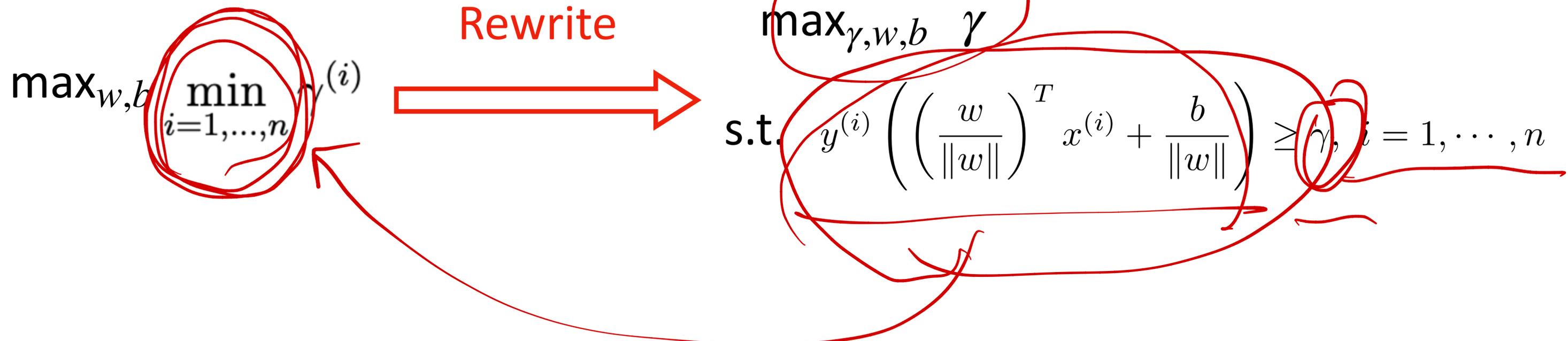
Recap: The Optimization Problem

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

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



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

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

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

Linear constraint



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

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

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

Linear constraint



$$\max_{\hat{\gamma},w,b} \frac{\hat{\gamma}}{\|w\|}$$

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

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

Recap: The Optimization Problem

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

Recap: The Optimization Problem

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

$$\hat{\gamma} = 2$$

$$\|w\| = 1$$

$$w, b, \gamma$$

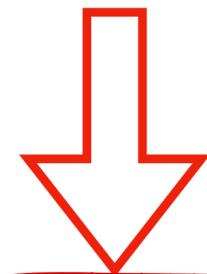
$$kw, kb$$

$$\Rightarrow k\gamma$$

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

Recap: The Optimization Problem

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



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

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

max $\|w\|$

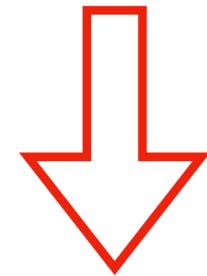
min $\frac{1}{2} \|w\|^2$

Recap: The Optimization Problem

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

Kernel

$\|w\|^2$



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

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

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

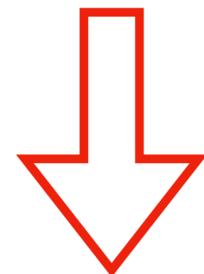
$\{x^{(i)}, y^{(i)}\}$

linear

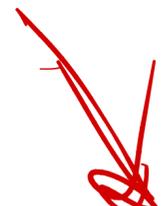
$\|w\|$

Recap: The Optimization Problem

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



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


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

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

Assumption: the training dataset is linearly separable

The Dual Problem in Optimization

In optimization, sometimes the primal optimization is hard to solve, then we may find a related alternative optimization problem that can be solved more easily, to solve the original problem in an indirect way

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

Quadratic Program

Quadratic Program

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

Quadratic Program

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

This is already a standard convex opt problem that is ready to be solved, why are we doing all the rest of things?

Generalized Lagrangian

Generalized Lagrangian

Primal optimization problem

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

Generalized Lagrangian

Primal optimization problem

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

Generalized Lagrangian

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

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

Generalized Lagrangian

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

Generalized Lagrangian

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

$\rho \rightarrow$ primal

$$\theta_{\mathcal{P}}(w) = \max_{\alpha, \beta : \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Generalized Lagrangian

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

$\min_w \theta_P(w)$

$$\theta_P(w) = \max_{\alpha, \beta: \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

$$g_i(w) \leq 0$$

$$h_i(w) = 0$$

$$\theta_P(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise.} \end{cases}$$

$$\min_w f(w) \\ \text{s.t. } g_i(w) \leq 0 \\ h_i(w) = 0$$

11

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

The Dual Problem in Optimization

In optimization, sometimes the primal optimization is hard to solve, then we may find a related alternative optimization problem that can be solved more easily, to solve the original problem in an indirect way

The Dual Problem

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

D → dual

The Dual Problem

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

The dual optimization problem

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

The Dual Problem

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

dual opt
problem

The dual optimization problem

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

α, β

The primal optimization problem

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

w

The Dual Problem

The Dual Problem

Lagrangian function

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

The Dual Problem

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

The dual optimization problem

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


The Dual Problem

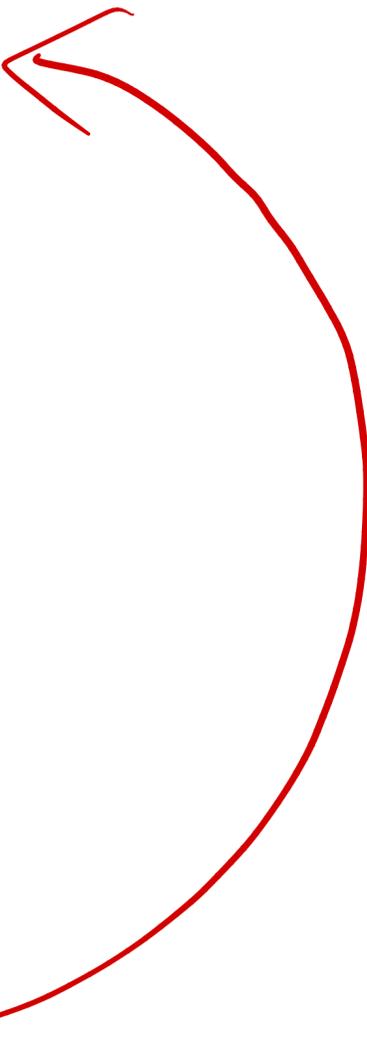
$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1]$$

The dual optimization problem

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

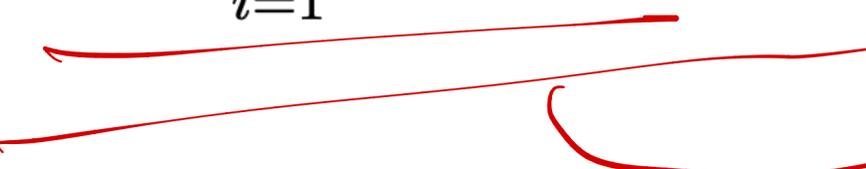
$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

The Dual Problem

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


The dual optimization problem

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

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


The Dual Problem

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

The dual optimization problem

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

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


The Dual Problem

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

The dual optimization problem

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

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

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

dual objective

The Dual Problem

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

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

kernel

convex
quadratic

program

$K(x^{(i)}, x^{(j)})$

The Dual Problem

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

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

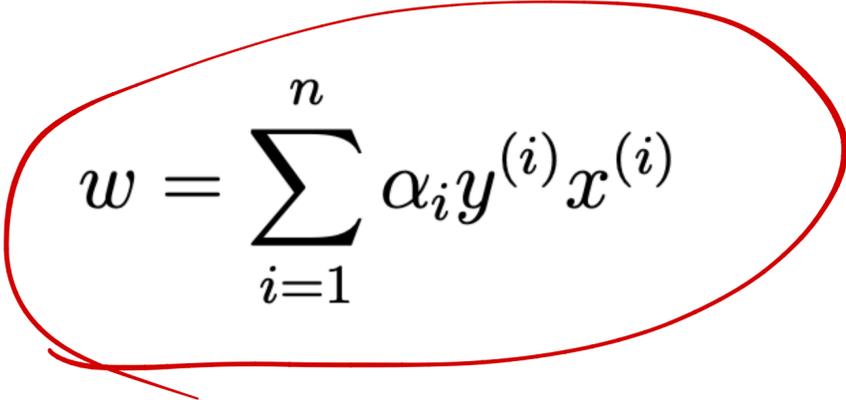
$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

The Dual Problem

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

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

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


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

The Dual Problem

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

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

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

The Dual Problem

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

s.t. $\alpha_i \geq 0, i = 1, \dots, n$

kernel

α

$n \alpha_1, \alpha_2, \dots, \alpha_n$

w^*

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

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

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

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

optimal w for the original?

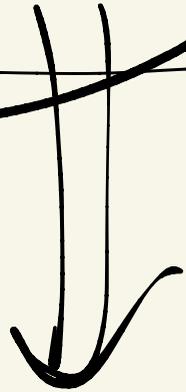
The Dual Problem

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

$$d^* \leq p^*$$

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

$$\Downarrow$$
$$\underline{f(x_0, y_0)}$$

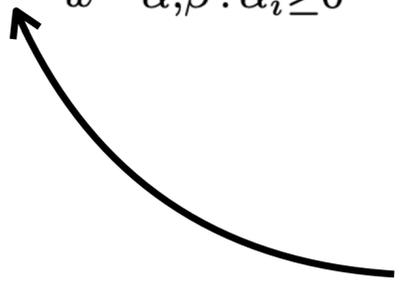


$$\Downarrow$$
$$f(x_1, y_1)$$

$$\left. \begin{array}{l} \underline{f(x, y_0)} \leq \underline{f(x, y)} \text{ for any } x, y \\ \underline{f(x_1, y)} \geq \underline{f(x, y)} \text{ for any } x, y \end{array} \right\} \leftarrow$$
$$\underline{f(x_0, y_0)} \leq \underline{f(x_1, y_0)} \leq \underline{f(x_1, y_1)}$$

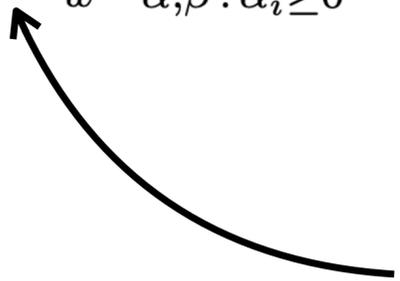
The Dual Problem

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

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


The Dual Problem

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

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


Under certain conditions: $d^* = p^*$

The Dual Problem

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

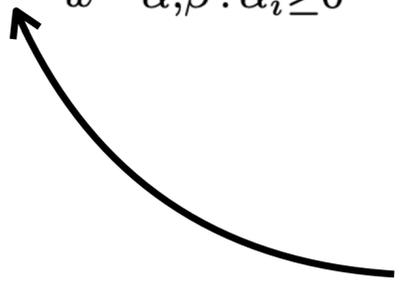
$$\max_x \min_y f(x, y) \leq \min_y \max_x f(x, y)$$

Under certain conditions: $d^* = p^*$

Zero-duality Gap

The Dual Problem

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

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


Under certain conditions: $d^* = p^*$

Zero-duality Gap



What are the conditions?

Slater's Condition

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

Slater's Condition

$$\min_w f(w)$$

$$\text{s.t. } g_i(w) \leq 0, \quad i = 1, \dots, k$$

$$h_i(w) = 0, \quad i = 1, \dots, l.$$

Handwritten note: $f(w) = \frac{1}{2} \|w\|^2$

- $f(w)$ and $g(w)$ are convex
- $h_i(w)$ is affine (i.e. linear) \rightarrow *no equality*
- $g_i(w)$ are strictly feasible for all i , which means there exists some w so that $g_i(w) < 0$ for all i

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

$$\boxed{1 - y^{(i)} (w^T x^{(i)} + b) < 0} \quad i = 1, 2, \dots, n$$

Slater's Condition

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

- $f(w)$ 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

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

Slater's Condition

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

- $f(w)$ 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

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

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

KKT Conditions

KKT Conditions

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

KKT Conditions

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

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

KKT Conditions

→ bridge

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

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

$w = (a, b)$
↓

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

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

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

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

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

KKT Conditions

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

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

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

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

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

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

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

Normal Lagrange multiplier equations

KKT Conditions

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

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

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

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

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

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

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

Normal Lagrange multiplier equations

The original constraints

And

KKT Conditions

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

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

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

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

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

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

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

KKT Conditions

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

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

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

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

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

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

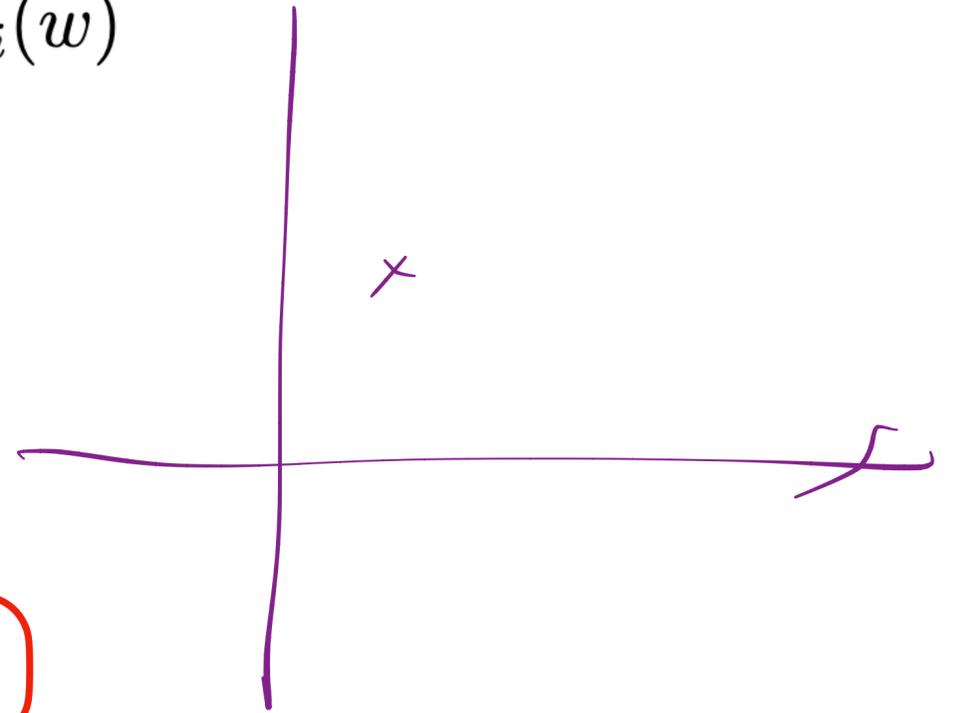
$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

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

α_i $\alpha \geq 0$

If $\alpha_i^* > 0$, then

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



Supporting Vectors

Supporting Vectors

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

Supporting Vectors

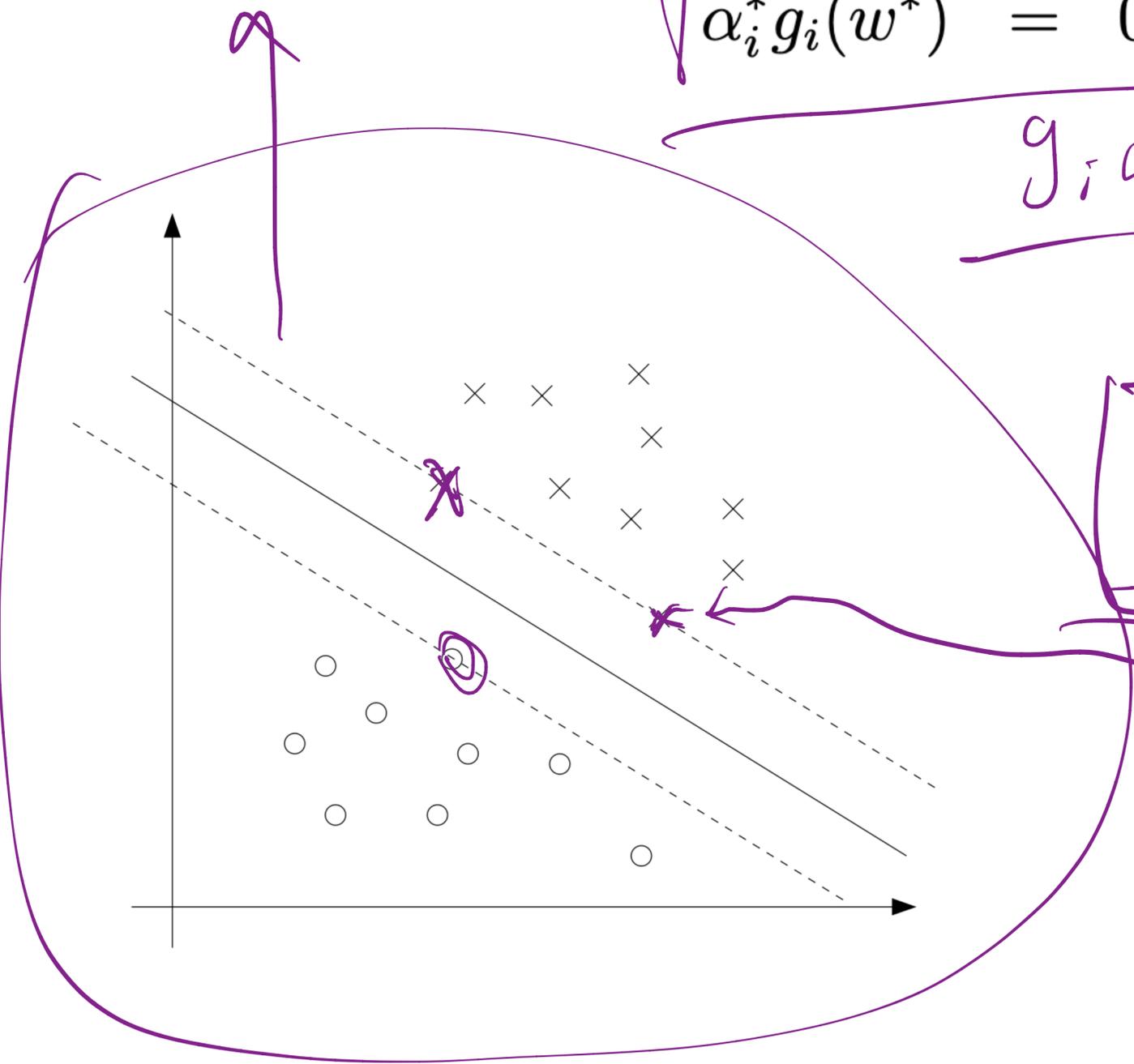
$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) = 0$$

$$i = 1$$

min
 $i = 1, 2, \dots, n$
 $y^{(i)}$

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



$x^{(i)}$

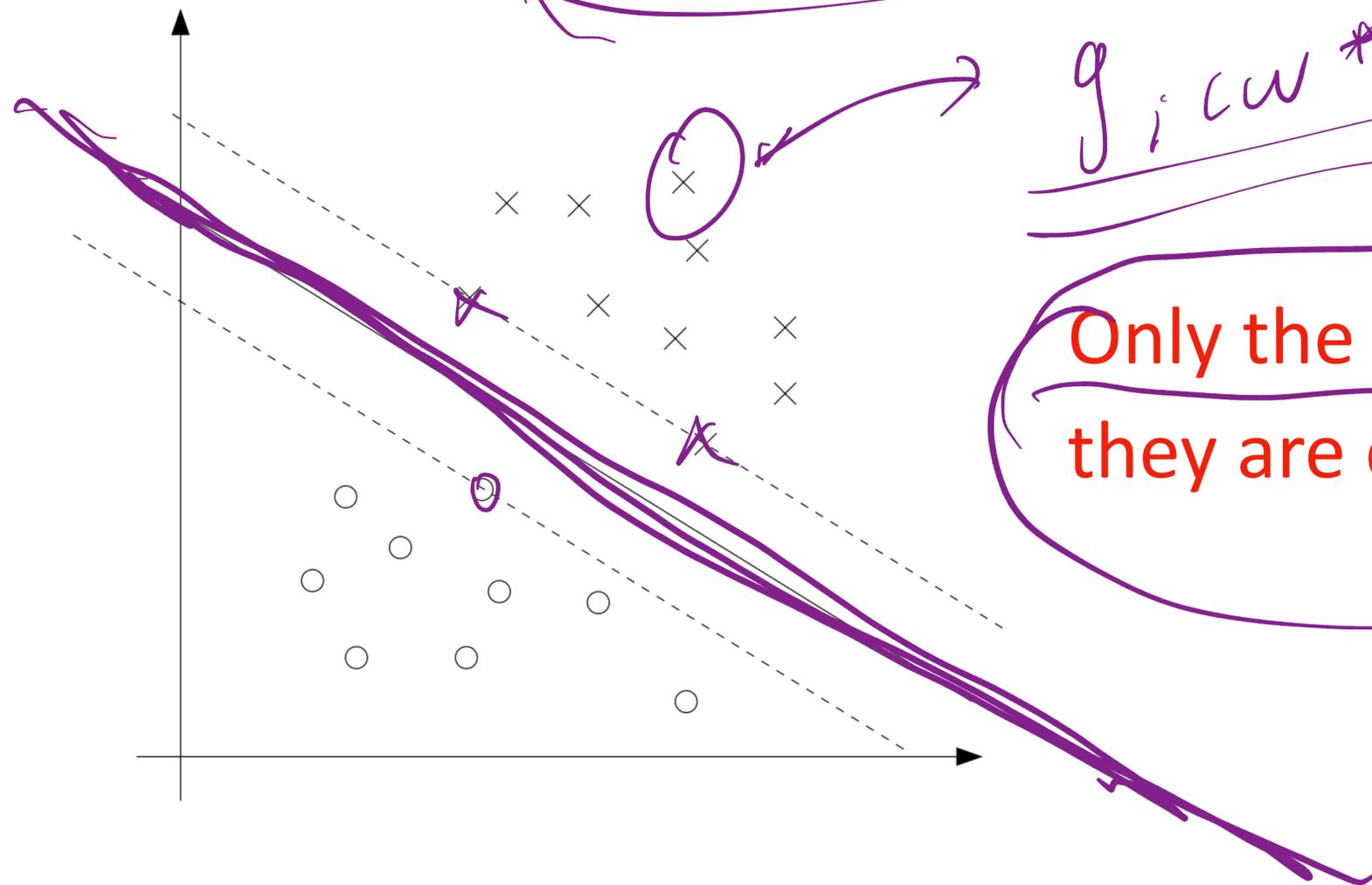
Supporting Vectors

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

Handwritten notes: $g_i(w^*) \neq 0$ (circled), α_i^* (circled), and arrows pointing from the equation to the plot below.

$$g_i(w^*) > 0 \Rightarrow \alpha_i^* = 0$$

Handwritten note: $g_i(w^*)$ (circled), α_i^* (circled).



Only the 3 points have non-zero α_i , and they are called supporting vectors

The Dual Problem of SVM

The Dual Problem of SVM

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

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

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

→ linear

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!

α_i, α_j

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!

After solving α (we'll talk about how later)

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!

After solving α (we'll talk about how later)

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

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!

After solving α (we'll talk about how later)

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

From KKT Conditions

The Dual Problem of SVM

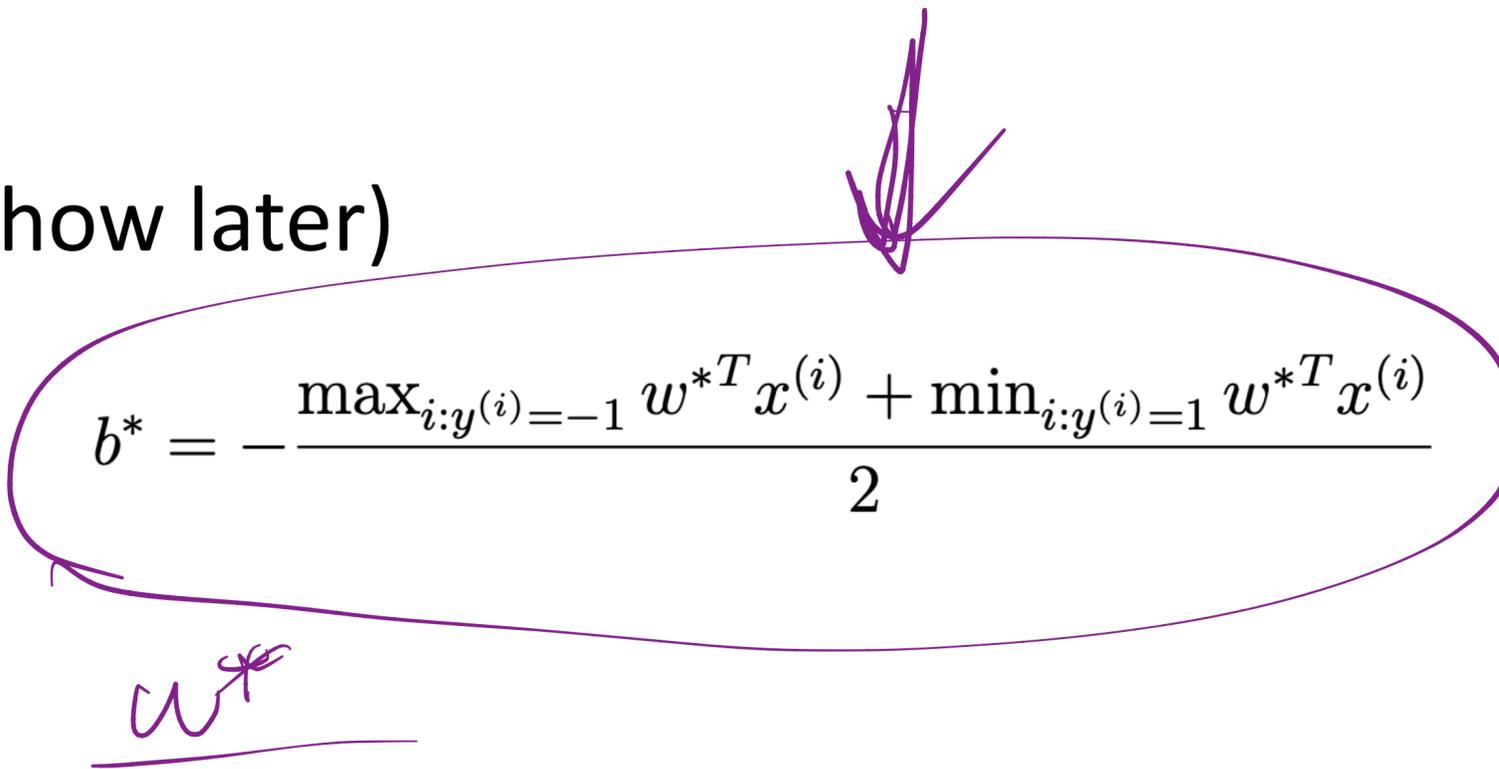
$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \end{aligned}$$

Kernel is all we need!

After solving α (we'll talk about how later)

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

From KKT Conditions

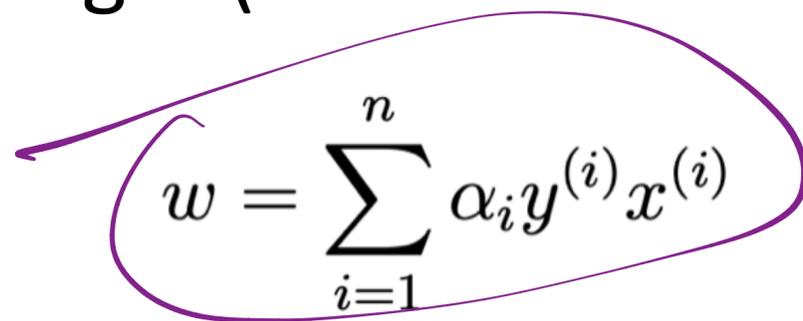

$$b^* = - \frac{\max_{i:y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i:y^{(i)}=1} w^{*T} x^{(i)}}{2}$$

The Dual Problem of SVM

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

Kernel is all we need!

After solving α (we'll talk about how later)


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

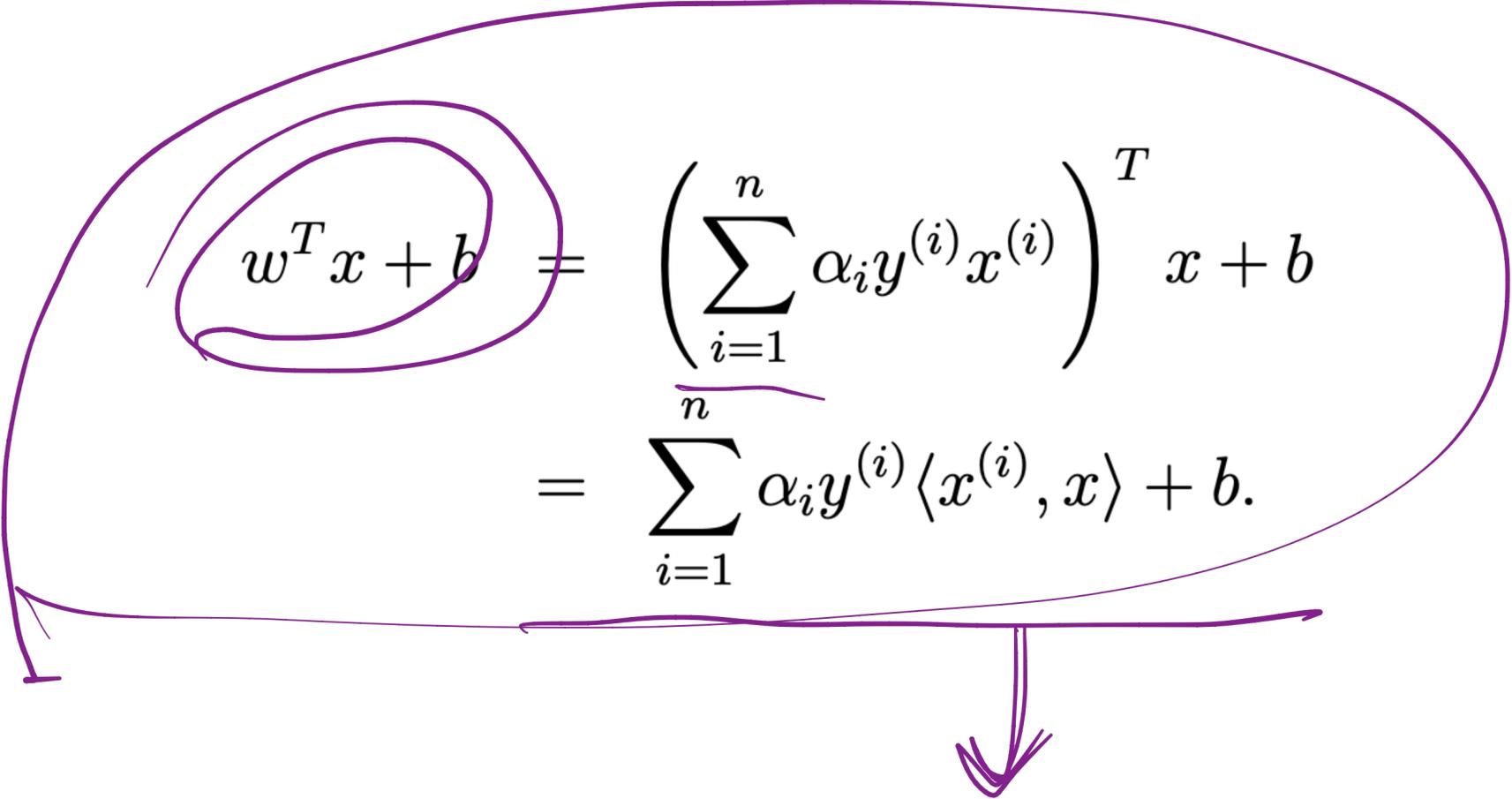
From KKT Conditions

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

From the original constraints

Inference

Inference

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


Inference

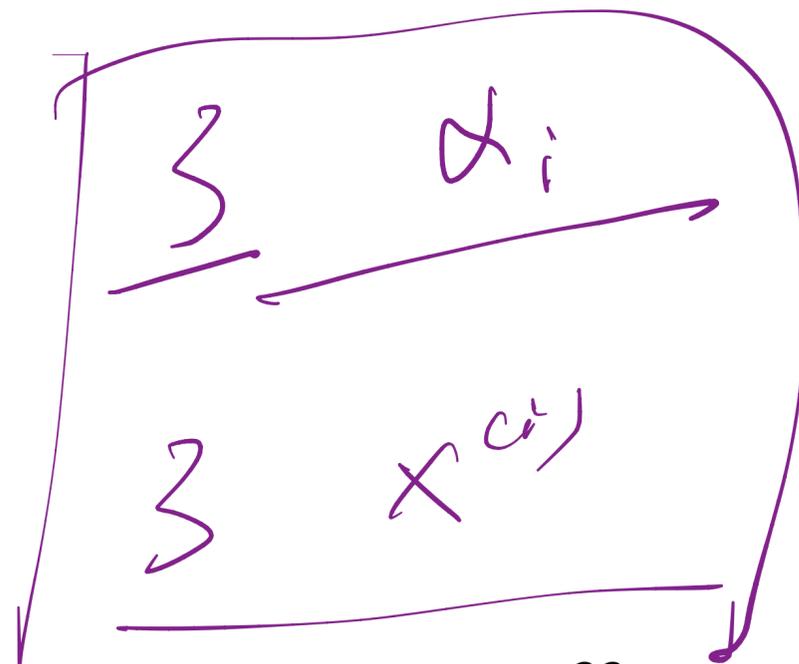
$$w^T x + b = \left(\sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right)^T x + b$$

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

kernel

Supporting
vector

We never need to really compute w



Inference

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

We never need to really compute w

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

Inference

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

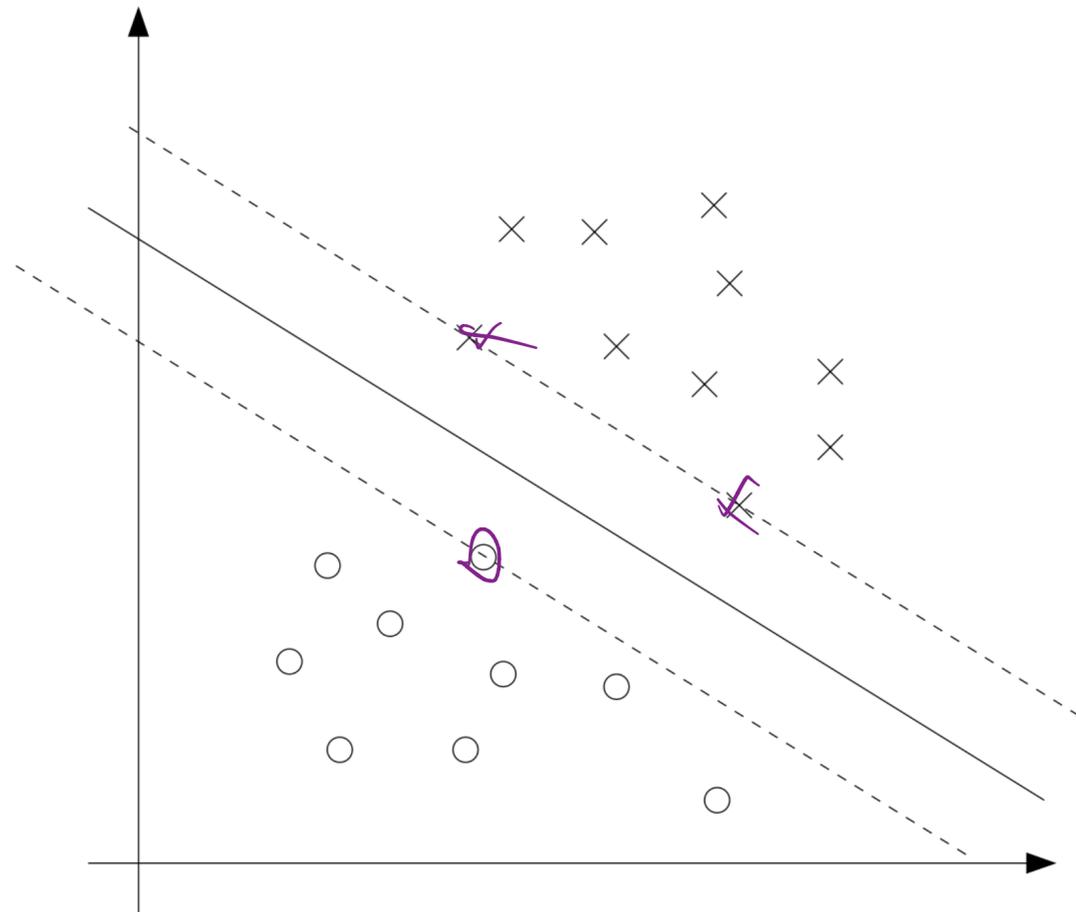
We never need to really compute w

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

Most α_i are 0, only the supporting examples will influence the final prediction

Inference

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



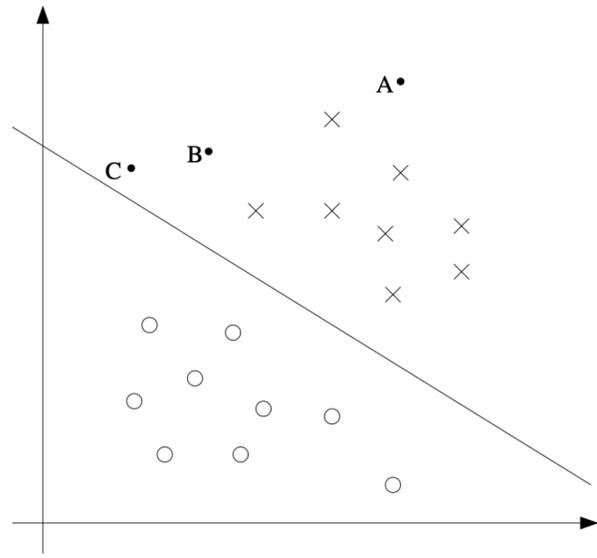
We never need to really compute w

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

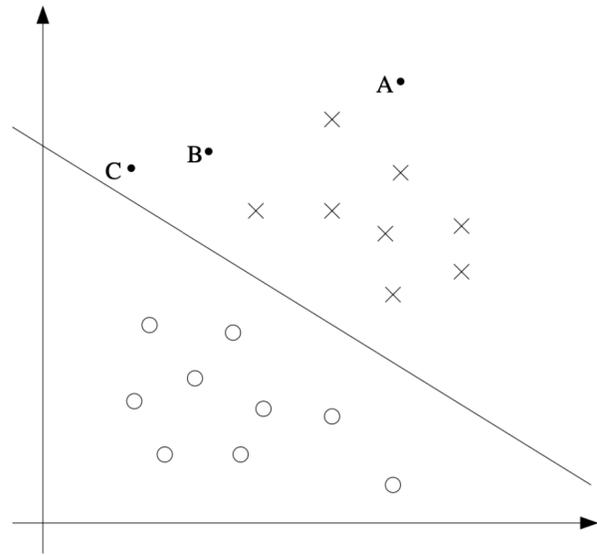
Most α_i are 0, only the supporting examples will influence the final prediction

Review of the High-Level Logic

Review of the High-Level Logic

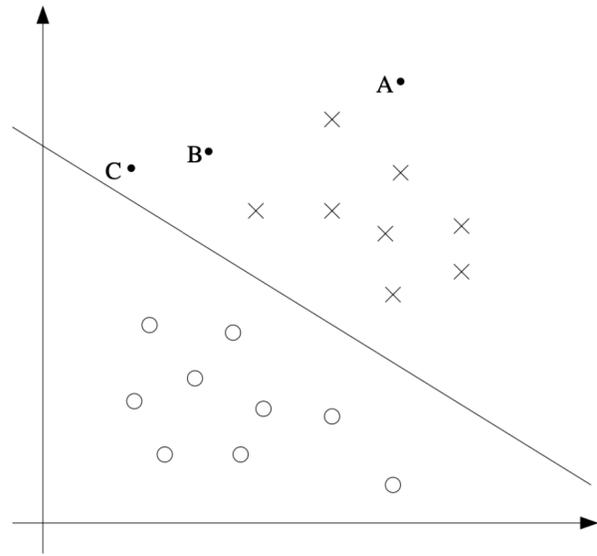


Review of the High-Level Logic



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

Review of the High-Level Logic

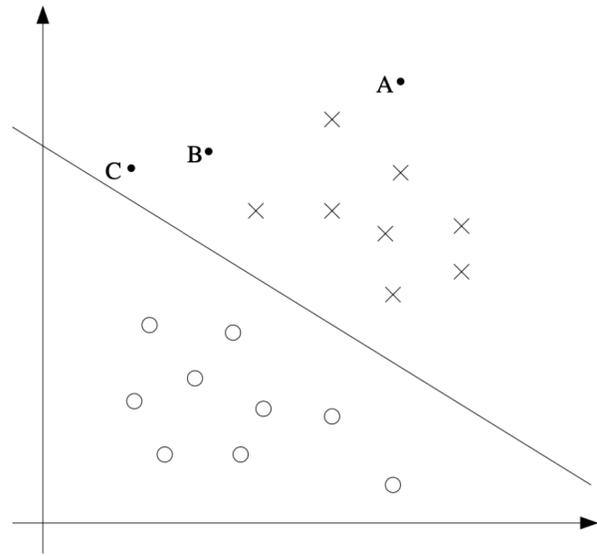


$$h_{w,b}(x) = g(w^T x + b).$$

Maximize
geometric
margin

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

Review of the High-Level Logic



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

Maximize
geometric
margin

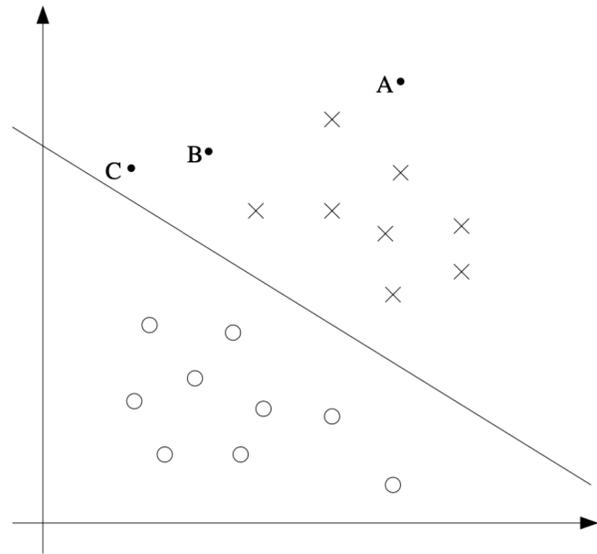
Problem
rewriting



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

$\|w\|$

Review of the High-Level Logic



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

Maximize
geometric
margin

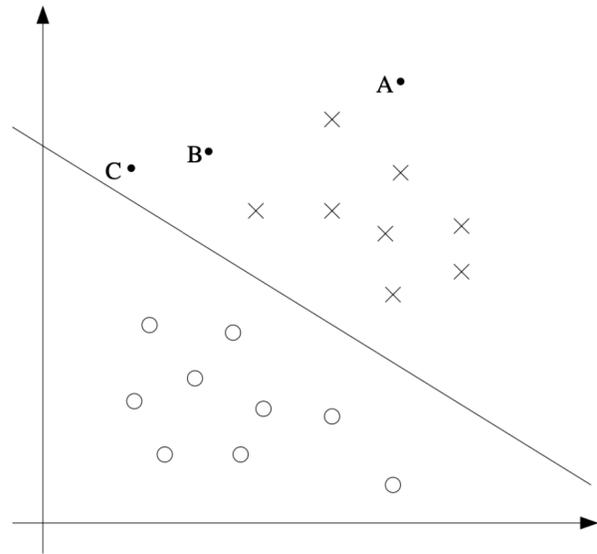
Problem
rewriting

Quadratic
Optimization
Problem

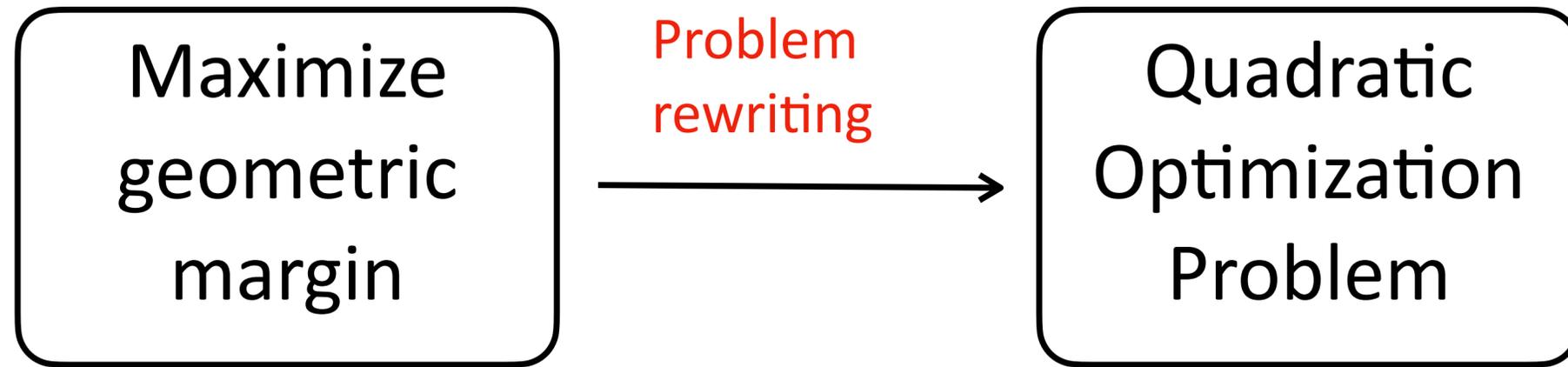
$$\gamma^{(i)} = y^{(i)} \left(\left(\frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right)$$

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

Review of the High-Level Logic



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

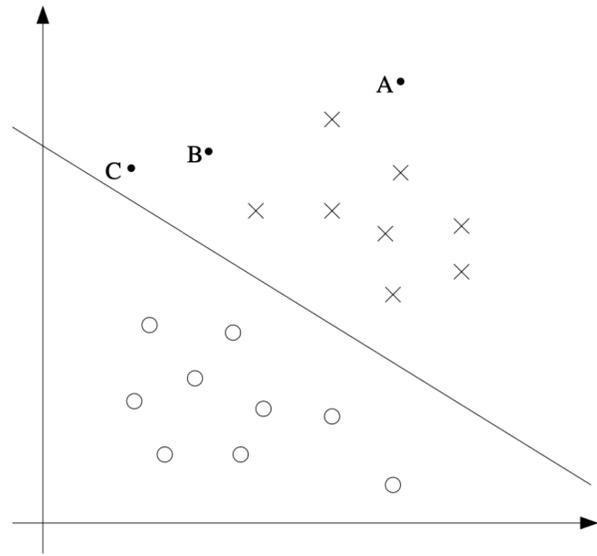


$$\gamma^{(i)} = y^{(i)} \left(\left(\frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right)$$

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

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

Review of the High-Level Logic



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

Maximize
geometric
margin

Problem
rewriting

Quadratic
Optimization
Problem

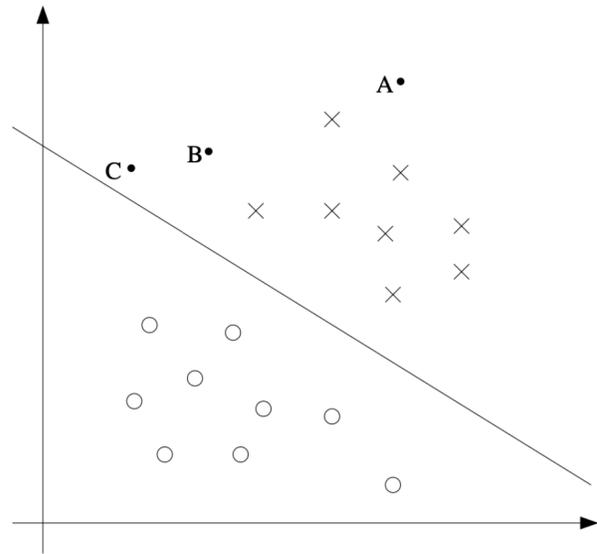
Finding a related
optimization problem
that is easier

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

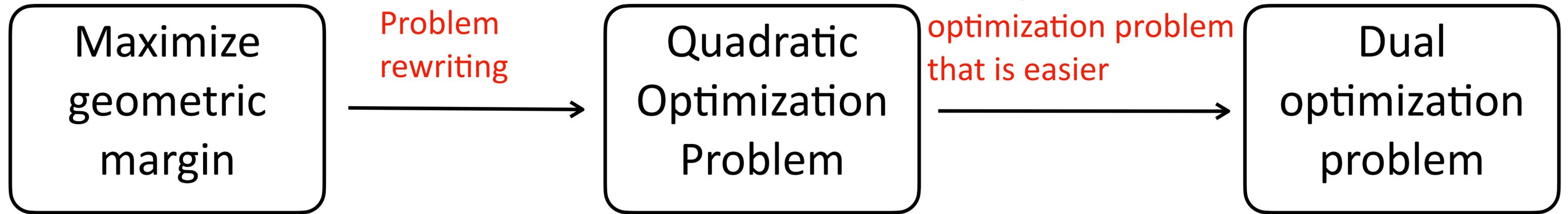
$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)} (w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

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

Review of the High-Level Logic



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



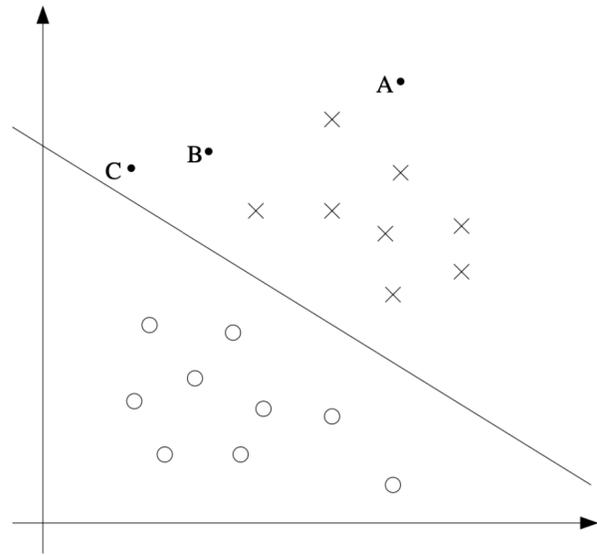
$$\gamma^{(i)} = y^{(i)} \left(\left(\frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right)$$

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

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

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

Review of the High-Level Logic



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

Maximize
geometric
margin

Problem
rewriting

Quadratic
Optimization
Problem

Finding a related
optimization problem
that is easier

Dual
optimization
problem

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

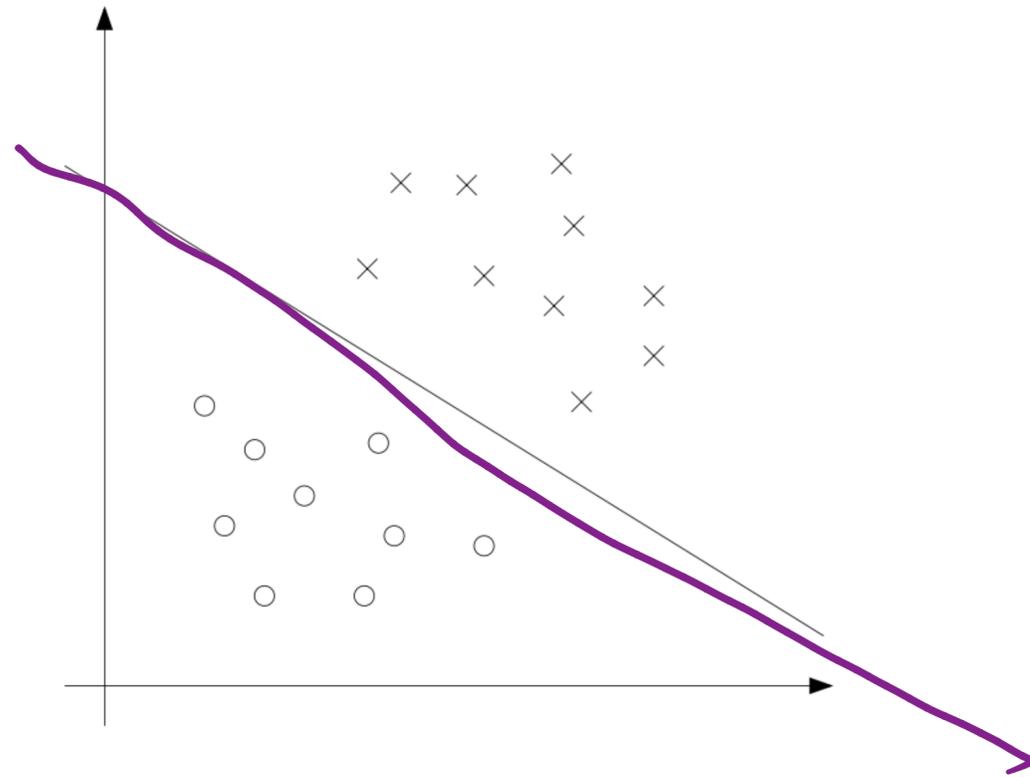
$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

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

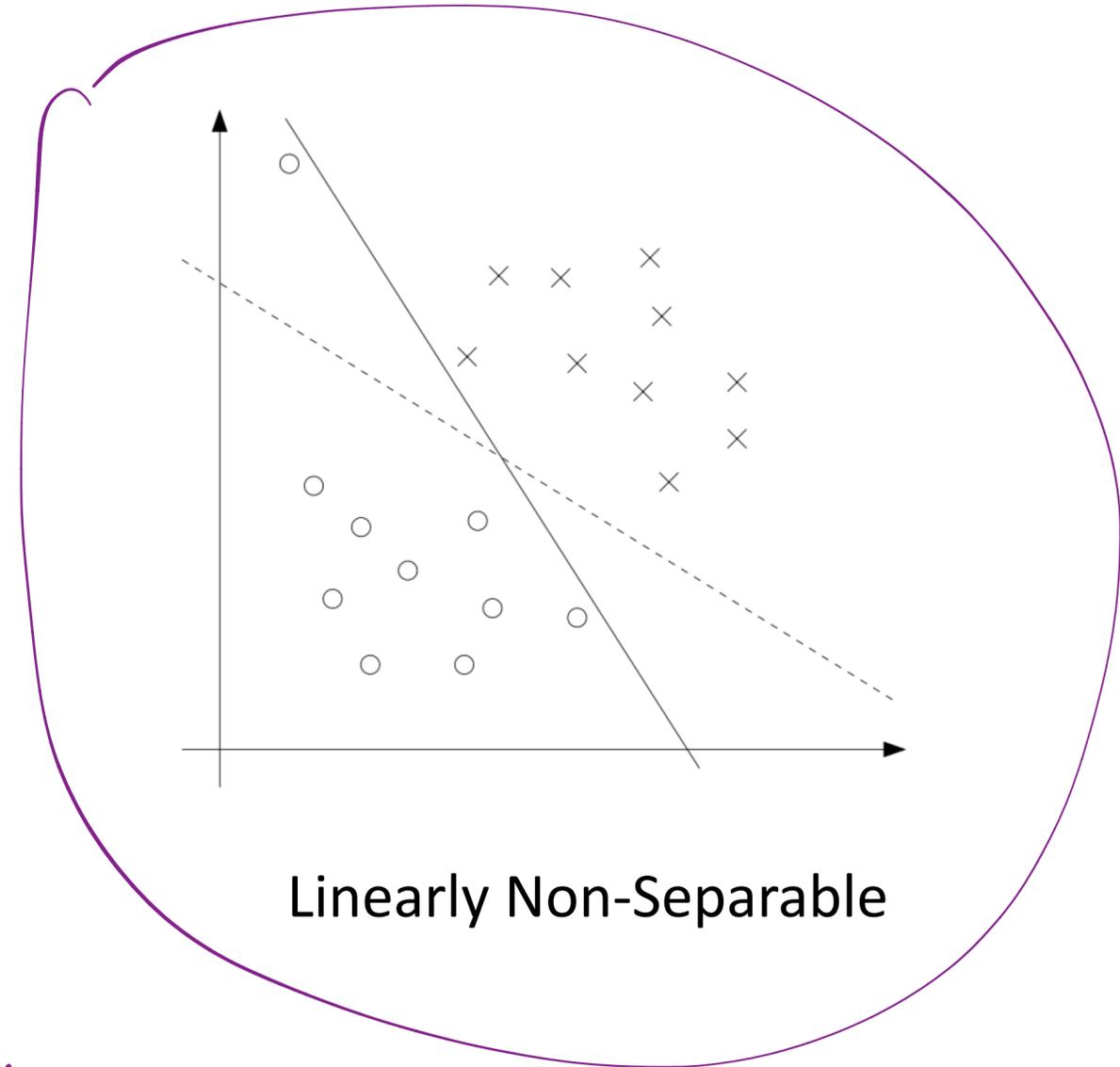
$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \end{aligned}$$

Kernel makes it very flexible in
non-linear cases!

The Non-Separable Case



Linearly Separable



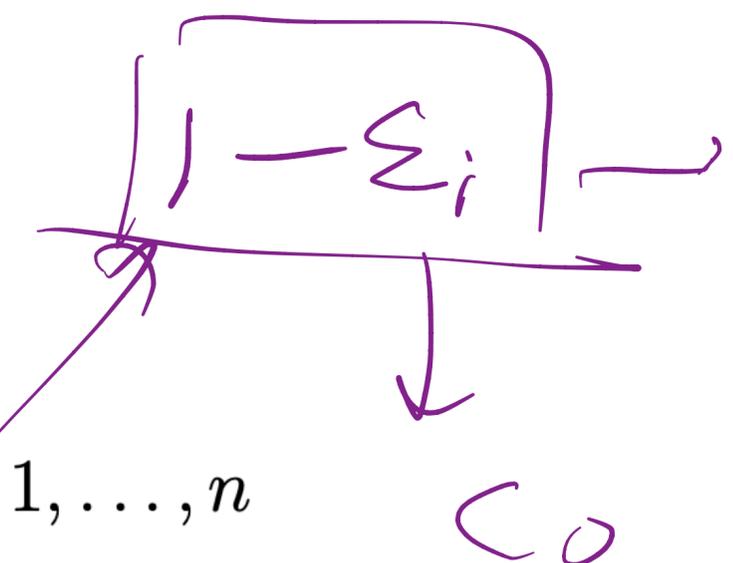
Linearly Non-Separable

$y^{(i)}$ $(w^T x + b)$ \rightarrow ~~ϕ~~

The Non-Separable Case

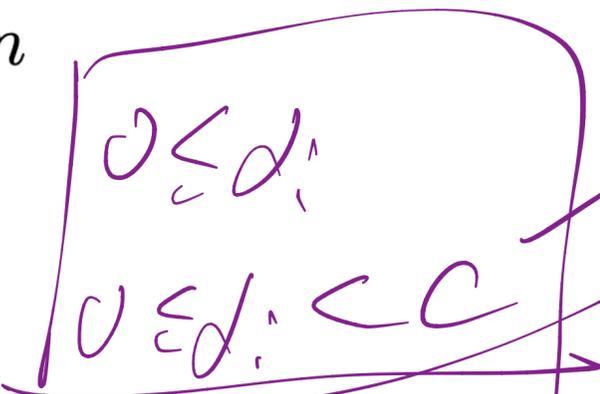
Primal opt problem:

$$\begin{aligned} \min_{\gamma, w, b} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$



Dual opt problem

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



Thank You!
Q & A