

# Support Vector Machines

Somak Aditya

Sudeshna Sarkar

CSE Department, IIT Kharagpur

Sep 1, 2023

These slides were assembled by Byron Boots, with only minor modifications from Eric Eaton's slides and grateful acknowledgement to the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

# Strengths of SVMs

- Good generalization
  - in theory
  - in practice
- Works well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick

# Notational Conventions

To better match notation used in SVMs  
...and to make matrix formulas simpler

We will use the following for the  $i^{\text{th}}$  instance

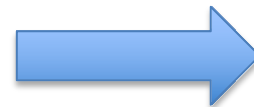
$i^{\text{th}}$  instance



$\mathbf{x}_i$

**Bold** denotes  
vector

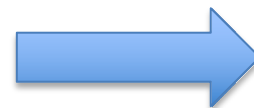
$i^{\text{th}}$  instance label



$y_i$

**Non-bold**  
denotes scalar

$j^{\text{th}}$  feature of  $i^{\text{th}}$  instance



$x_{ij}$

# Linear Separators

- Training instances

$$\mathbf{x} \in \mathbb{R}^{d+1}, \mathbf{x}_0 = \mathbf{1}$$

$$y \in \{-1, 1\}$$

- Model parameters

$$\theta \in \mathbb{R}^{d+1}$$

- Hyperplane

$$\theta^\top \mathbf{x} = \langle \theta, \mathbf{x} \rangle = 0$$

- Decision function

$$\begin{aligned} h(\mathbf{x}) &= \text{sign}(\theta^\top \mathbf{x}) \\ &= \text{sign}(\langle \theta, \mathbf{x} \rangle) \end{aligned}$$

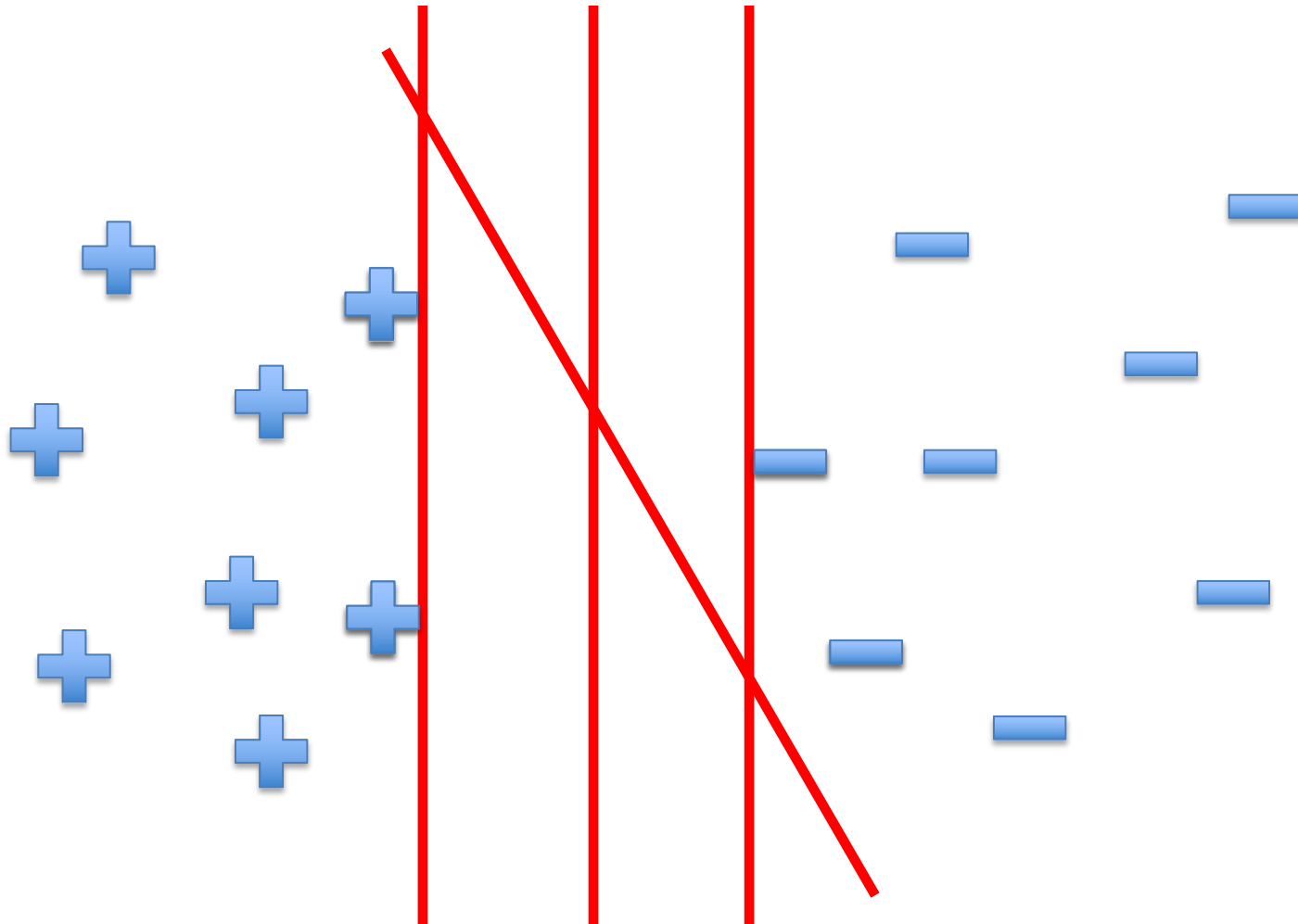
Recall:

Inner (dot) product:

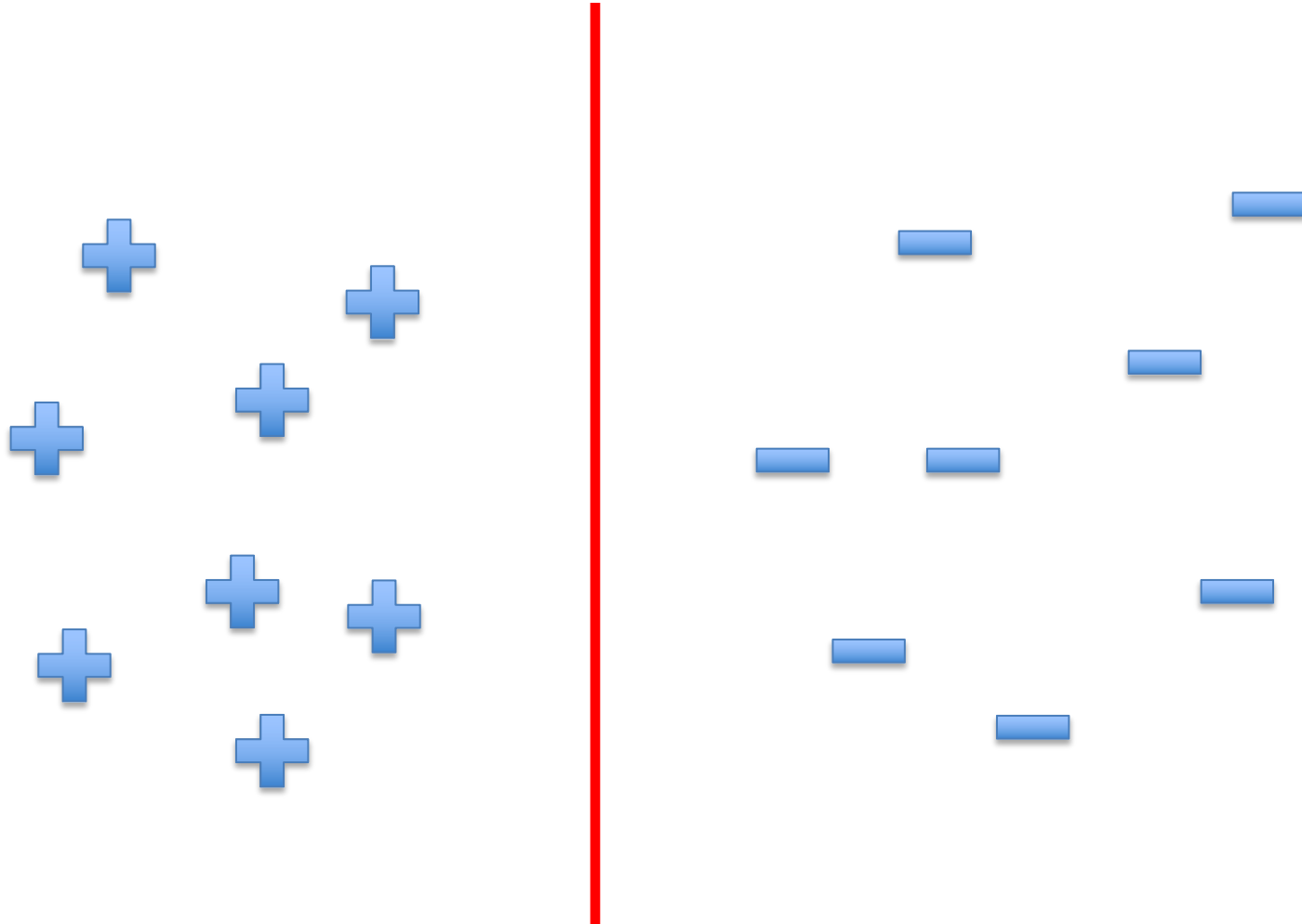
$$\begin{aligned} \langle \mathbf{u}, \mathbf{v} \rangle &= \mathbf{u} \cdot \mathbf{v} \\ &= \mathbf{u}^\top \mathbf{v} \end{aligned}$$

$$= \sum_i u_i v_i$$

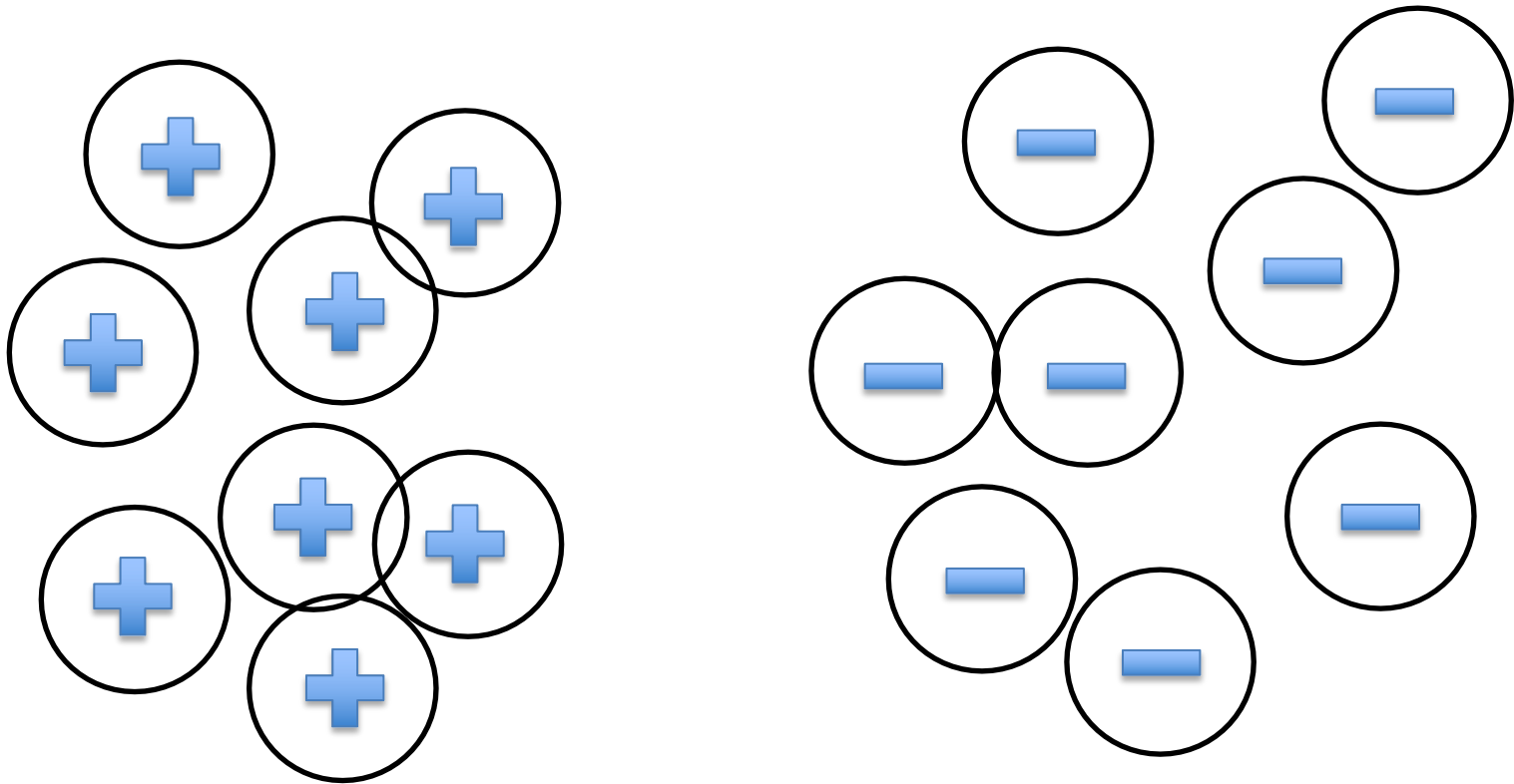
# Intuitions



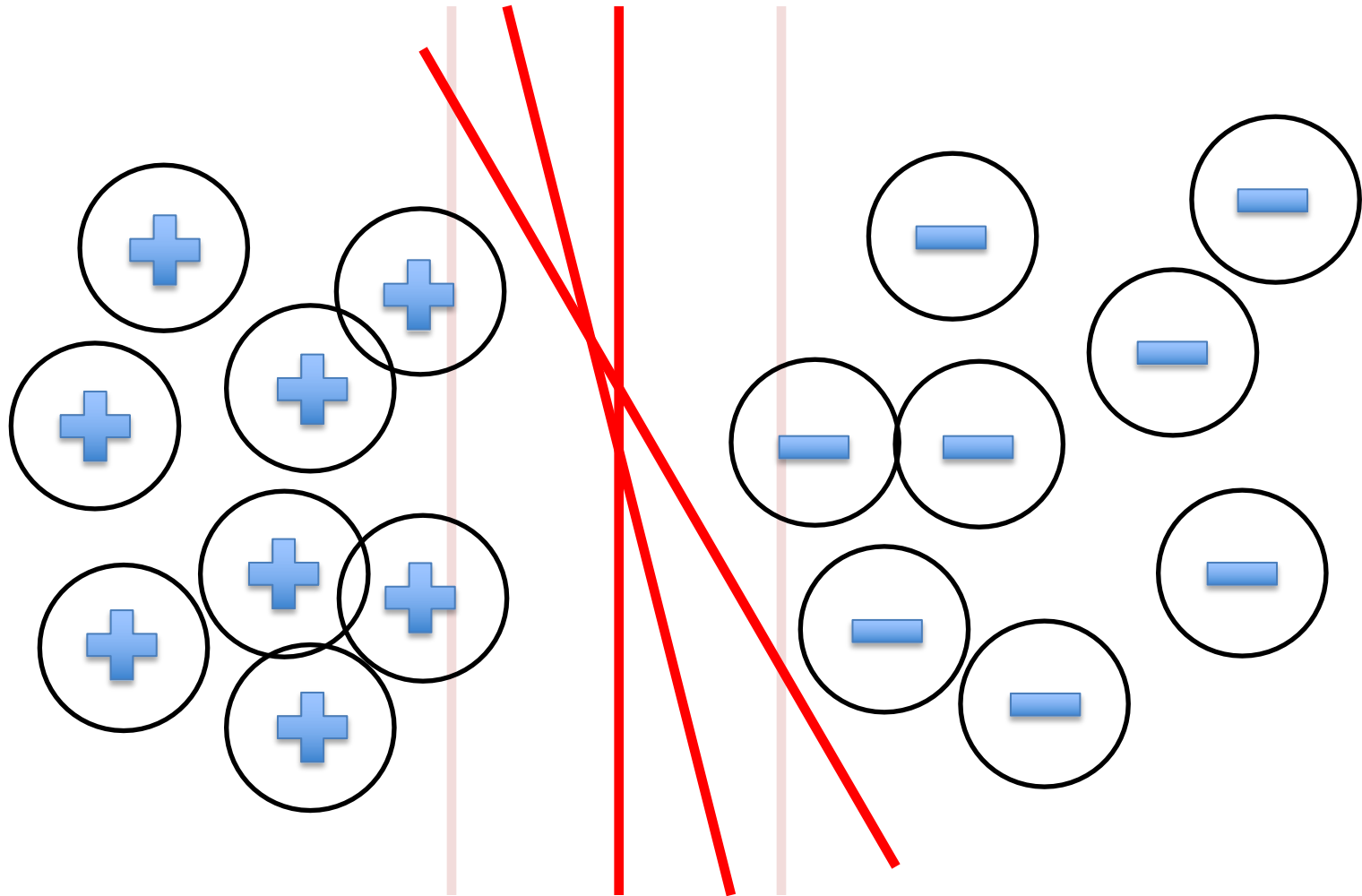
# A "Good" Separator



# Noise in the Observations

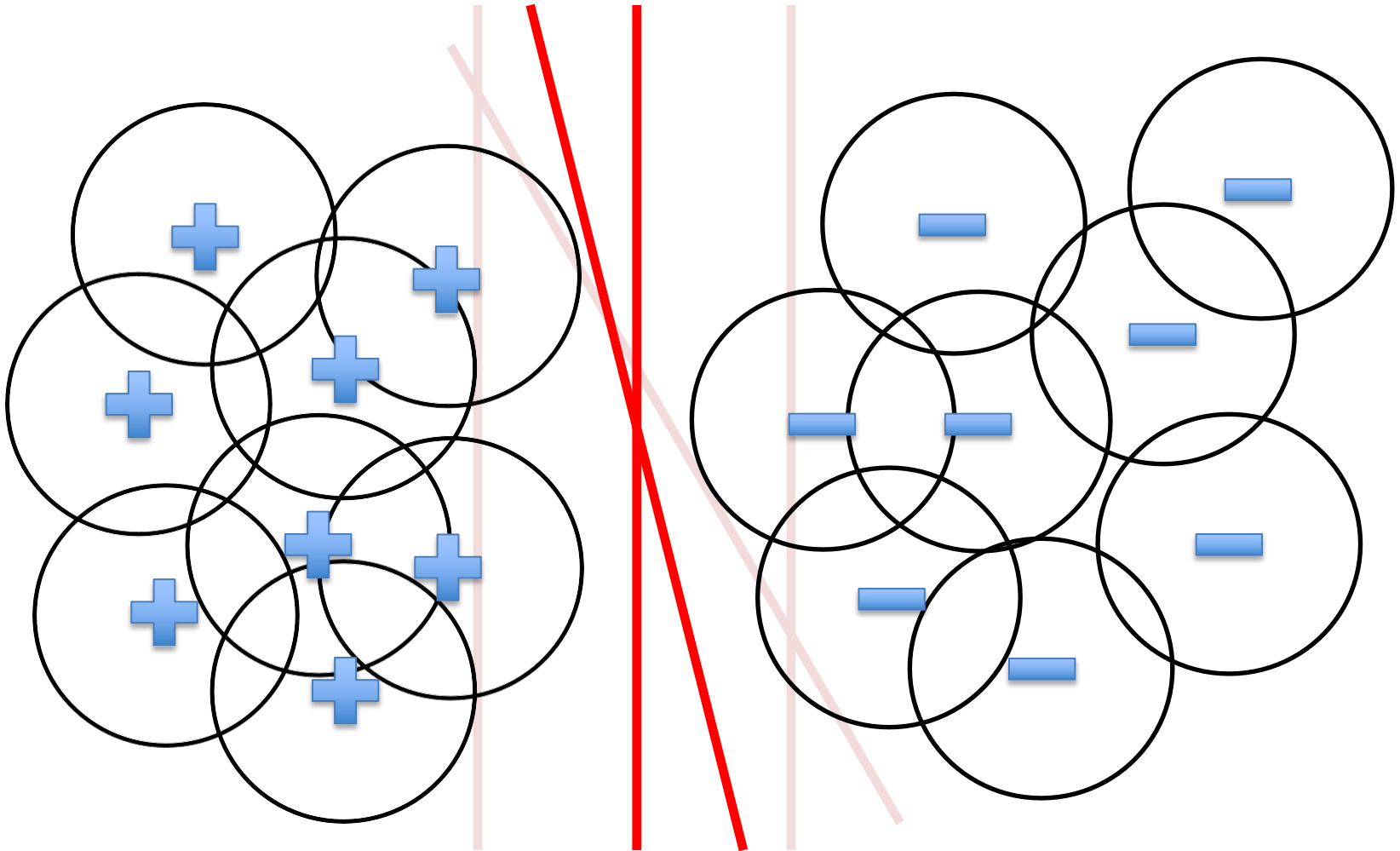


# Ruling Out Some Separators

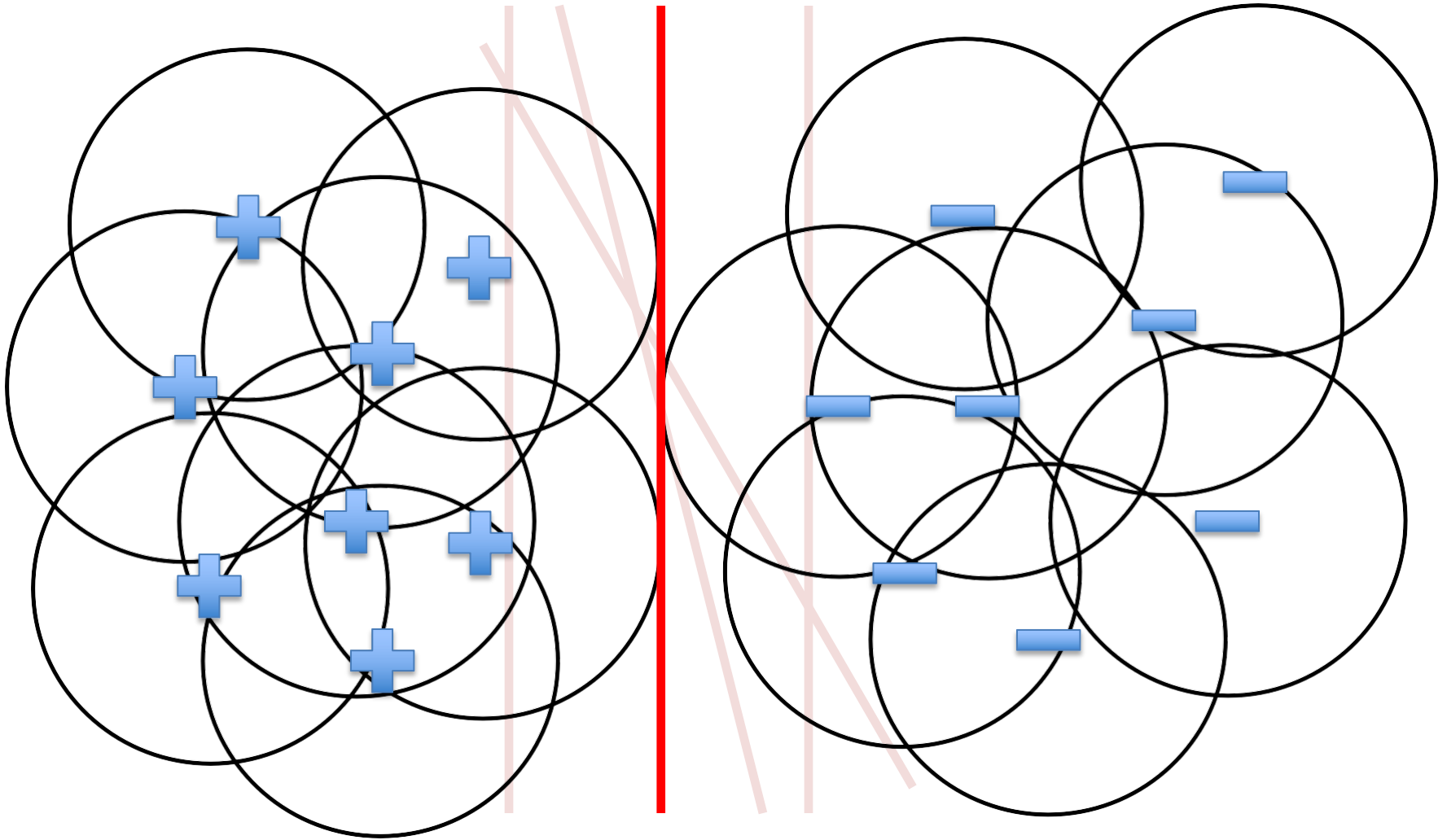




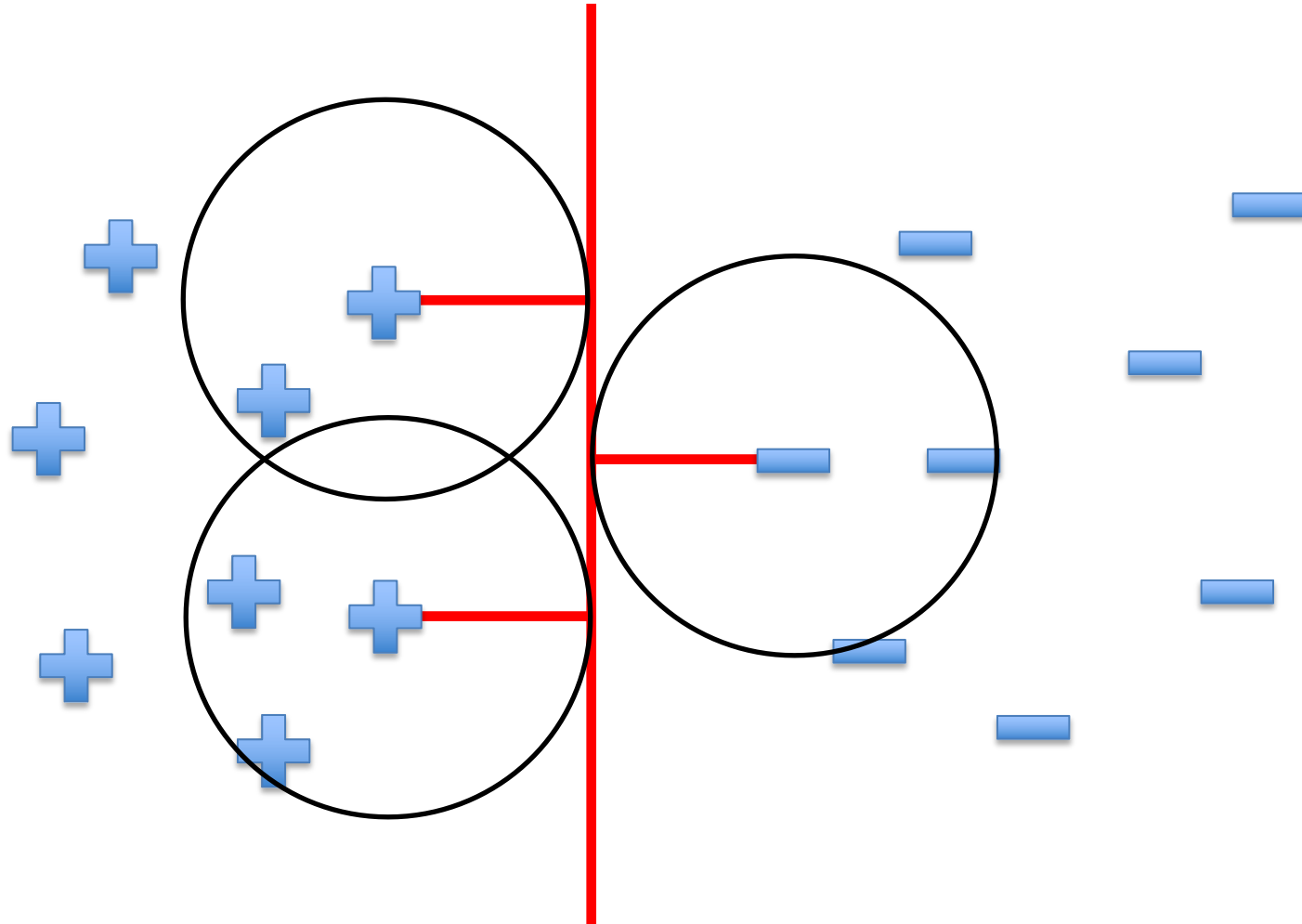
# Lots of Noise



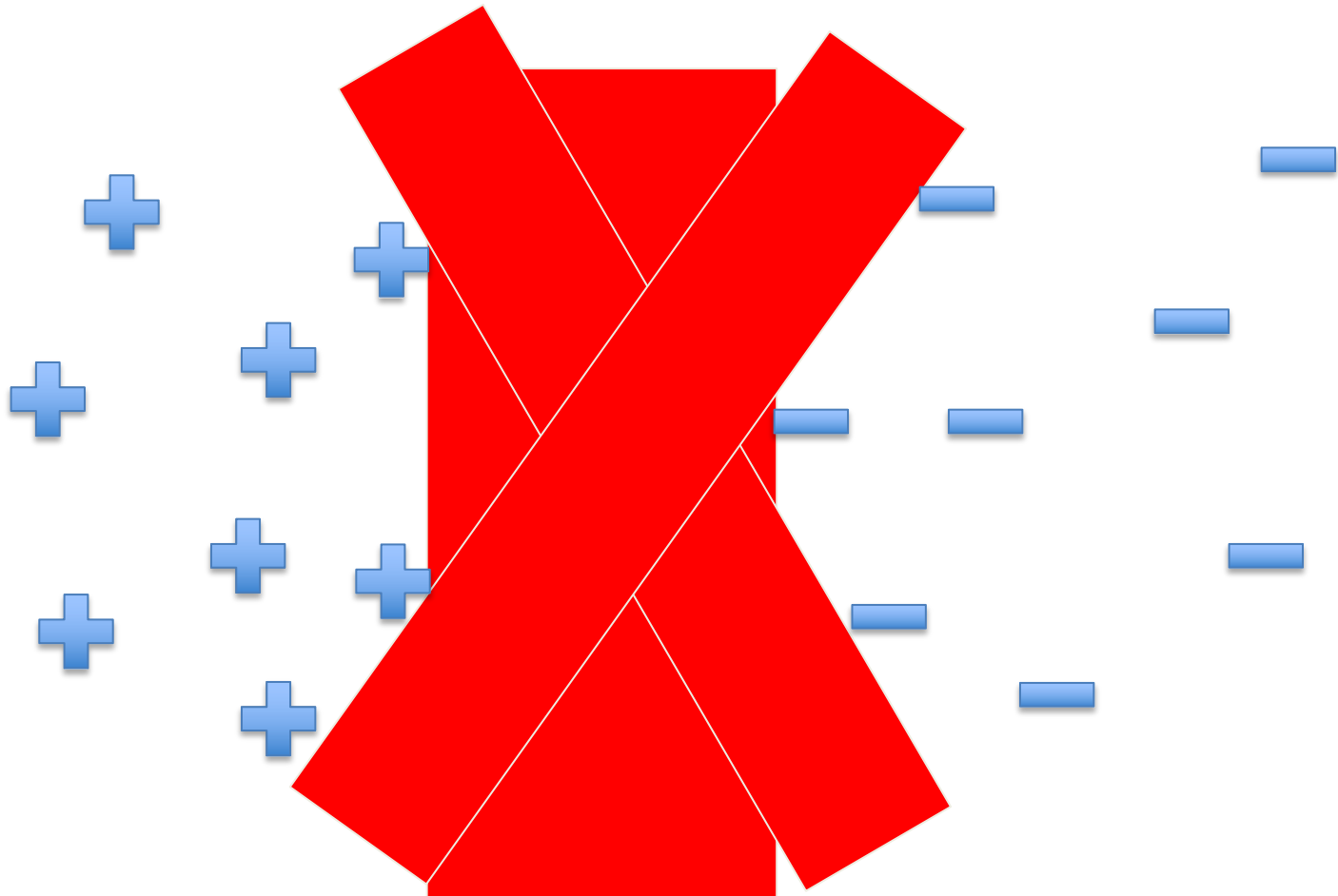
# Only One Separator Remains



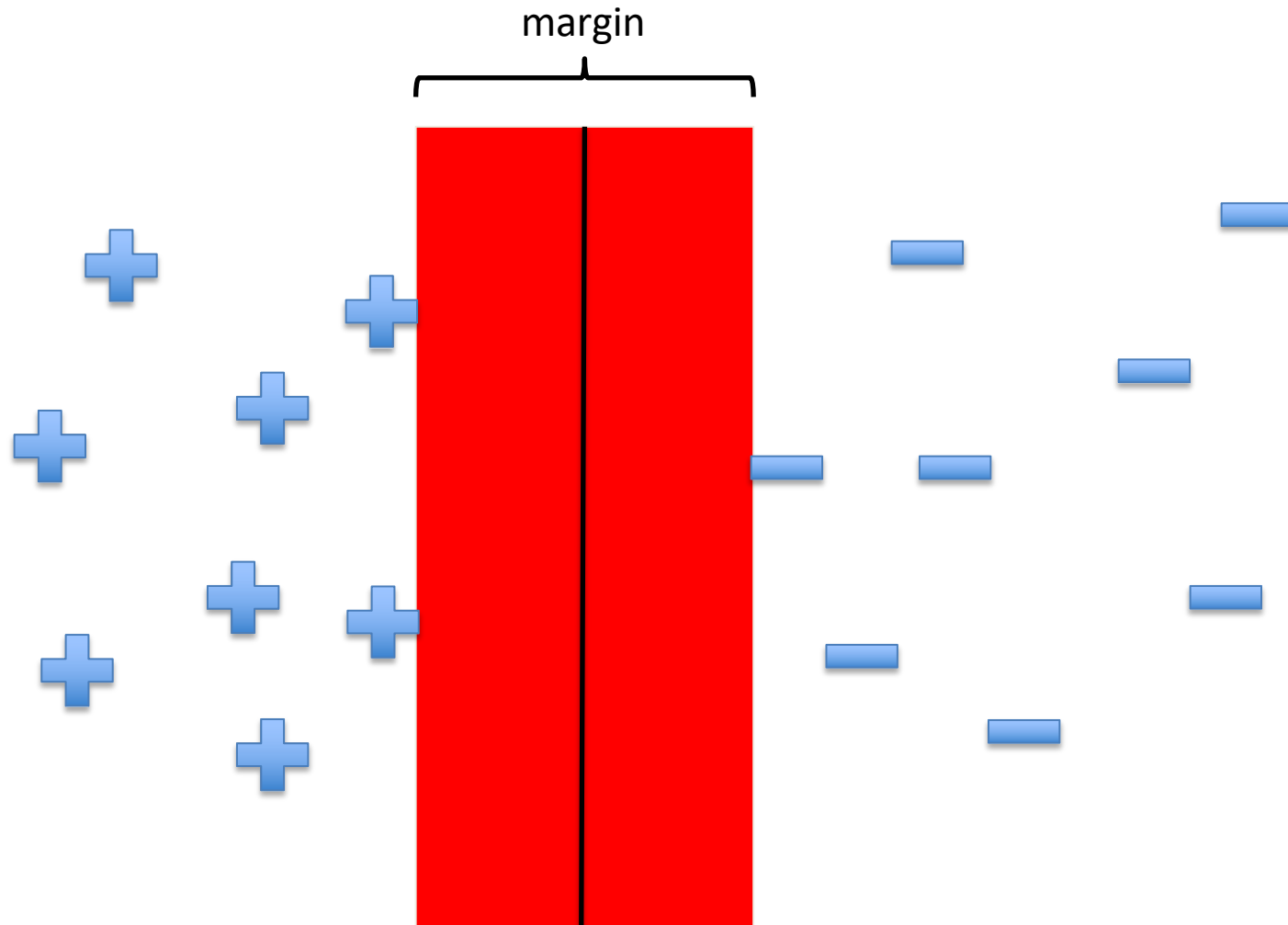
# Maximizing the Margin



# "Fat" Separators



# "Fat" Separators



# Why Maximize Margin

Increasing margin reduces *capacity*

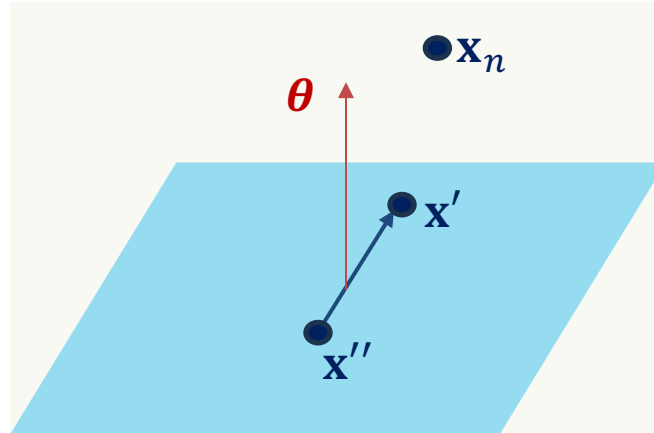
- i.e., fewer possible models
- *What about bias? Variance?*

Lesson from Learning Theory:

- If the following holds:
  - $H$  is sufficiently constrained in size
  - and/or the size of the training data set  $n$  is large,then low training error is likely to be evidence of low generalization error

# Computing the margin

Margin = The distance between  $\mathbf{x}_n$  and the plane  $\theta^\top \mathbf{x} = 0$ ,  
where  $|\theta^\top \mathbf{x}_n| = 1$



# Computing the margin

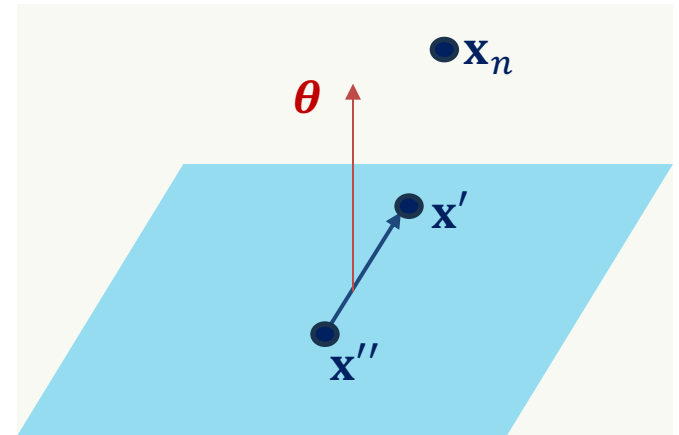
Proposition:

The vector  $\theta$  is orthogonal to the plane in the  $X$  space

Take any two points  $\mathbf{x}'$  and  $\mathbf{x}''$  on the plane.

$$\theta^\top \mathbf{x}' = 0 \text{ and } \theta^\top \mathbf{x}'' = 0 \\ \Rightarrow \theta^\top (\mathbf{x}' - \mathbf{x}'') = 0$$

Hence  $\theta$  is orthogonal to any vector that lies on the plane  $\Rightarrow \theta$  is orthogonal to the plane





# Margin: distance between $\mathbf{x}_n$ and the plane

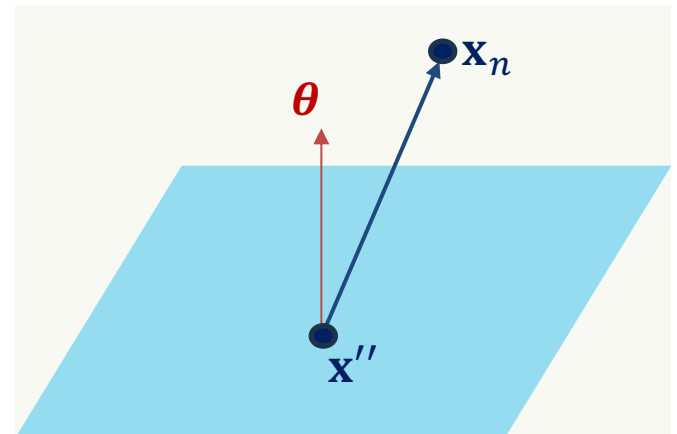
Take any point  $\mathbf{x}$  on the plane

Projection of  $\mathbf{x}_n - \mathbf{x}$  on  $\boldsymbol{\theta}$  (direction orthogonal to the plane)

$$\hat{\boldsymbol{\theta}} = \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|} \Rightarrow \text{distance} = |\hat{\boldsymbol{\theta}}^T (\mathbf{x}_n - \mathbf{x})|$$

Projection of the vector  $(\mathbf{x}_n - \mathbf{x})$  along  $\boldsymbol{\theta}$   
computed by taking the vector product of  $(\mathbf{x}_n - \mathbf{x})$  with the unit vector in the direction of  $\boldsymbol{\theta}$

$\|\boldsymbol{\theta}\|$  is the norm of  $\boldsymbol{\theta}$



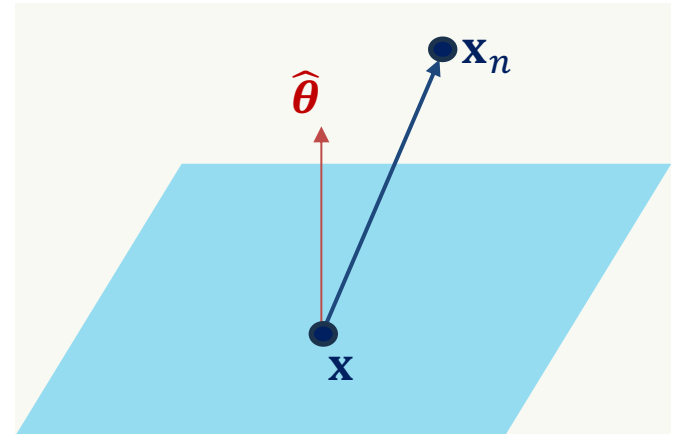
Margin: distance between  $\mathbf{x}_n$  and the plane

$$\text{distance} = \frac{1}{\|\boldsymbol{\theta}\|} |\boldsymbol{\theta}^\top \mathbf{x}_n - \boldsymbol{\theta}^\top \mathbf{x}|$$
$$\frac{1}{\|\boldsymbol{\theta}\|} |\boldsymbol{\theta}^\top \mathbf{x}_n - \boldsymbol{\theta}^\top \mathbf{x}| = \frac{1}{\|\boldsymbol{\theta}\|}$$

$\mathbf{x}$  is a point on the plane.

Hence  $\Rightarrow \boldsymbol{\theta}^\top \mathbf{x} = 0$

$|\boldsymbol{\theta}^\top \mathbf{x}_n| = 1$  for the nearest point  $\mathbf{x}_n$  (due to our normalization)



# The optimization problem

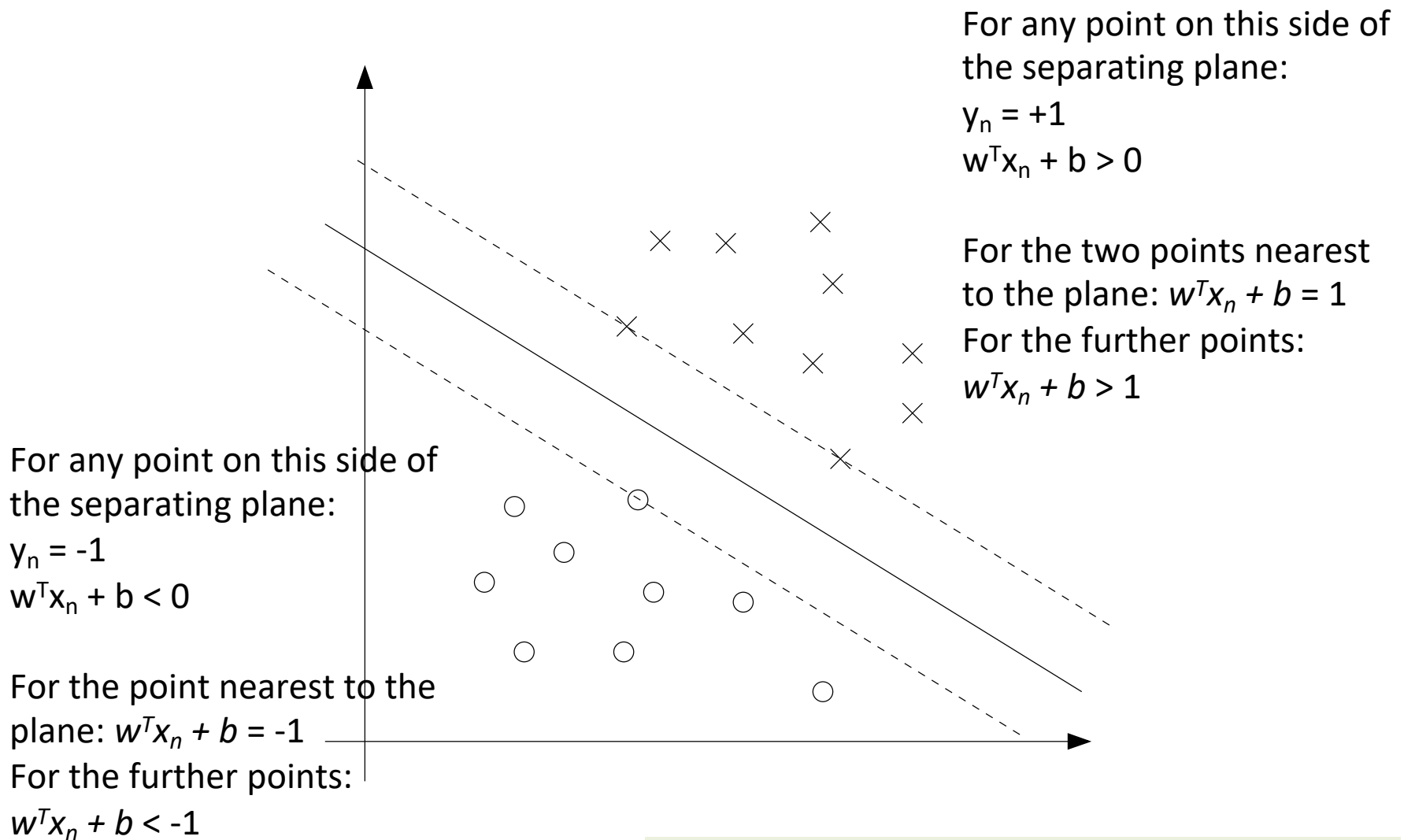
$$\begin{aligned} &\text{Maximize } \frac{1}{\|\boldsymbol{\theta}\|} \\ &\text{subject to } \min_{n=1,2,\dots,N} \|\boldsymbol{\theta}^T \mathbf{x}_n\| = 1 \end{aligned}$$

This optimization problem is too complex, because of

- (i) the norm in the objective function, and
- (ii) the minimum term in the constraints

Can we find an equivalent optimization problem that is easier to tackle?

# The geometry



Notice:  $|\theta^T \mathbf{x}_n| = y_n(\theta^T \mathbf{x}_n)$

# Equivalent optimization problem

Maximize  $\frac{1}{\|\boldsymbol{\theta}\|}$

subject to  $\min_{n=1,2,\dots,N} |\boldsymbol{\theta}^\top \mathbf{x}_n| = 1$

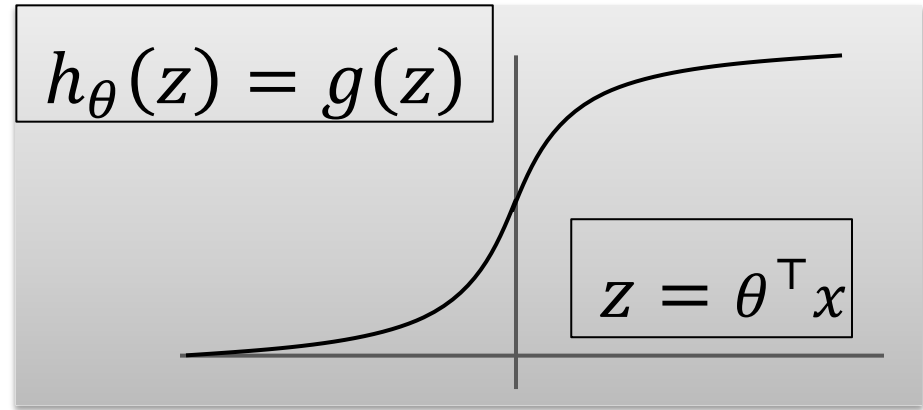
Notice:  $|\boldsymbol{\theta}^\top \mathbf{x}_n| = y_n(\boldsymbol{\theta}^\top \mathbf{x}_n)$

Minimize  $\frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta}$

subject to  $y_n(\boldsymbol{\theta}^\top \mathbf{x}_n) \geq 1$ , for  $n = 1, 2, \dots, N$

# Alternative View of Logistic Regression

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^{\top}x}}$$



If  $y = 1$ , we want  $h_{\theta}(x) \approx 1, \theta^{\top}x \gg 0$

If  $y = 0$ , we want  $h_{\theta}(x) \approx 0, \theta^{\top}x \ll 0$

$$J(\theta) = - \sum_{i=1}^n [y_i \underbrace{\log h_{\theta}(x_i)}_{cost_1(\theta^{\top}x_i)} + (1 - y_i) \underbrace{\log(1 - h_{\theta}(x_i))}_{cost_0(\theta^{\top}x_i)}]$$

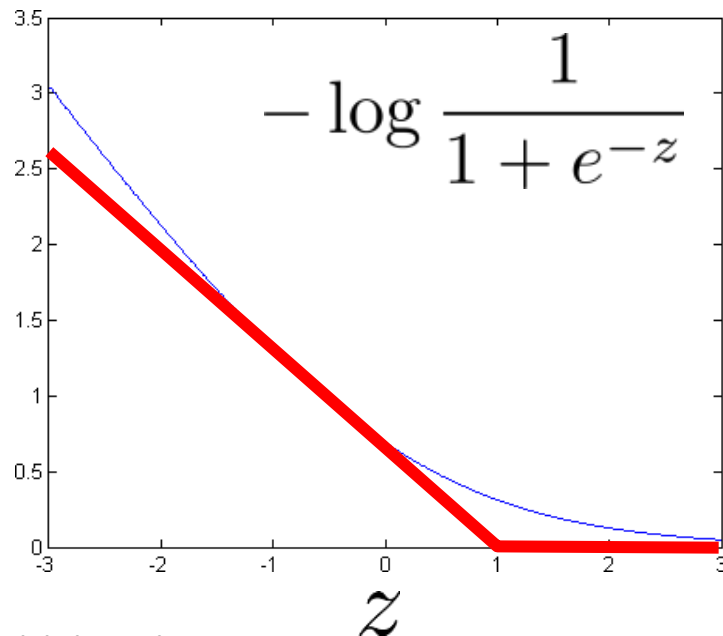
$\min_{\theta} J(\theta)$

# Alternative View of Logistic Regression

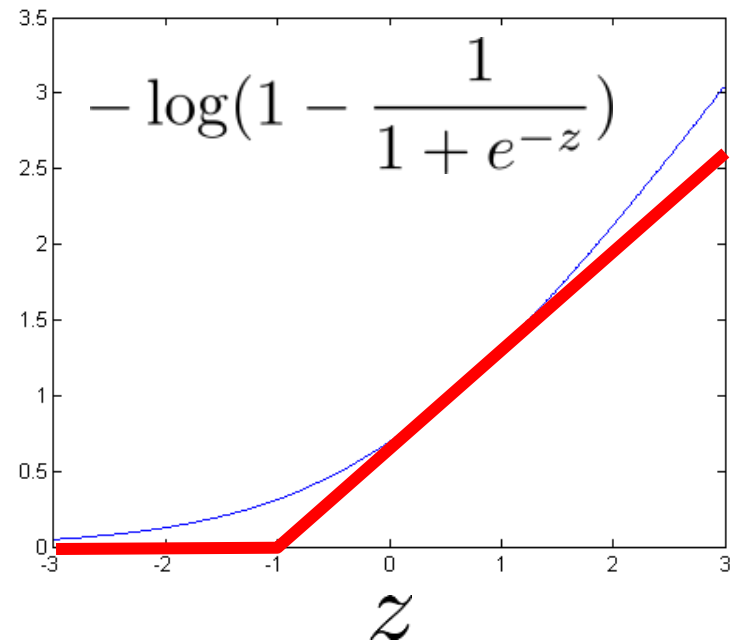
Cost of example:  $-y_i \log h_{\theta}(\mathbf{x}_i) - (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))$

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \quad z = \theta^T \mathbf{x}$$

If  $y = 1$  (want  $\theta^T \mathbf{x} \gg 0$ ):

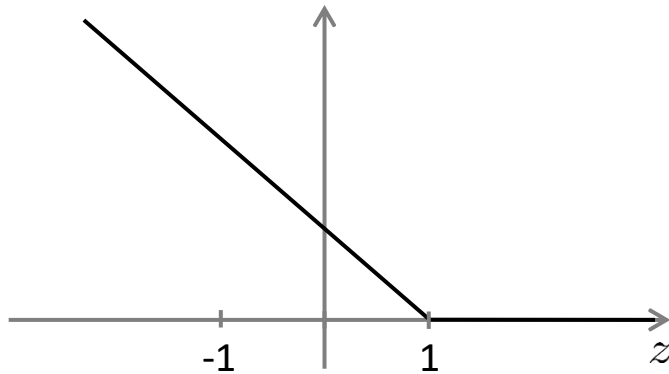


If  $y = 0$  (want  $\theta^T \mathbf{x} \ll 0$ ):

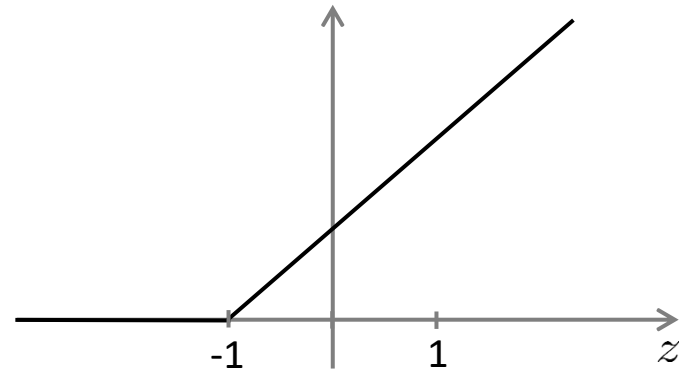


# Support Vector Machine

If  $y = 1$  (want  $\theta^\top x \geq 1$ )



If  $y = -1$  (want  $\theta^\top x \leq -1$ )



$$l_{hinge}(h(\mathbf{x})) = \max(0, 1 - y \cdot h(\mathbf{x}))$$



# Support Vector Machine

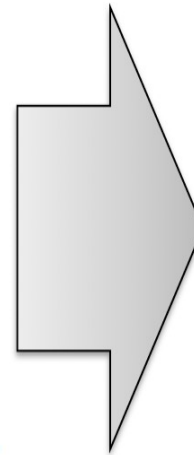
$$\min_{\theta} C \sum_{i=1}^n [y_i \text{cost}_1(\theta^\top x_i) + (1 - y_i) \text{cost}_0(\theta^\top x_i)] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

$y = 1 / 0$

with  $C = 1$

$y = +1 / -1$

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{2} \sum_{j=1}^d \theta_j^2 \\ \text{s.t.} \quad & \theta^\top \mathbf{x}_i \geq 1 \quad \text{if } y_i = 1 \\ & \theta^\top \mathbf{x}_i \leq -1 \quad \text{if } y_i = -1 \end{aligned}$$



$$\begin{aligned} \min_{\theta} \quad & \frac{1}{2} \sum_{j=1}^d \theta_j^2 \\ \text{s.t.} \quad & y_i (\theta^\top \mathbf{x}_i) \geq 1 \end{aligned}$$