# Support Vector Machine (SVM)

Shusen Wang

# Project a Point onto a Hyperplane

# Project a Point onto a Hyperplane

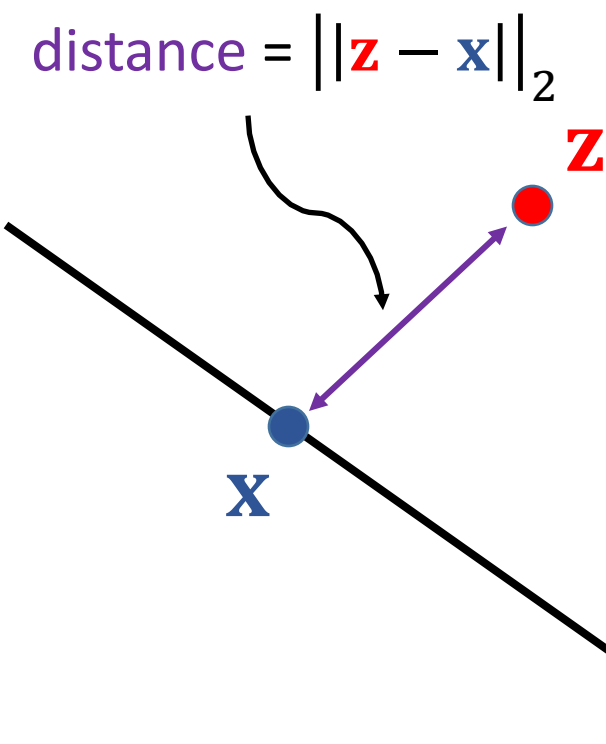**Question**: how to project **z** onto the hyperplane?

**z**
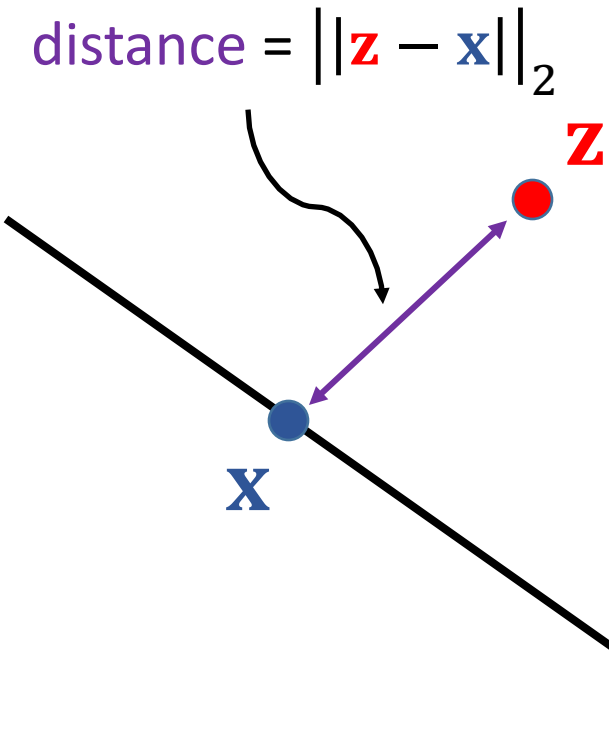
Hyperplane $\mathbf{w}^T\mathbf{x} + b = 0$

# Project a Point onto a Hyperplane

**Question**: how to project $\mathbf{z}$ onto the hyperplane?

**Solution**: find $\mathbf{x}$ on the hyperplane such that $\left\|\mathbf{z}-\mathbf{x}\right\|_2^2$ is minimized.

- $\min\limits_{\mathbf{x}}\left\|\mathbf{z}-\mathbf{x}\right\|_2^2;$    s.t. $\mathbf{w}^T\mathbf{x}+b=0$

distance $=\left\|\mathbf{z}-\mathbf{x}\right\|_2$

$\mathbf{Z}$

$\mathbf{X}$

Hyperplane $\mathbf{w}^T\mathbf{x}+b=0$

# Project a Point onto a Hyperplane

**Solution**: find $\mathbf{x}$ on the hyperplane such that $\left|\left|\mathbf{z} - \mathbf{x}\right|\right|_2^2$ is minimized.

distance = $\left|\left|\mathbf{z} - \mathbf{x}\right|\right|_2$

$\mathbf{Z}$

$\mathbf{X}$

Hyperplane $\mathbf{w}^T\mathbf{x} + b = 0$

- $\min\limits_{\mathbf{x}}\left|\left|\mathbf{z} - \mathbf{x}\right|\right|_2^2;$    s.t. $\mathbf{w}^T\mathbf{x} + b = 0$

- Solve the problem using the KKT conditions:

$$\begin{cases} \dfrac{\partial\left|\left|\mathbf{z} - \mathbf{x}\right|\right|_2^2}{\partial\mathbf{x}} + \lambda\dfrac{\partial\left(\mathbf{w}^T\mathbf{x} + b\right)}{\partial\mathbf{x}} = 0; \\ \mathbf{w}^T\mathbf{x} + b = 0. \end{cases}$$

- Solution: $\mathbf{x} = \mathbf{z} - \dfrac{\mathbf{w}^T\mathbf{z} + b}{\left|\left|\mathbf{w}\right|\right|_2^2}\mathbf{w}$

# Project a Point onto a Hyperplane

**Question**: how to project $\mathbf{z}$ onto the hyperplane?

**Solution**: find $\mathbf{x}$ on the hyperplane such that $\left\|\mathbf{z} - \mathbf{x}\right\|_2^2$ is minimized.

distance $= \left\|\mathbf{z} - \mathbf{x}\right\|_2$

$\mathbf{z}$

$\mathbf{x}$

- Solution: $\mathbf{x} = \mathbf{z} - \dfrac{\mathbf{w}^T\mathbf{z} + b}{\left\|\mathbf{w}\right\|_2^2}\mathbf{w}$

- The $\ell_2$ distance between $\mathbf{z}$ and the hyperplane is

$$\left\|\mathbf{z} - \mathbf{x}\right\|_2 = \frac{\left|\mathbf{w}^T\mathbf{z} + b\right|}{\left\|\mathbf{w}\right\|_2}.$$

z
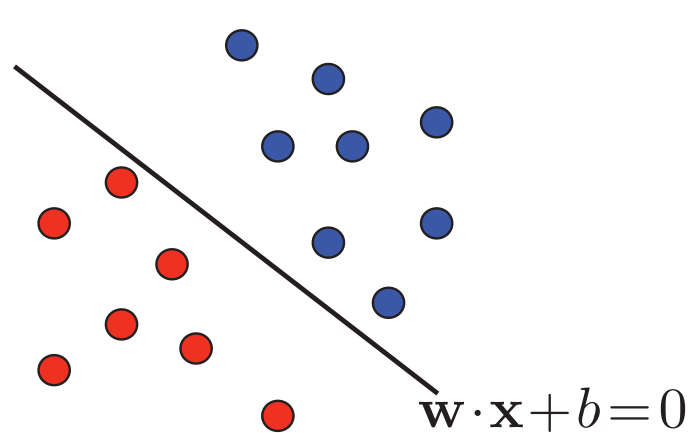
Hyperplane $\mathbf{w}^T\mathbf{x} + b = 0$

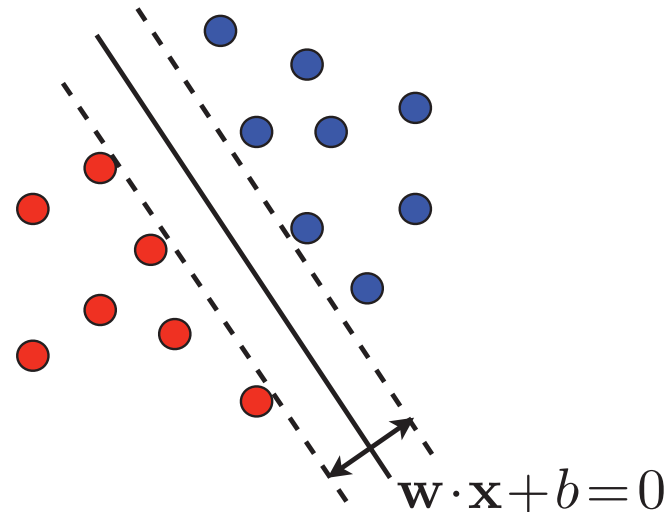# Support Vector Machine (SVM)

# Support Vector Machine (SVM)

Separate data by a hyperplane (assume the data are separable)



$\mathbf{w}\cdot\mathbf{x}+b=0$

$\mathbf{w}\cdot\mathbf{x}+b=0$

An arbitrary hyperplane.

The hyperplane that maximizes the margin.
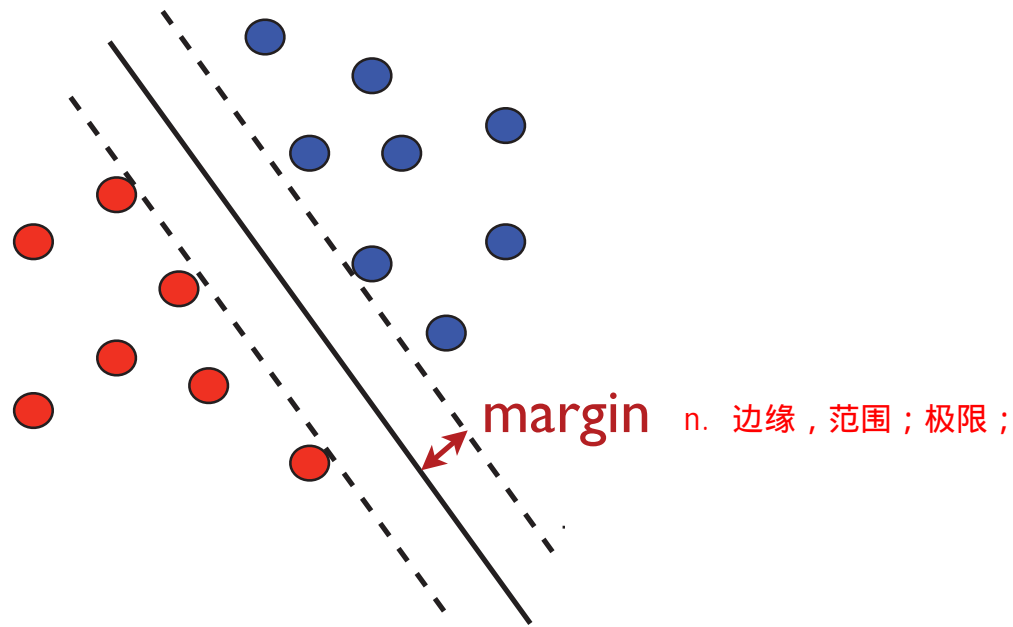
The figure is from the book "*Foundations of Machine Learning*"
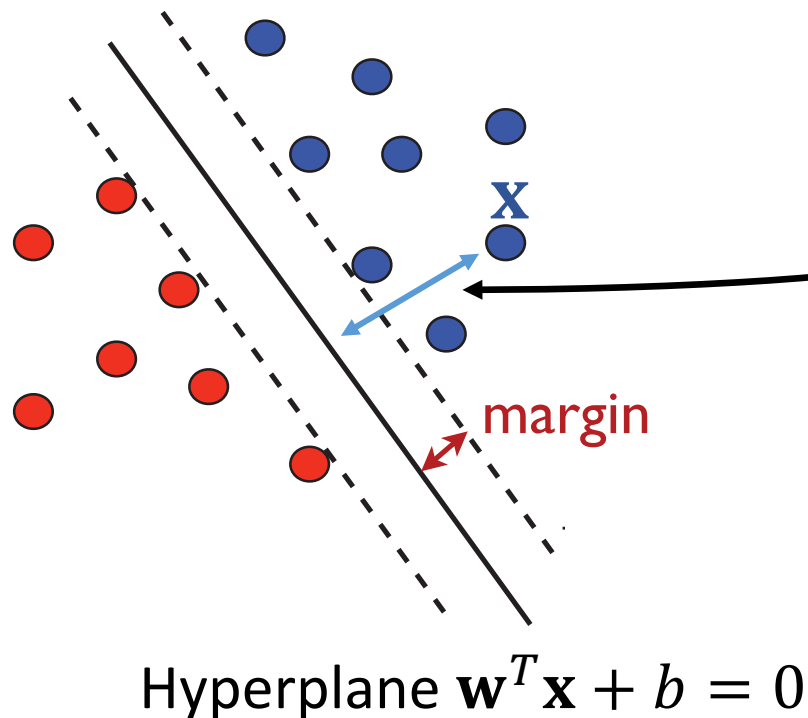
# Support Vector Machine (SVM)

Separate data by a hyperplane (assume the data are separable)



margin  n.

Hyperplane $\mathbf{w}^T\mathbf{x} + b = 0$

# Support Vector Machine (SVM)

Separate data by a hyperplane (assume the data are separable)



- The distance between any feature vector, **x**, and the hyperplane is

$$\text{dist} = \frac{\left| \mathbf{w}^T \mathbf{x} + b \right|}{\|\mathbf{w}\|_2}.$$

z                                x

margin

Hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$

# Support Vector Machine (SVM)

Separate data by a hyperplane (assume the data are separable)



Hyperplane $\mathbf{w}^T\mathbf{x} + b = 0$
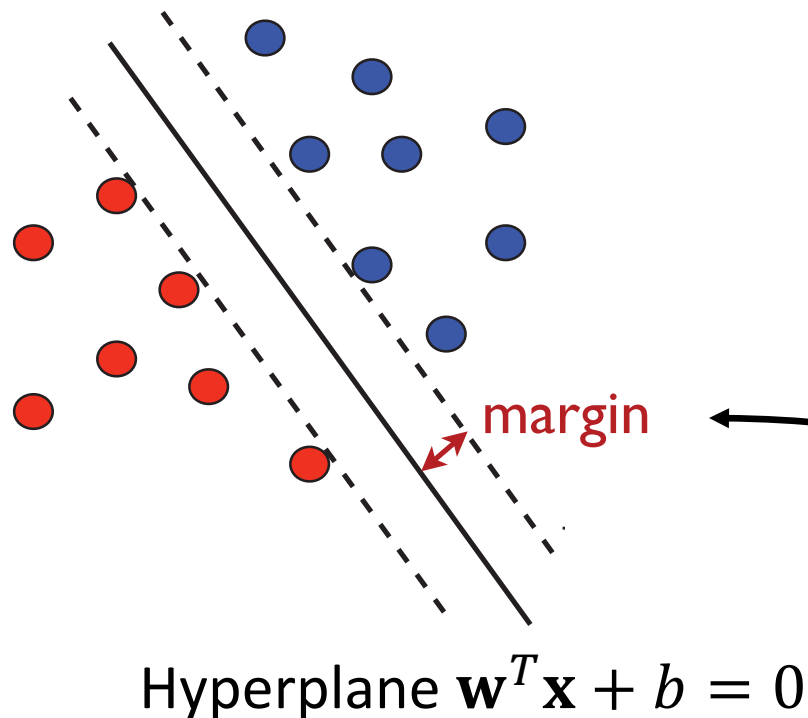
- The distance between any feature vector, $\mathbf{x}$, and the hyperplane is
$$\text{dist} = \frac{|\mathbf{w}^T\mathbf{x}+b|}{||\mathbf{w}||_2}.$$

- The margin is the smallest distance:
$$\min_j \frac{|\mathbf{w}^T\mathbf{x}_j+b|}{||\mathbf{w}||_2}$$

*The figure is from the book "Foundations of Machine Learning"*

# Support Vector Machine (SVM)

Separate data by a hyperplane (assume the data are separable)

Positive samples
$(y_j = +1)$

margin

Negative samples
$(y_j = -1)$

Hyperplane $\mathbf{w}^T\mathbf{x} + b = 0$

- The distance between any feature vector, $\mathbf{x}$, and the hyperplane is
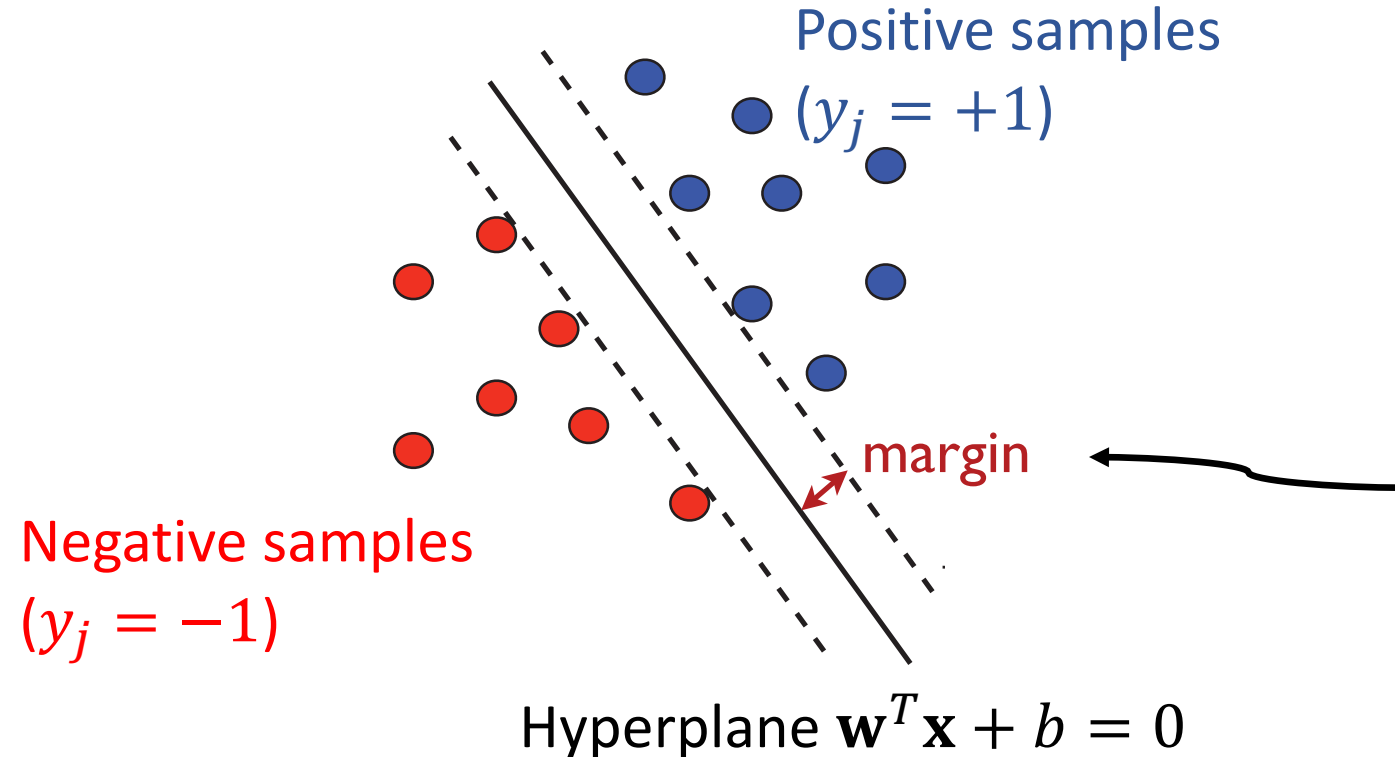$$\text{dist} = \frac{|\mathbf{w}^T\mathbf{x}+b|}{\|\mathbf{w}\|_2}.$$

- The margin is the smallest distance:
$$\min_j \frac{|\mathbf{w}^T\mathbf{x}_j+b|}{\|\mathbf{w}\|_2} = \min_j \frac{y_j(\mathbf{w}^T\mathbf{x}_j+b)}{\|\mathbf{w}\|_2}$$

# Support Vector Machine (SVM)

Margin $= \min_{j} \dfrac{y_j(\mathbf{w}^T\mathbf{x}_j + b)}{||\mathbf{w}||_2}$ ; we want to maximize the margin.

# Support Vector Machine (SVM)

$$\text{Margin} = \min_j \frac{y_j \left( \mathbf{w}^T \mathbf{x}_j + b \right)}{\|\mathbf{w}\|_2} \; ; \text{ we want to maximize the margin.}$$

Define $\bar{\mathbf{x}}_j = \left[ \mathbf{x}_j ; 1 \right] \in \mathbb{R}^{d+1}$

Define $\bar{\mathbf{w}} = [\mathbf{w}, b] \in \mathbb{R}^{d+1}$

➔ $\mathbf{x}_j^T \mathbf{w} + b = \bar{\mathbf{x}}_j^T \bar{\mathbf{w}}$

# Support Vector Machine (SVM)

Margin $= \min\limits_{j} \dfrac{y_j \mathbf{w}^T \mathbf{x}_j}{||\mathbf{w}||_2}$ ;    we want to maximize the margin.

Support Vector Machine (SVM):    $\max\limits_{\mathbf{w}} \min\limits_{j} \dfrac{y_j \mathbf{w}^T \mathbf{x}_j}{||\mathbf{w}||_2}$

# Support Vector Machine (SVM)

Support Vector Machine (SVM): $\max\limits_{\mathbf{w}} \min\limits_{j} \dfrac{y_j \mathbf{w}^T \mathbf{x}_j}{\left\|\mathbf{w}\right\|_2}$

# Support Vector Machine (SVM)

Support Vector Machine (SVM): $\quad \max_{\mathbf{w}} \min_{j} \dfrac{y_j \mathbf{w}^T \mathbf{x}_j}{||\mathbf{w}||_2}$

$$\operatorname*{argmax}_{\mathbf{w}} \min_{j} \frac{y_j \mathbf{w}^T \mathbf{x}_j}{||\mathbf{w}||_2} = \operatorname*{argmax}_{\mathbf{w}} \frac{\min_{j} y_j \mathbf{w}^T \mathbf{x}_j}{||\mathbf{w}||_2}$$

$$= \operatorname*{argmax}_{\mathbf{w}} \frac{1}{||\mathbf{w}||_2}, \qquad \text{s.t.} \quad \left( \min_{j} \ y_j \mathbf{w}^T \mathbf{x}_j \right) = 1$$

$$= \operatorname*{argmin}_{\mathbf{w}} ||\mathbf{w}||_2^2, \qquad \text{s.t.} \quad \left( \min_{j} \ y_j \mathbf{w}^T \mathbf{x}_j \right) = 1$$

$$= \operatorname*{argmin}_{\mathbf{w}} ||\mathbf{w}||_2^2, \qquad \text{s.t.} \quad y_j \mathbf{w}^T \mathbf{x}_j \geq 1 \ \text{ for all } j$$

# Support Vector Machine (SVM)

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2^2, \quad \text{s.t.} \quad y_j \mathbf{w}^T \mathbf{x}_j \geq 1 \text{ for all } j \in \{1, \cdots, n\}.$$
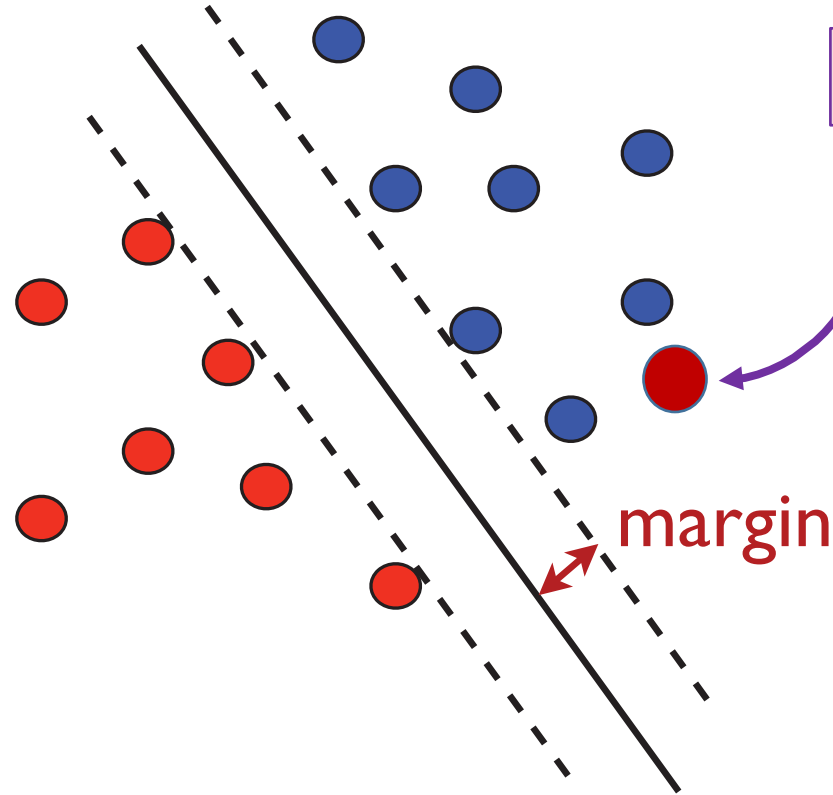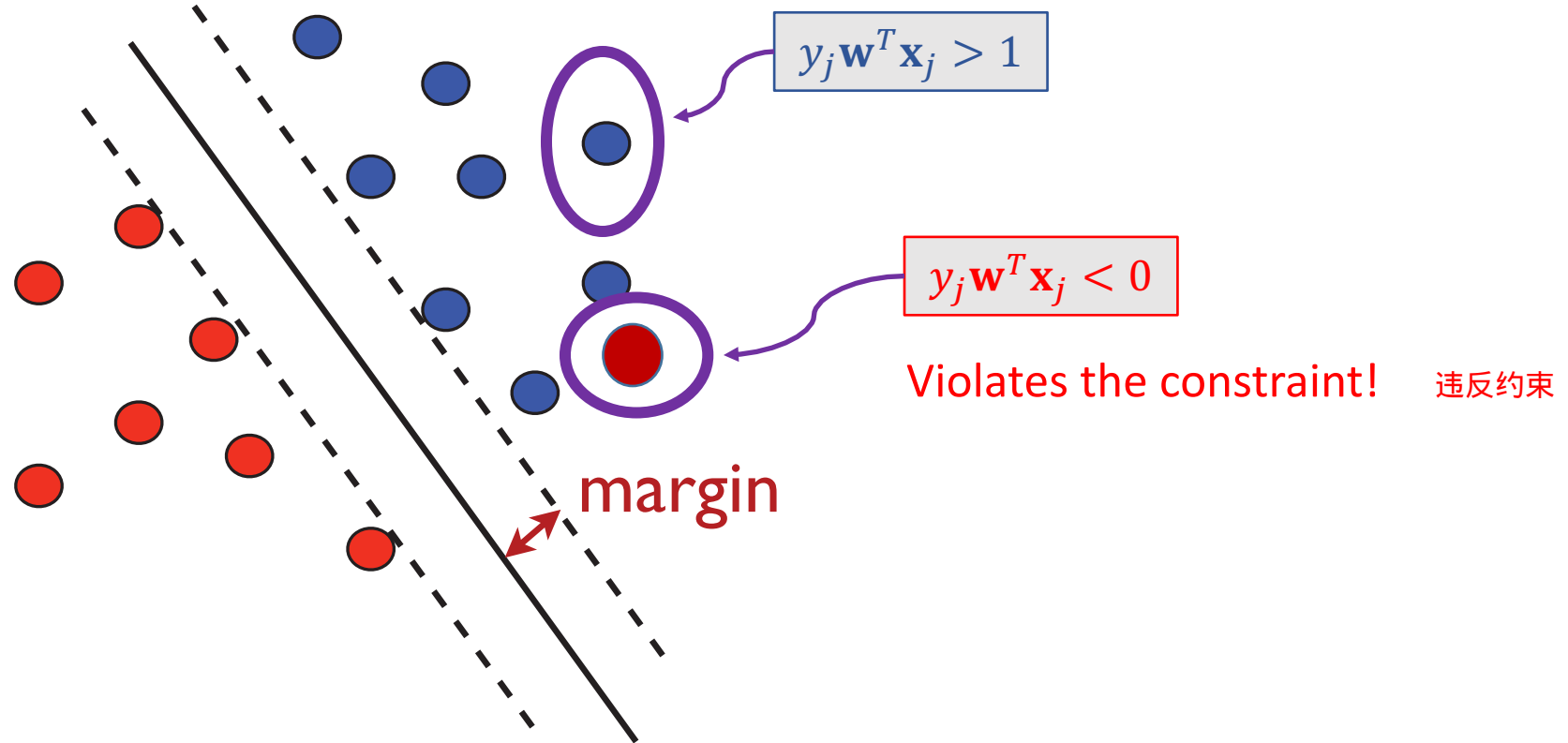
**Equivalent form of SVM** SVM

# Support Vector Machine (SVM)

$$\min_{\mathbf{w}} \left|\left|\mathbf{w}\right|\right|_2^2, \quad \text{s.t.} \quad y_j \mathbf{w}^T \mathbf{x}_j \geq 1 \text{ for all } j \in \{1, \cdots, n\}.$$



What if the data is inseparable?

margin

# Support Vector Machine (SVM)

$$\min_{\mathbf{w}} \left\|\mathbf{w}\right\|_2^2, \qquad \text{s.t.} \qquad y_j \mathbf{w}^T \mathbf{x}_j \geq 1 \ \text{ for all } j \in \{1, \cdots, n\}.$$



$y_j \mathbf{w}^T \mathbf{x}_j > 1$

$y_j \mathbf{w}^T \mathbf{x}_j < 0$

Violates the constraint!

margin

# Support Vector Machine (SVM)

$$\min_{\mathbf{w}} ||\mathbf{w}||_2^2, \quad \text{s.t.} \quad 1 - y_j \mathbf{w}^T \mathbf{x}_j \leq 0 \text{ for all } j \in \{1, \cdots, n\}.$$

**Relax**

$$\min_{\mathbf{w}, \xi_j} ||\mathbf{w}||_2^2 + \lambda \sum_j [\xi_j]_+, \quad \text{s.t.} \quad 1 - y_j \mathbf{w}^T \mathbf{x}_j = \xi_j \text{ for all } j \in \{1, \cdots, n\}.$$

- $[\xi_j]_+ = \max\{\xi_j, 0\}$

# Support Vector Machine (SVM)

$$\min_{\mathbf{w}} \left\|\mathbf{w}\right\|_2^2, \quad \text{s.t.} \quad 1 - y_j \mathbf{w}^T \mathbf{x}_j \leq 0 \text{ for all } j \in \{1, \cdots, n\}.$$

**Relax**

$$\min_{\mathbf{w}, \xi_j} \left\|\mathbf{w}\right\|_2^2 + \lambda \sum_j [\xi_j]_+, \quad \text{s.t.} \quad 1 - y_j \mathbf{w}^T \mathbf{x}_j = \xi_j \text{ for all } j \in \{1, \cdots, n\}.$$

- $[\xi_j]_+ = \max\{\xi_j, 0\}$
- $\xi_j \leq 0$ means the constraint $1 - y_j \mathbf{w}^T \mathbf{x}_j \leq 0$ is satisfied
  - ➔ no penalty!
- $\xi_j > 0$ means the constraint is violated (because the data is inseparable)
  - ➔ penalize the violation $\xi_j$.        min

# Support Vector Machine (SVM)

$$\min_{\mathbf{w}} \left\|\mathbf{w}\right\|_2^2, \quad \text{s.t.} \quad 1 - y_j \mathbf{w}^T \mathbf{x}_j \leq 0 \ \text{ for all } j \in \{1, \cdots, n\}.$$

**Relax**

$$\min_{\mathbf{w}, \xi_j} \left\|\mathbf{w}\right\|_2^2 + \lambda \sum_j [\xi_j]_+, \quad \text{s.t.} \quad 1 - y_j \mathbf{w}^T \mathbf{x}_j = \xi_j \ \text{ for all } j \in \{1, \cdots, n\}.$$

**Equivalent**

$$\min_{\mathbf{w}, b} \left\|\mathbf{w}\right\|_2^2 + \lambda \sum_j \left[1 - y_j \mathbf{w}^T \mathbf{x}_j\right]_+.$$

# Comparisons

SVM: $\min_{\mathbf{w}} \left\| \mathbf{w} \right\|_2^2 + \lambda \sum_j g(y_j \mathbf{w}^T \mathbf{x}_j).$

Hinge loss: $g(z) = [1 - z]_+.$



Hinge loss
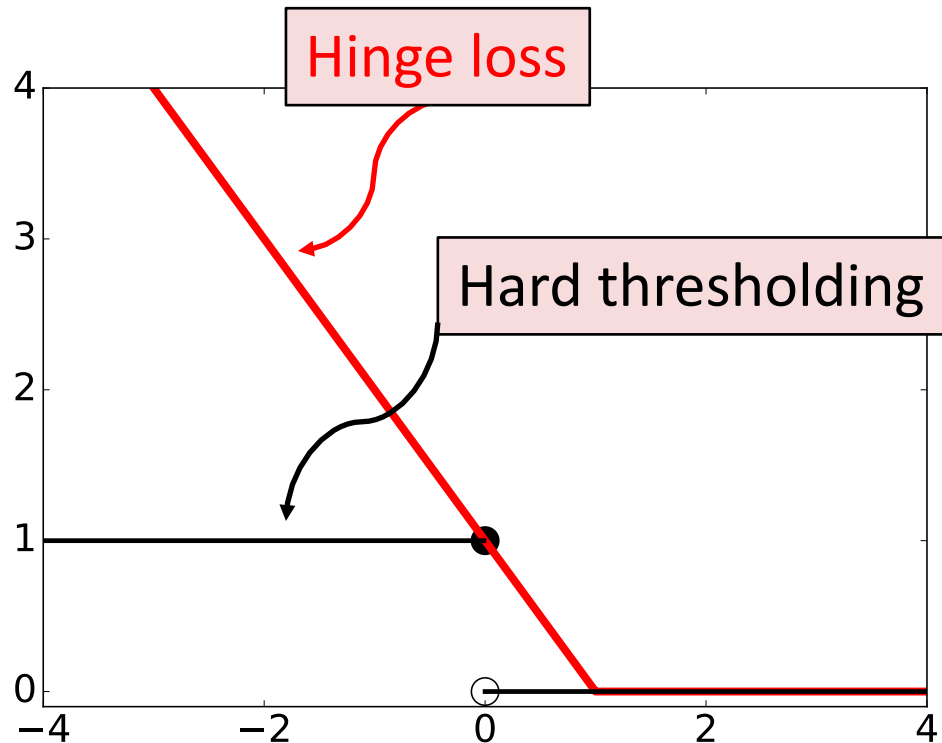
# Comparisons

SVM: $\min_{\mathbf{w}} \left\| \mathbf{w} \right\|_2^2 + \lambda \sum_j g(y_j \mathbf{w}^T \mathbf{x}_j).$
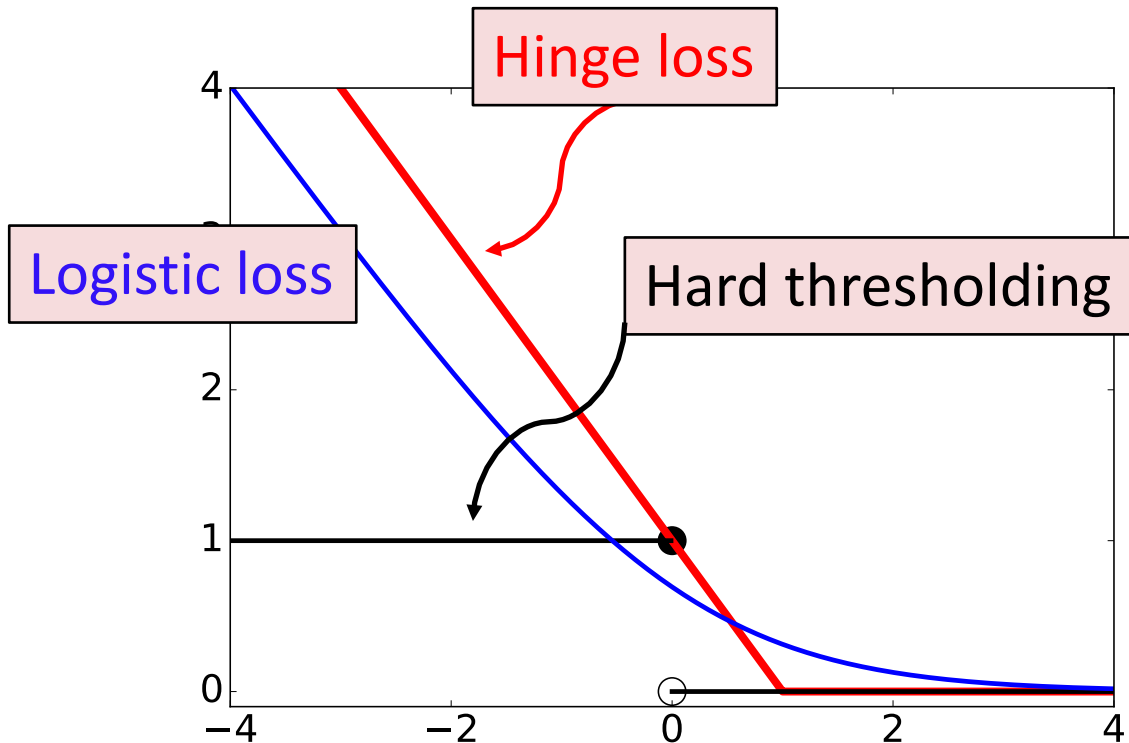
Hinge loss: $g(z) = [1 - z]_+.$

Hard thresholding: $h(\mathrm{z}) = \begin{cases} 1, & \text{if } z < 0; \\ 0, & \text{if } z \geq 0. \end{cases}$

# Comparisons

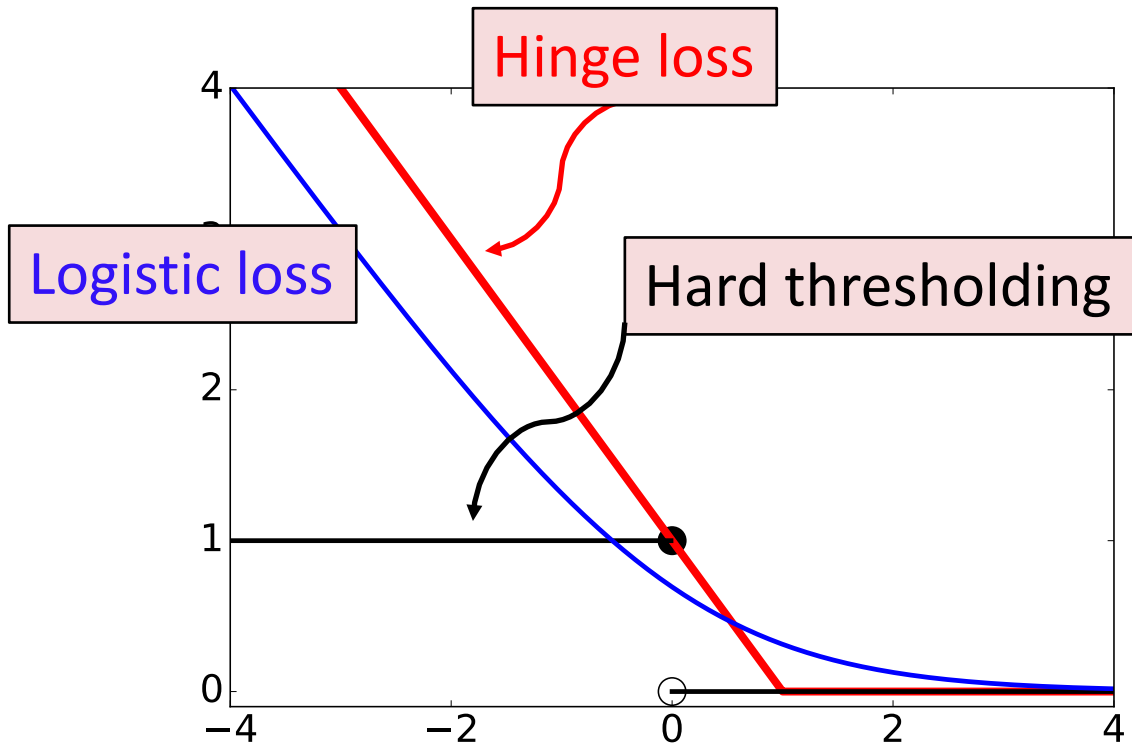SVM: $\min_{\mathbf{w}} \big||\mathbf{w}|\big|_2^2 + \lambda \sum_j g(y_j \mathbf{w}^T \mathbf{x}_j).$

Hinge loss: $g(z) = [1 - z]_+.$

Hard thresholding: $h(z) = \begin{cases} 1, & \text{if } z < 0; \\ 0, & \text{if } z \geq 0. \end{cases}$

Logistic loss: $l(z) = \log(1 + e^{-z}).$

# Comparisons



- Convexity
  - Hinge loss and logistic loss are convex.
  - Global optima can be efficiently found.

- Smoothness
  - Hinge loss is non-smooth.
  - Logistic loss is smooth.

- Logistic regression is easier to solve than SVM.
  - GD for logistic regression has linear convergence.
  - Algorithms for SVM have sub-linear convergence.