# LECTURE 10

**SPRING 2021**
**APPLIED MACHINE LEARNING**
**CIHANG XIE**

1

# TODAY

- Support Vector Machine

  **--** review

  **--** Lagrangian duality

  **--** kernel trick

$$\theta^T x + b = +1$$

$$\theta^T x + b = 0$$

$$\theta^T x + b = -1$$

$$\max_{\theta} \gamma$$

$$\text{s.t.} \quad y^{(i)}(\theta^T x^{(i)} + b) \geq 1, \forall i$$

margin $\gamma$

$$\max_{\boldsymbol{\theta}} \boldsymbol{\gamma}$$

$$\text{s.t. } \boldsymbol{y^{(i)}(\theta^T x^{(i)} + b) \geq 1, \forall i}$$

$$\theta^T x + b = +1$$

$$\boldsymbol{\theta^T x + b = 0}$$

$$\theta^T x + b = -1$$

$\theta$

$x^+$

$x^-$

relationship between $\gamma$ & $\boldsymbol{\theta}$ ?

margin $\gamma$

$$\max_{\boldsymbol{\theta}} \boldsymbol{\gamma}$$
$$\text{s.t.} \quad \boldsymbol{y}^{(i)}(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}) \geq \boldsymbol{1}, \forall i$$

$$\min_{\boldsymbol{\theta}} \frac{\boldsymbol{1}}{\boldsymbol{2}} \boldsymbol{\theta}^T \boldsymbol{\theta}$$
$$\text{s.t.} \quad \boldsymbol{y}^{(i)}(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}) \geqq \boldsymbol{1}, \forall i$$

$\boldsymbol{\theta}$

$\boldsymbol{x}^+$

$\boldsymbol{x}^-$

relationship between $\boldsymbol{\gamma}$ & $\boldsymbol{\theta}$ ?

margin $\gamma$

# SUPPORT VECTOR MACHINE (SVM)

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

$$\text{s.t.} \quad y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b) \geq 1, \forall i$$

$$\theta^T x + b = +1$$

$$\theta^T x + b = -1$$

outlier

**NOT linearly separable**

$\theta^T x + b = +1$

outlier

$\theta^T x + b = -1$

Add some flexibilities?

# SOFT MARGIN SVM

$$\min_{\boldsymbol{\theta},\boldsymbol{\xi},b} \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + C \sum_{i=1}^{N} \xi_i$$

$$\text{s.t.} \quad y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b) + \xi_i \geq 1,$$

$$\xi_i \geq 0, \forall i$$

**$\xi_i$ is the "slack" variable**

- for $0 < \xi_i \leq 1$ point is between margin and correct side of hyperplane. This is a margin violation

- for $\xi_i > 1$ point is misclassified

# SOFT MARGIN SVM

$$\min_{\boldsymbol{\theta}, \boldsymbol{\xi}, \boldsymbol{b}} \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + C \sum_{i=1}^{N} \xi_i$$

$$\text{s.t.} \quad y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b) + \xi_i \geq 1,$$

$$\xi_i \geq 0, \forall i$$

**C is a regularization parameter:**

• small C allows constraints to be easily ignored → large margin

• large C makes constraints hard to ignore → narrow margin

• C = ∞ enforces all constraints: hard margin

# GRADIENT DESCENT FOR SVM

$$y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b) + \xi_i \geq 1 \,\&\, \xi_i \geq 0$$

$$\xi_i = \max\{0, 1 - y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b)\}$$

$$\min_{\theta,b} \frac{1}{2}\boldsymbol{\theta}^T\boldsymbol{\theta} + C\sum_{i=1}^{N}\max\{0, 1 - y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b)\}$$

# GRADIENT DESCENT FOR SVM

$$\text{COST}(\theta, b) = \frac{1}{2}\theta^T\theta + C\sum_{i=1}^{N}\max\{0, 1 - y^{(i)}(\theta^T x^{(i)} + b)\}$$

$$= \sum_{i=1}^{N}(\frac{1}{2N}\theta^T\theta + C\max\{0, 1 - y^{(i)}(\theta^T x^{(i)} + b)\})$$

**For each data point $x^{(i)}$**

$$\frac{\partial Cost(\theta, b)}{\partial \theta_j} = \begin{cases} \frac{1}{N}\theta_j - C\, y^{(i)}x_j^{(i)} & , \text{if } 1 - y^{(i)}(\theta^T x^{(i)} + b) > 0 \\ \frac{1}{N}\theta_j, & \text{otherwise} \end{cases}$$
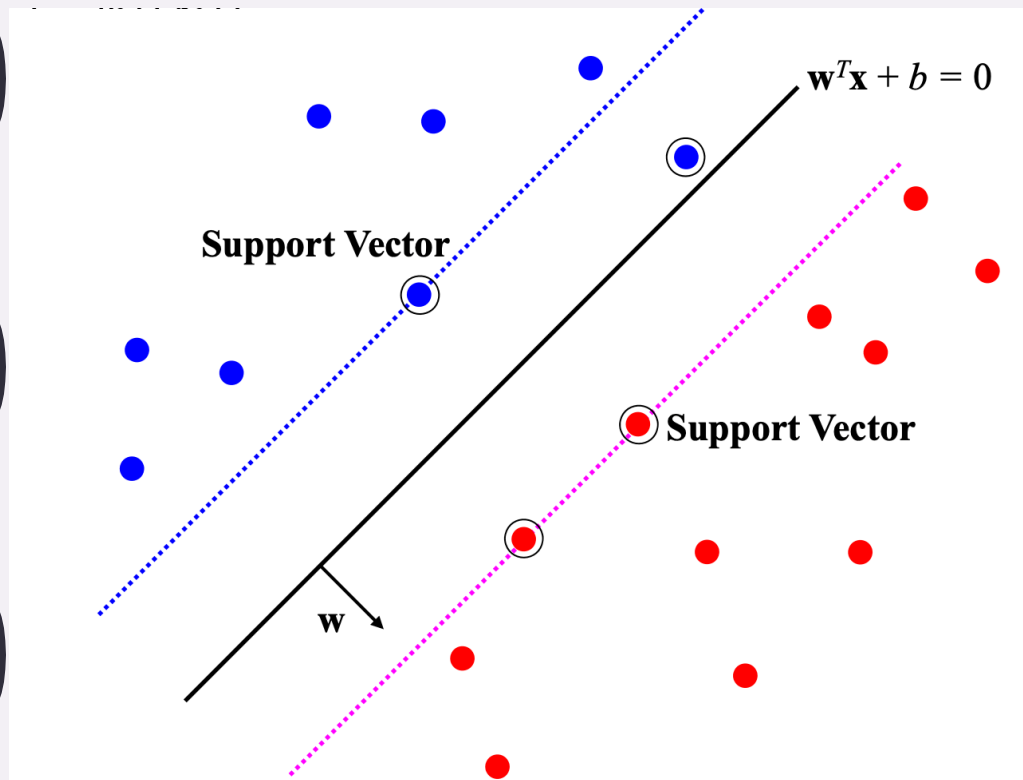
$$\frac{\partial Cost(\theta, b)}{\partial b} = \begin{cases} -C\, y^{(i)}, & \text{if } 1 - y^{(i)}(\theta^T x^{(i)} + b) > 0 \\ 0, & \text{otherwise} \end{cases}$$

# OPTIMIZATION

$$\min_{\theta, b} \left\{ \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} \right\} + \left\{ C \sum_{i=1}^{N} \max\left\{ 0, 1 - \boldsymbol{y}^{(i)}(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}) \right\} \right\}$$

Regularization          Model fit to data



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

**Support Vector**

**Support Vector**

**w**

1. $\boldsymbol{y}^{(i)}(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}) > 1$ => Point is outside margin. No contribution to loss

2. $\boldsymbol{y}^{(i)}(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}) = 1$ => Point is on margin. No contribution to loss.

3. $\boldsymbol{y}^{(i)}(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}) < 1$ => Point violates margin constraint. Contributes to loss

# RECALL: LOGISTIC REGRESSION

- Maximum likelihood estimation:

$$\max_{\theta} \quad ll(w) = \max_{\theta} \quad \sum_{i} \log P(y^{(i)}|x^{(i)};\theta)$$

with:

$$P(y^{(i)} = +1|x^{(i)};\theta) = \frac{1}{1+e^{-\theta^T x^{(i)}}}$$

$$P(y^{(i)} = -1|x^{(i)};\theta) = 1 - \frac{1}{1+e^{-\theta^T x^{(i)}}}$$

$$\boldsymbol{P(y^{(i)}|x^{(i)};\theta) = \frac{1}{1 + e^{-y^{(i)}(\theta^T x^{(i)})}}}$$

# RECALL: LOGISTIC REGRESSION

$$\max_{\boldsymbol{\theta}} \sum_i \log P\big(\boldsymbol{y}^{(i)} \big| \boldsymbol{x}^{(i)}; \boldsymbol{\theta}\big) = \max_{\boldsymbol{\theta}} \sum_i \log \frac{\mathbf{1}}{\mathbf{1} + \boldsymbol{e}^{-\boldsymbol{y}^{(i)}\left(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}\right)}}$$

$$= \max_{\boldsymbol{\theta}} \sum_{\boldsymbol{i}} \left(\log 1 - \log\left(1 + \boldsymbol{e}^{-\boldsymbol{y}^{(i)}\left(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}\right)}\right)\right)$$

$$= \max_{\boldsymbol{\theta}} \sum_{\boldsymbol{i}} -\log\left(1 + \boldsymbol{e}^{-\boldsymbol{y}^{(i)}\left(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}\right)}\right)$$

$$= \min_{\boldsymbol{\theta}} \sum_{\boldsymbol{i}} \log\left(1 + \boldsymbol{e}^{-\boldsymbol{y}^{(i)}\left(\boldsymbol{\theta}^T \boldsymbol{x}^{(i)} + \boldsymbol{b}\right)}\right)$$

# RELATIONSHIP TO LOGISTIC REGRESSION

$$\min_{\boldsymbol{\theta},\boldsymbol{b}} \boldsymbol{\lambda}\boldsymbol{\theta}^{T}\boldsymbol{\theta} - \sum_{i} \log \boldsymbol{P}\left(\boldsymbol{y}^{(i)}\big|\boldsymbol{x}^{(i)};\boldsymbol{\theta},\boldsymbol{b}\right)$$

$$\min_{\boldsymbol{\theta},\boldsymbol{b}} \left\{\boldsymbol{\lambda}\boldsymbol{\theta}^{T}\boldsymbol{\theta}\right\} + \sum_{i} \left\{\log(1 + \boldsymbol{e}^{-\boldsymbol{y}^{(i)}\left(\boldsymbol{\theta}^{T}\boldsymbol{x}^{(i)}+\boldsymbol{b}\right)})\right\}$$

Regularization          Logistics Loss

$$\min_{\boldsymbol{\theta},\boldsymbol{b}} \left\{\frac{1}{2}\boldsymbol{\theta}^{T}\boldsymbol{\theta}\right\} + C\sum_{i=1}^{N} \left\{\max\left\{0, 1 - \boldsymbol{y}^{(i)}\left(\boldsymbol{\theta}^{T}\boldsymbol{x}^{(i)} + \boldsymbol{b}\right)\right\}\right\}$$

Regularization   Applied Machine Learning Hinge Loss

16

# RELATIONSHIP TO LOGISTIC REGRESSION



Hinge Loss
$$\max\{0, 1 - y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b)\}$$

Logistic Loss
$$log(1 + e^{-y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b)})$$

$$y^{(i)}(\boldsymbol{\theta}^T x^{(i)} + b)$$

Logistic loss is sometime viewed as the **smooth version** of the Hinge loss.

# HW2 (DUE MAY 9)

# QUESTIONS?