

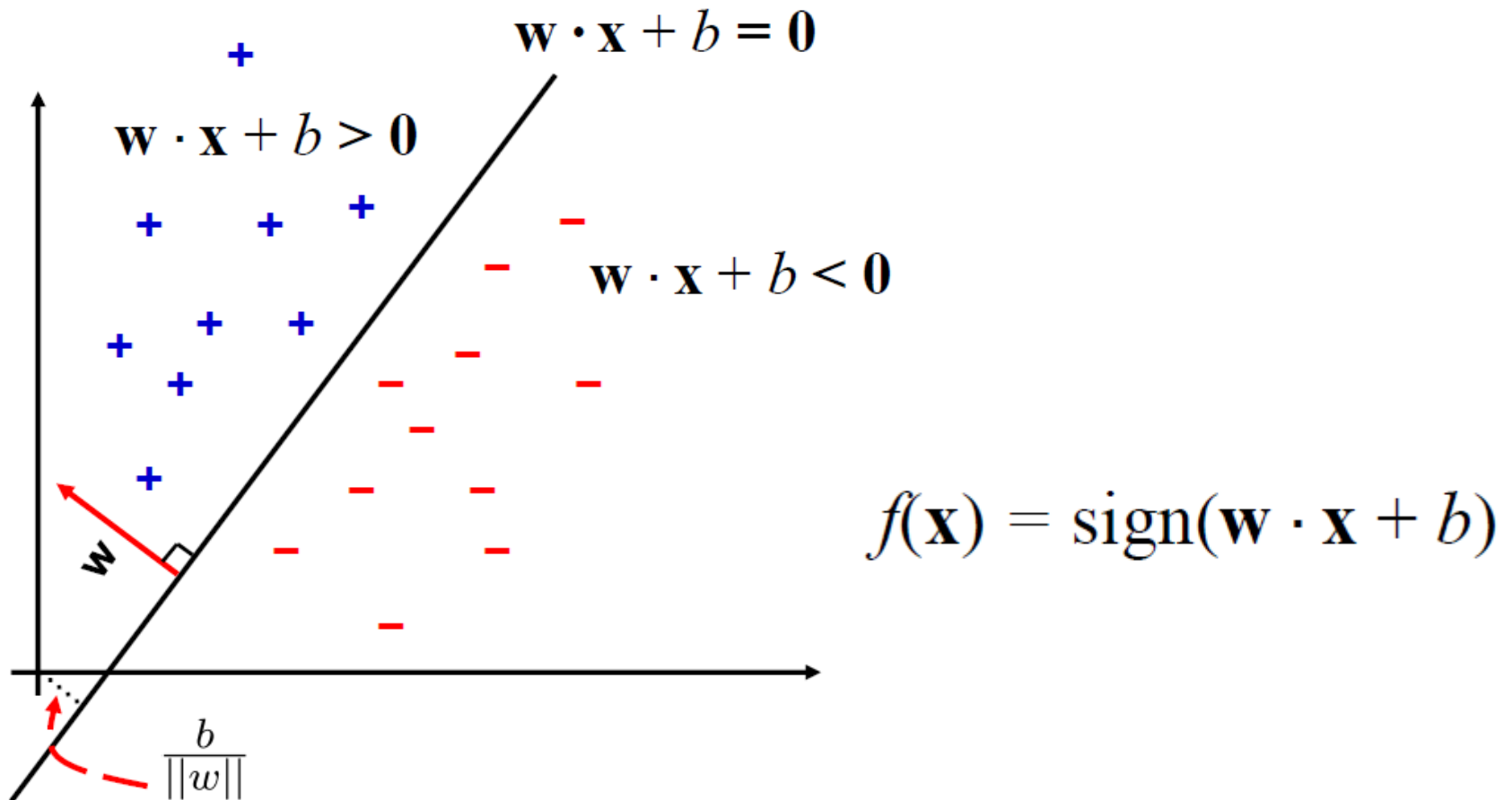
# Lecture #4: Max-Margin Classification and Support Vector Machines

**Janardhan Rao (Jana) Doppa**

School of EECS, Washington State University

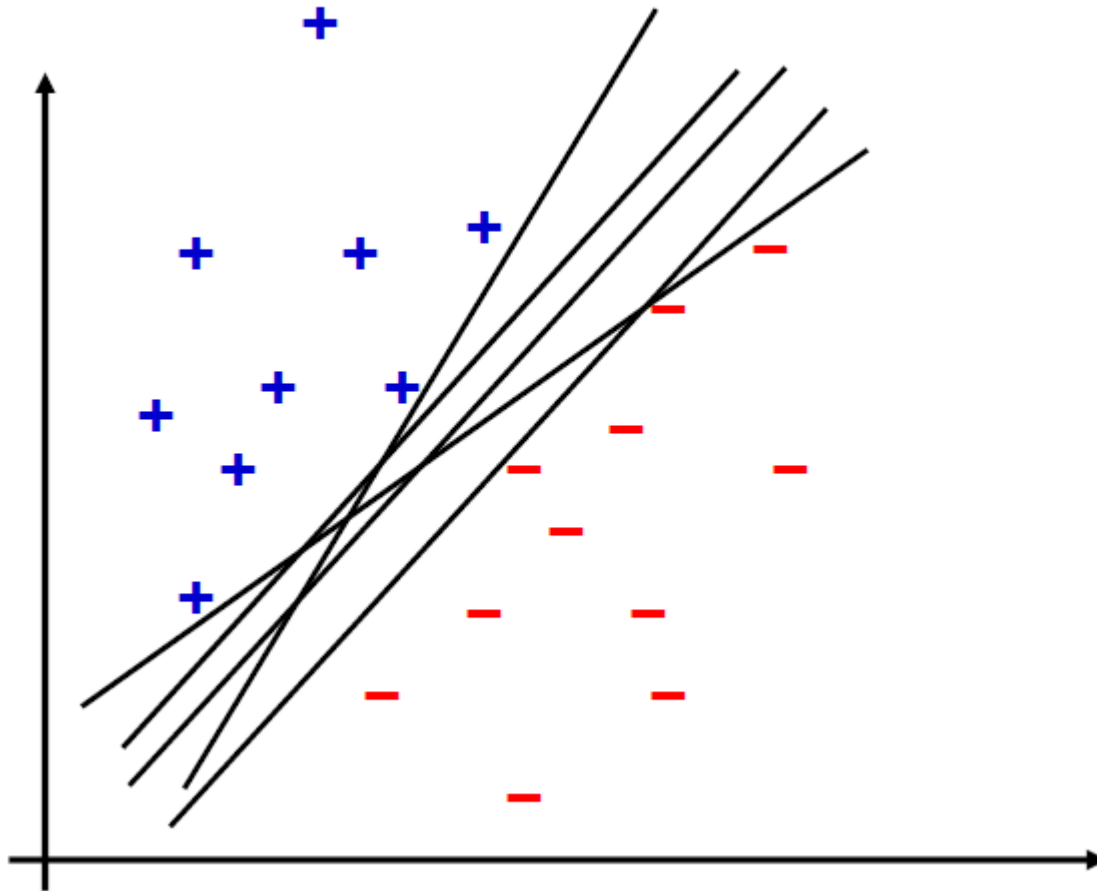
# Perceptron Revisited: Linear Separator

- Binary classification can be viewed as the task of separating classes in a given feature space



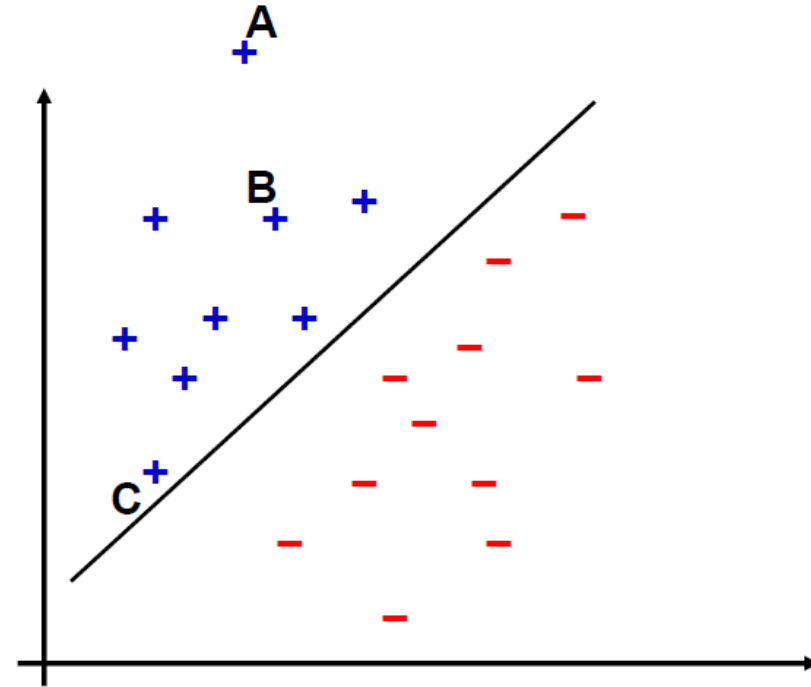
# Linear Separators

- Which of the linear separators is optimal?



# Intuition of Margin

- Consider points A, B, and C
- We are quite confident in our prediction for A because it is far from the decision boundary
- In contrast, we are not so confident in our prediction for C because a slight change in the decision boundary may flip the decision

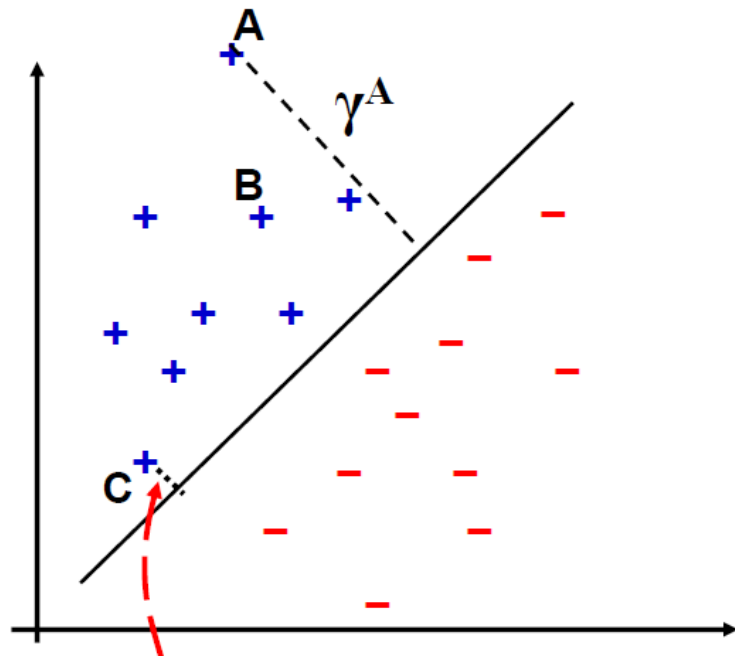


Given a training set, we would like to make all predictions correct and confident! This leads to the concept of margin.

# Geometric Margin

- The geometric margin of  $(\mathbf{w}, b)$  w.r.t. example  $(x_i, y_i)$  is the distance from  $x_i$  to the decision surface, which can be computed as:

$$\gamma_i = \frac{y_i(\mathbf{w} \cdot \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$



- Given a training set  $S = (x_i, y_i): i = 1, 2, \dots, N$  the geometric margin of the classifier w.r.t.  $S$  is

$$\gamma = \min_{i=1,2,\dots,N} \gamma_i$$

# Maximum Margin Classifier

- Given a linearly separable training set  $(x_i, y_i): i = 1, 2, \dots, N$ , we would like to find a linear classifier with maximum margin
- This can be represented as an optimization problem

$$\max_{\mathbf{w}, b, \gamma} \gamma$$

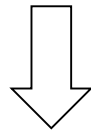
$$\text{subject to: } y^{(i)} \frac{(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|} \geq \gamma, \quad i = 1, \dots, N$$

# Maximum Margin Classifier

- Let  $\gamma' = \gamma \|w\|$ , we can rewrite the optimization problem as follows:

$$\max_{\mathbf{w}, b, \gamma}$$

$$\text{subject to: } y^{(i)} \frac{(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|} \geq \gamma, \quad i = 1, \dots, N$$



$$\max_{\mathbf{w}, b, \gamma'} \frac{\gamma'}{\|\mathbf{w}\|}$$

$$\text{subject to: } y^i (\mathbf{w} \cdot \mathbf{x}^i + b) \geq \gamma', \quad i = 1, \dots, N$$

# Maximum Margin Classifier

- Note that rescaling  $w$  and  $b$  by  $1/\gamma'$  will not change the classifier -- we can thus further reformulate the optimization problem

$$\begin{aligned} & \max_{\mathbf{w}, b} \frac{\gamma'}{\|\mathbf{w}\|} \\ & \text{subject to : } y^i (\mathbf{w} \cdot \mathbf{x}^i + b) \geq \gamma', \quad i = 1, \dots, N \end{aligned}$$



$$\begin{aligned} & \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \quad (\text{or equivalently } \min_{\mathbf{w}, b} \|\mathbf{w}\|^2) \\ & \text{subject to : } y^i (\mathbf{w} \cdot \mathbf{x}^i + b) \geq 1, \quad i = 1, \dots, N \end{aligned}$$



# Maximum Margin Classifier

- Maximizing the geometric margin is equivalent to minimizing the magnitude of  $\mathbf{w}$  subject to maintaining a functional margin of at least 1

$$\max_{\mathbf{w}, b} \frac{\gamma'}{\|\mathbf{w}\|}$$

$$\text{subject to : } y^i (\mathbf{w} \cdot \mathbf{x}^i + b) \geq \gamma', \quad i = 1, \dots, N$$



$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \quad (\text{or equivalently } \min_{\mathbf{w}, b} \|\mathbf{w}\|^2)$$

$$\text{subject to : } y^i (\mathbf{w} \cdot \mathbf{x}^i + b) \geq 1, \quad i = 1, \dots, N$$

# Maximum Margin Classifier: Formulation

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to:  $y^i (\mathbf{w} \cdot \mathbf{x}^i + b) \geq 1, \quad i = 1, \dots, N$

- This results in a **quadratic optimization (QP) problem** with linear inequality constraints
- This is a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist
  - ▲ One could solve for  $\mathbf{w}$  using any of these methods

# Characteristics of Solution

- Weights  $\mathbf{w}$  can be represented as a linear combination of the training examples

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

# Characteristics of Solution

- Many of the  $\alpha_i$ 's are zero
  - ▲ Weights  $\mathbf{w}$  is a linear combination of small number of data points
- $\mathbf{x}_i$  with non-zero  $\alpha_i$  are called support vectors (SVs)
  - ▲ The decision boundary is determined only by the SVs
  - ▲ Let  $t_j (j = 1, 2, \dots, s)$  be the indices of the  $s$  support vectors. We can write

$$\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}$$

# Characteristics of Solution

- For classifying a new input example  $\mathbf{z}$ , compute

$$\mathbf{w}^T \mathbf{z} + b = \sum_{j=1}^s \alpha_{t_j} y_{t_j} (\mathbf{x}_{t_j}^T \mathbf{z}) + b$$

- ▲ Classify  $\mathbf{z}$  as positive if the sum is positive, and negative otherwise
- ▲ **Note:**  $\mathbf{w}$  need not be formed explicitly, rather we can classify  $\mathbf{z}$  by taking inner products with the support vectors (useful when we generalize the notion of inner product later)