

The Maximum Margin Principle

Announcements

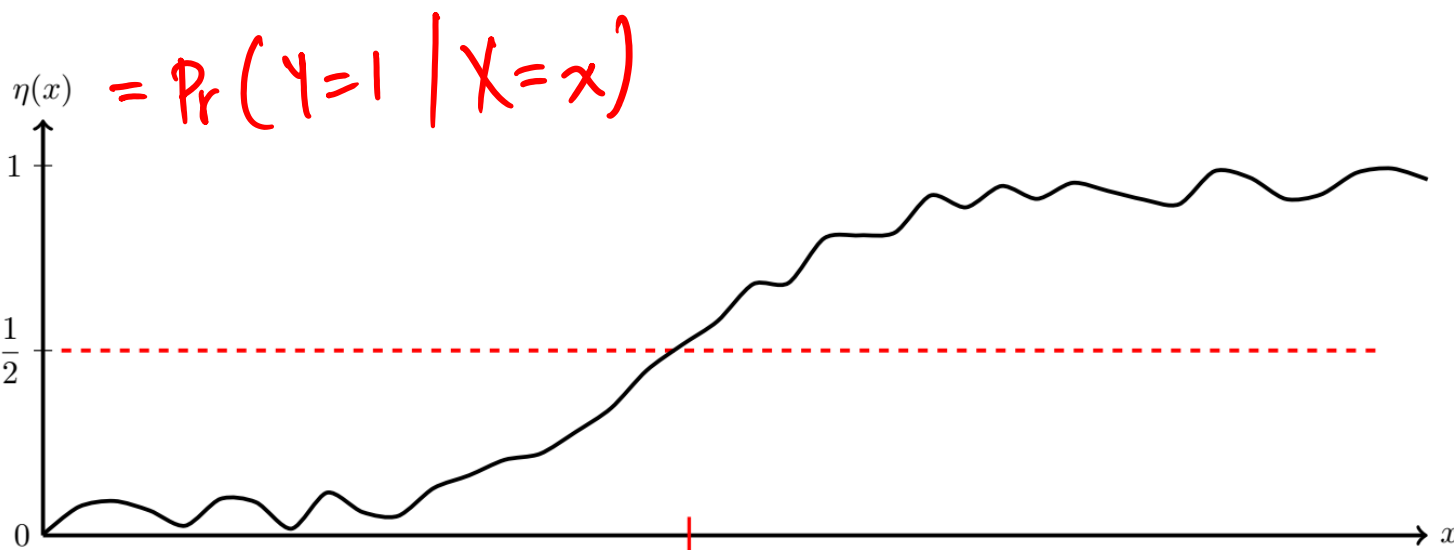
- None

Outline

- Hyperplanes
- Max-margin hyperplanes
- Optimal soft-margin hyperplanes

Drawback of Plug-in Classifiers

- Plug-in methods require estimation of (conditional) densities or mass functions, which can be more difficult than estimating a decision boundary



$\eta(x)$ is quite complicated but the decision regions are simple and η is smooth near $1/2$

Linear Classifiers

- Binary classification
- Training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$.
- Assume the labels are -1 and 1
- Recall a linear classifier has the form

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$
$$\text{sign}(t) = \begin{cases} 1 & t \geq 0 \\ -1 & t < 0 \end{cases}$$

- How can we use the training data to directly optimize for \mathbf{w} and b ?

$$\mathbf{x}_i \in \mathbb{R}^d$$

$$\mathbf{w} \in \mathbb{R}^d$$

$$b \in \mathbb{R}$$

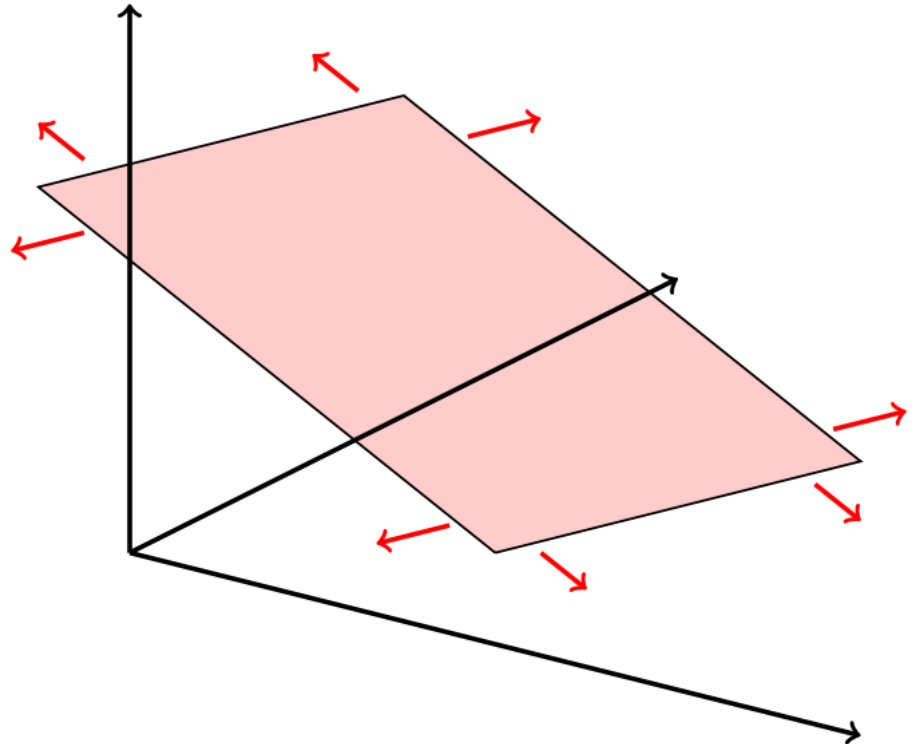
Hyperplanes

- A *hyperplane* is a subset of \mathbb{R}^d of the form

$$H = \{x : w^T x + b = 0\}$$

for some $w \in \mathbb{R}^d$, $b \in \mathbb{R}$.

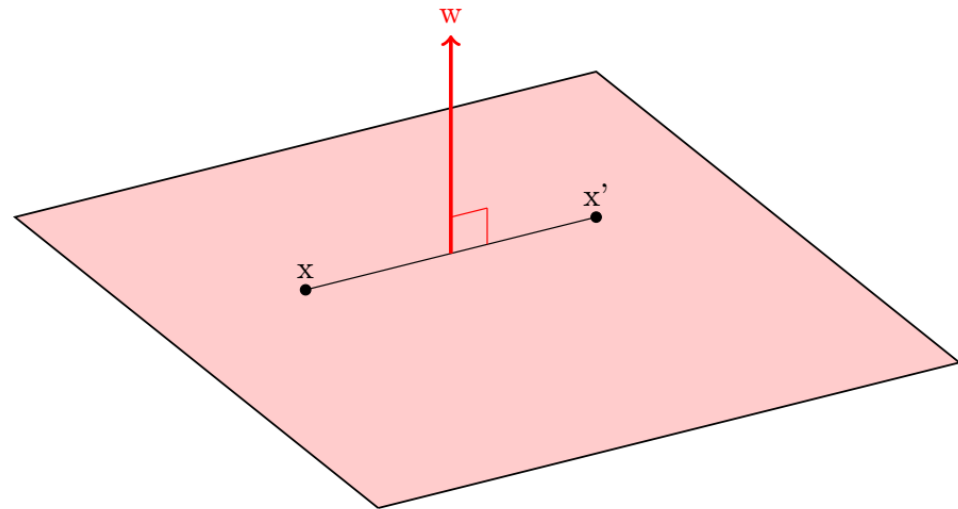
- In general, a hyperplane is an *affine subspace* of dimension $d - 1$



Normal Vectors

- The vector w is orthogonal to the hyperplane, and for this reason is called a *normal vector*.
- To say that w is ~~orthogonal to a hyperplane means~~ *a normal vector* that it is orthogonal to every vector that lies in the hyperplane. Every such vector can be written as the difference of two points x and x' in the hyperplane.

$$\left. \begin{array}{l} w^T x + b = 0 \\ w^T x' + b = 0 \end{array} \right\} \Rightarrow w^T (x - x') = 0 \Rightarrow w \perp x - x'$$



Distance from Point to Hyperplane

- Let $\mathcal{H} = \{x \mid w^T x + b = 0\}$. The distance from $z \notin \mathcal{H}$ to \mathcal{H} is

$$\frac{|w^T z + b|}{\|w\|}$$

- Let $\mathcal{H} = \{x \mid x_1 - 5x_2 + 5x_3 - 7 = 0\} \subseteq \mathbb{R}^3$. The distance from \mathcal{H} to

$$z = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

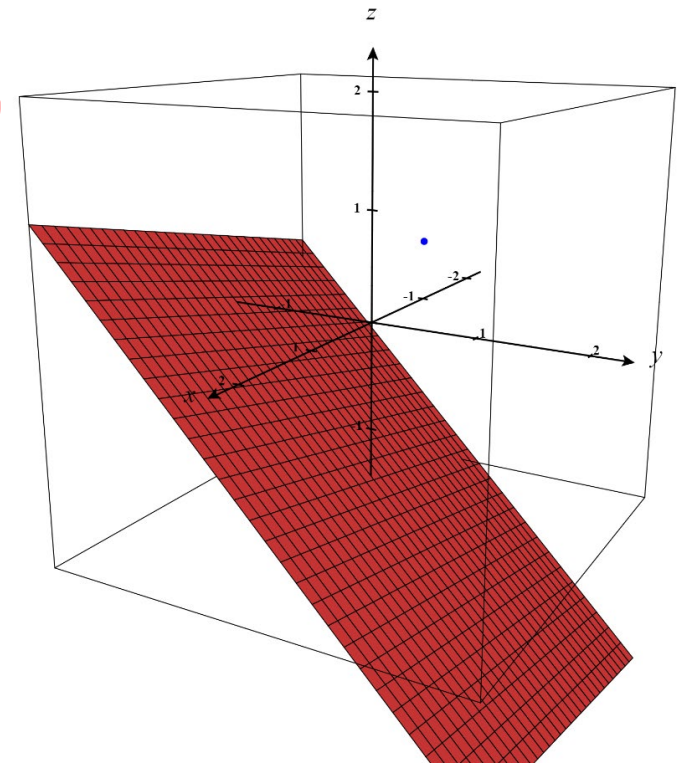
$$w = \begin{bmatrix} 1 \\ -5 \\ 5 \end{bmatrix}, b = -7$$

is

$$|w^T z + b| = |1 - 5 + 5 - 7| = 6$$

$$\|w\| = \sqrt{1^2 + (-5)^2 + 5^2} = \sqrt{51}$$

$$\text{distance} : \frac{6}{\sqrt{51}}$$



Separating Hyperplanes

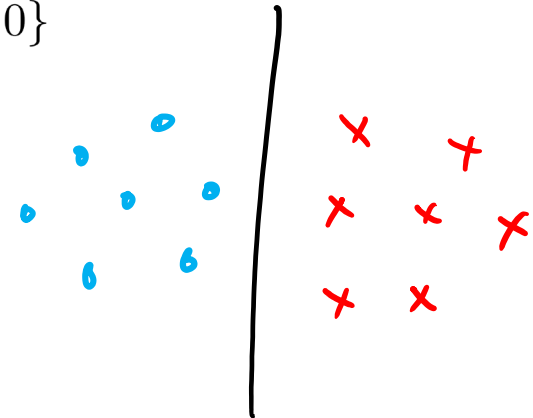
- Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ be training data for a binary classification problem
- Assume $y_i \in \{-1, 1\}$.
- We say the training data are *linearly separable* if there exist $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ such that

$$\forall i \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) > 0$$

- In this case we refer to $\underbrace{\mathbf{w}^T \mathbf{x}_i + b}_{\text{take sign to get the predicted label}}$

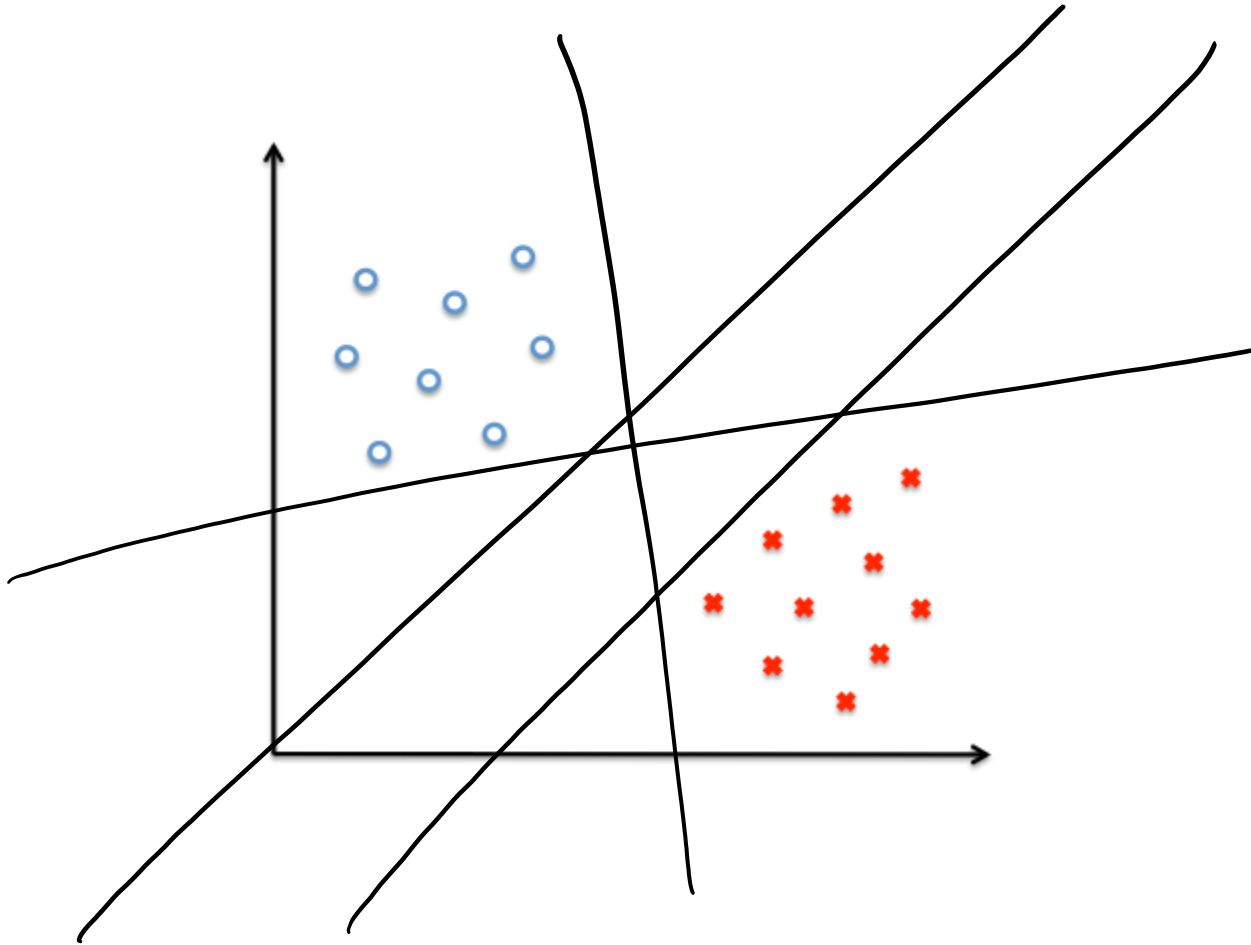
$$\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$$

as a *separating hyperplane*.



Separating Hyperplanes

- Are all separating hyperplanes equally good?



Poll: Margin of a Hyperplane

- Let $\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$ be a separating hyperplane.
- The *margin* ρ of a \mathcal{H} is the distance from \mathcal{H} to the nearest training point \mathbf{x}_i .
- **Poll:** A formula for ρ is

(A) $\rho(\mathbf{w}, b) = \min_{\mathbf{z} \in \mathbb{R}^d} \frac{|\mathbf{w}^T \mathbf{z} + b|}{\|\mathbf{w}\|}$

(B) $\rho(\mathbf{w}, b) = \max_{\mathbf{z} \in \mathbb{R}^d} \frac{|\mathbf{w}^T \mathbf{z} + b|}{\|\mathbf{w}\|}$

(C) $\rho(\mathbf{w}, b) = \min_{i=1, \dots, n} \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$

(D) $\rho(\mathbf{w}, b) = \max_{i=1, \dots, n} \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$



Max-Margin Hyperplane

- The *margin* ρ of a separating hyperplane is the distance from the hyperplane to the nearest training point:

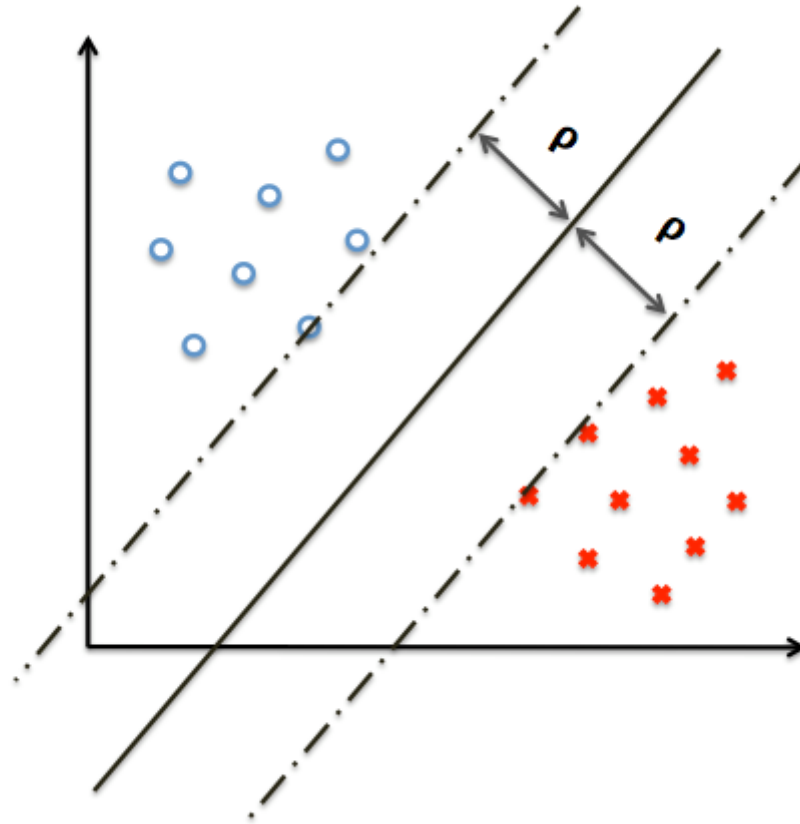
$$\rho(\mathbf{w}, b) := \min_{i=1, \dots, n} \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

- The *maximum margin* or *optimal* separating hyperplane is the solution of

$$\max_{\mathbf{w}, b} \left(\min_{i=1, \dots, n} \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|} \right)$$

$$\text{s.t. } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0 \quad \forall i=1, \dots, n$$

Max-Margin Hyperplane



Canonical Form

- A separating hyperplane is said to be in *canonical form* if \mathbf{w} and b are such that

$$\begin{aligned} y_i(\mathbf{w}^T \mathbf{x}_i + b) &\geq 1 \quad \forall i \\ y_i(\mathbf{w}^T \mathbf{x}_i + b) &= 1 \quad \text{for some } i \end{aligned}$$

- Every separating hyperplane can be represented in canonical form. If $\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$ and $\alpha > 0$, then

$$\mathcal{H} = \{\mathbf{x} : (\alpha \mathbf{w})^T \mathbf{x} + (\alpha b) = 0\}$$

and

$$\forall \mathbf{x} \quad \text{sign} \{ (\alpha \mathbf{w})^T \mathbf{x} + (\alpha b) \} = \text{sign} \{ \mathbf{w}^T \mathbf{x} + b \}$$

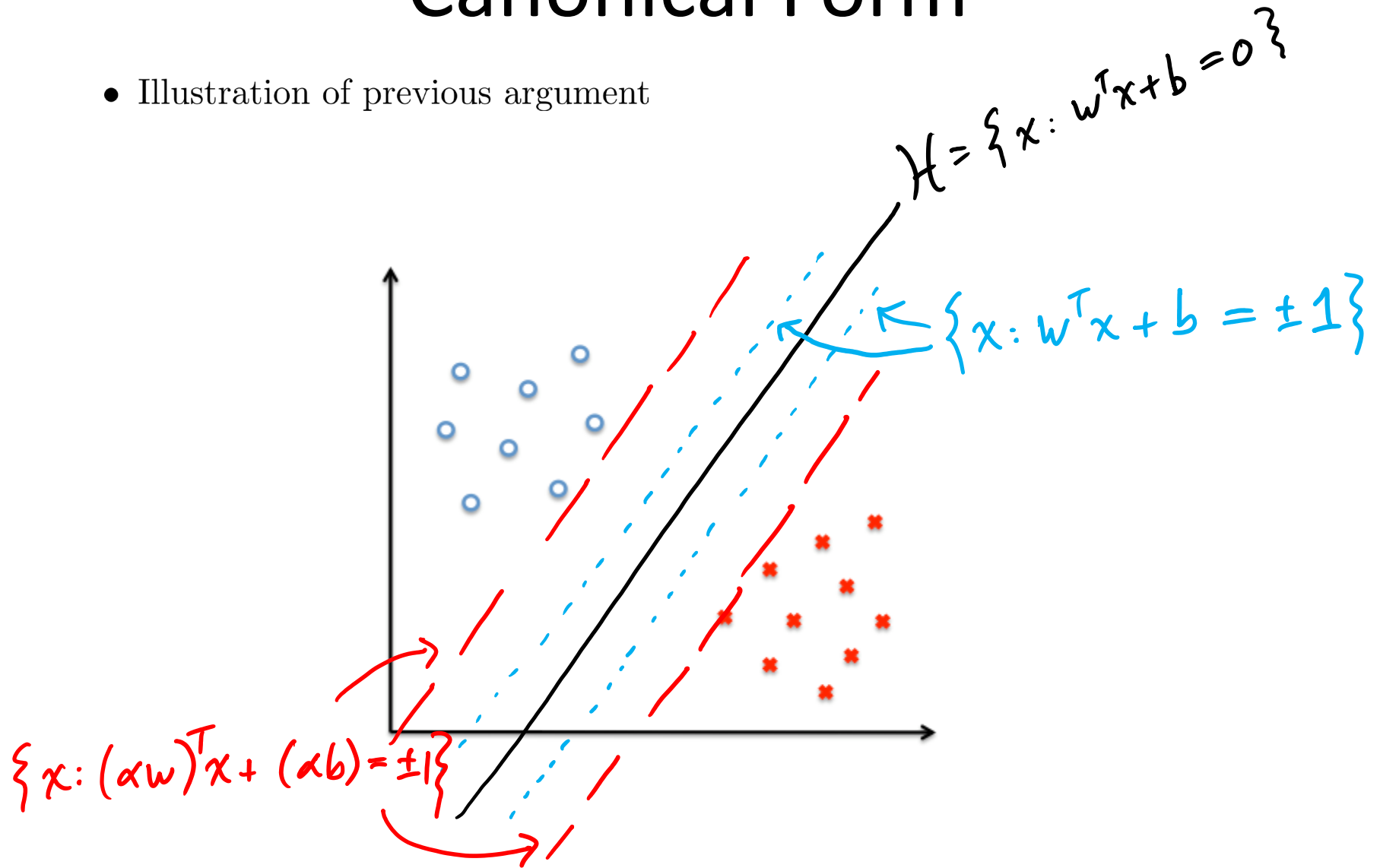
- Thus we can always scale \mathbf{w} and b such that the smallest value of

$$y_i((\alpha \mathbf{w})^T \mathbf{x}_i + (\alpha b))$$

is 1

Canonical Form

- Illustration of previous argument



Max-Margin Hyperplane

- This allows us to write the max-margin hyperplane

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

- Equivalently,

$$\begin{aligned} \max_{\mathbf{w}, b} \quad & \frac{1}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & \forall i \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ & \exists i \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 \end{aligned}$$

\Rightarrow

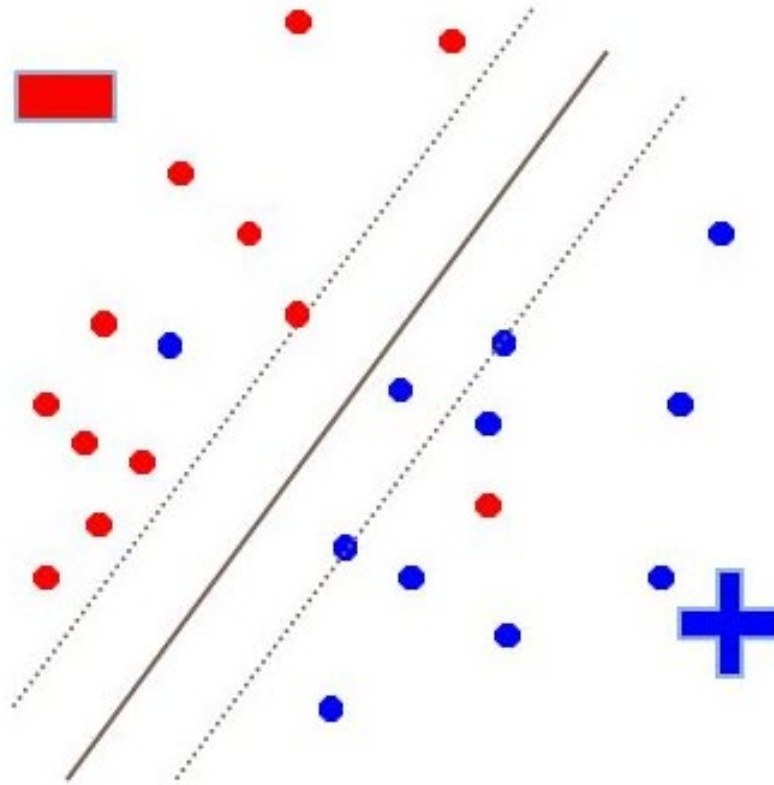
$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \forall i \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ & \exists i \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 \end{aligned}$$

↑ satisfied even if omitted

- This is an example of a quadratic program

Non-Separable Data

- What if the training data are not linearly separable?



ksee

Optimal Soft-Margin Hyperplane

- Introduce *slack variables* $\xi_1, \dots, \xi_n \geq 0$.
- The *optimal soft-margin hyperplane* is the solution of

$$\xi = \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_n \end{bmatrix}$$

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

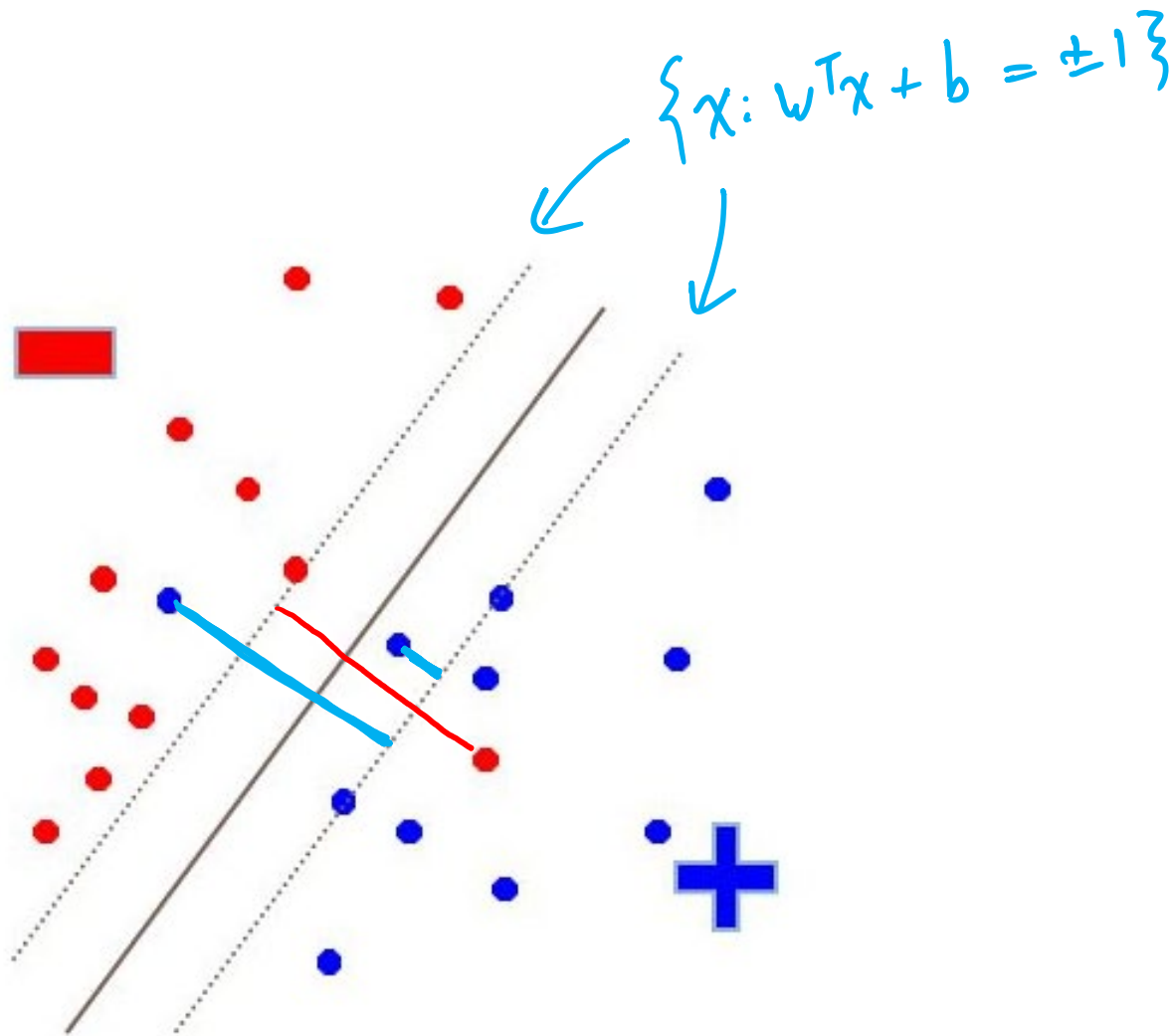
$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$$

$$\xi_i \geq 0$$

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

- C is a user-defined parameter
- This is another quadratic program

Optimal Soft-Margin Hyperplane



Poll 2

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

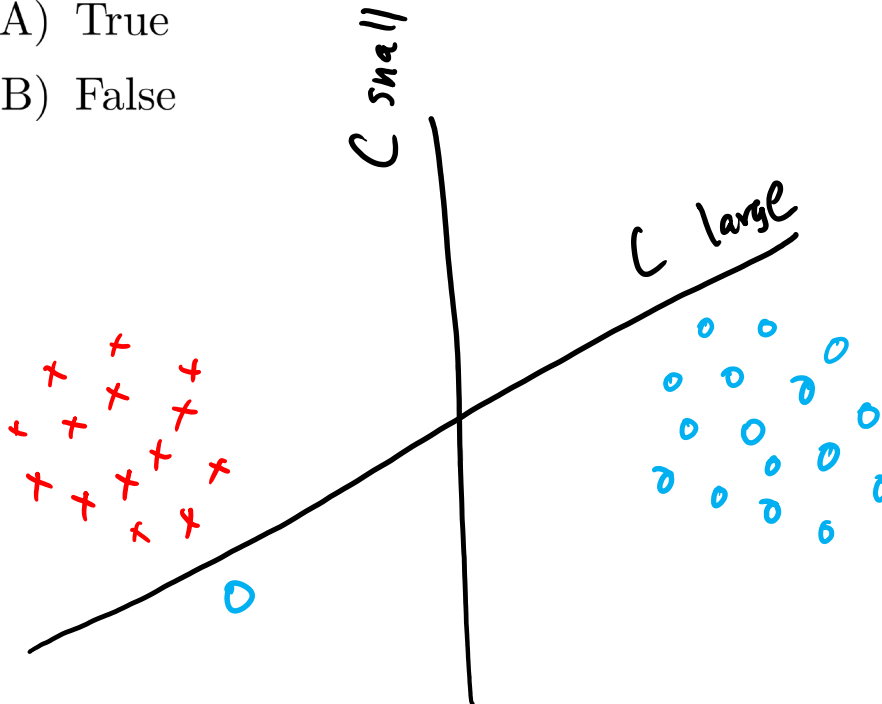
$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$$

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

- True or False: As C increases, the solution becomes more sensitive to outliers like the one shown

(A) True

(B) False



Poll 3

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

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

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

- True or False: If $C = 0$, the OSM hyperplane recovers the max-margin hyperplane in the case of linearly separable data.

(A) True

(B) False ✓

$\mathbf{w} = \mathbf{0}$ is optimal

Closing Thoughts

- The optimal soft margin hyperplane classifier is a special case of a much more general classifier that we will study soon: the support vector machine