

# Machine Learning from Data

Lecture 23: Spring 2021

# Today's Lecture

- Support Vector Machines (SVMs)
  - Maximizing the margins

## RECAP: Linear Models, RBFs, Neural Networks

### Linear Model with Nonlinear Transform

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1}^{\bar{d}} w_j \Phi_j(\mathbf{x}) \right)$$



### Neural Network

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1}^m w_j \theta(\mathbf{v}_j^T \mathbf{x}) \right)$$

gradient descent



### k-RBF-Network

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1}^k w_j \phi(\|\mathbf{x} - \boldsymbol{\mu}_j\|) \right)$$

k-means

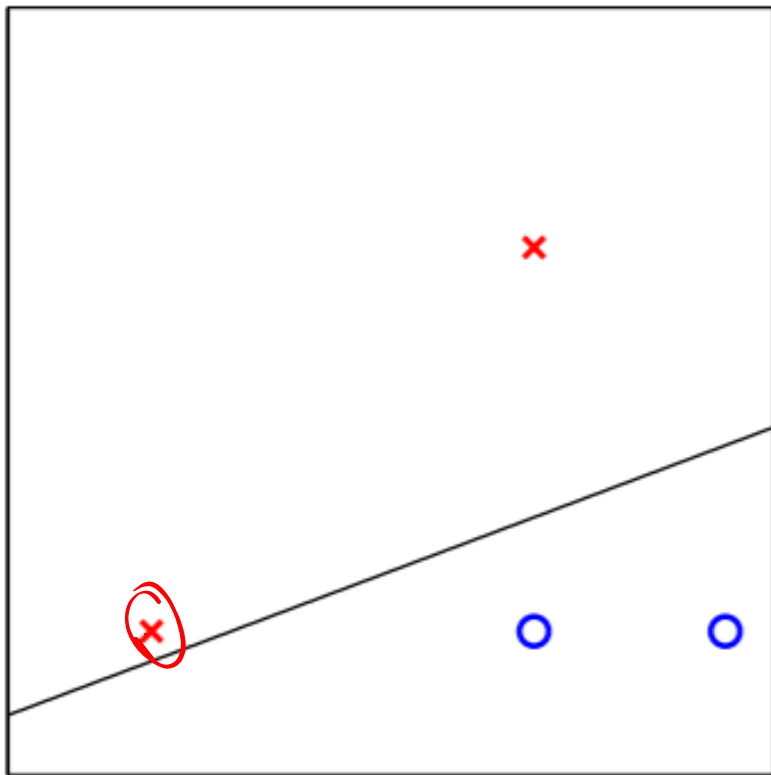


Neural Network: generalization of linear model by adding layers.

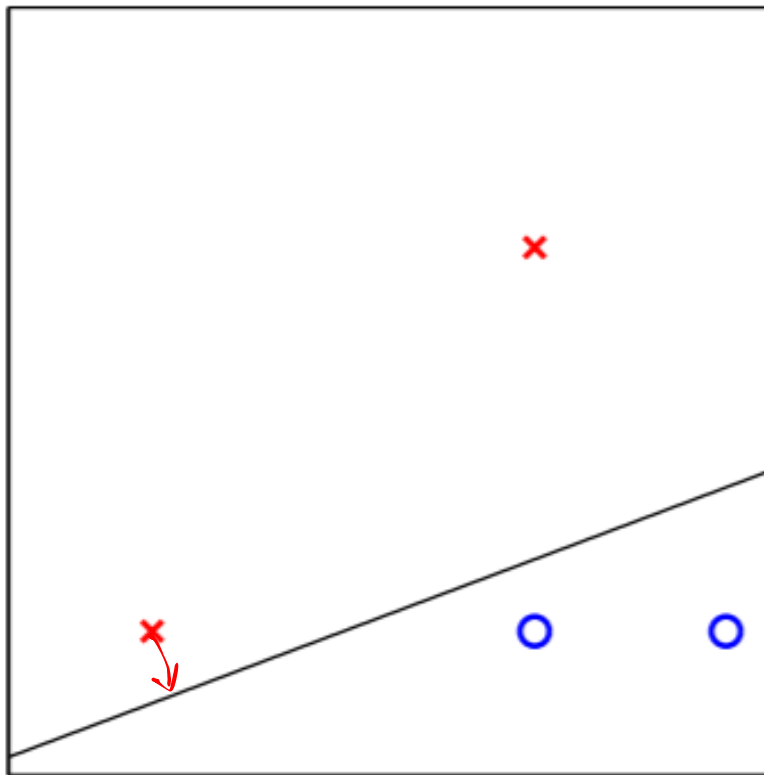
**Support Vector Machine:** more 'robust' linear model



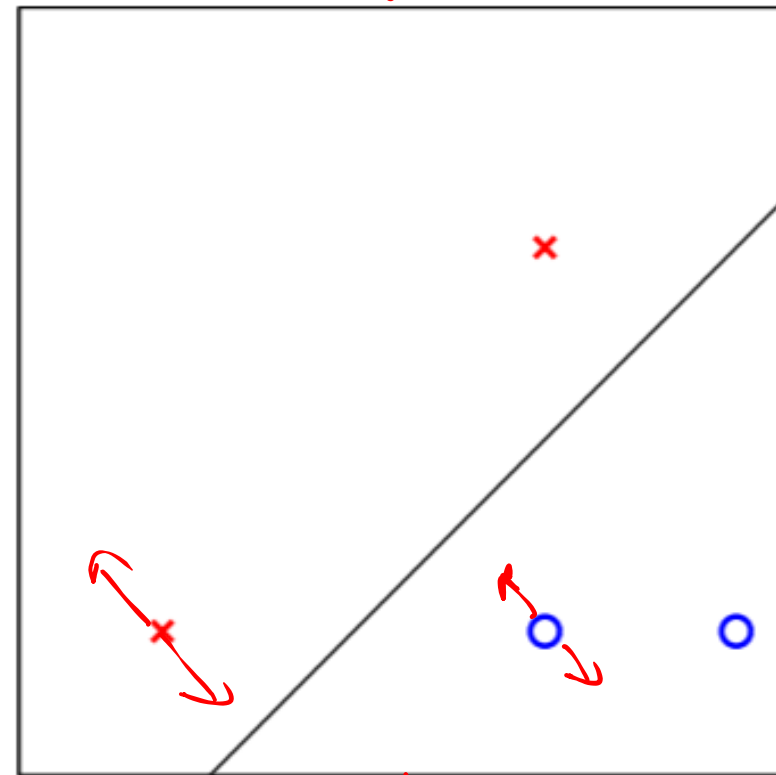
# Which Separator Do You Pick?



(1)



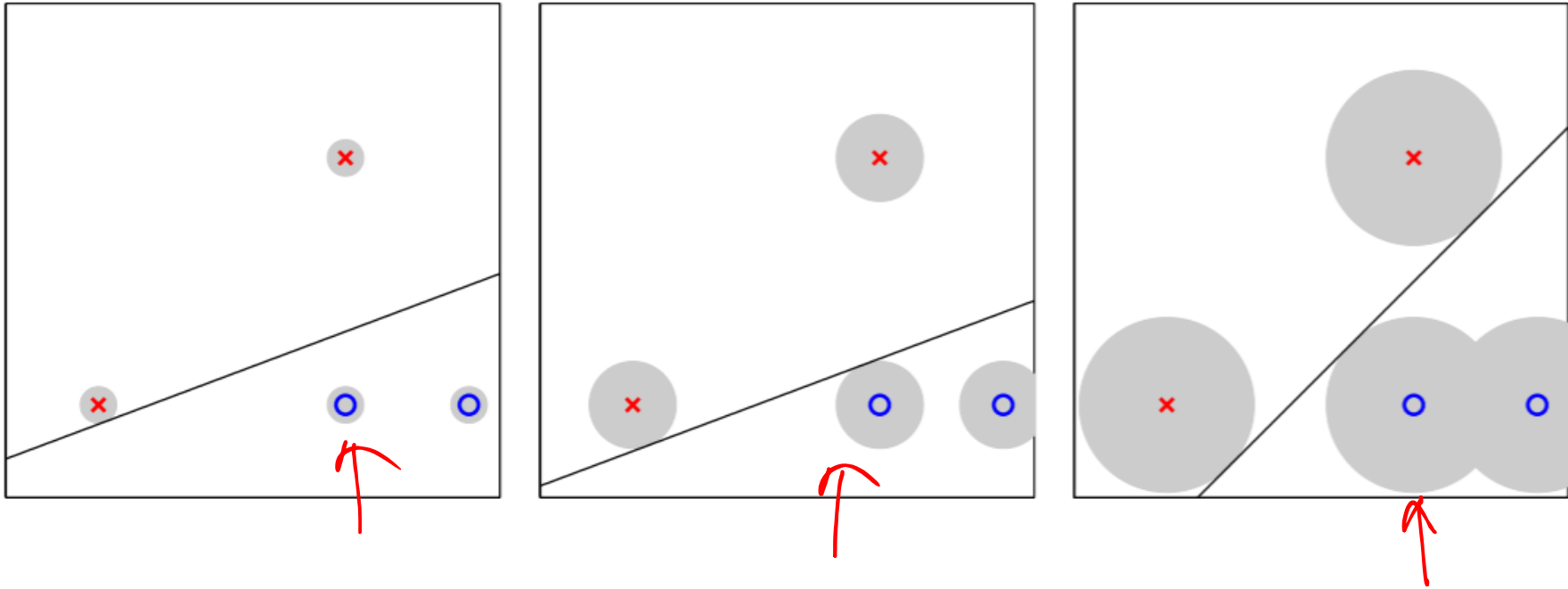
(2)



(3)

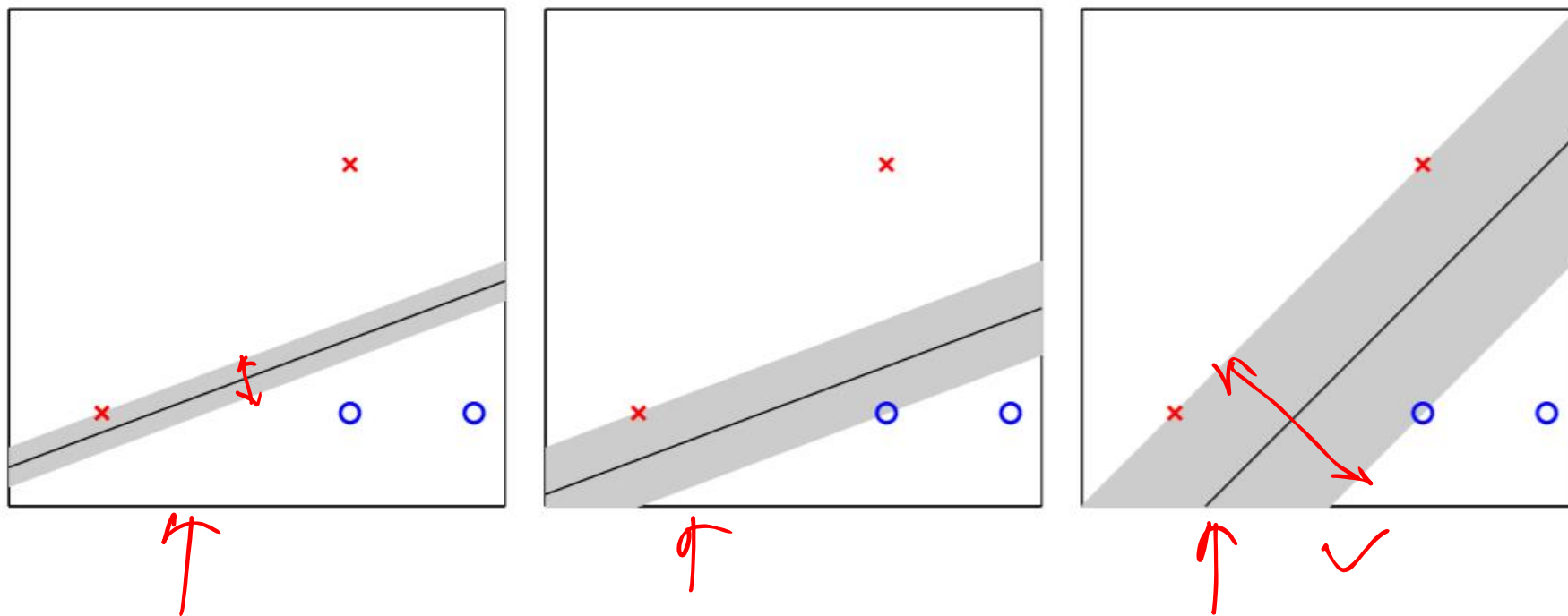


## Robustness to Noisy Data



Being robust to noise (measurement error) is good (remember regularization).

## Thicker Cushion Means More Robustness



We call such hyperplanes **fat**

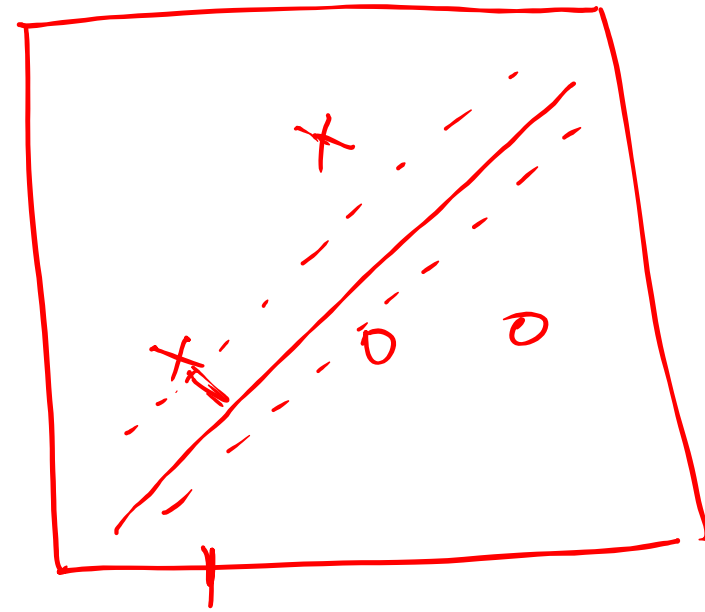
## Two Crucial Questions

1. Can we efficiently find the fattest separating hyperplane?
2. Is a fatter hyperplane better than a thin one?

Optimal Hyperplane.  $\rightarrow$  Separates the data.

Thickness: Distance to nearest datapoint

$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad y = \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix}$$





Hyperplane  $\rightarrow$  set of weights.

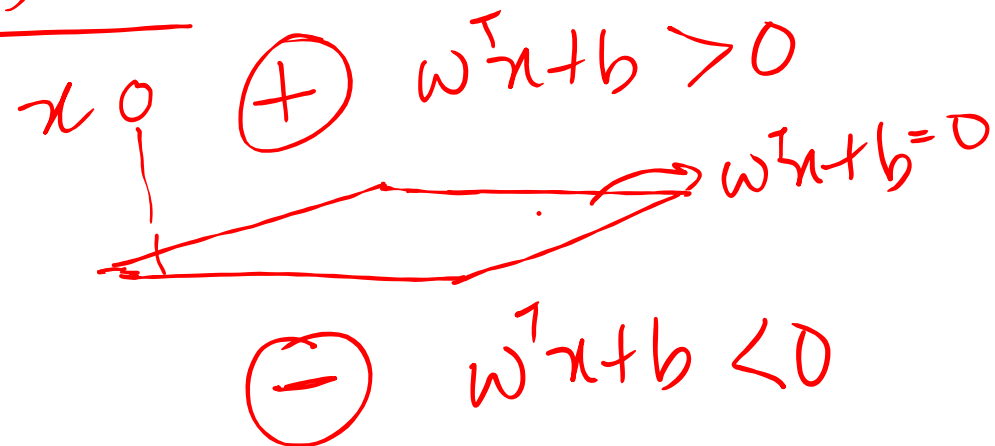
$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$w_0$  (bias)  $\rightarrow b$

earlier  $h(x) = \text{sign}(w^T x)$

$$h(x) = \text{sign}(w^T x + b)$$

3-d space



Separating

$$\text{sign}(\omega^T x_n + b) = y_n$$

$$\rightarrow y_n(\omega^T x_n + b) > 0 \checkmark, n = 1, 2, \dots, N$$

$$\min_n y_n(\omega^T x_n + b) = \rho > 0 \checkmark$$

$$\text{OR } \min_n y_n \left( \frac{\omega^T x_n}{\rho} + \frac{b}{\rho} \right) = 1 \checkmark$$

Definition:  $h = (\omega, b)$  separates the data iff

$$\min_n y_n(\omega^T x_n + b) = 1$$

$$h = (\omega, b)$$

$$\text{dist}(x, h) = \left| \underset{\substack{\uparrow \\ \text{direction of } u}}{u} \cdot (x - x_1) \right|$$

$\omega \rightarrow$  normal (in the direction of  $u$ )

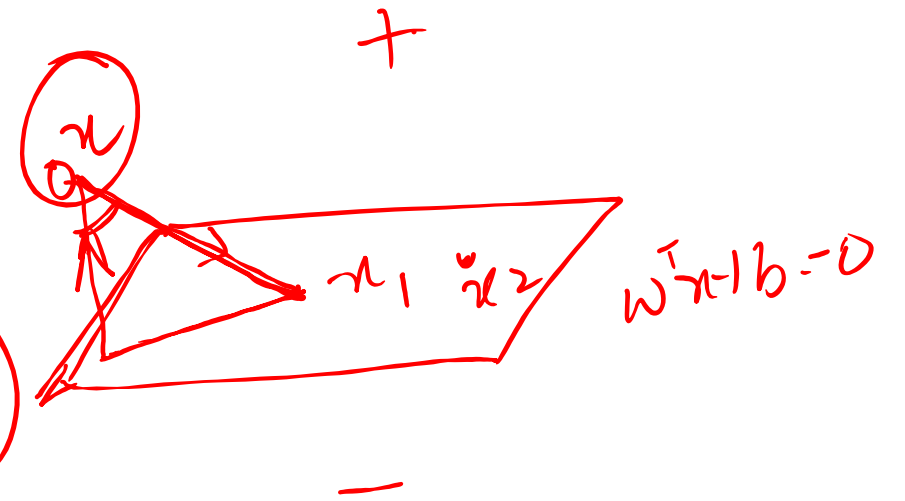
$$\text{or } u = \frac{\omega}{\|\omega\|}$$

$x_2 - x_1 \rightarrow$  vector

$$\left. \begin{aligned} \omega^T x_1 + b &= 0 \\ \omega^T x_2 + b &= 0 \end{aligned} \right\}$$

$$\omega^T x_1 - \omega^T x_2 = 0 \Rightarrow \omega \cdot (x_1 - x_2) = 0$$

$\omega / \|\omega\|$



$$\text{dist}(x, h) = \left| \frac{\omega}{\|\omega\|} \cdot (x - x_1) \right|$$

$$= \frac{1}{\|\omega\|} |\omega^T x - \omega^T x_1|$$

$$= \frac{1}{\|\omega\|} |\omega^T x + b - (\omega^T x_1 + b)|$$

simpler

↓  
0

$$\text{dist}(x, h) = \frac{1}{\|\omega\|} |\omega^T x + b|$$

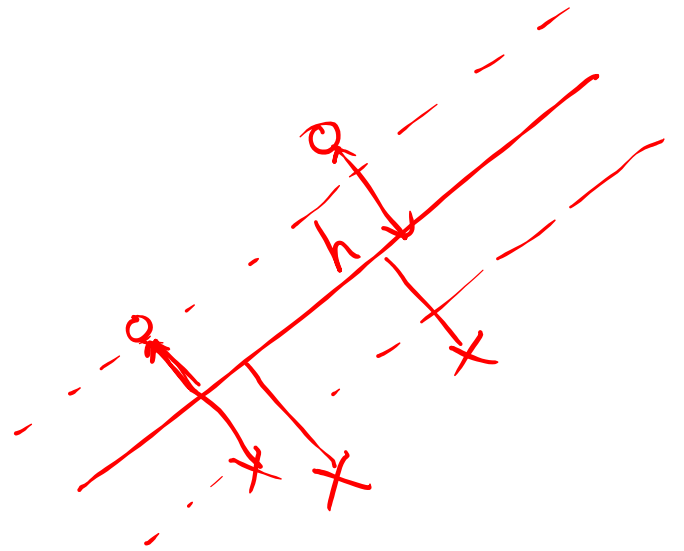
# Determine Thickness of a Separating Hyperplane

The thickness - margin ( $\rho(h)$ )

$$\rho(h) = \min_n \text{dist}(x_n, h)$$

$$= \min_n \frac{|\omega^T x_n + b|}{\|\omega\|}$$

$$\rho(h) = \frac{1}{\|\omega\|} \min_n |\omega^T x_n + b|$$



$$\min_n y_n (w^T x_n + b) = 1$$

$$\min_n |y_n (w^T x_n + b)| = 1$$

Separating  
hyperplane

$$\min_n |w^T x_n + b| = 1$$

Substitute

$$r(h) = \frac{1}{\|w\|}$$

Thickness or largest.  
Margin

$$\text{dist}(x, h) = \frac{|w^T x + b|}{\|w\|}$$

$$\left. \begin{array}{l} \text{Maximize } \frac{1}{\|w\|} \\ \text{s.t. } \min_n y_n (w^T x_n + b) = 1 \end{array} \right\}$$

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

$$\text{s.t. } \min_n y_n (w^T x_n + b) = 1$$

$$\Leftrightarrow$$

$$y_n (w^T x_n + b) \geq 1$$

$$\boxed{\begin{array}{l} \min \frac{1}{2} w^T w \\ \text{s.t. } y_n (w^T x_n + b) \geq 1 \end{array}}$$

- i) Variables in optimization ( $w, b$ )
  - ii) Objective function is quadratic.
  - iii) Linear inequality constraints
- Quadratic Program (QP).



$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad y = \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix} \Rightarrow$$

$$\omega_1, \omega_2, b$$

$$\min \frac{1}{2} (\omega_1^2 + \omega_2^2)$$

Support Vectors

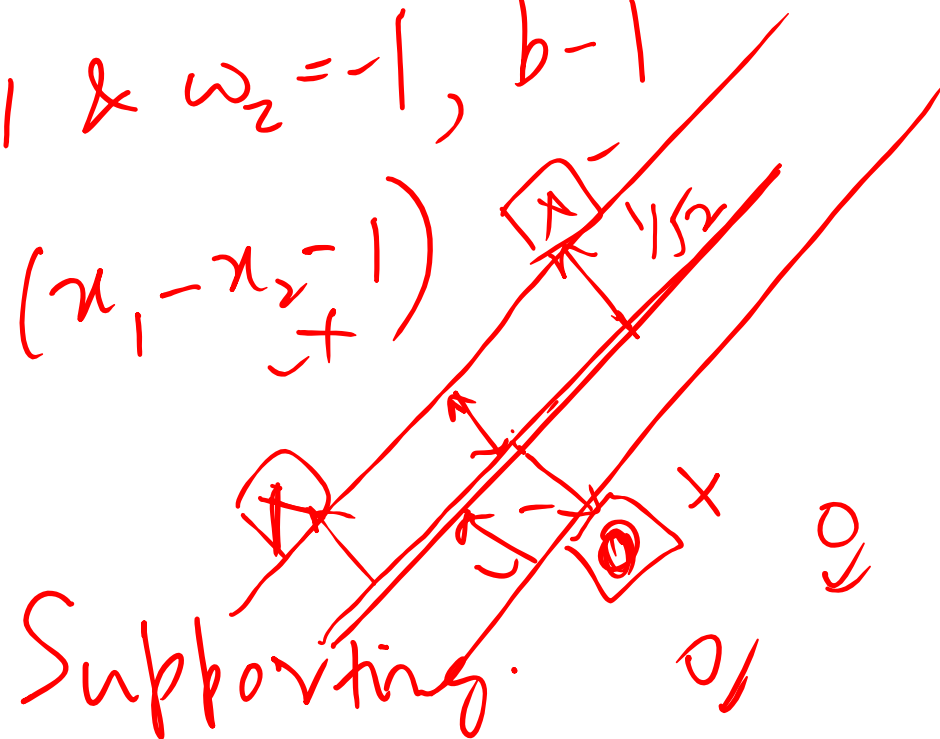
$$\begin{cases} -b \geq 1 \\ -1(2\omega_1 + 2\omega_2 + b) \geq 1 \\ 1(2\omega_1 + b) \geq 1 \\ 1(3\omega_1 + b) \geq 1 \end{cases} \Rightarrow$$

$$\omega_1 \geq 1, \omega_2 \leq -1$$

$$\omega_1 = 1 \& \omega_2 = -1, b = 1$$

$$\therefore h_{\text{thick}} = \text{sign}(x_1 - x_2 - 1)$$

$$\text{Thickness} = \frac{1}{\|\omega\|} = \frac{1}{\sqrt{2}}$$

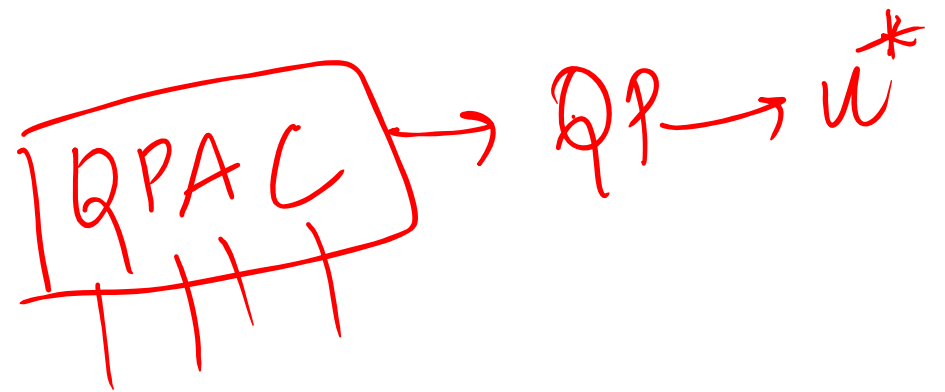


# Quadratic Program in standard form

$$\begin{array}{ll} \min_u & \frac{1}{2} u^T Q u + \underbrace{p^T u}_{\uparrow} \\ \text{s.t.} & A u \geq C \end{array}$$

$\left\{ \begin{array}{l} Q: \text{ is a } q \times q \text{ matrix} \\ A: \text{ is a } N \times q \text{ matrix} \\ p: \text{ is a vector } q \times 1 \\ C: \text{ is a vector } N \times 1 \end{array} \right.$

$$u \in \mathbb{R}^q$$



find  $\rightarrow b, w$   
 $\rightarrow u = \begin{bmatrix} b \\ w \end{bmatrix} \in \mathbb{R}^{d+1}$ ,  $\min \frac{1}{2} w^T w$

$$\underbrace{w^T w}_{\text{minimize}} = \begin{bmatrix} b & w^T \end{bmatrix} \begin{bmatrix} 0 & 0_d^T \\ 0_d & I_d \end{bmatrix} \begin{bmatrix} b \\ w^T \end{bmatrix}$$

$I_d$ :  $d \times d$   
 identity matrix  
 $0_d$  -  $d$  dimensional  
 zero vector.

$$= u^T \begin{bmatrix} 0 & 0_d^T \\ 0_d & I_d \end{bmatrix} u$$

$$\rightarrow Q, \quad p = 0_{d+1}$$

$$\checkmark C = 1$$

A

$$y_n (w^T u_n + b) \geq 1 \Rightarrow$$

$$\rightarrow a_n \Rightarrow y_n [1$$

$$\begin{bmatrix} y_n & y_n u_n^T \\ u_n^T \end{bmatrix} u \geq 1$$

$$\rightarrow C_n = 1$$

We have identified QPAC  $\rightarrow$  QP  $\rightarrow u^* = \begin{bmatrix} b^* \\ w^* \end{bmatrix}$

## Linear Hard-Margin SVM with QP

- 1: Let  $\mathbf{p} = \mathbf{0}_{d+1}$  ( $(d+1)$ -dimensional zero vector) and  $\mathbf{c} = \mathbf{1}_N$  ( $N$ -dimensional vector of ones). Construct matrices  $\mathbf{Q}$  and  $\mathbf{A}$ , where

$$\mathbf{Q} = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix}, \quad \mathbf{A} = \underbrace{\begin{bmatrix} y_1 & -y_1 \mathbf{x}_1^T \\ \vdots & \vdots \\ y_N & -y_N \mathbf{x}_N^T \end{bmatrix}}_{\text{signed data matrix}}.$$

- 2: Calculate  $\begin{bmatrix} b^* \\ \mathbf{w}^* \end{bmatrix} = \mathbf{u}^* \leftarrow \text{QP}(\mathbf{Q}, \mathbf{p}, \mathbf{A}, \mathbf{c})$ .

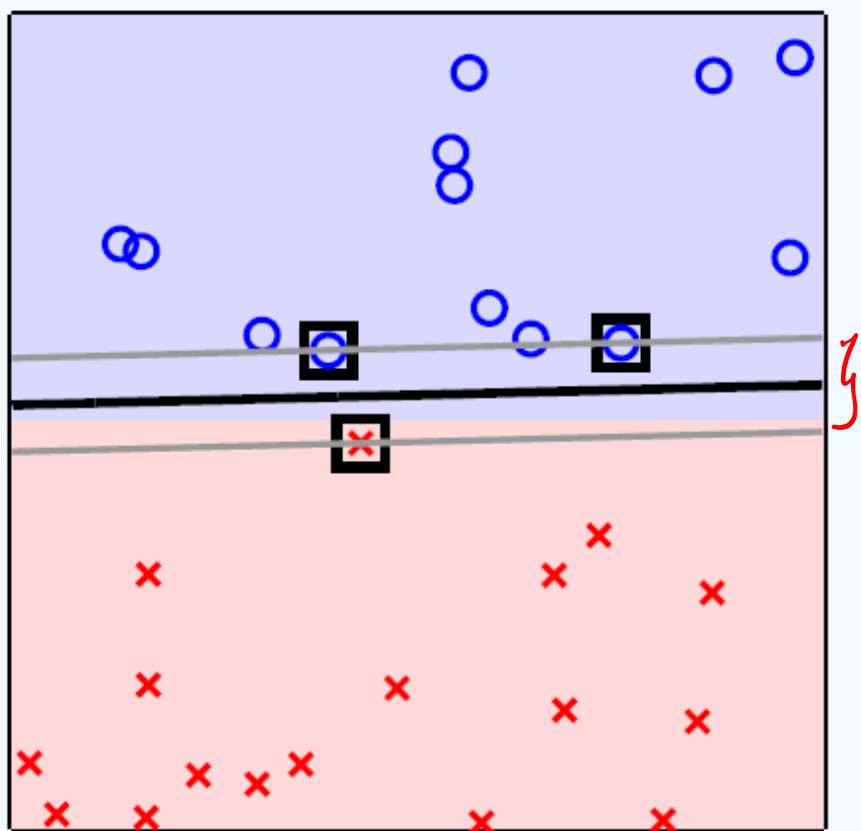
- 3: Return the hypothesis  $g(\mathbf{x}) = \text{sign}(\mathbf{w}^{*T} \mathbf{x} + b^*)$ .

## Toy Example

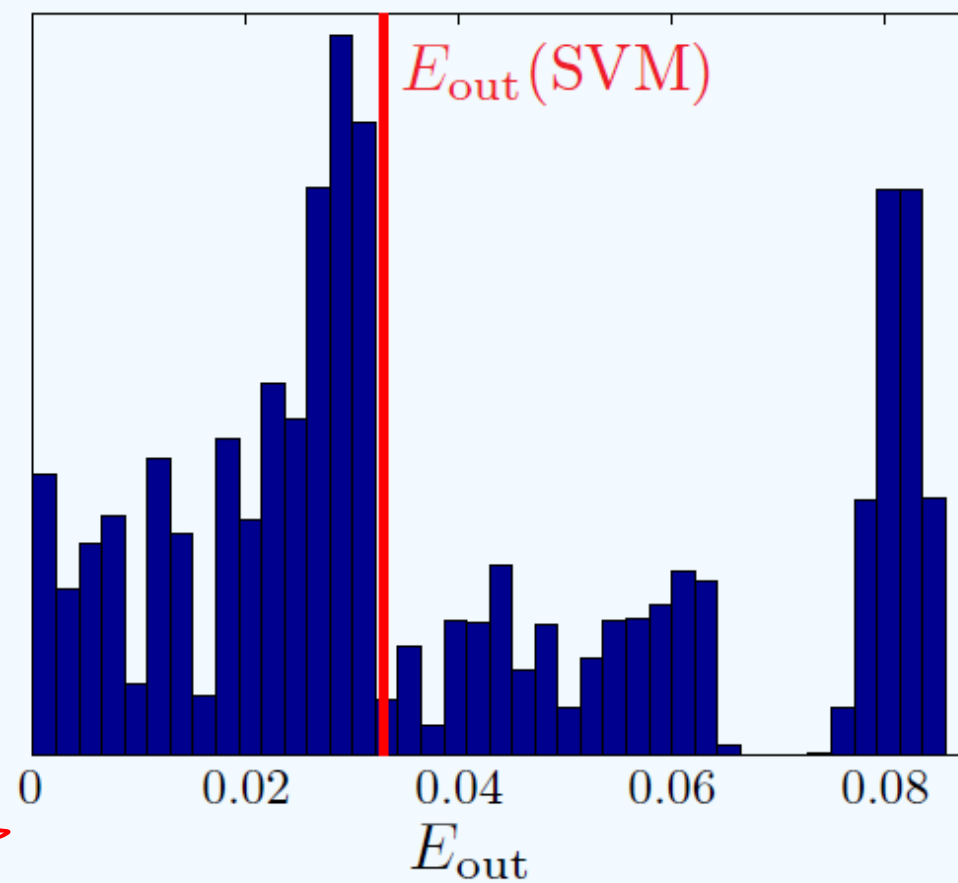
11-1

$$\checkmark Q = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \checkmark \mathbf{p} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad A = \begin{bmatrix} -1 & 0 & 0 \\ -1 & -2 & -2 \\ 1 & 2 & 0 \\ 1 & 3 & 0 \end{bmatrix} \quad \checkmark \mathbf{c} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

→ A standard QP-solver gives  $(b^*, w_1^*, w_2^*) = (-1, 1, -1)$ , the same solution we computed manually, but obtained in less than a millisecond.  $\square$



(a) Data and SVM separator



(b) Histogram of  $E_{\text{out}}(\text{PLA})$

optimal hyperplane

regularization

minimize:  
subject to:

$$\mathbf{w}^T \mathbf{w}$$

$$E_{\text{in}} = 0$$

$$E_{\text{in}} \checkmark$$

$$\mathbf{w}^T \mathbf{w} \leq C$$



Thanks!