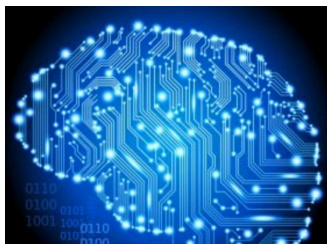




Lecture 15: Learning 4

http://cs.nju.edu.cn/yuy/course_ai17.ashx



Previously...



Learning

Decision tree learning

Neural networks

Why we can learn

Linear model



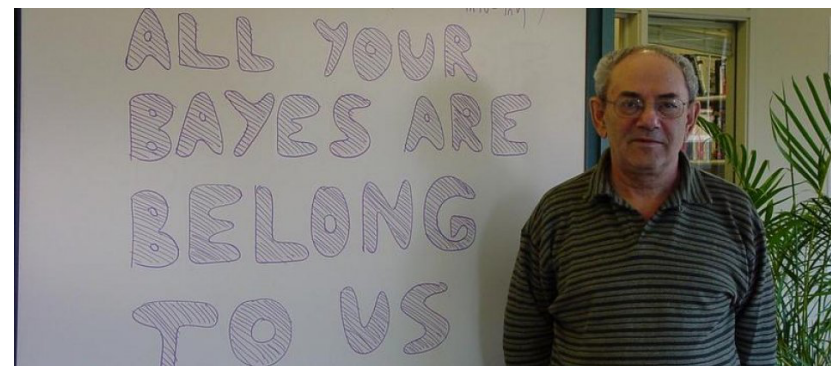
$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

$$\mathbf{w} = w_1, w_2, \dots, w_n \quad b$$



$$w_1 \cdot x_1 + w_2 \cdot x_2 + \dots + w_n \cdot x_n + b$$

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

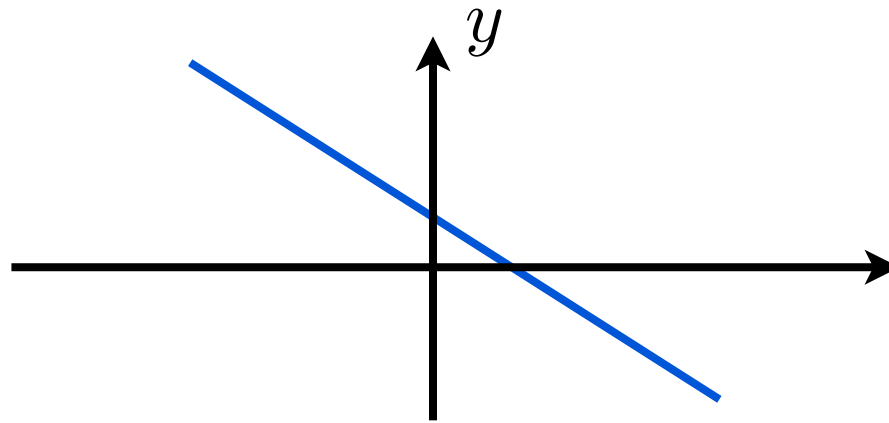


Vladimir Vapnik

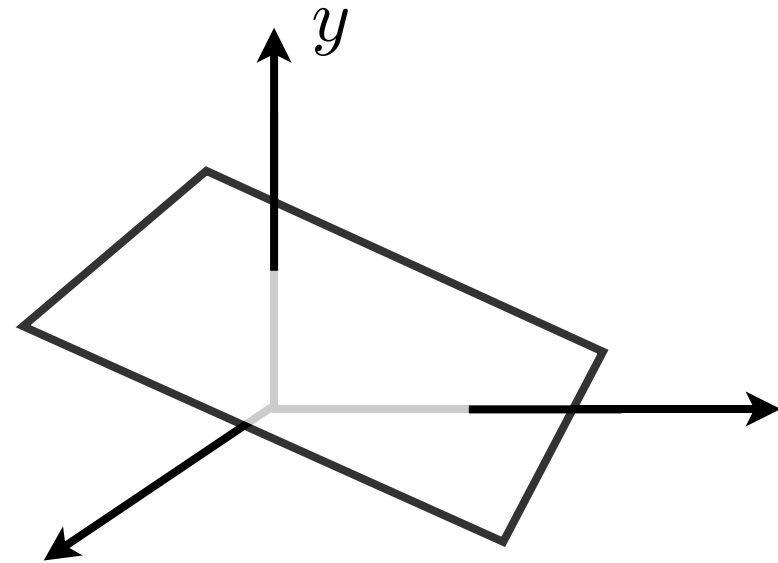
Linear model



$$y = ax + b$$



$$y = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$



is the following a linear model?

$$y = w_1 \cdot x + w_2 \cdot x^2 + b$$

Least square regression



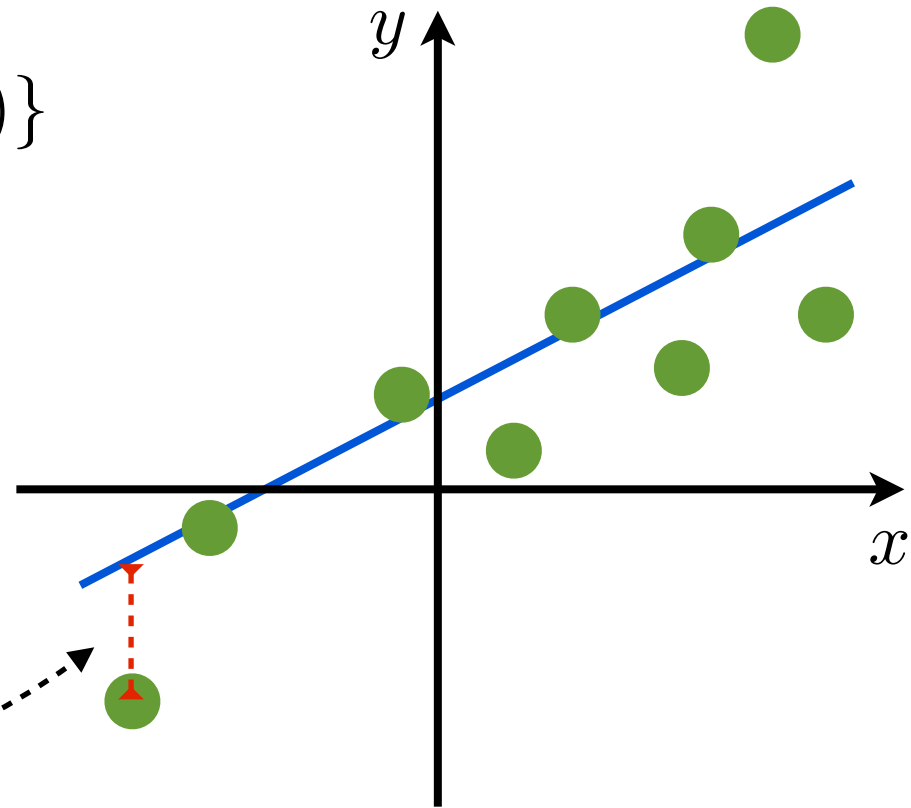
Regression: $y \in \mathbb{R}$

Training data:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_m, y_m)\}$$

Least square loss:

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$$





Least square regression

$$L(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$$

$$\frac{\partial L(\mathbf{w}, b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m 2(\mathbf{w}^\top \mathbf{x}_i + b - y_i) = 0$$

$$\frac{\partial L(\mathbf{w}, b)}{\partial \mathbf{w}} = \frac{1}{m} \sum_{i=1}^m 2(\mathbf{w}^\top \mathbf{x}_i + b - y_i) \mathbf{x}_i^\top = 0$$

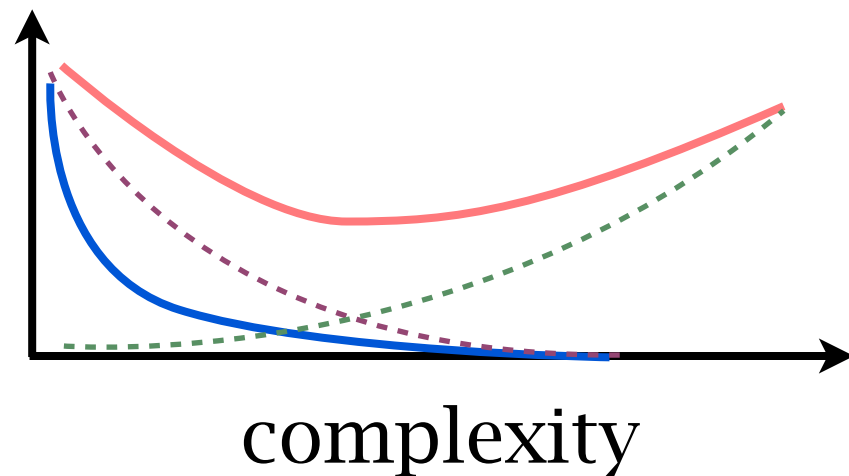
$$b = \frac{1}{m} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i) = \bar{y} - \mathbf{w}^\top \bar{\mathbf{x}}$$

$$\mathbf{w} = \left(\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top - \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \right)^{-1} \left(\frac{1}{m} \sum_{i=1}^m (y_i \mathbf{x}_i) - \bar{y} \bar{\mathbf{x}} \right)$$

$$= \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y) = (X^\top X)^{-1} X^\top Y$$

*closed
form
solution*

Complexity of linear models



$$f(x) = w^T x$$

possibility of w

Regularization



make hypothesis space small

→ better generalization ability

make numerical analysis stable

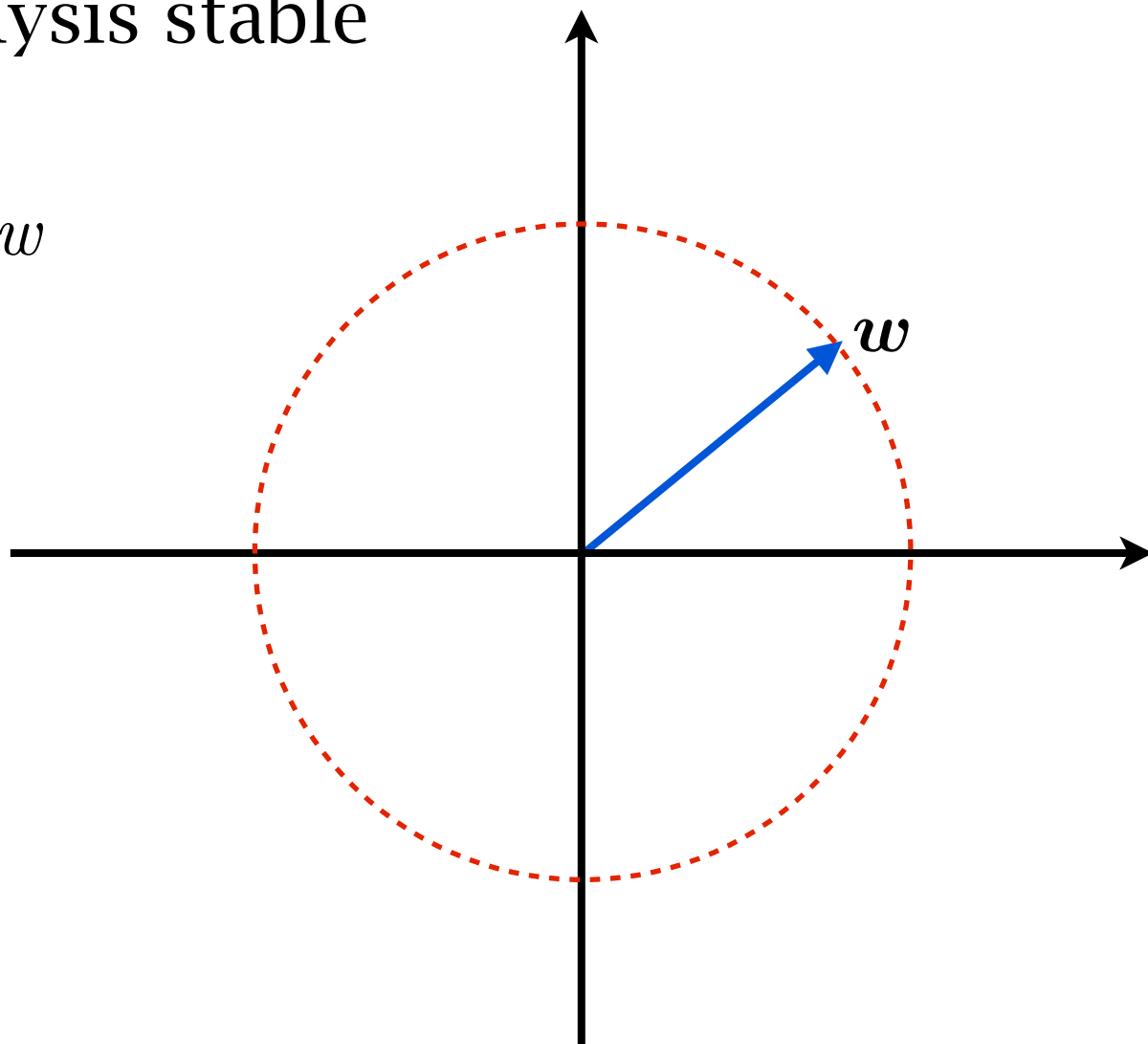
restrict the norm of w

$$\|w\|_p = \left(\sum_{i=1}^n |w_i|^p \right)^{1/p}$$

$$\|w\|_2 = \sqrt{\sum_{i=1}^n w_i^2}$$

$$\|w\|_1 = \sum_{i=1}^n |w_i|$$

$$\|w\|_\infty = \max_{i=1, \dots, n} |w_i|$$



Ridge regression



Regression: $y \in \mathbb{R}$

Training data:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_m, y_m)\}$$

objective:

$$\arg \min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$$

$$s.t. \quad \|\mathbf{w}\|_2 \leq \theta$$

or:

$$\arg \min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 + \lambda \|\mathbf{w}\|_2$$

Ridge regression



centered data, no bias:

$$\arg \min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2$$

closed form solution:

$$\mathbf{w} = \left(\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top - \bar{\mathbf{x}} \bar{\mathbf{x}}^\top + \lambda \mathbf{I} \right)^{-1} \left(\frac{1}{m} \sum_{i=1}^m (y_i \mathbf{x}_i) - \bar{y} \bar{\mathbf{x}} \right)$$

$$= (\text{var}(\mathbf{x}) + \lambda \mathbf{I})^{-1} \text{cov}(\mathbf{x}, y)$$

$$= (X^\top X + \lambda I)^{-1} X^\top Y$$

\mathbf{I} is the identity matrix



Least square v.s. ridge regression



$$\begin{aligned}\mathbf{w} &= \left(\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top - \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \right)^{-1} \left(\frac{1}{m} \sum_{i=1}^m (y_i \mathbf{x}_i) - \bar{y} \bar{\mathbf{x}} \right) \\ &= \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y) = (X^\top X)^{-1} X^\top Y\end{aligned}$$

$$\begin{aligned}\mathbf{w} &= \left(\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top - \bar{\mathbf{x}} \bar{\mathbf{x}}^\top + \lambda \mathbf{I} \right)^{-1} \left(\frac{1}{m} \sum_{i=1}^m (y_i \mathbf{x}_i) - \bar{y} \bar{\mathbf{x}} \right) \\ &= (\text{var}(\mathbf{x}) + \lambda \mathbf{I})^{-1} \text{cov}(\mathbf{x}, y) \\ &= (X^\top X + \lambda \mathbf{I})^{-1} X^\top Y\end{aligned}$$

 stable solution

Least absolute shrinkage and selection operator (LASSO)



Regression: $y \in \mathbb{R}$

Training data:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_m, y_m)\}$$

objective:

$$\arg \min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$$

$$s.t. \quad \|\mathbf{w}\|_1 \leq \theta$$

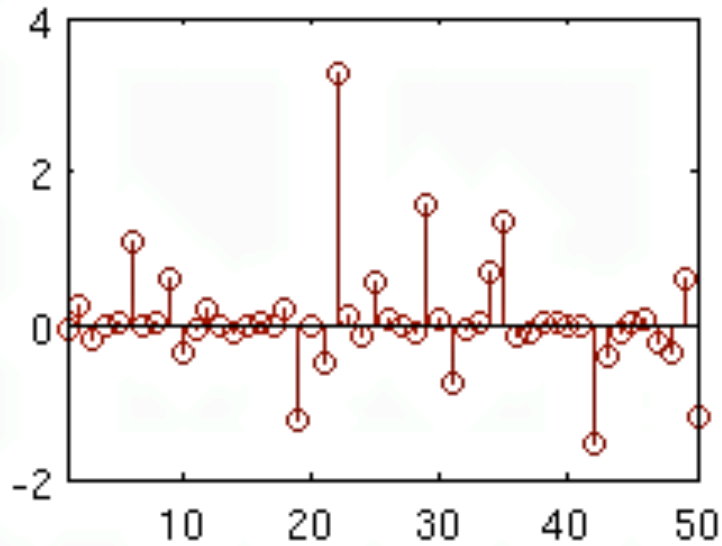
or:

$$\arg \min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

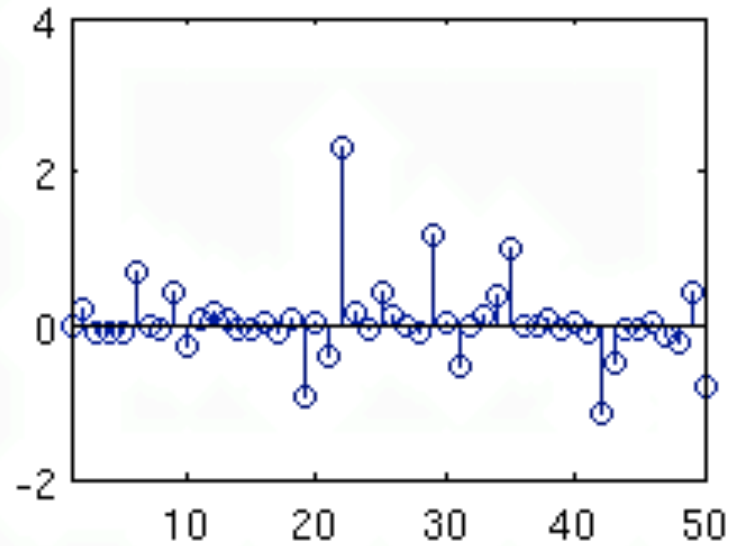
Comparing different regressions



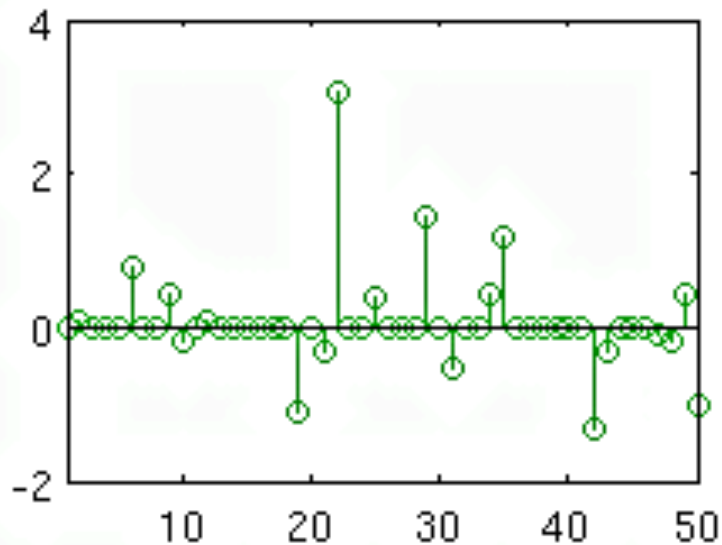
Least Squares



Ridge Regression



LASSO



[Pictures from www.cs.ubc.ca/~schmidtm/Software/L1General/examples.html]

A general framework



objective function:

$$\arg \min_{\mathbf{w}, b} L(\mathbf{w}, b) + \|\mathbf{w}\|_p$$

general optimization: gradient descent

$$(\mathbf{w}, b)_{-} = \eta \frac{\partial(L(\mathbf{w}, b) + \|\mathbf{w}\|_p)}{\partial(\mathbf{w}, b)}$$

good for convex objective functions

$$f(\alpha \mathbf{w}_1 + (1 - \alpha) \mathbf{w}_2) \geq \alpha f(\mathbf{w}_1) + (1 - \alpha) f(\mathbf{w}_2)$$

linear, quadratic

convex + convex \rightarrow convex

Linear classifier



model space: \mathbb{R}^{n+1}

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

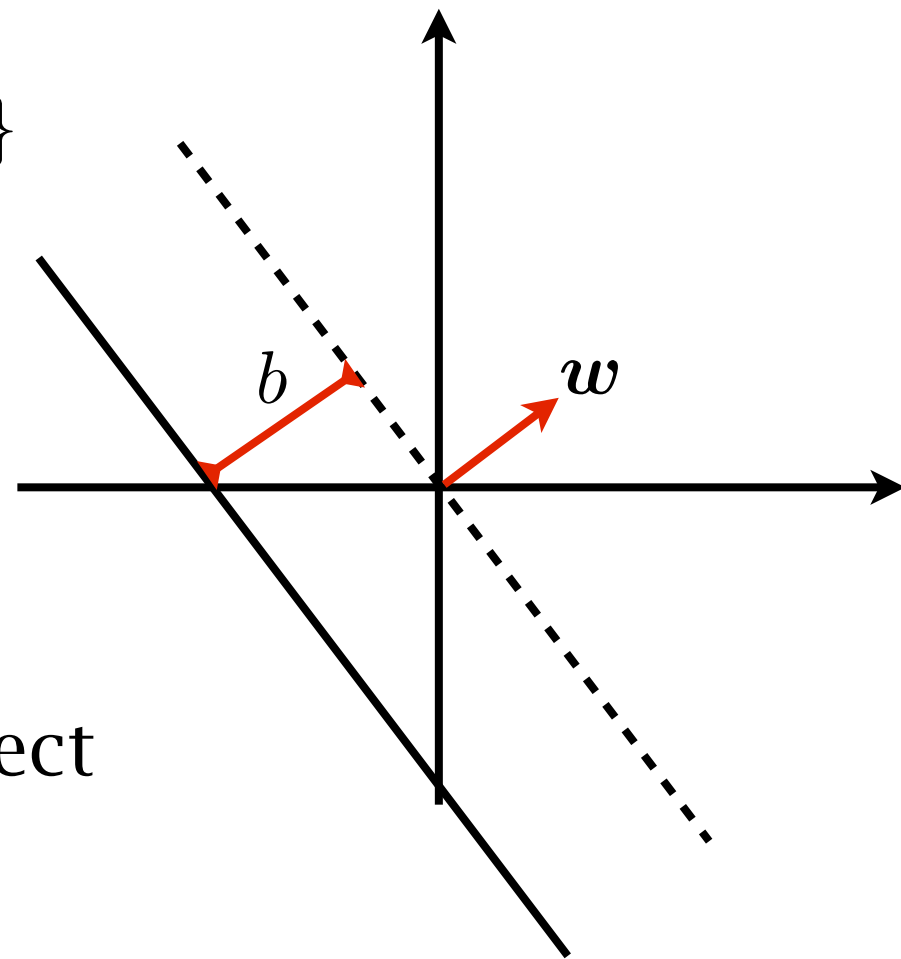
for classification $y \in \{-1, +1\}$

we predict an instance by

$$\text{sign}(\mathbf{w}^\top \mathbf{x} + b) = \begin{cases} +1, & \mathbf{w}^\top \mathbf{x} + b > 0 \\ -1, & \mathbf{w}^\top \mathbf{x} + b < 0 \\ \text{random}, & \text{otherwise} \end{cases}$$

for an example (\mathbf{x}, y) , a correct prediction means

$$y(\mathbf{w}^\top \mathbf{x} + b) > 0$$



Idea classifier



$$\arg \min_{\mathbf{w}, b} \sum_i I(y(\mathbf{w}^\top \mathbf{x} + b) \leq 0)$$

not convex
hard to solve

Prototype

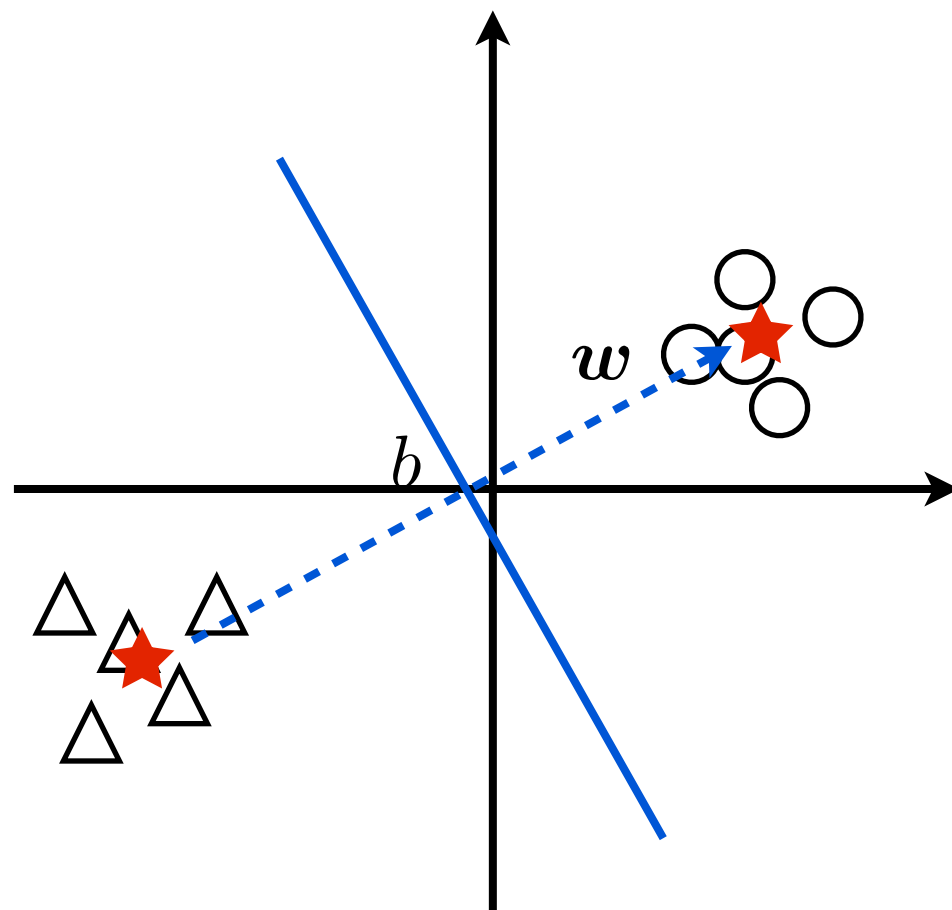
simple, but too restricted

$$\bar{\mathbf{x}}^+ = \frac{1}{\sum_{i:y_i=+1} 1} \sum_{i:y_i=+1} \mathbf{x}_i$$

$$\bar{\mathbf{x}}^- = \frac{1}{\sum_{i:y_i=-1} 1} \sum_{i:y_i=-1} \mathbf{x}_i$$

$$\mathbf{w} = \bar{\mathbf{x}}^+ - \bar{\mathbf{x}}^-$$

$$b = -\mathbf{w}^\top \cdot \frac{\bar{\mathbf{x}}^+ + \bar{\mathbf{x}}^-}{2}$$



Perceptron



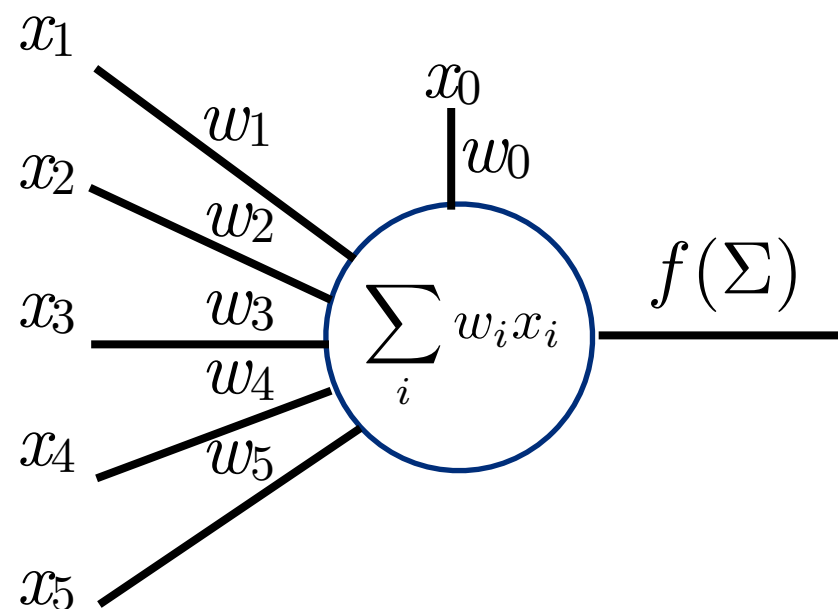
feed training examples one by one

1. $\mathbf{w} = 0$
2. for each example (\mathbf{x}, y)
if $\text{sign}(y\mathbf{w}^\top \mathbf{x}) < 0$

$$\mathbf{w} = \mathbf{w} + y\mathbf{x}$$

gradient ascent

$$\frac{\partial y\mathbf{w}^\top \mathbf{x}}{\partial \mathbf{w}} = y\mathbf{x}$$



$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

when all examples are with length 1 and are linearly separable by \mathbf{w}^* , perceptron algorithm makes at most $\left(1 / \min_{\mathbf{x}} \frac{|\mathbf{w}^{*\top} \mathbf{x}|}{\|\mathbf{x}\|_2}\right)^2$ mistakes

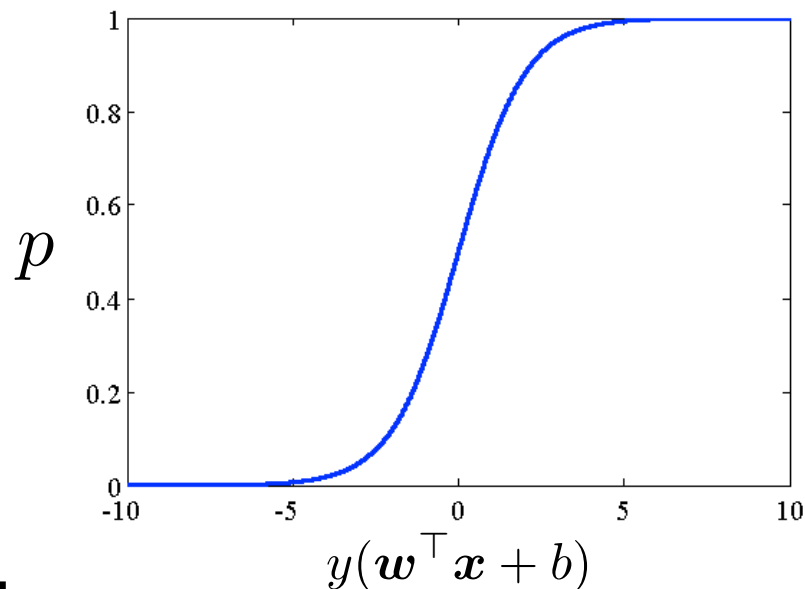
Logistic regression



assume logit model: for a positive example

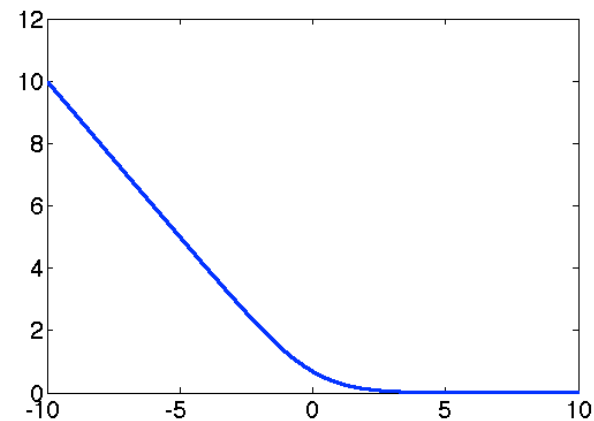
$$\mathbf{w}^\top \mathbf{x} = \log \frac{p(+1 | \mathbf{x})}{1 - p(+1 | \mathbf{x})}$$

$$\text{so that } p(y | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x})}}$$



minimize negative log-likelihood:

$$\begin{aligned} \arg \min_{\mathbf{w}, b} -\log \prod_{i=1}^m p(y_i | \mathbf{x}_i, \mathbf{w}) &= -\sum_i \log p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \sum_i \log \left(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i)} \right) \end{aligned}$$



convex

Linear classifier revisit



model space: \mathbb{R}^{n+1}

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

for classification $y \in \{-1, +1\}$

Original objective:

$$\arg \min_{\mathbf{w}, b} \sum_i I(y(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0)$$

0-1 loss
hard to optimize

Surrogate objective:

$$\arg \min_{\mathbf{w}, b} \sum_i \log \left(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)} \right)$$

logistic regression

$$\arg \min_{\mathbf{w}, b} \sum_i \max\{-y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0\}$$

perceptron

Linear classifier revisit



0-1 loss

$$I(y(\mathbf{w}^\top \mathbf{x} + b) \leq 0)$$

logistic regression

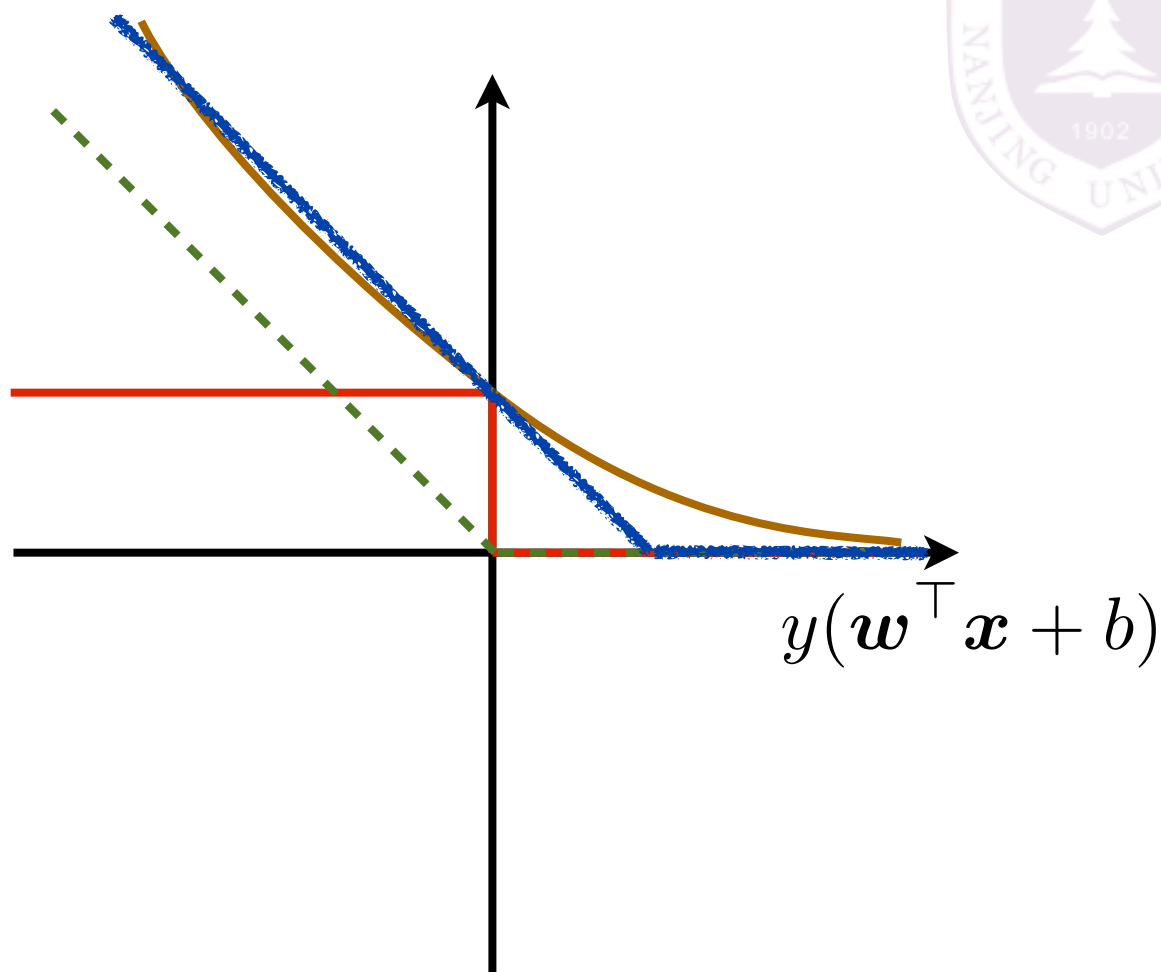
$$\log_2(1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)})$$

perceptron

$$\max\{-y(\mathbf{w}^\top \mathbf{x} + b), 0\}$$

hinge loss

$$\max\{1 - y(\mathbf{w}^\top \mathbf{x} + b), 0\}$$



Support vector machines (SVM)



hinge loss + L2-norm

$$\arg \min_{\mathbf{w}, b} \sum_i \max(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0) + \lambda \|\mathbf{w}\|_2$$

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2 + C \sum_i \xi_i$$

$$\begin{aligned} s.t. \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

$$\begin{aligned} \max(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0) &= \xi_i \\ \xi_i &\geq 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \\ \xi_i &\geq 0 \end{aligned}$$

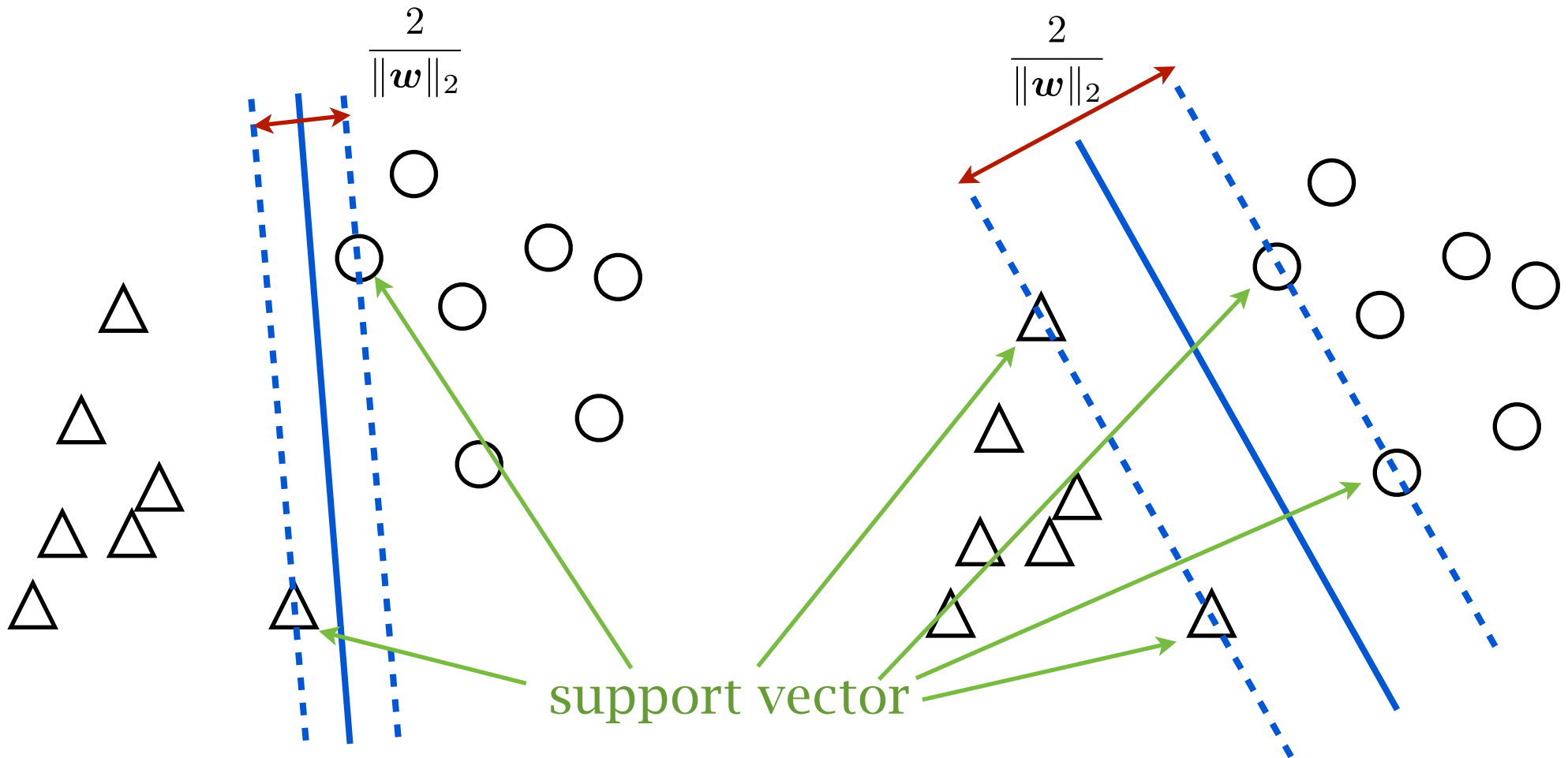
quadratic

Support vector machines (SVM)



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

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



Scoring functions



$$\frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 \quad \text{least square regression}$$

$$\frac{1}{m} \sum_{i=1}^m |\mathbf{w}^\top \mathbf{x}_i + b - y_i| \quad \text{LAD regression}$$

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 + \lambda \|\mathbf{w}\|_2 \quad \text{ridge regression}$$

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 + \lambda \|\mathbf{w}\|_1 \quad \text{LASSO}$$

Scoring functions



$$\sum_i I(y(\mathbf{w}^\top \mathbf{x} + b) > 0)$$

0-1 loss

$$\sum_i \max\{-y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0\}$$

perceptron

$$\sum_i \log\left(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}\right)$$

logistic regression

$$\sum_i \log\left(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)}\right) + \lambda \|\mathbf{w}\|_2$$

regularized LR

$$\sum_i \max(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0) + \lambda \|\mathbf{w}\|_2$$

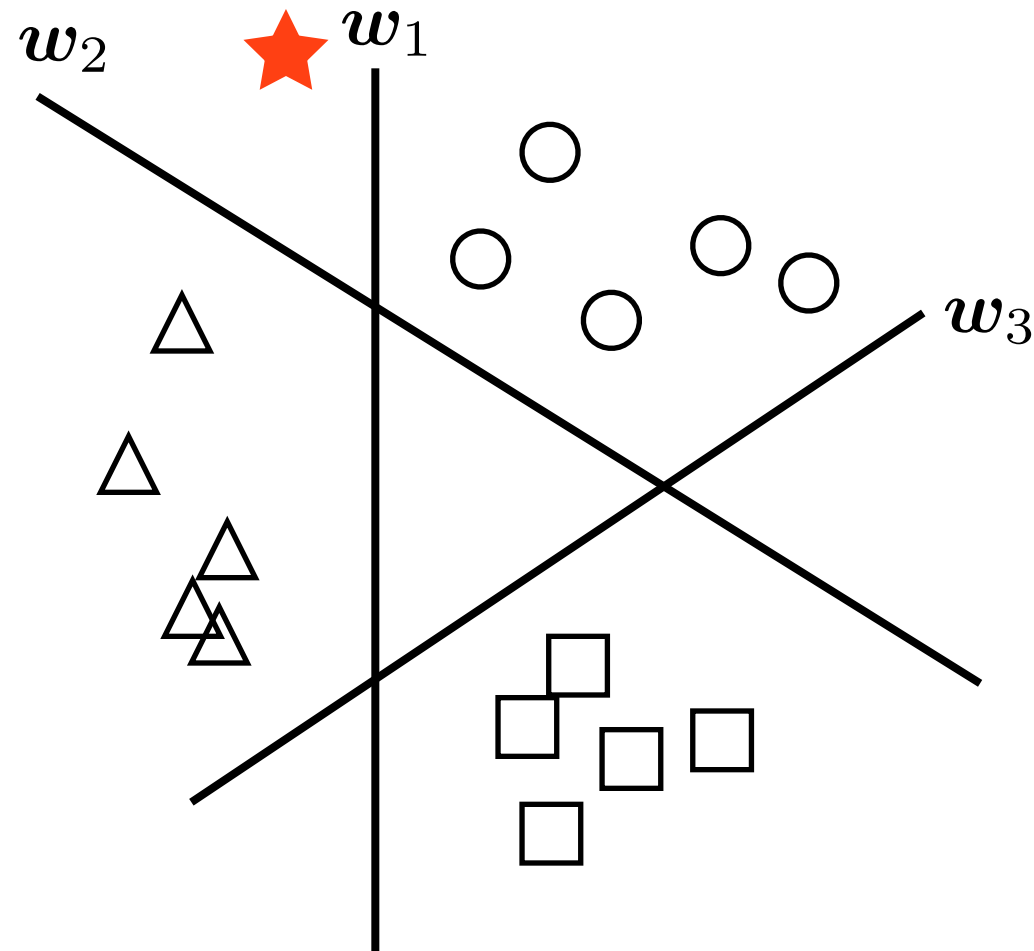
SVM

minimize loss + regularization

Multi-class classification



one-vs-rest

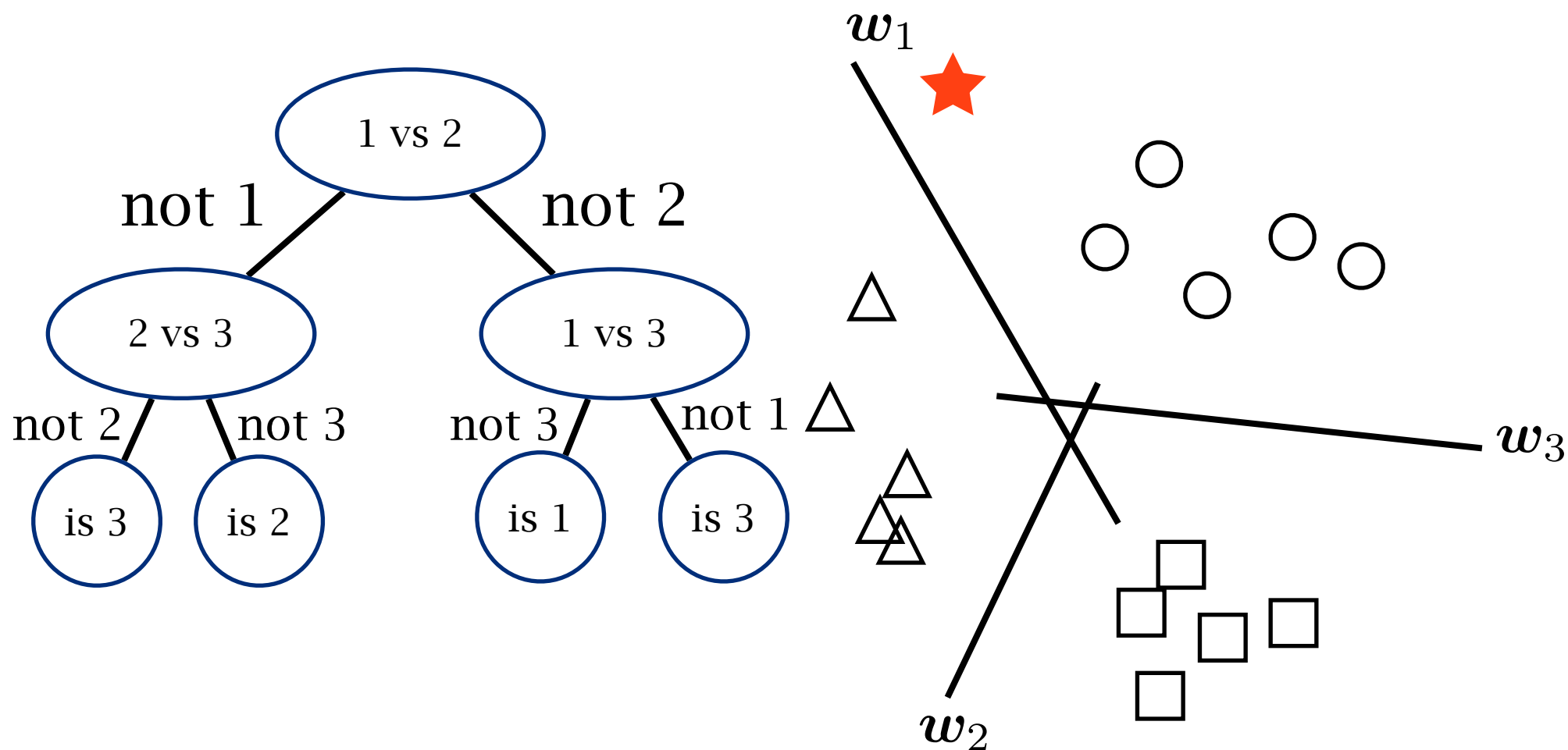


for C classes, need to train C binary classifiers

Multi-class classification

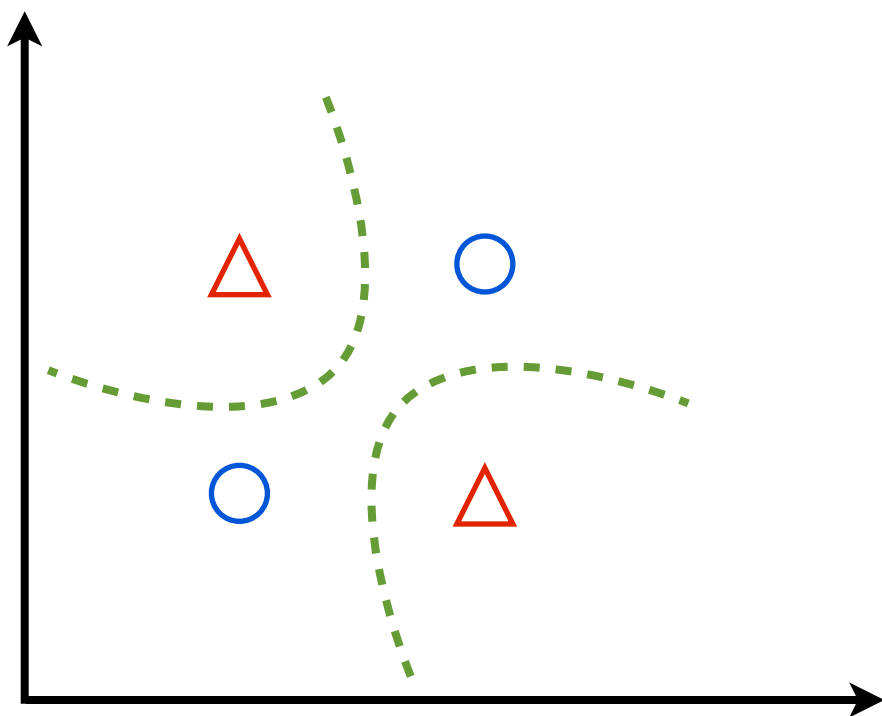


one-vs-one

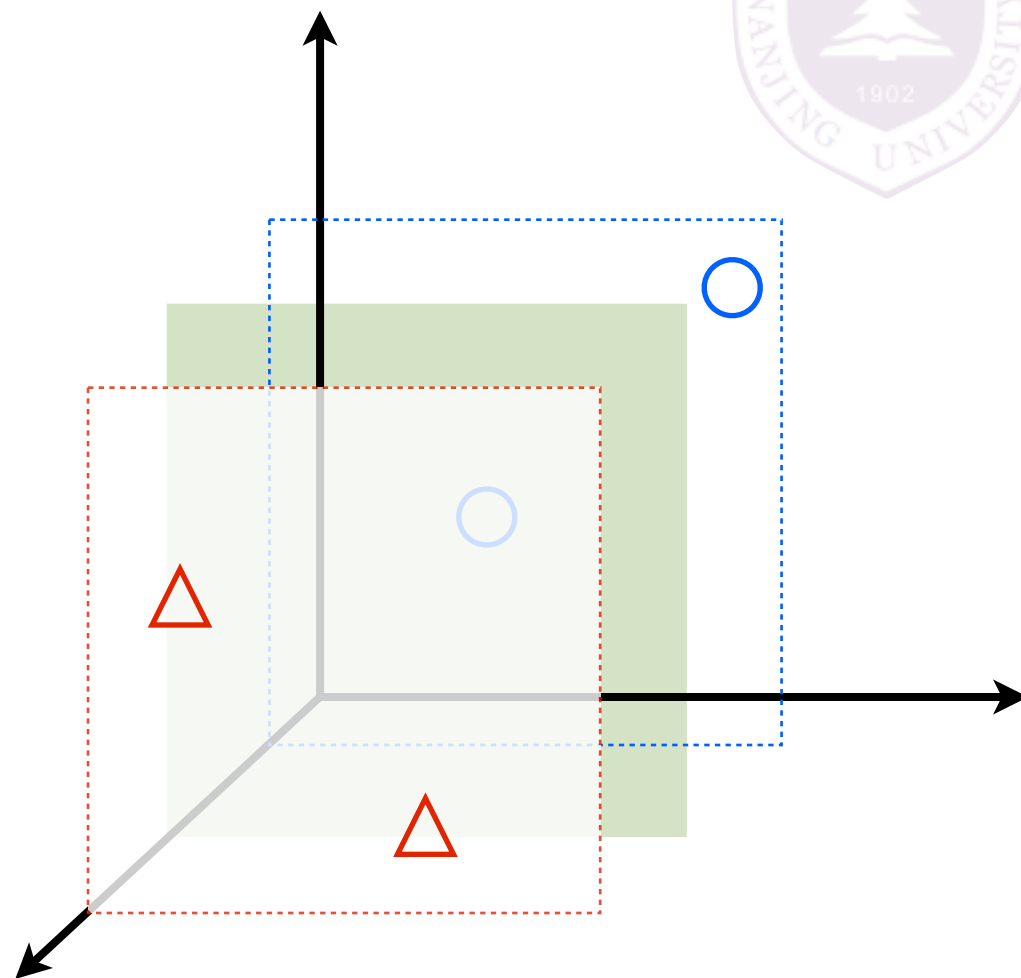


for C classes, need to train $C(C-1)/2$ binary classifiers

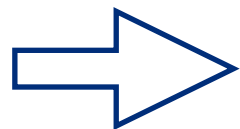
Linearity v.s. dimensionality



XOR in 2D



x_1	x_2	y
0	0	+1
0	1	-1
1	0	-1
1	1	+1



x_1	x_2	x_1x_2	y
0	0	0	+1
0	1	0	-1
1	0	0	-1
1	1	1	+1

$$w = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix}, b = -0.5$$