

# CS480/680: Intro to ML

## Lecture 08: Soft-margin SVM



# Outline

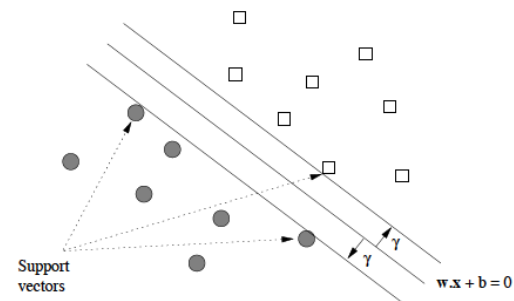
- Formulation
- Dual
- Optimization
- Extension

# Hard-margin SVM

Primal

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } \forall i, y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$$



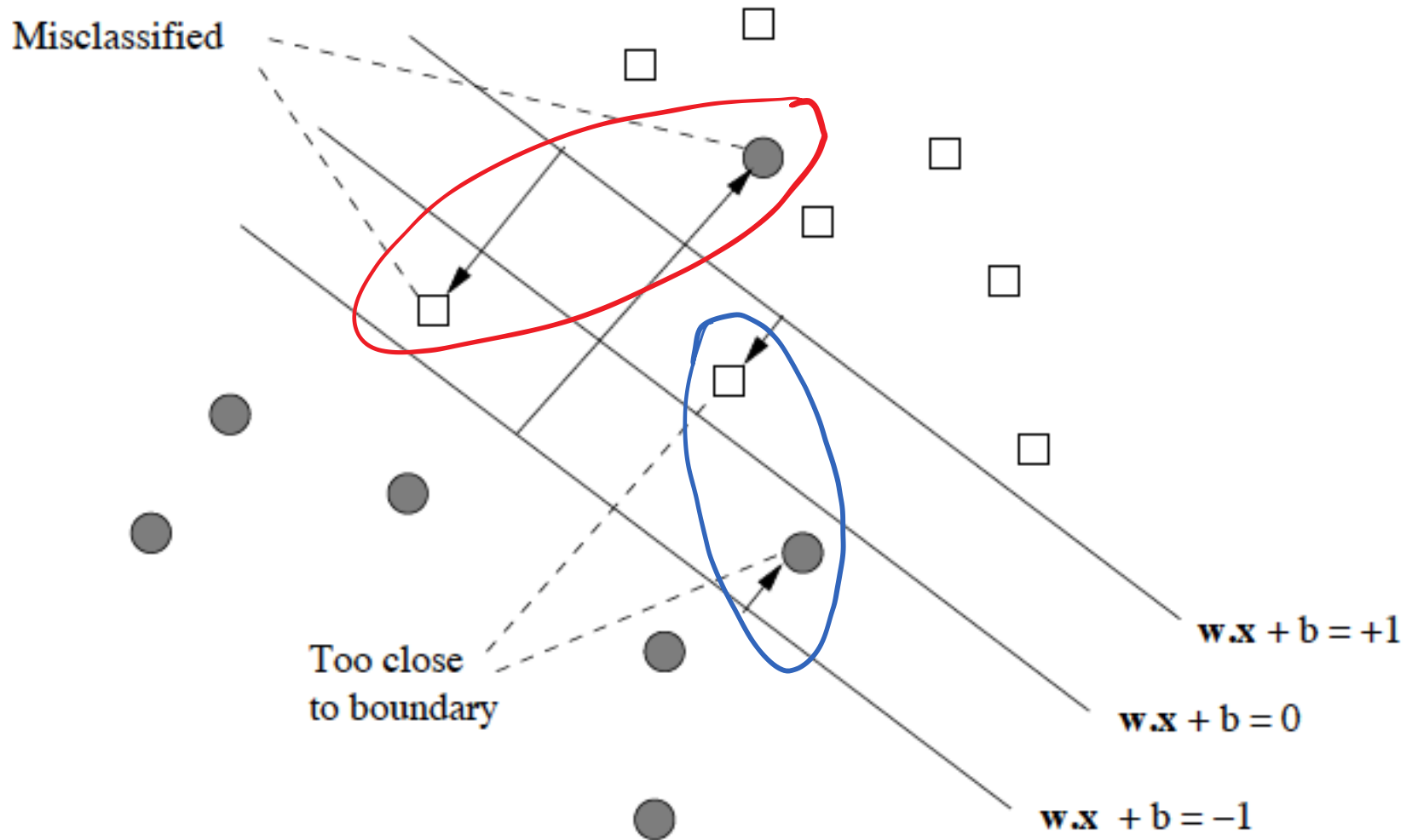
Dual

$$\min_{\alpha \geq 0} \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_k \alpha_k$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

hard constraint

# What if **in**separable?



# Soft-margin (Cortes & Vapnik'95)

Primal

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } \forall i, y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$$

$y_i \hat{y}_i \geq 1$   
 $1 - y_i \hat{y}_i \leq 0$   
 hard constraint

propto 1/margin

hyper-parameter

training error

Primal

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (1 - y_i \hat{y}_i)_+$$

$$\forall i, \hat{y}_i = \mathbf{w}^\top \mathbf{x}_i + b$$

$\begin{cases} 0, & 1 - y_i \hat{y}_i \leq 0 \\ 1 - y_i \hat{y}_i, & \text{o.w.} \end{cases}$

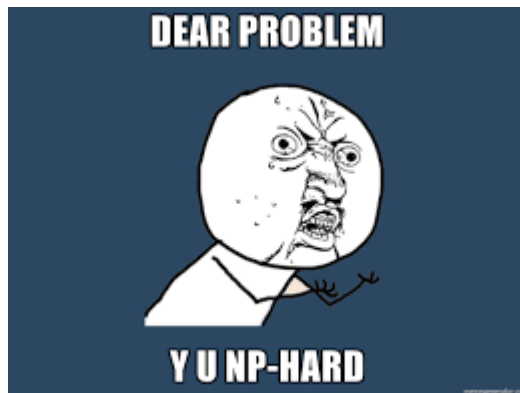
soft constraint

prediction  
(no sign)

# Zero-one loss

$$\Pr(\underbrace{\text{sign}(f(X))}_{\text{your prediction}} \neq Y) = \mathbf{E}[1 - Y f(X) \geq 0]$$

your prediction



$$\frac{1}{n} \sum_{i=1}^n 1 - Y_i \hat{Y}_i \geq 0$$

- Find prediction rule  $f$  so that on an unseen random  $X$ , our prediction  $\text{sign}(f(X))$  has small chance to be different from the true label  $Y$

# The hinge loss



$$(1 - y\hat{y})_+ = \max\{1 - y\hat{y}, 0\}$$

zero-one  $1 - y\hat{y} \geq 0$

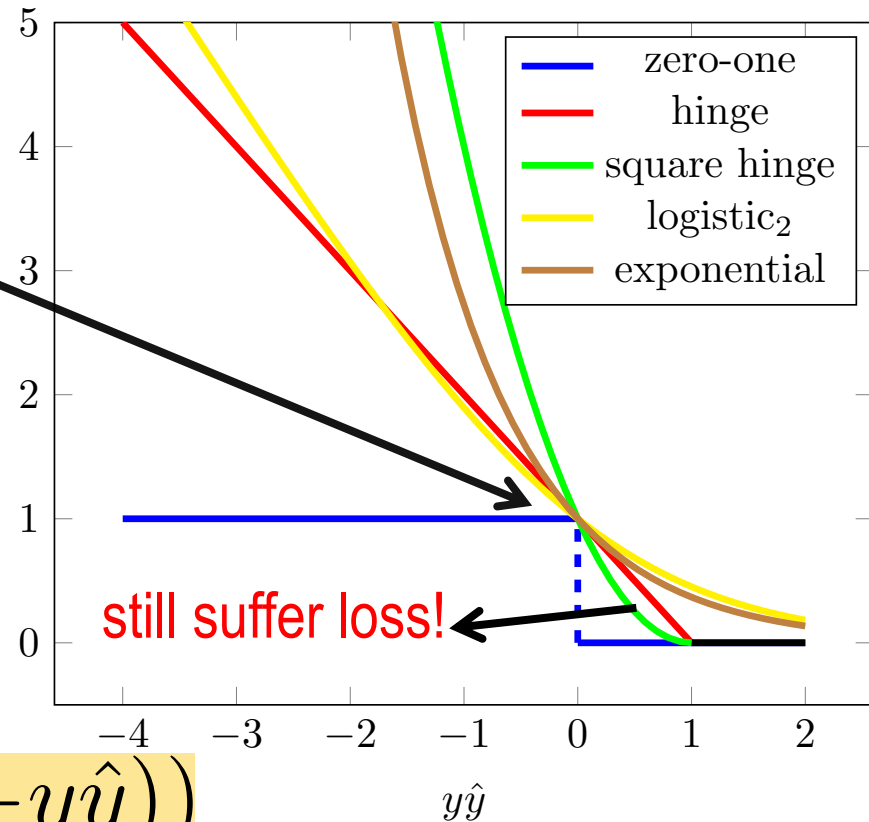
squared hinge  $(1 - y\hat{y})_+^2$

exponential loss  $\exp(-y\hat{y})$

logistic loss  $\log_2(1 + \exp(-y\hat{y}))$

upper  
bound

loss



# Classification-calibration

- Want to minimize zero-one loss
- End up with minimizing some **other** loss

**Theorem (Bartlett, Jordan, McAuliffe'06).** Any convex margin loss  $\ell$  is classification-calibrated **iff**  $\ell$  is differentiable at 0 and  $\ell'(0) < 0$ .

**Classification calibration.**  $\arg \min_a \mathbf{E}[\ell(Y a) | X = x]$  has the same sign as  $2\eta(x) - 1$ , i.e., the Bayes rule.

$$\eta(x)\ell(a) + (1 - \eta(x))\ell(-a)$$



# Outline

- Formulation
- Dual
- Optimization
- Extension

# Important optimization trick

$$\min_x f(x)$$



joint over  
 $x$  and  $t$

$$\min_{x, t} t$$

$$\text{s.t. } f(x) \leq t$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (1 - y_i \hat{y}_i)_+ \\ \forall i, \hat{y}_i = \mathbf{w}^\top \mathbf{x}_i + b$$



↓

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

Slack for "wrong" prediction

$$\text{s.t. } \forall i, (1 - y_i \hat{y}_i)_+ \leq \xi_i \longrightarrow \begin{cases} 1 - y_i \hat{y}_i \leq \xi_i \\ 0 \leq \xi_i \end{cases}$$

# Lagrangian

$$\min_{\mathbf{w}, b, \boldsymbol{\xi}} \max_{\alpha \geq 0, \beta \leq 0} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i C \xi_i + \alpha_i (1 - y_i \hat{y}_i - \xi_i) + \beta_i \xi_i$$
$$\max_{\alpha \geq 0, \beta \leq 0} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i C \xi_i + \alpha_i (1 - y_i \hat{y}_i - \xi_i) + \beta_i \xi_i$$

$$\frac{\partial}{\partial b} = \sum_i \alpha_i y_i = 0$$

$$\frac{\partial}{\partial \xi_i} = C - \alpha_i + \beta_i = 0$$

$$\frac{\partial}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0$$

# Dual problem

$$\max_{\substack{C \geq \\ \alpha \geq 0}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^{\top} \mathbf{x}_j$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

only dot product is needed!

$$\min_{\substack{C \geq \\ \alpha \geq 0}} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^{\top} \mathbf{x}_j - \sum_{i=1}^n \alpha_i$$

# The effect of C

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (1 - y_i \hat{y}_i)_+$$

$$\mathbf{R}^d \times \mathbf{R} \quad \forall i, \hat{y}_i = \mathbf{w}^\top \mathbf{x}_i + b$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } \forall i, y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$$

- $C \rightarrow 0$ ?
- $C \rightarrow \text{inf}$ ?

$$\min_{\alpha \geq 0} \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_i \alpha_i$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

$$\mathbf{R}^n \quad \min_{C \geq \alpha \geq 0} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^n \alpha_i$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

# Karush-Kuhn-Tucker conditions

- Primal constraints on  $\mathbf{w}$ ,  $b$  and  $\xi$ :  $(1 - y_i \hat{y}_i)_+ \leq \xi_i$
- Dual constraints on  $\alpha$  and  $\beta$ :  $\alpha \geq 0 \quad \beta \leq 0$
- **Complementary slackness**

$$\alpha_i (1 - y_i \hat{y}_i - \xi_i) = 0$$
$$\beta_i \xi_i = 0$$

- **Stationarity**

$$\frac{\partial}{\partial b} = \sum_i \alpha_i y_i = 0$$
$$\frac{\partial}{\partial \xi_i} = C - \alpha_i + \beta_i = 0$$
$$\frac{\partial}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0$$

$$\min_{\mathbf{w}, b, \xi} \max_{\alpha \geq 0, \beta \leq 0} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i C \xi_i + \alpha_i (1 - y_i \hat{y}_i - \xi_i) + \beta_i \xi_i$$

# Parsing the equations

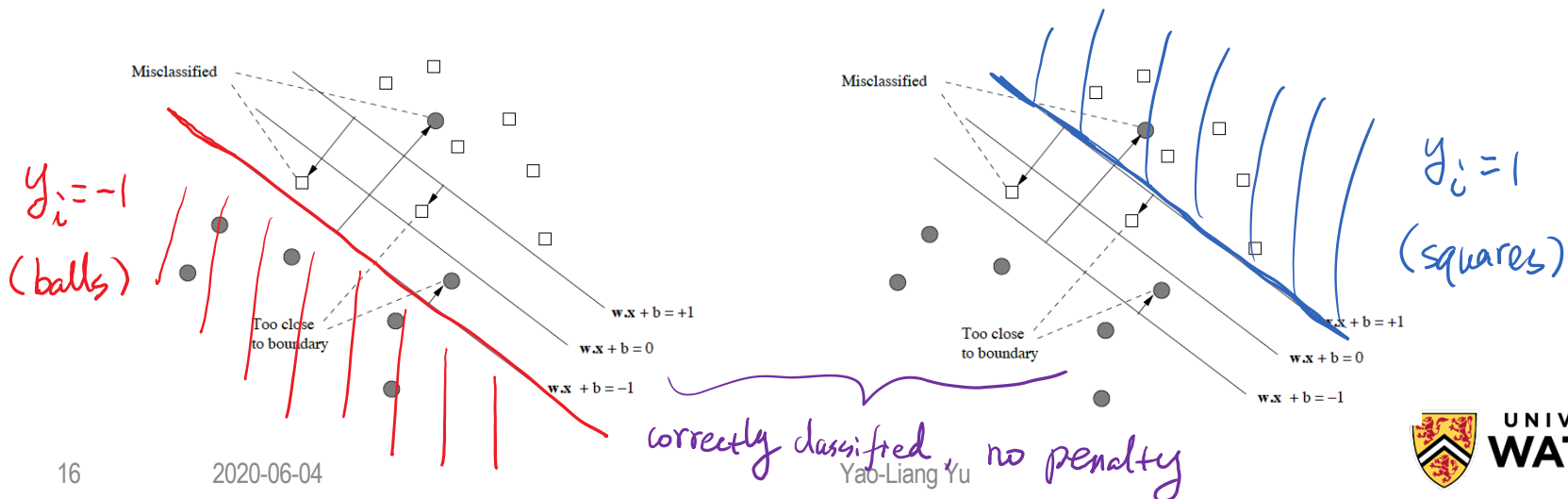
$$\alpha_i(1 - y_i\hat{y}_i - \xi_i) = 0$$

$$(C - \alpha_i)\xi_i = 0$$

$$(1 - y_i\hat{y}_i)_+ \leq \xi_i$$

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

- $\alpha_i < C \implies \xi_i = 0 \implies y_i\hat{y}_i \geq 1$





# Parsing the equations

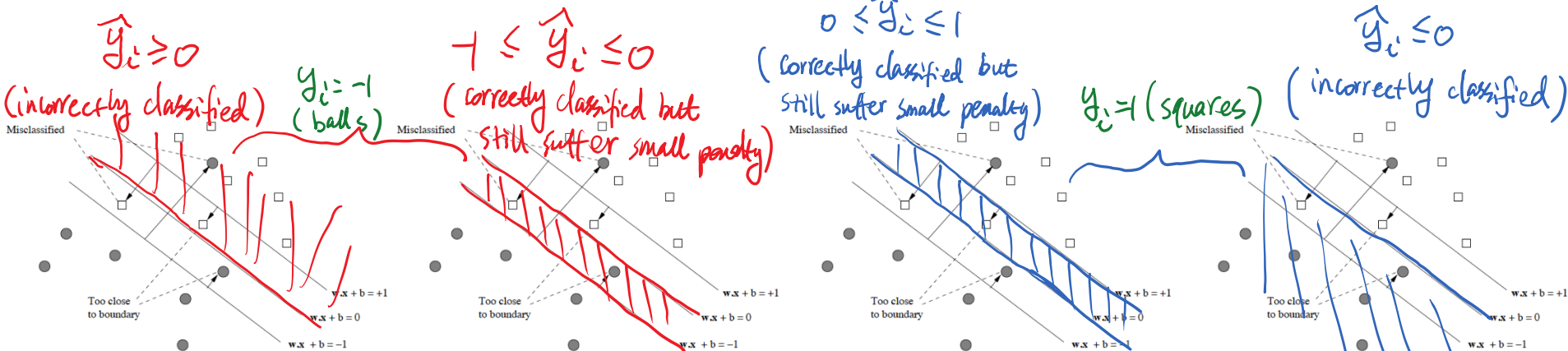
$$\alpha_i(1 - y_i\hat{y}_i - \xi_i) = 0$$

$$(C - \alpha_i)\xi_i = 0$$

$$(1 - y_i\hat{y}_i)_+ \leq \xi_i$$

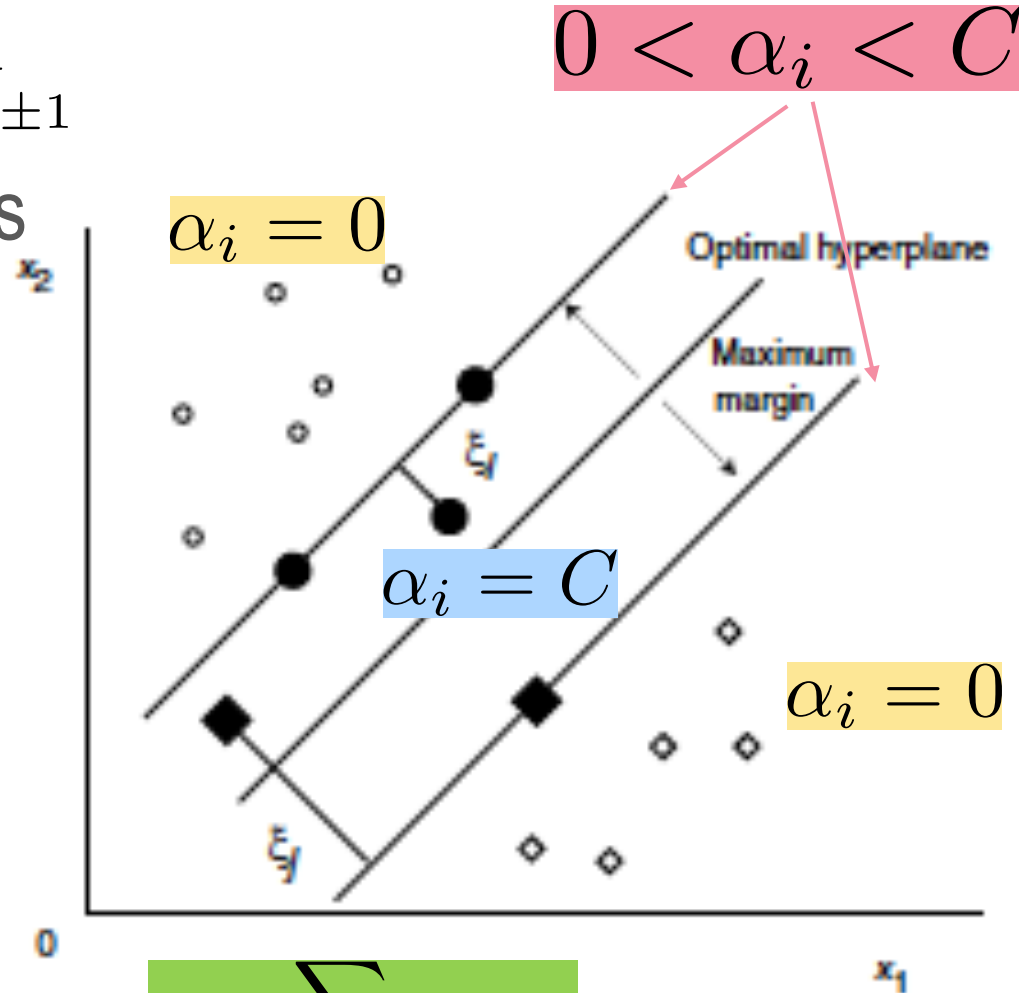
$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

•  $\alpha_i > 0 \implies 1 - y_i\hat{y}_i - \xi_i = 0 \implies y_i\hat{y}_i \leq 1$



# Support Vectors

- $\alpha_i = 0$  : on **or** beyond  $H_{\pm 1}$
- $\alpha_i = C$  : “incorrect” points
- $0 < \alpha_i < C$  : on  $H_{\pm 1}$
- $x_i$  is on  $H_{\pm 1}$  : what do we know about its  $\alpha_i$ ?
- $x_i$  is **strictly** beyond  $H_{\pm 1}$  : what is  $\alpha_i$ ?
- $x_i$  is **strictly** within  $H_{\pm 1}$  : what is  $\alpha_i$ ?



$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

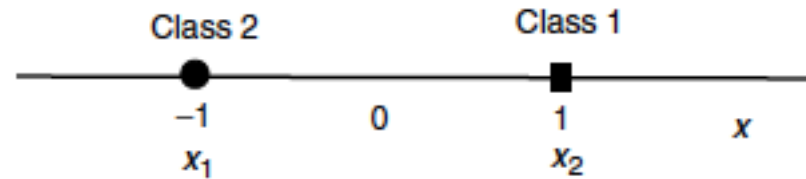
# Recover $b$

- Take any  $i$  such that  $C > \alpha_i > 0$
- Then  $\mathbf{x}_i$  is on the hyperplane:

$$1 - y_i \hat{y}_i = 0$$

- How to recover  $\xi$  ?

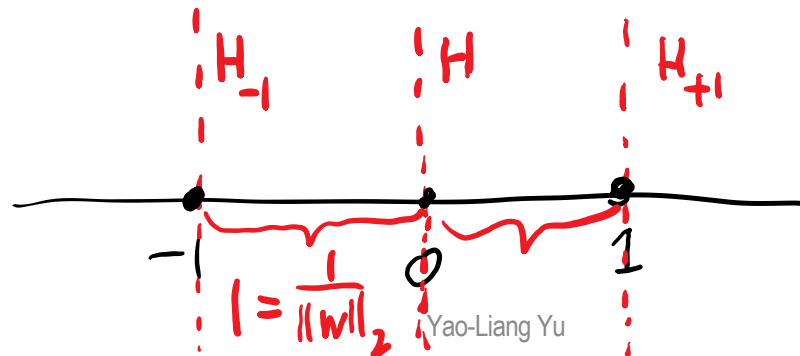
# A simple example



$$\min_{C \geq \alpha_1 \geq 0} 2\alpha_1^2 - 2\alpha_1 \quad \longrightarrow \quad \alpha_1 = \alpha_2 = \begin{cases} 1/2, & C \geq 1/2 \\ C, & C \leq 1/2 \end{cases}$$

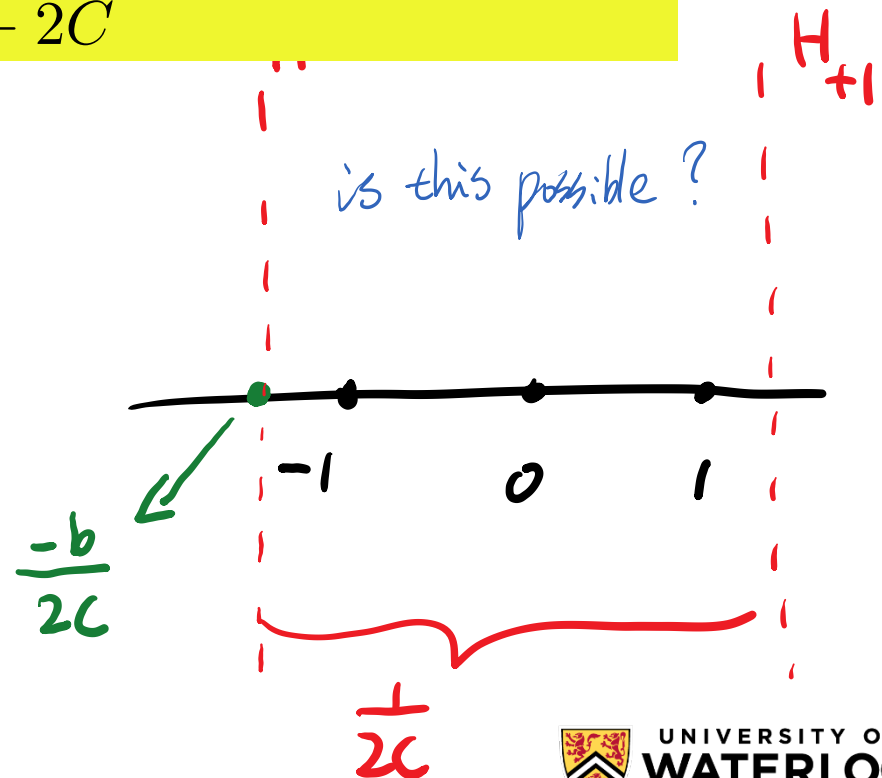
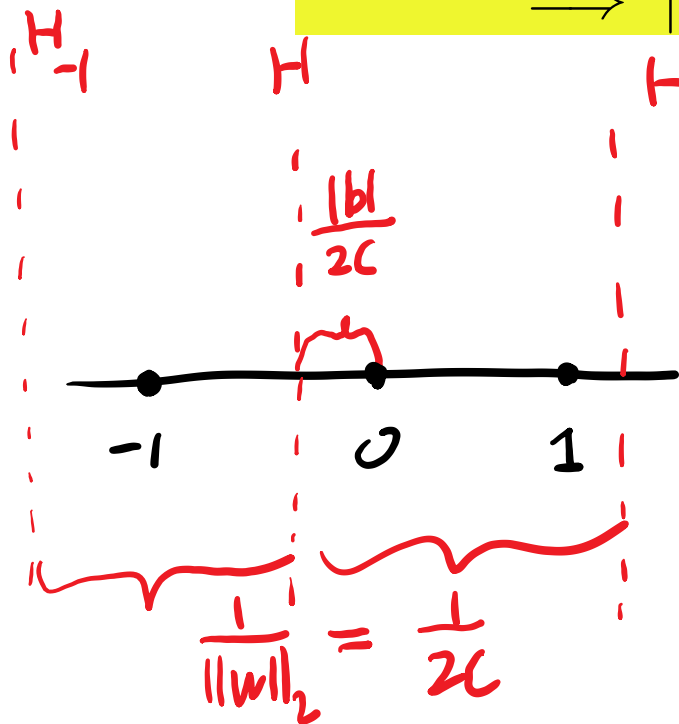
- $C > 1/2 \rightarrow \mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i = 2\alpha_1 = 1$

$0 < \alpha_1 = \frac{1}{2} < C \rightarrow x_2$  is on  $H_1 \rightarrow b=0 \rightarrow$  hard-margin



# A simple example cont'

- $C \leq \frac{1}{2} \rightarrow w=2C$ ,  $H_t = \{x : w^\top x + b = t\}$
- $\alpha_1 = C \implies (x_1, y_1 = -1)$  is on or above  $H_{-1}$
- $\implies y_i(w^\top x_i + b) \leq 1$
- $\implies |b| \leq 1 - 2C$



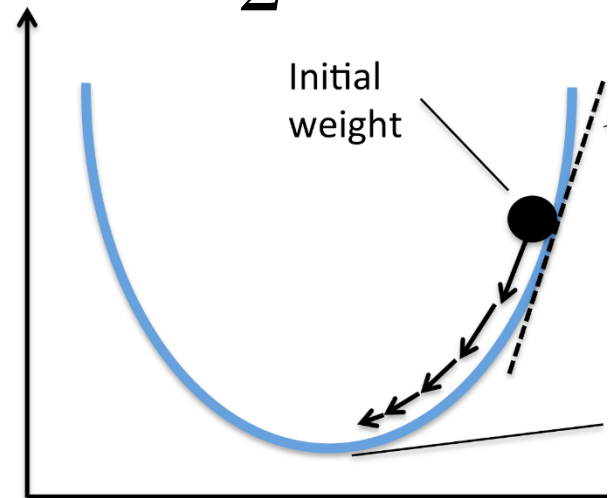
# Outline

- Formulation
- Dual
- Optimization
- Extension

# Gradient Descent

$$\min_{\mathbf{w}, b} L(\mathbf{w}) := \frac{C}{n} \sum_{i=1}^n \ell(y_i \hat{y}_i) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla L(\mathbf{w}_t)$$



$$\frac{C}{n} \sum_i \ell'(y_i \hat{y}_i) y_i \mathbf{x}_i + \mathbf{w}_t$$

(Generalized) gradient

Step size (learning rate)

- const., if  $L$  is smooth
- diminishing, otherwise

$O(nd)$  !

# Stochastic Gradient Descent (SGD)

$$\mathbf{w}_{t+1} = (1 - \eta_t)\mathbf{w}_t - \eta_t C \frac{1}{n} \sum_i \ell'(y_i \hat{y}_i) y_i \mathbf{x}_i$$

↪ average over  $n$  samples

a **random** sample suffices

$$\mathbf{w}_{t+1} = (1 - \eta_t)\mathbf{w}_t - \eta_t C \ell'(y_{i_t} \hat{y}_{i_t}) y_{i_t} \mathbf{x}_{i_t}$$

$O(d)$

- diminishing step size, e.g.,  $1/\sqrt{t}$  or  $1/t$
- averaging, momentum, variance-reduction, etc.
- sample w/o replacement; cycle; permute in each pass



# The derivative

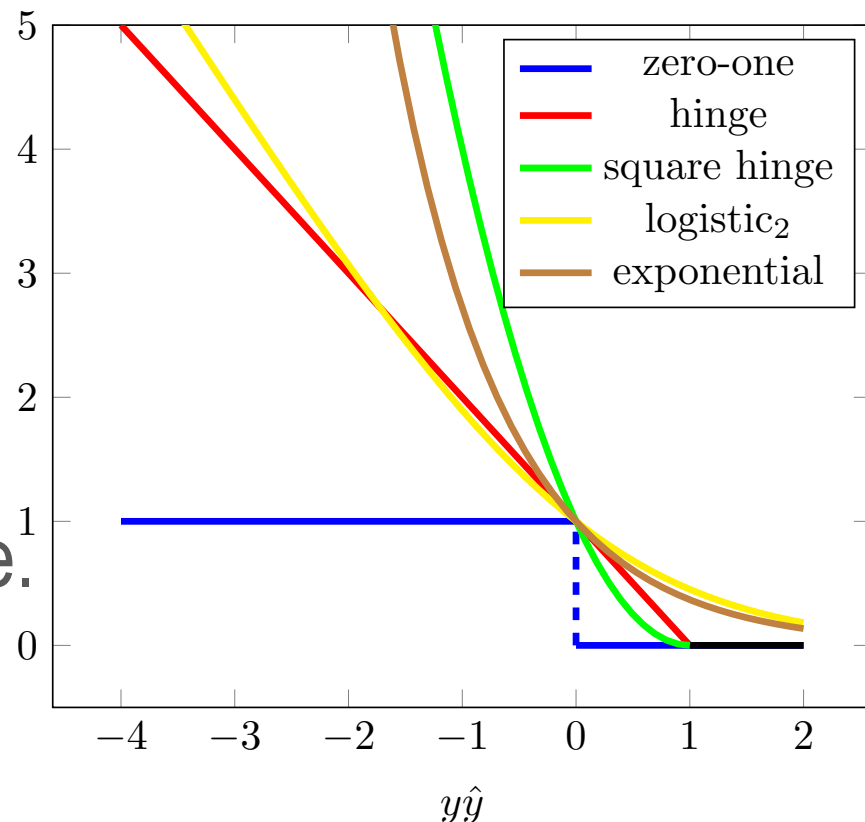
$$\ell'_{\text{hinge}}(t) = \begin{cases} -1, & t \leq 1 \\ 0, & t \geq 1 \end{cases}$$

What about zero-one loss? <sup>loss</sup>

All other losses are differentiable.

What about perceptron?

$$\mathbf{w}_{t+1} = (1 - \cancel{\eta_t}) \mathbf{w}_t - \cancel{\eta_t} C \ell'(y_{i_t} \hat{y}_{i_t}) y_{i_t} \mathbf{x}_{i_t}$$



# Solving the dual

$$\begin{array}{l}
 \min_{\alpha \geq 0} \quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^n \alpha_i \\
 \text{s.t.} \quad \sum_i \alpha_i y_i = 0
 \end{array}$$

$$\begin{array}{l}
 \alpha_{t+1} \leftarrow \alpha_t - \eta_t (K \odot \mathbf{y}\mathbf{y}^\top) \alpha_t + \eta_t \mathbf{1} \\
 \alpha_{t+1} \leftarrow \text{prox}(\alpha_{t+1}) \quad O(n^2)
 \end{array}$$

- Can choose constant step size  $\eta_t = \eta$
- Faster algorithms exist: e.g., choose a pair of  $\alpha_p$  and  $\alpha_q$  and derive a closed-form update

# Outline

- Formulation
- Dual
- Optimization
- Extension

# Multiclass (Crammer & Singer'01)

$$\min_W \frac{1}{2} \|W\|_F^2$$

$$\text{s.t. } \forall i, \forall k \neq y_i,$$

$$\mathbf{x}_i^\top \mathbf{w}_{y_i} \geq 1 +$$

Prediction for  
correct class

$$\geq 1 +$$

$$\mathbf{x}_i^\top \mathbf{w}_k$$

Prediction for  
wrong classes

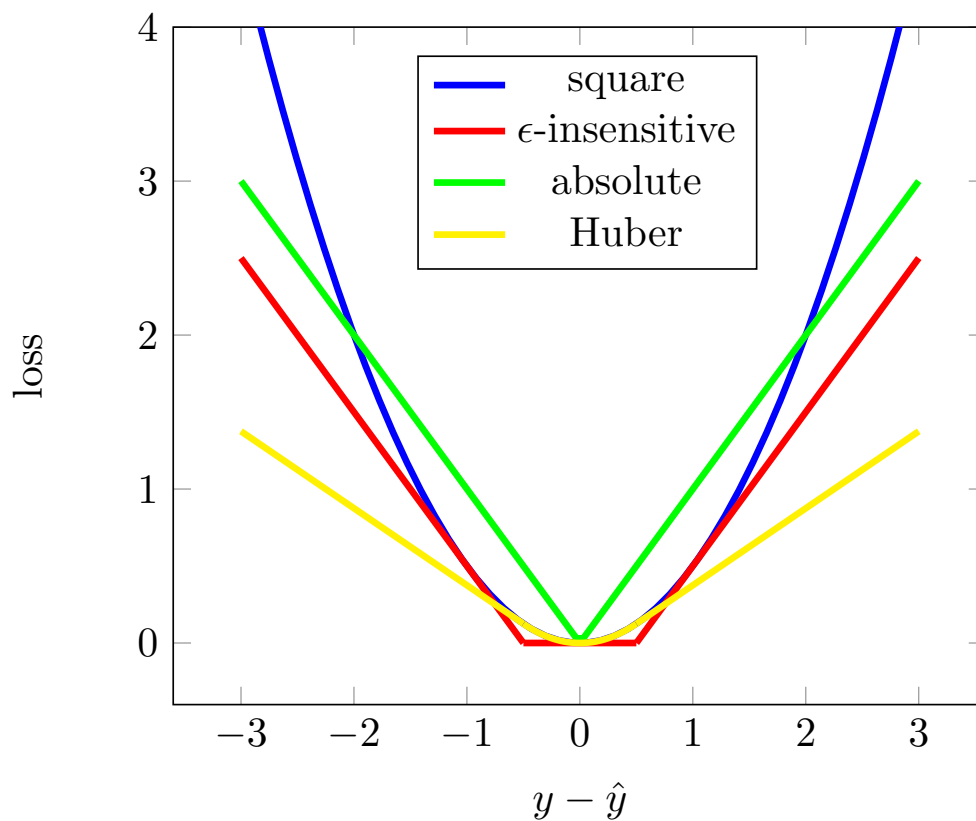
separate by a “safety margin”



- Soft-margin is similar
- Many other variants

# Regression (Drucker et al.'97)

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (|y - \hat{y}| - \epsilon)_+$$



# Questions?

