

Lecture 8: Machine Learning VI

Linear Models

http://cs.nju.edu.cn/yuy/course_dm13ms.ashx



Linear model

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



Linear model

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

$$w_1, w_2, \dots, w_n \quad b$$



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



Linear model



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

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



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

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



Linear model

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

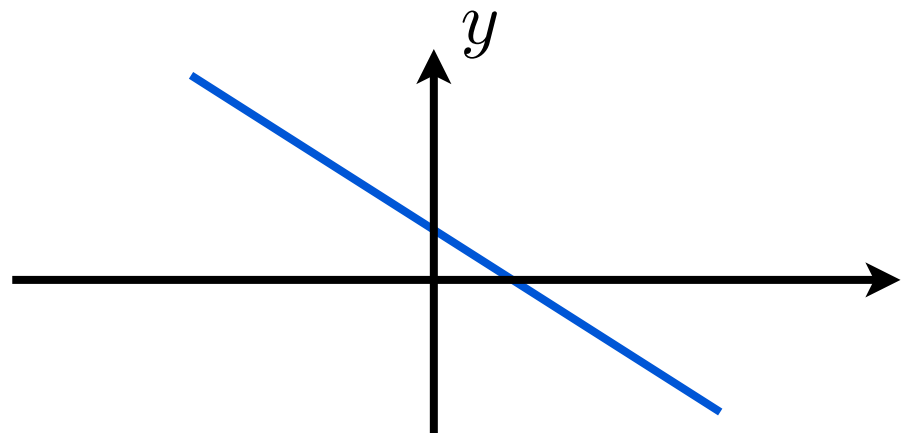
$$\mathbf{w} = w_1, w_2, \dots, w_n \quad b$$



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

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

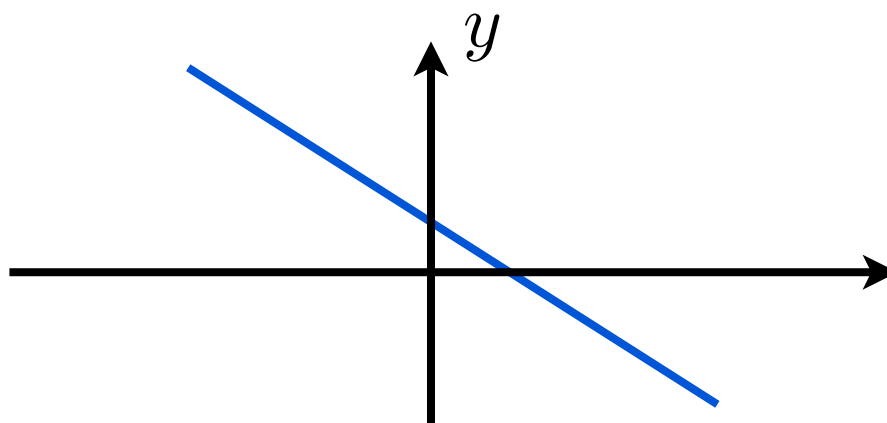
$$y = ax + b$$



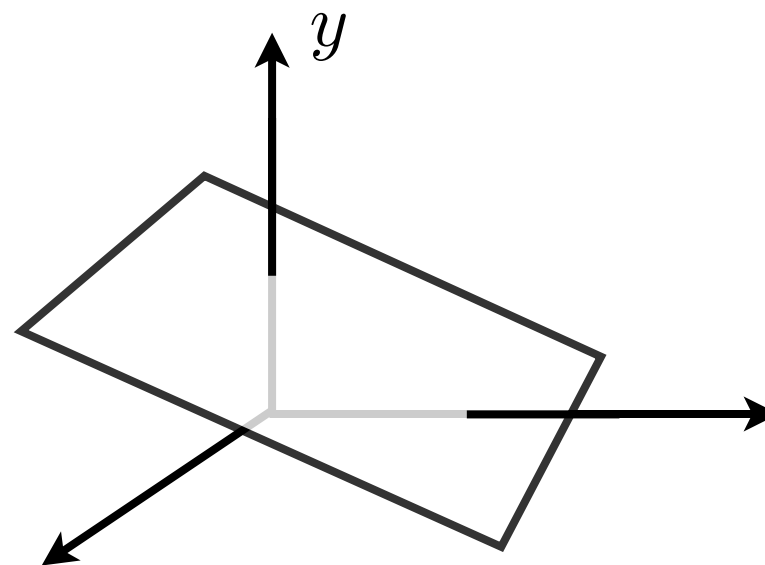
Linear model



$$y = ax + b$$



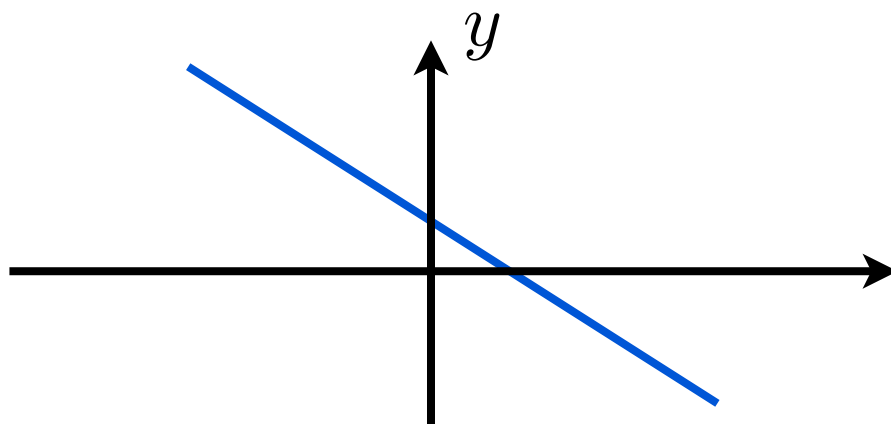
$$y = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$



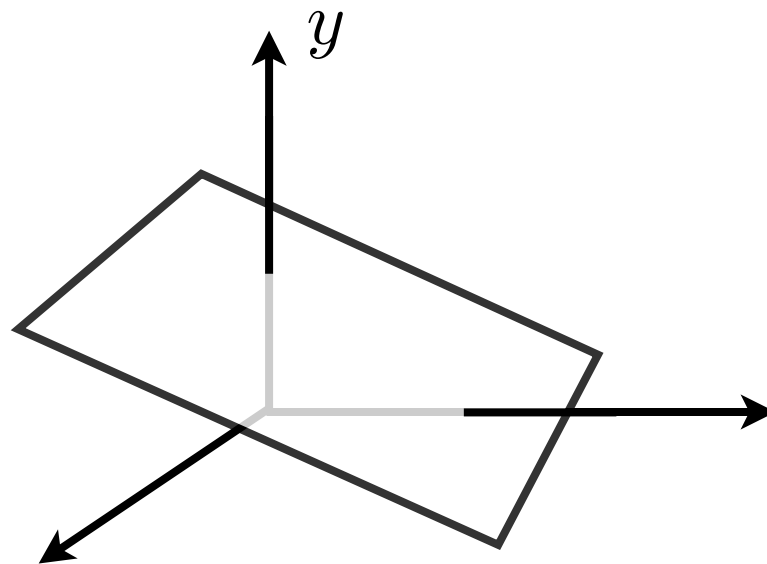
Linear model



$$y = ax + b$$



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



is the following a linear model?

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

Linear model



y



x_1

x_2

...

x_n

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

output/response
variable

linear relationship
independent parameters

basis

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

we sometimes omit the bias

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

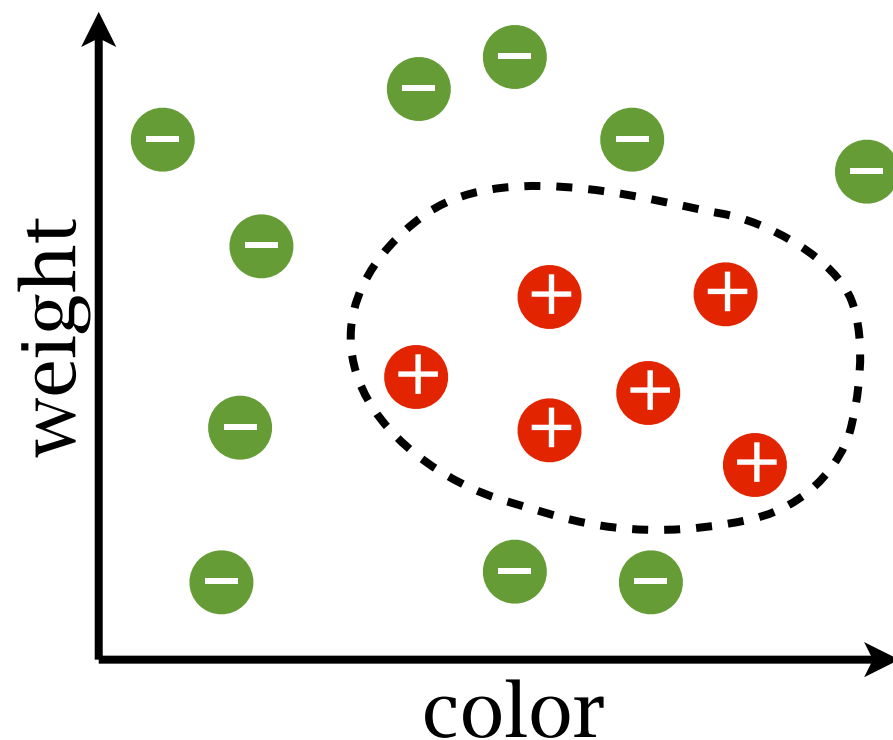
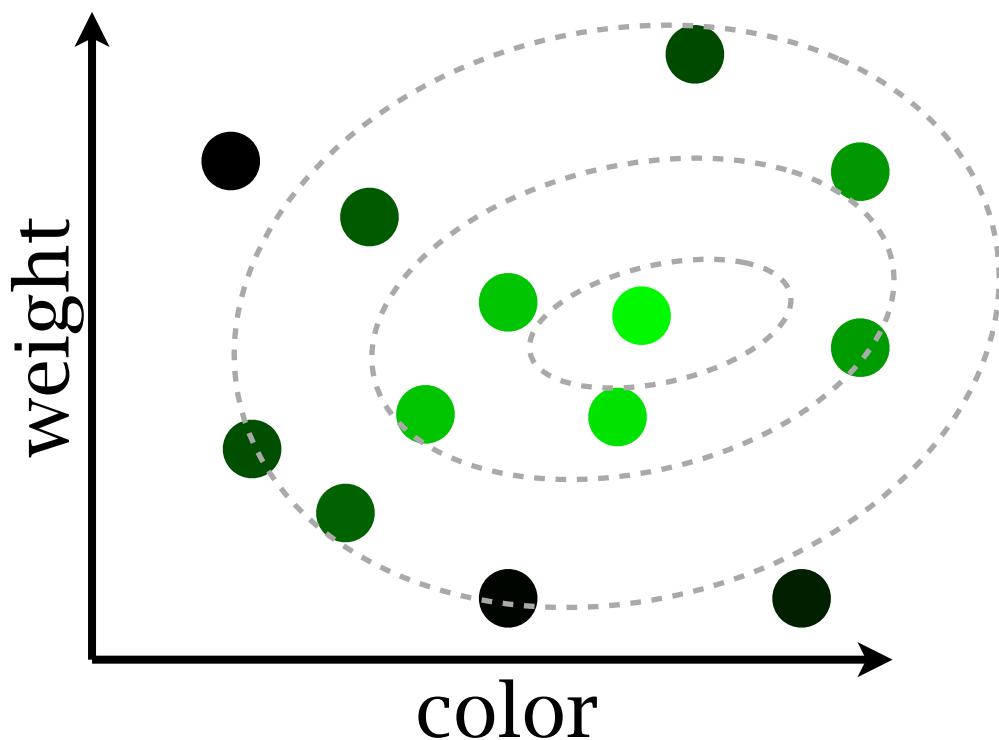
1. \mathbf{x} is with a constant element
2. practically as good as with bias (centered data)

Learning with linear models



linear models can be used for both regression and classification tasks

suitable for continuous features



Least square regression



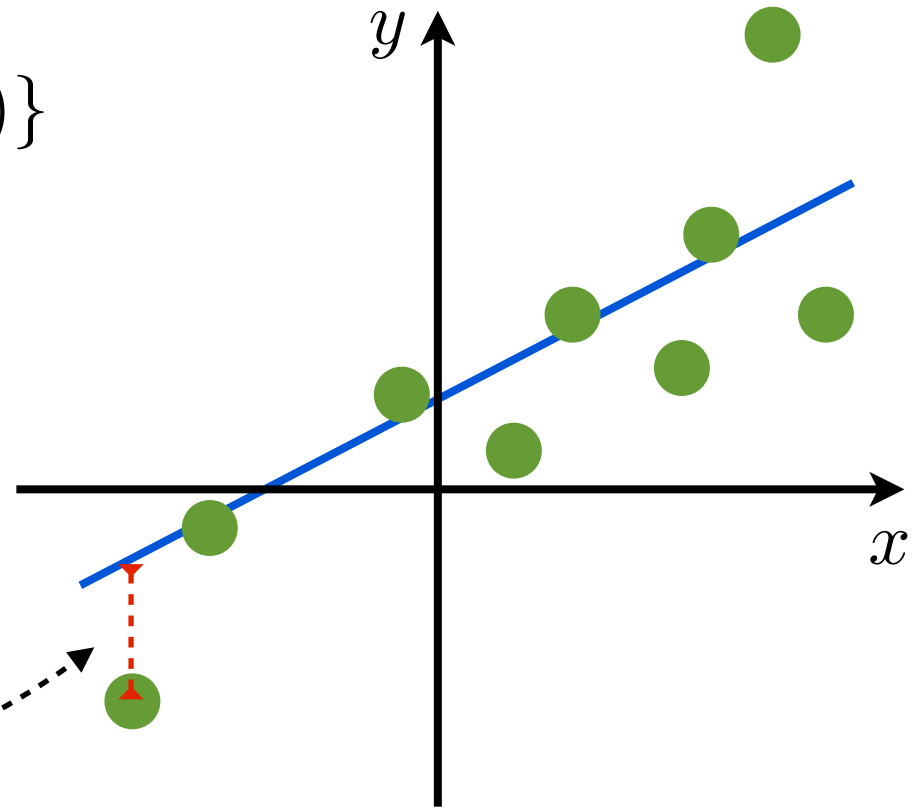
Regression: $y \in \mathbb{R}$

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

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



Least square regression

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

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

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



Least square regression



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

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

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

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

Least square regression



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

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

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

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

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

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



Least square regression

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

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

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

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

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

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

*closed
form
solution*

Least absolute deviation regression



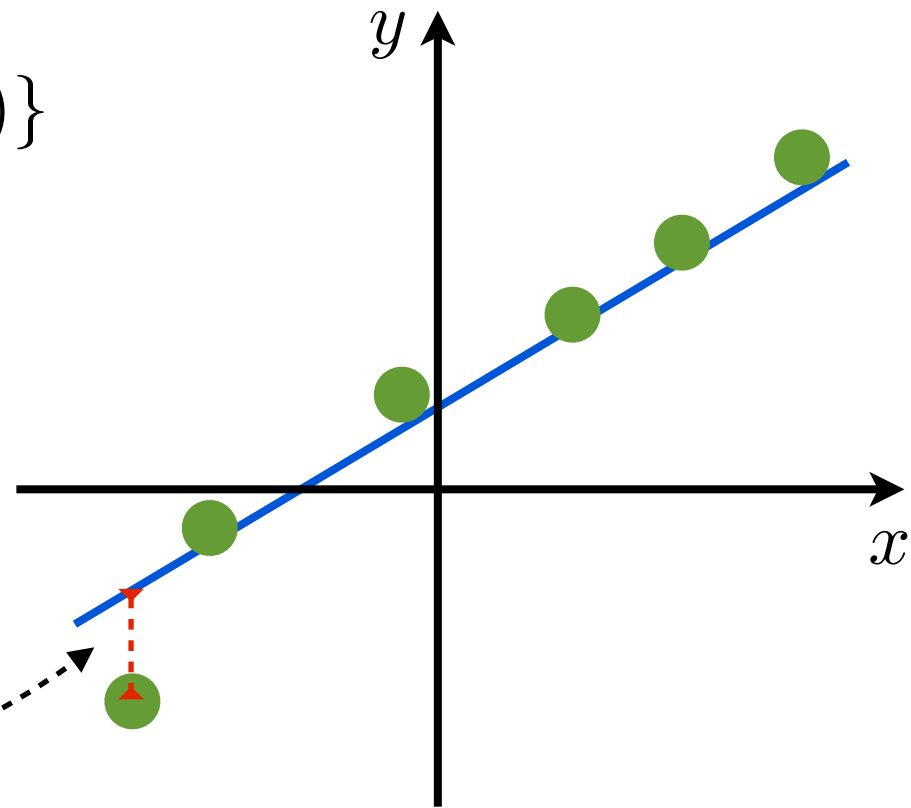
Regression: $y \in \mathbb{R}$

Training data:

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

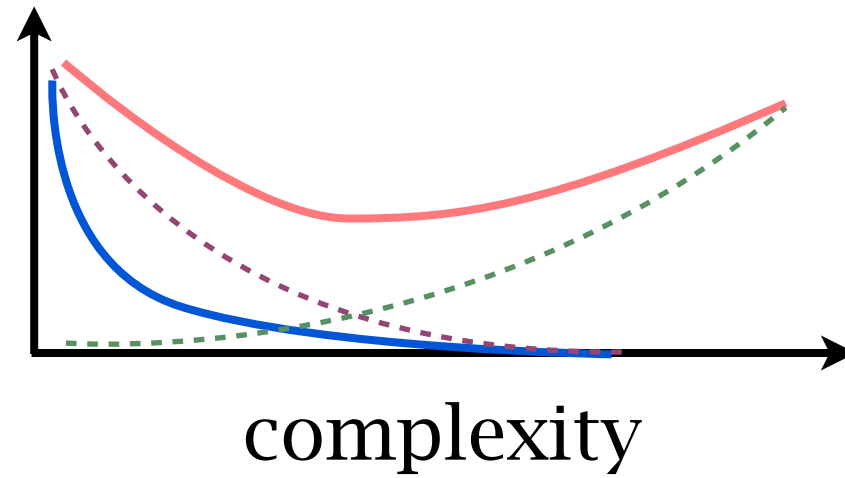
LAD loss:

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

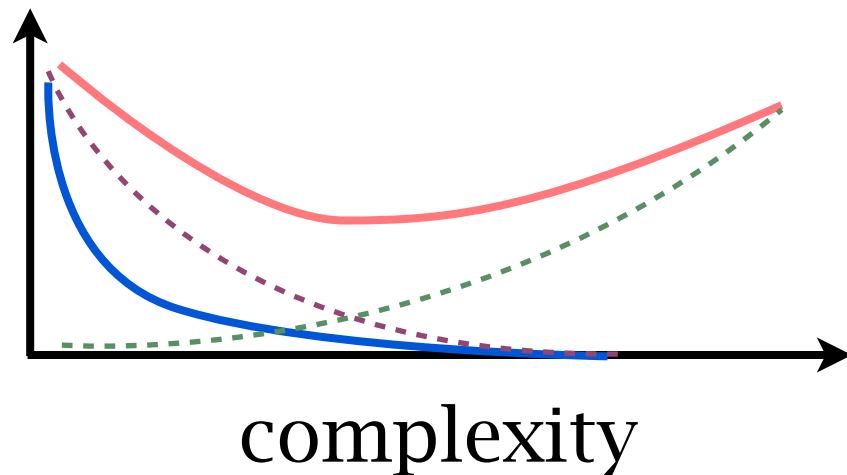


compare with least square regression:
robust to noise
unstable solution

Complexity of linear models



Complexity of linear models



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

↑
possibility of w

Regularization



make hypothesis space small

→ better generalization ability

make numerical analysis stable

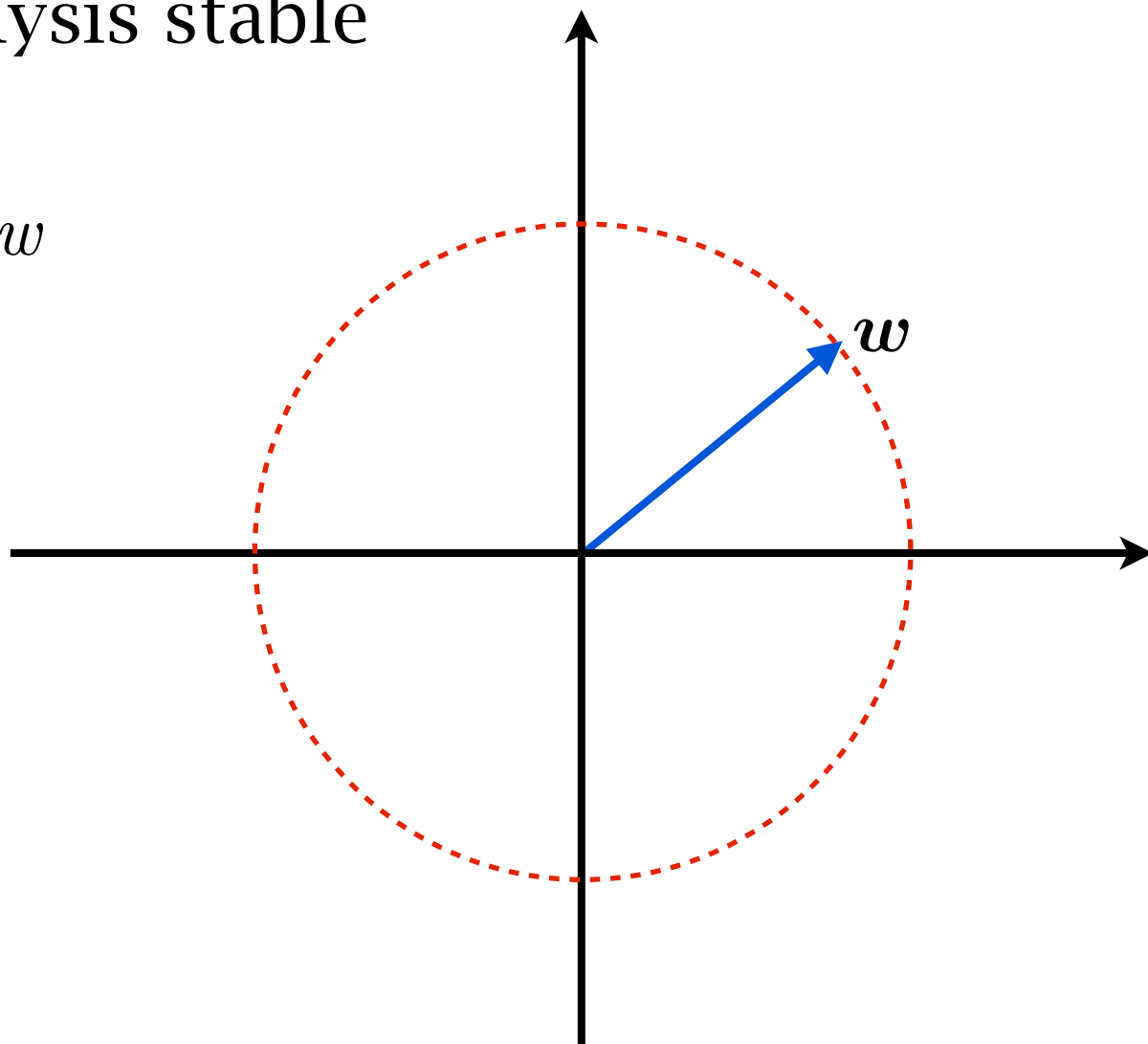
restrict the norm of w

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

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

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

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



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:

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

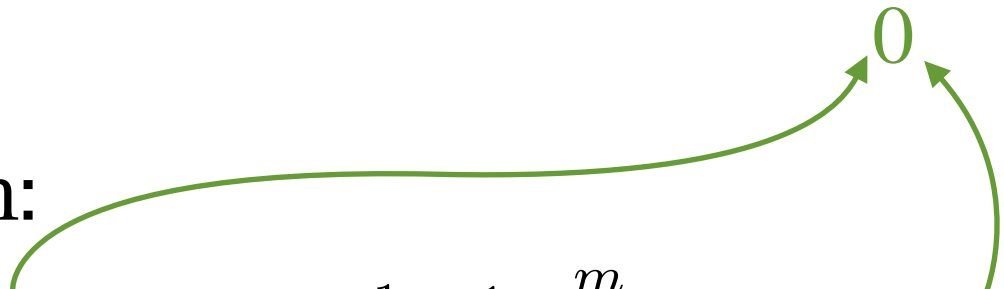
closed form solution:

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

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

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

\mathbf{I} is the identity matrix



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:

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

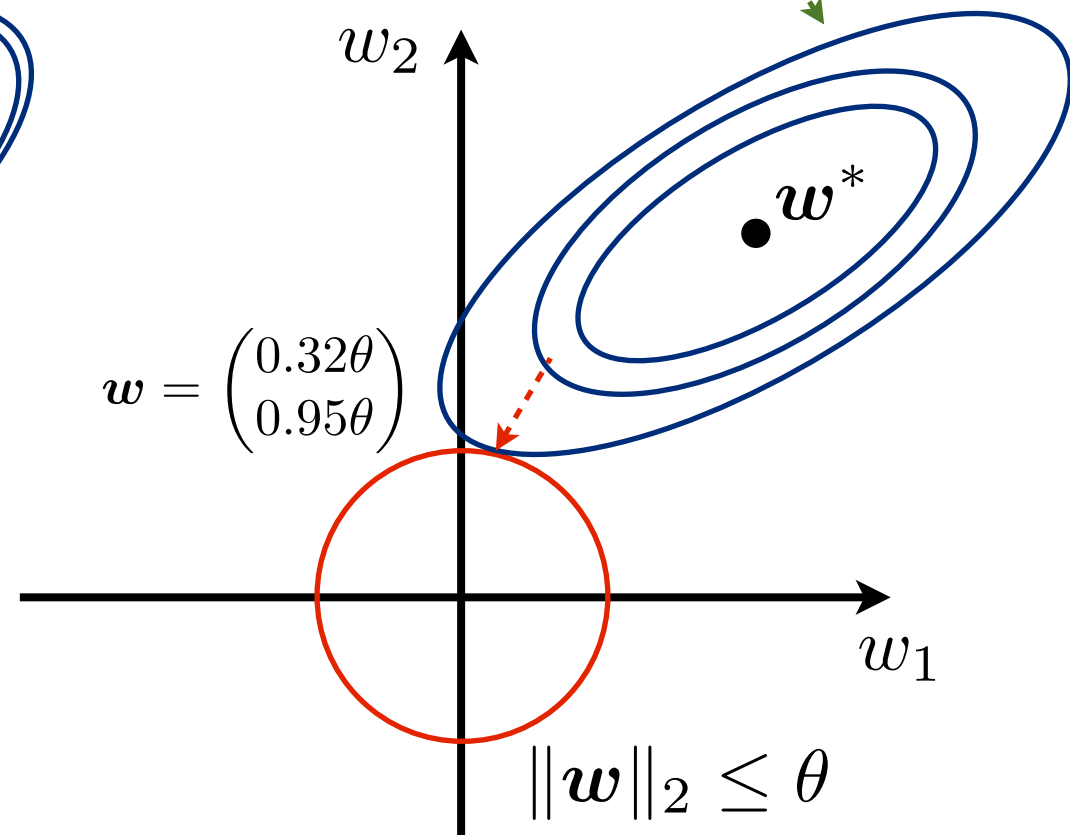
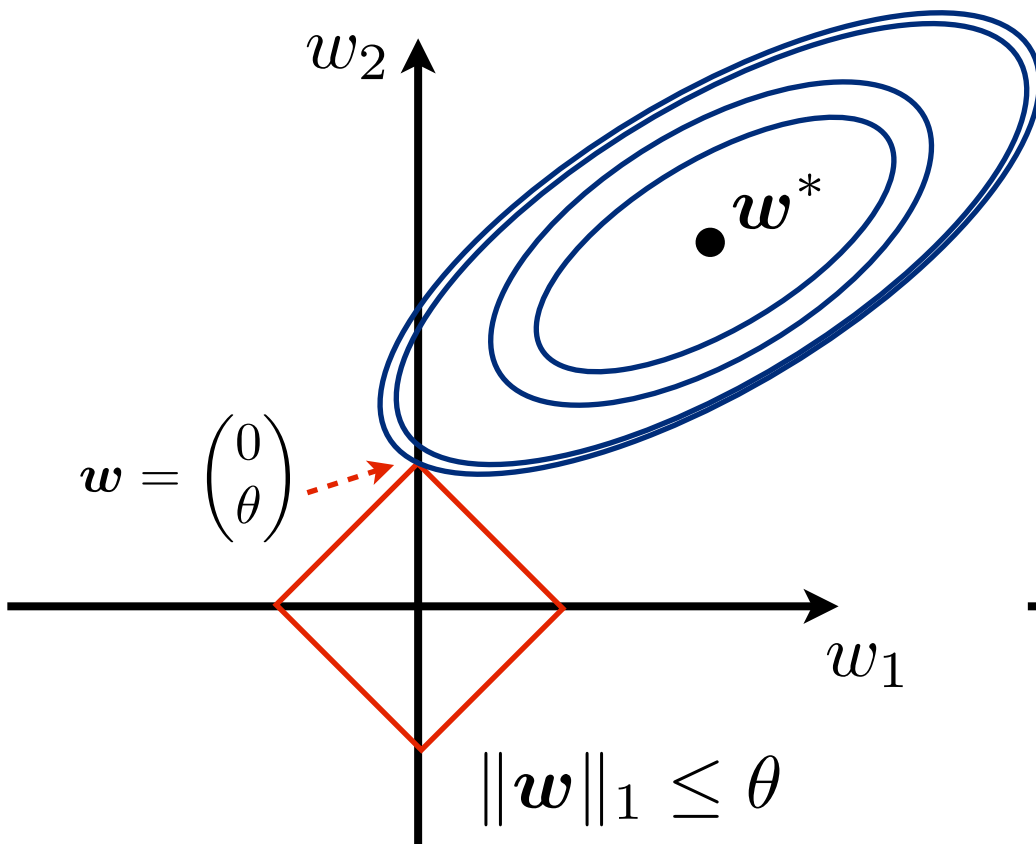
Comparing ridge regression with lasso



L1-norm leads to sparser solution, but worse empirical loss

sparse: many zero elements

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



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}

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

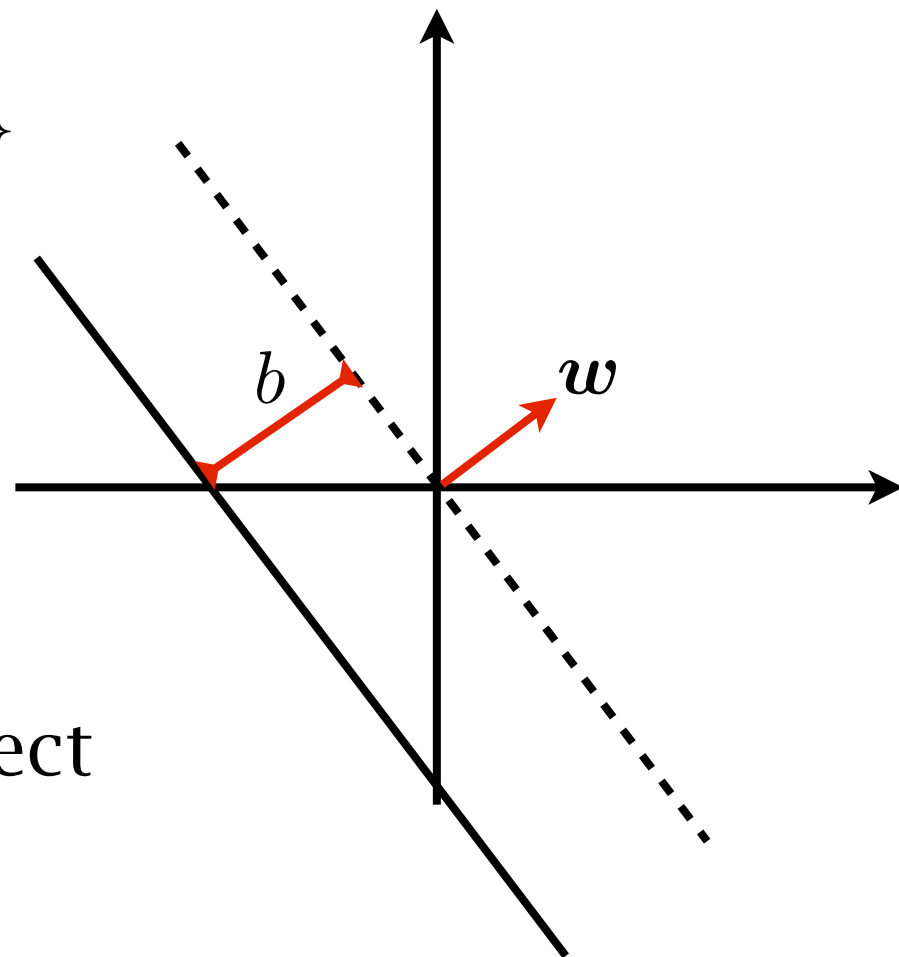
for classification $y \in \{-1, +1\}$

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

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



Prototype

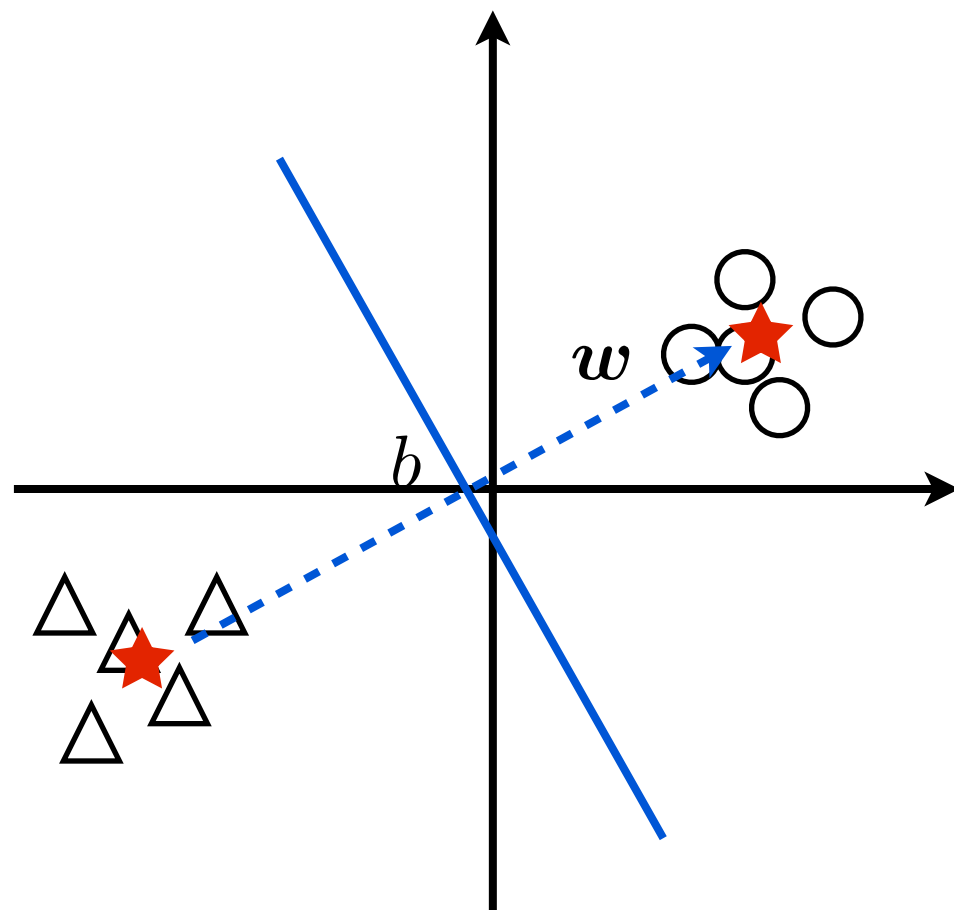
simple, but too restricted

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

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

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

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



Perceptron

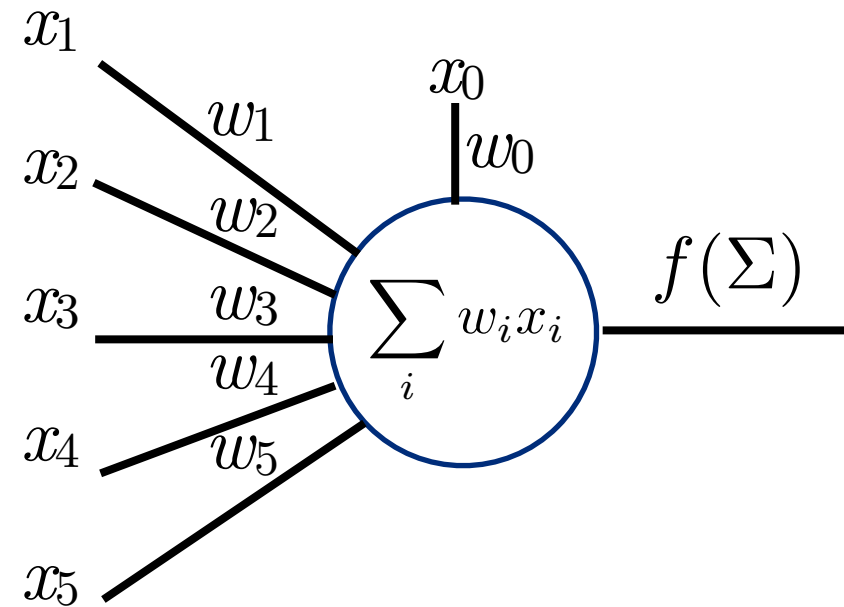


feed training examples one by one

1. $\mathbf{w} = 0$

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

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



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

Perceptron

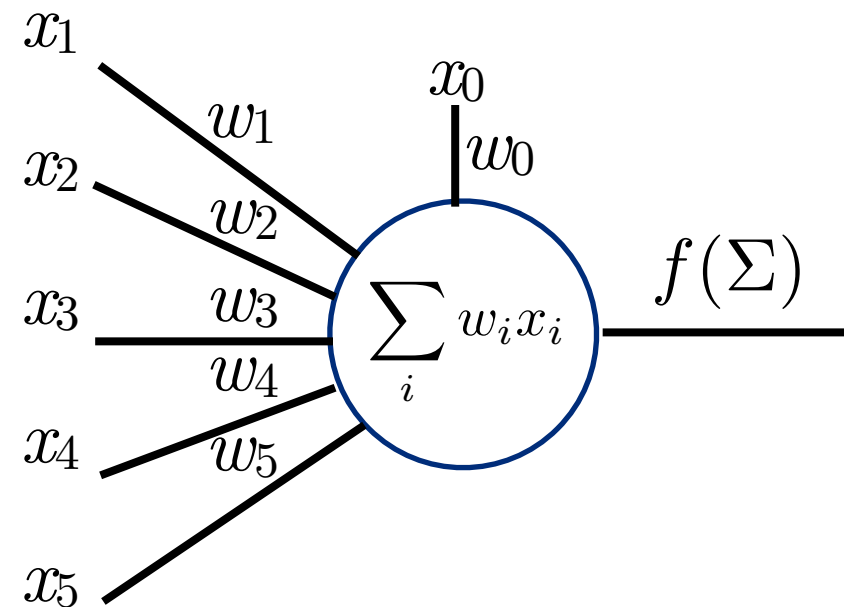


feed training examples one by one

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

gradient ascent

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



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

Perceptron



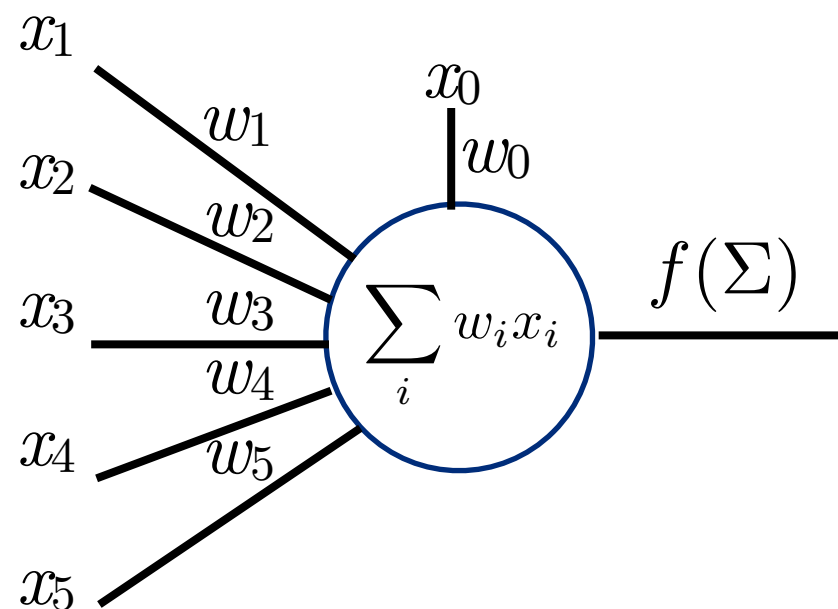
feed training examples one by one

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

$$w = w + yx$$

gradient ascent

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



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

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

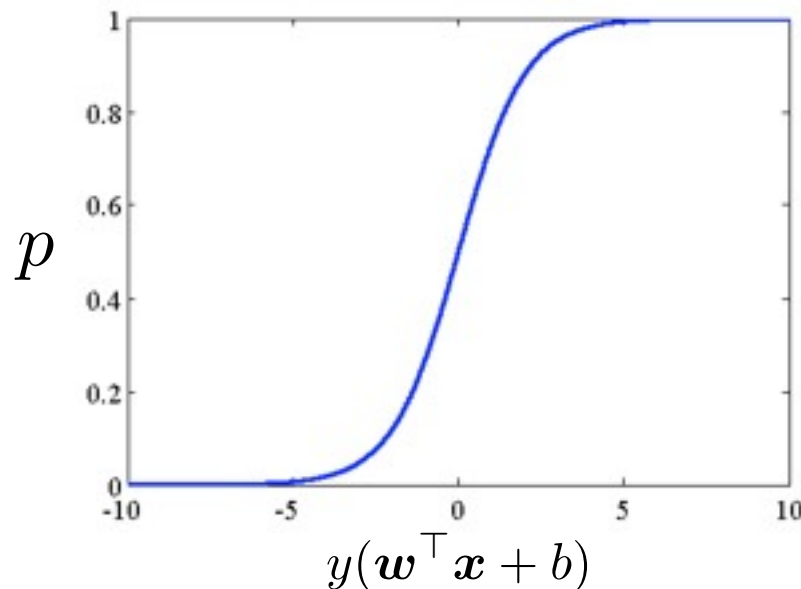
Logistic regression



assume logit model: for a positive example

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

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



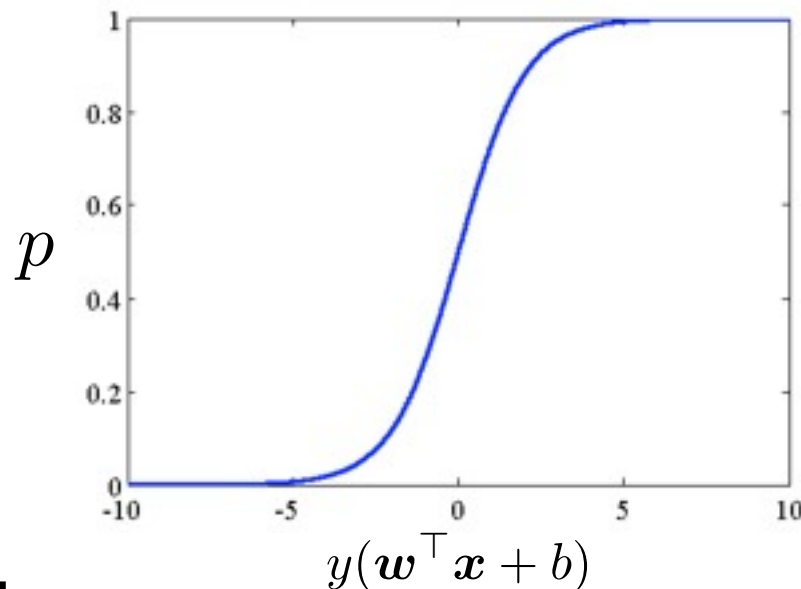
Logistic regression



assume logit model: for a positive example

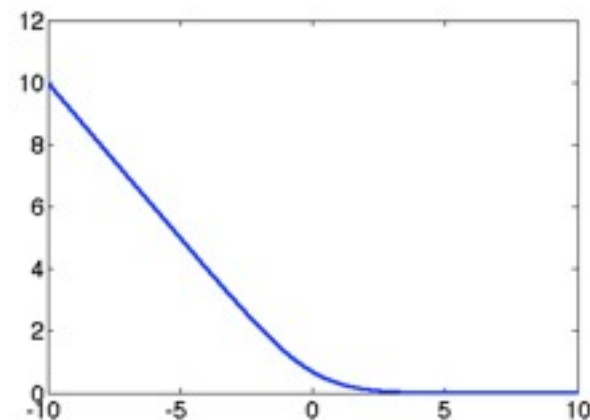
$$\mathbf{w}^\top \mathbf{x} = \log \frac{p(+1 | \mathbf{x})}{1 - p(+1 | \mathbf{x})}$$

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

Logistic regression



Maximize a posterior (minimize negative a posterior)

$$\arg \min_{\mathbf{w}, b} - \log \left(\prod_{i=1}^m p(y_i | \mathbf{x}_i, \mathbf{w}) \right) p(\mathbf{w})$$

a prior: $\mathbf{w} \sim \mathcal{N}(0, \delta \mathbf{I})$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; 0, \delta \mathbf{I}) = \frac{1}{\delta \sqrt{2\pi}} e^{-\frac{\|\mathbf{w}-0\|_2^2}{2\delta^2}}$$

$$= - \sum_i \log p(y_i | \mathbf{x}_i, \mathbf{w}) - \log p(\mathbf{w})$$

$$= \sum_i \log \left(1 + e^{-y_i (\mathbf{w}^\top \mathbf{x}_i)} \right) + \frac{1}{2\delta^2} \|\mathbf{w}\|_2^2 + \text{const}$$

Logistic regression



Maximize a posterior (minimize negative a posterior)

$$\arg \min_{\mathbf{w}, b} - \log \left(\prod_{i=1}^m p(y_i | \mathbf{x}_i, \mathbf{w}) \right) p(\mathbf{w})$$

a prior: $\mathbf{w} \sim \mathcal{N}(0, \delta \mathbf{I})$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; 0, \delta \mathbf{I}) = \frac{1}{\delta \sqrt{2\pi}} e^{-\frac{\|\mathbf{w}-0\|_2^2}{2\delta^2}}$$

$$= - \sum_i \log p(y_i | \mathbf{x}_i, \mathbf{w}) - \log p(\mathbf{w})$$

$$= \sum_i \log \left(1 + e^{-y_i (\mathbf{w}^\top \mathbf{x}_i)} \right) + \frac{1}{2\delta^2} \|\mathbf{w}\|_2^2 + \text{const}$$

convex

regularized logistic regression

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:

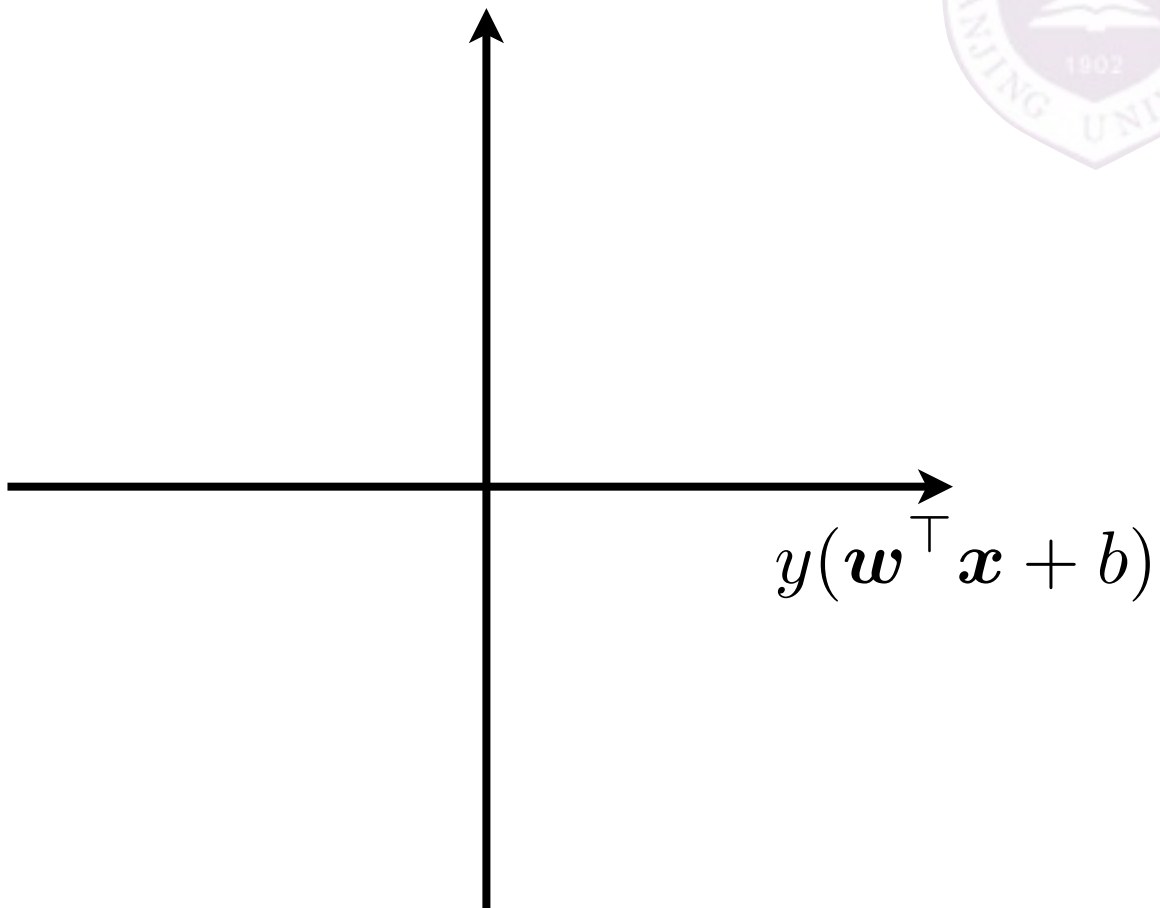
$$\arg \min_{\mathbf{w}, b} \sum_i \log \left(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)} \right)$$

logistic regression

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

perceptron

Linear classifier revisit



Linear classifier revisit

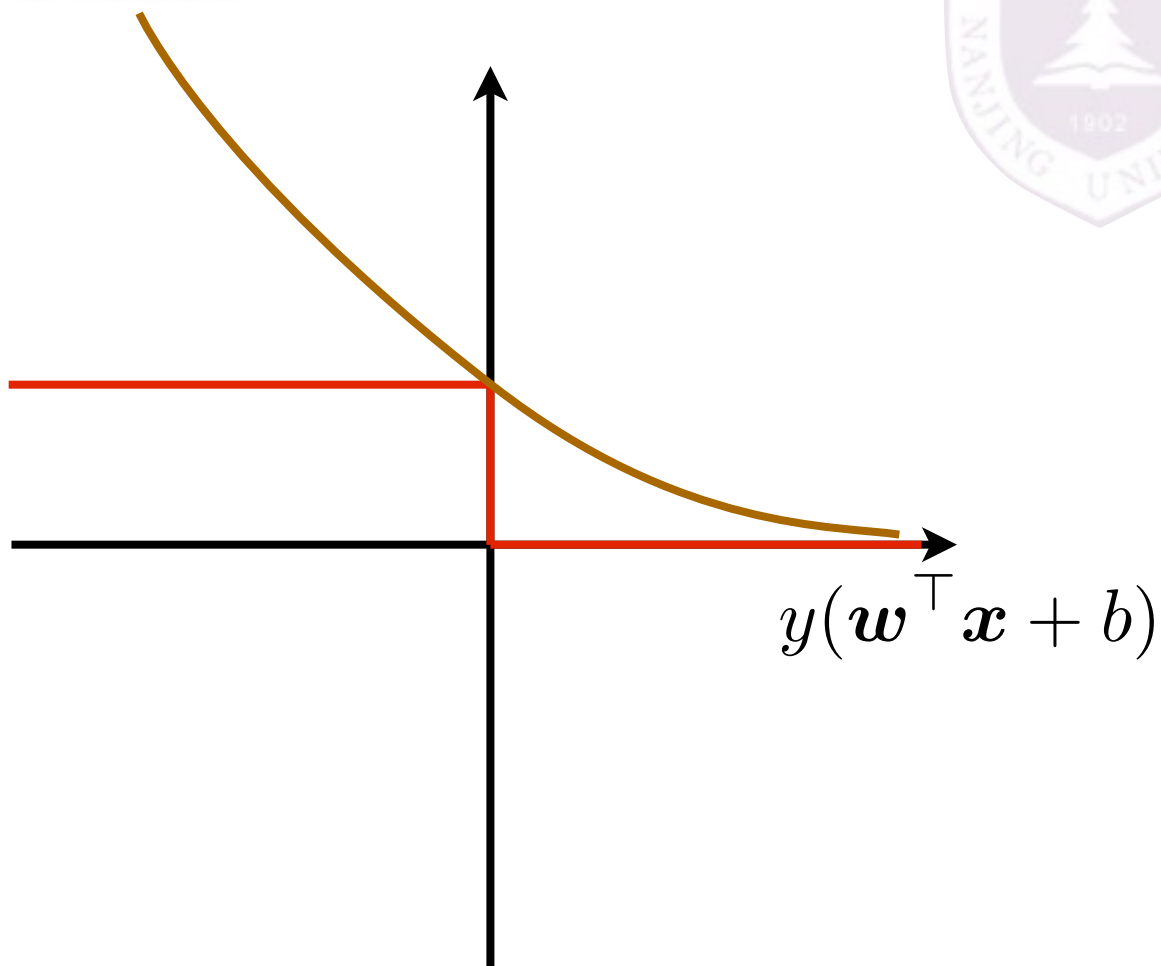


0-1 loss

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

logistic regression

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



Linear classifier revisit



0-1 loss

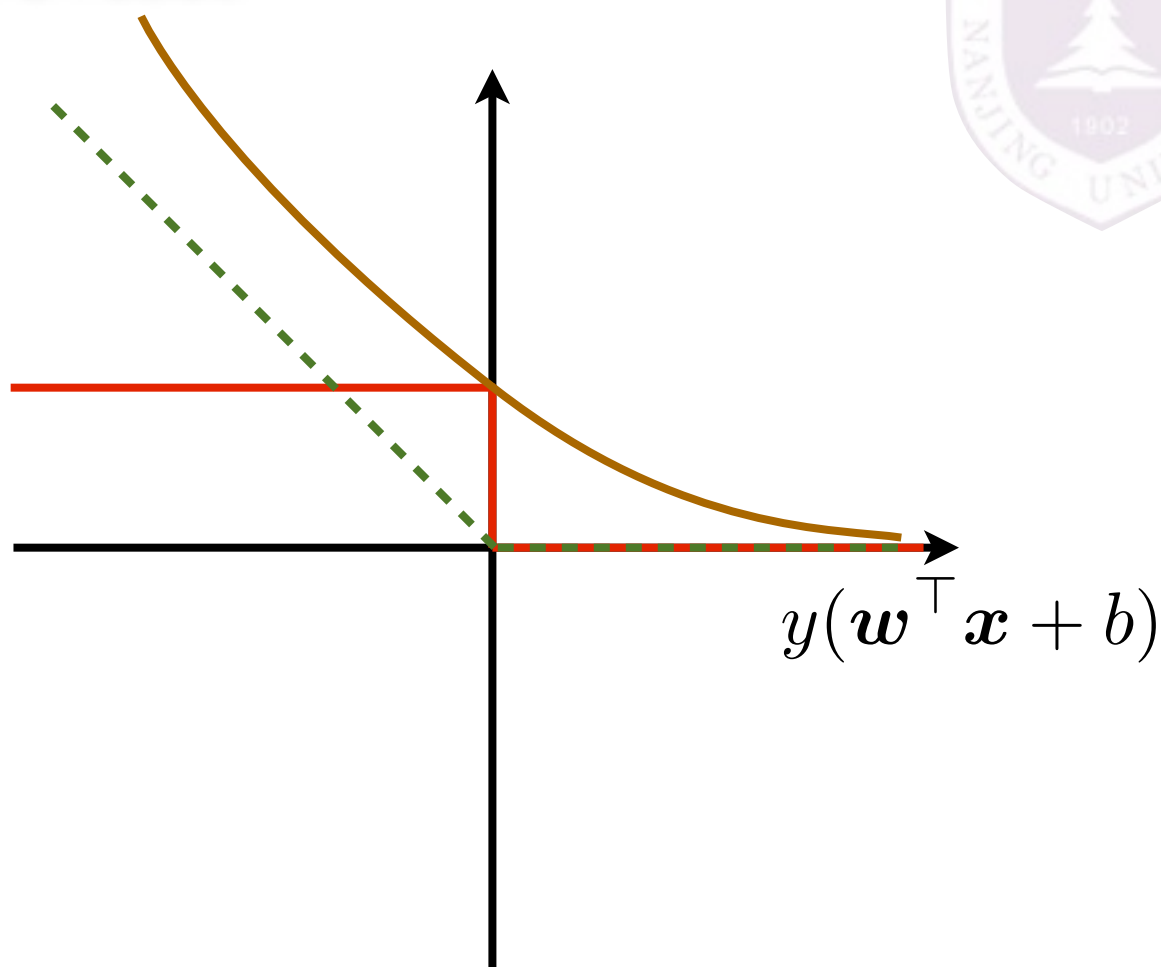
$$I(y(\mathbf{w}^\top \mathbf{x} + b) \leq 0)$$

logistic regression

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

perceptron

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



Linear classifier revisit



0-1 loss

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

logistic regression

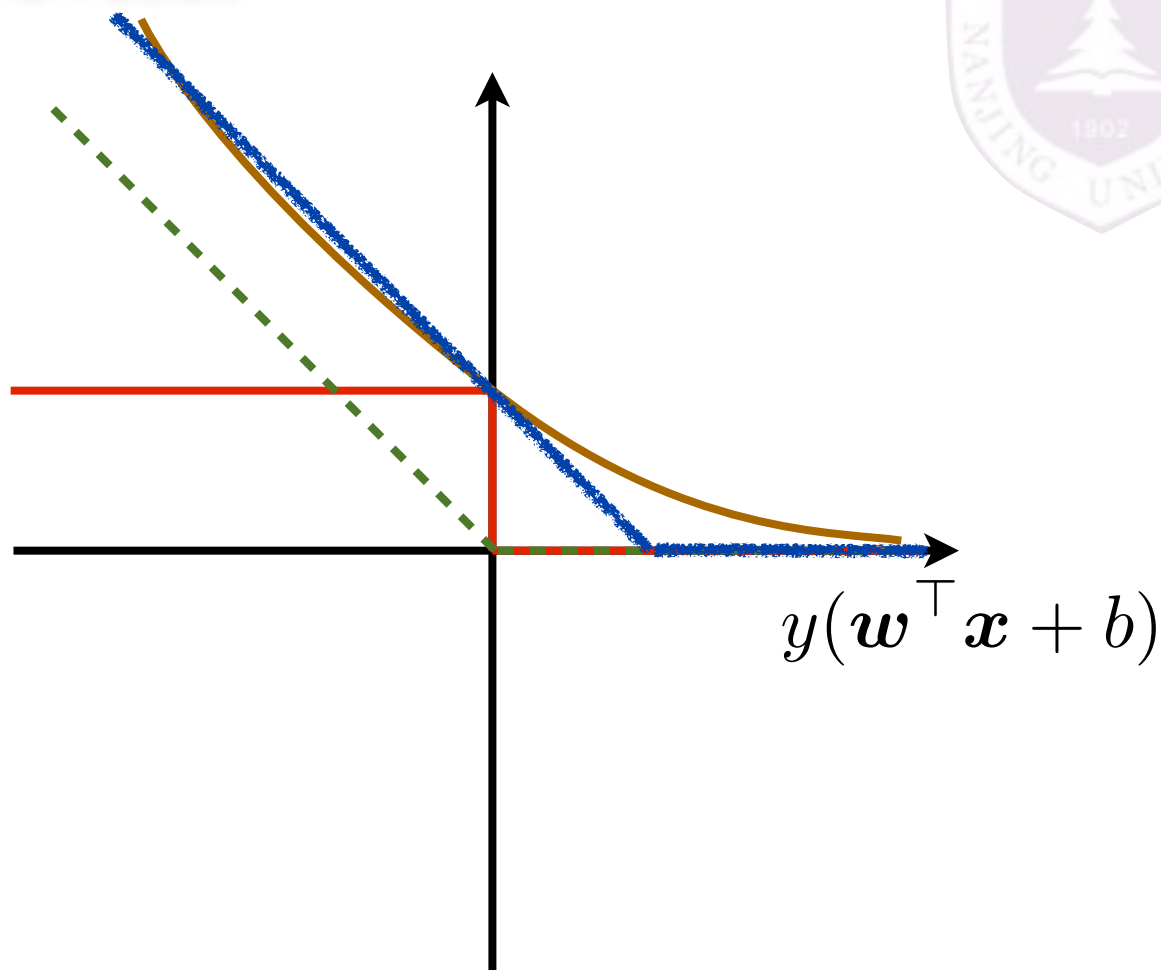
$$\log_2(1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)})$$

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

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

Support vector machines (SVM)



hinge loss + L2-norm

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

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

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

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

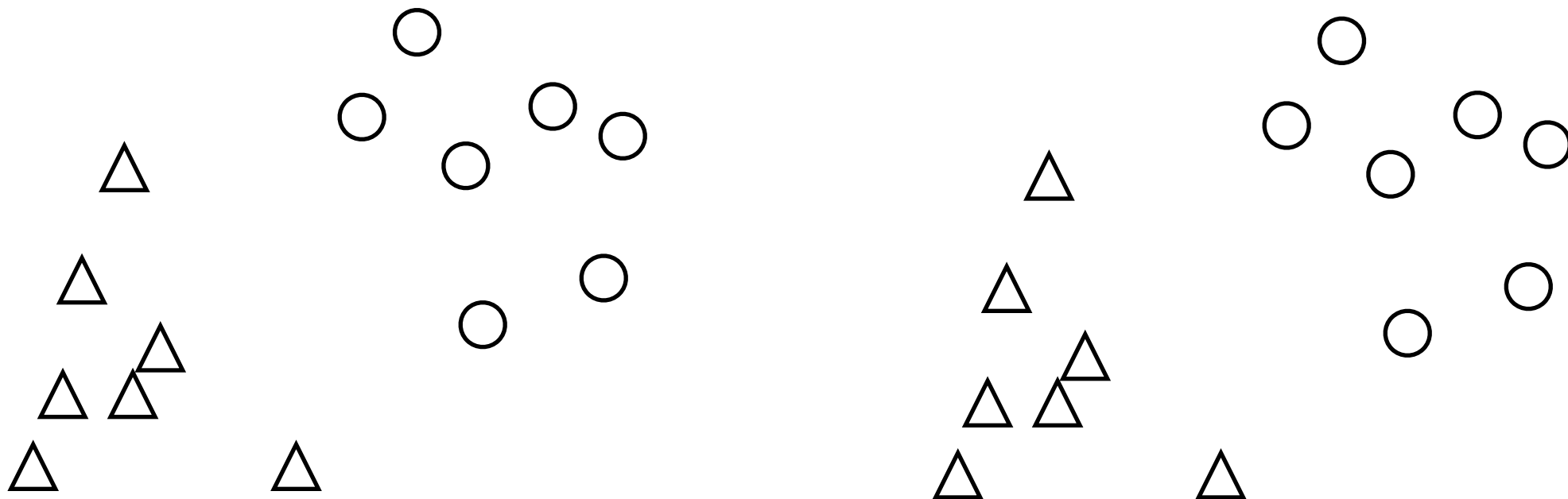
quadratic

Support vector machines (SVM)



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

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

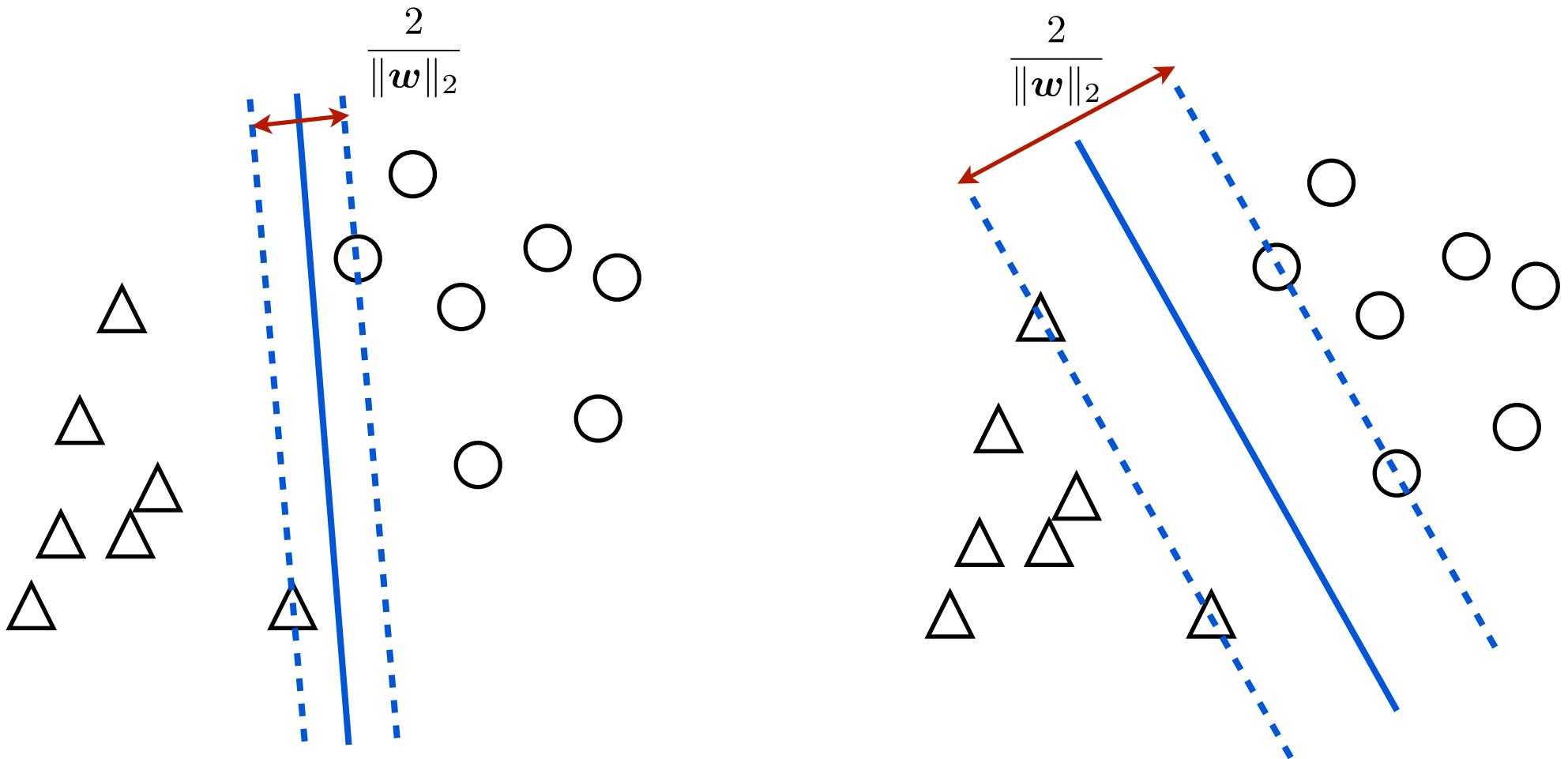


Support vector machines (SVM)



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

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

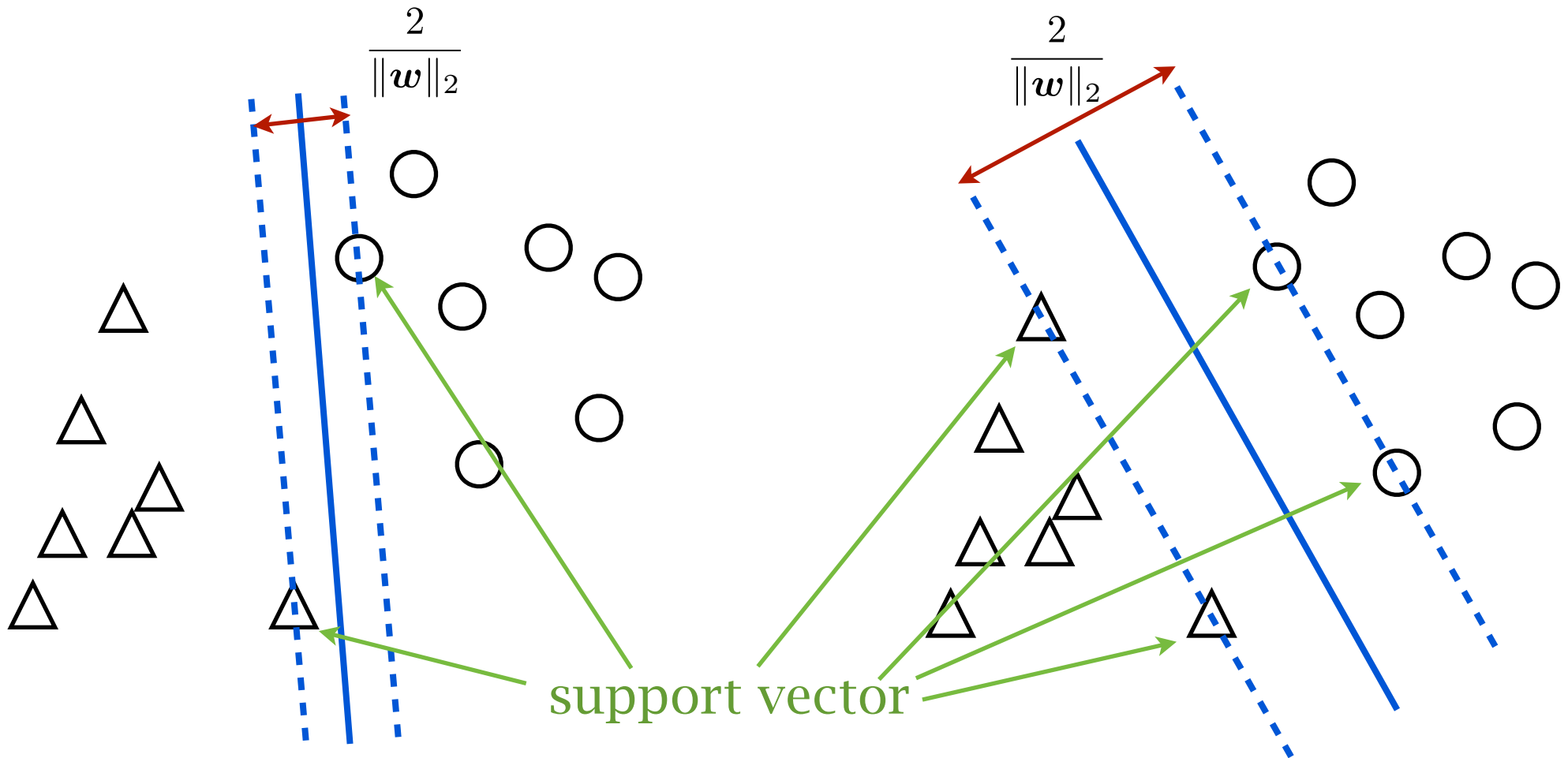


Support vector machines (SVM)



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

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





Support vector machines (SVM)

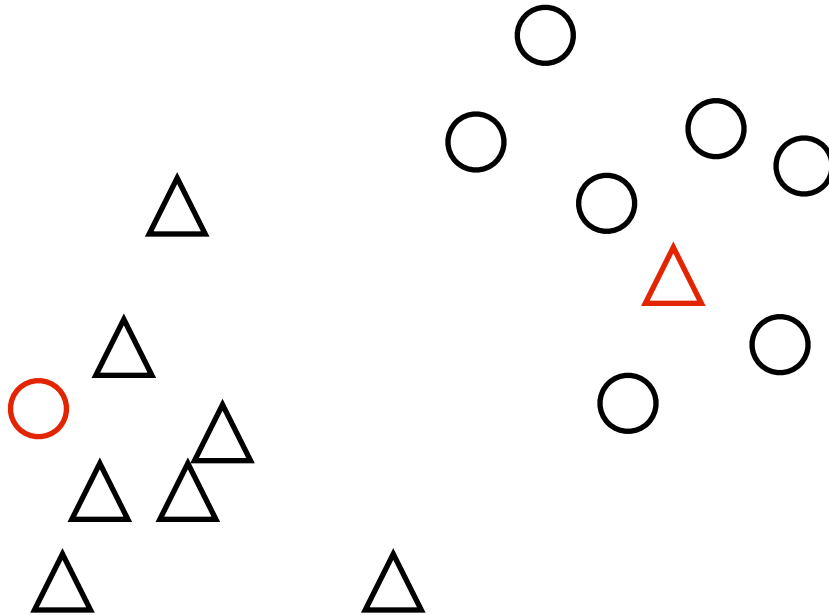
$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2$$

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

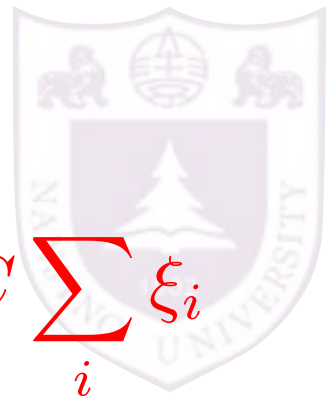
$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2 + C \sum_i \xi_i$$

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

$$\xi_i \geq 0$$



Support vector machines (SVM)



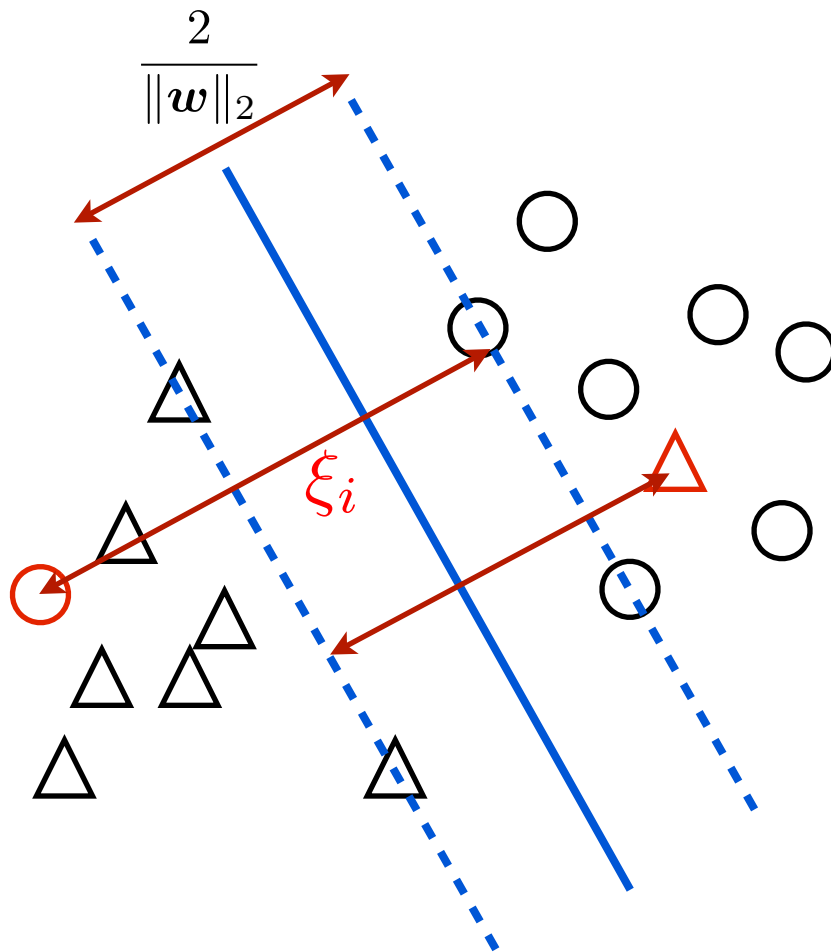
$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2$$

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

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

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

$$\xi_i \geq 0$$



slack variables

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

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

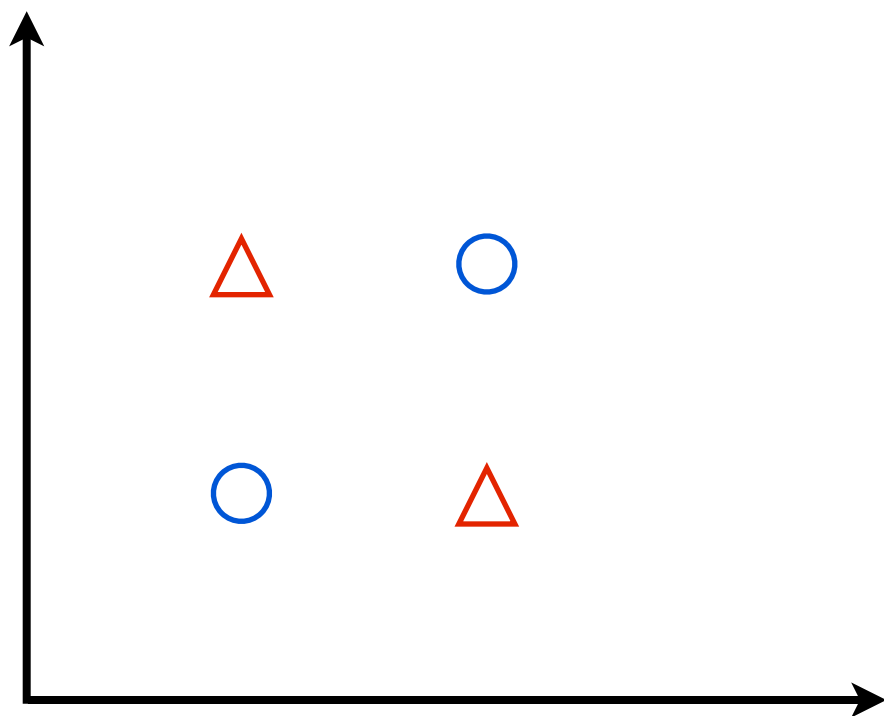
regularized LR

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

SVM

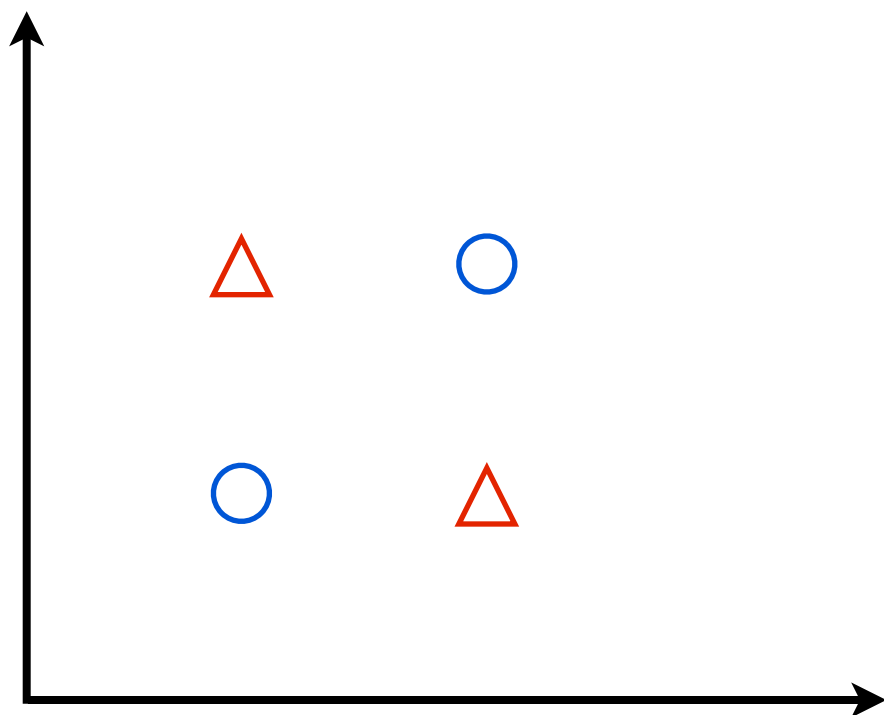
minimize loss + regularization

Linearity v.s. dimensionality

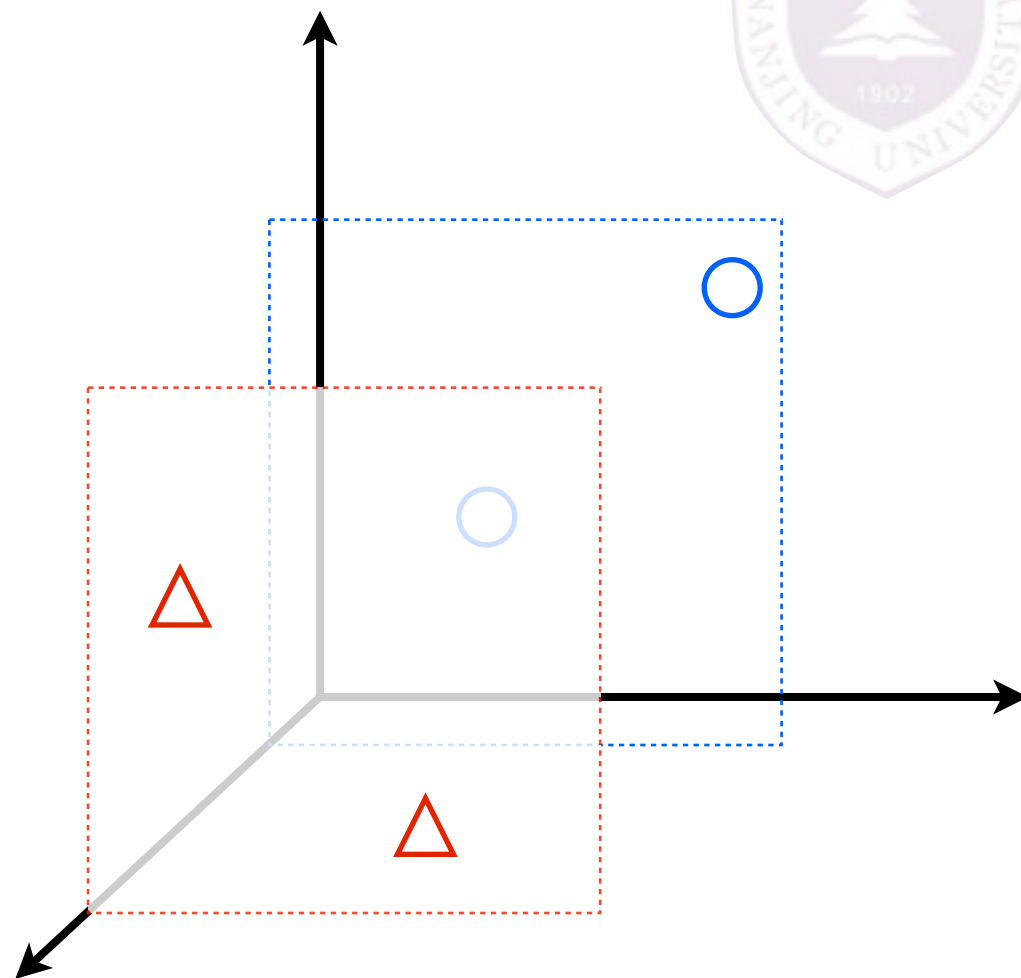


XOR in 2D

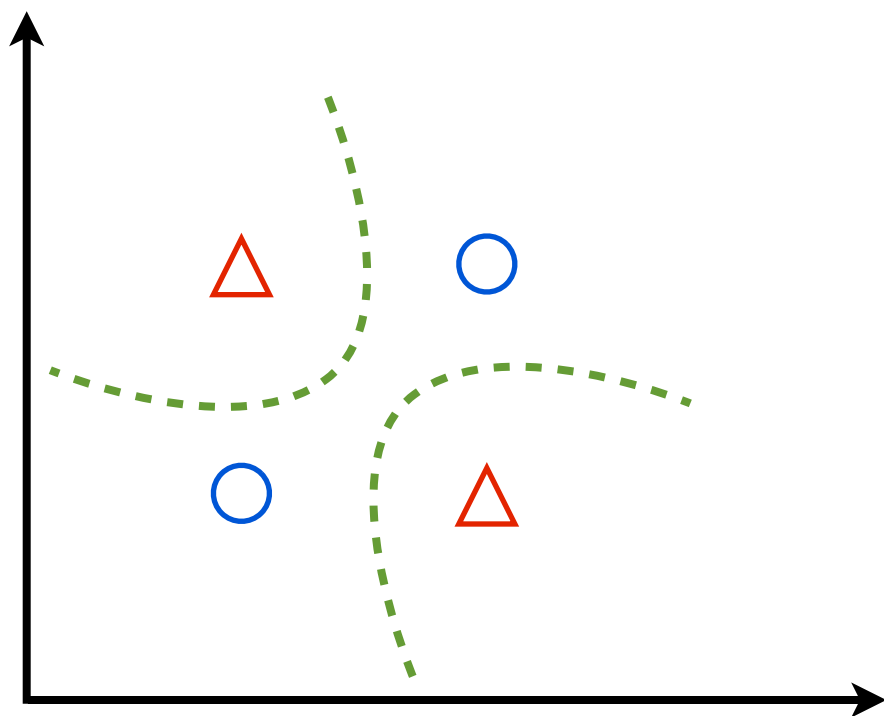
Linearity v.s. dimensionality



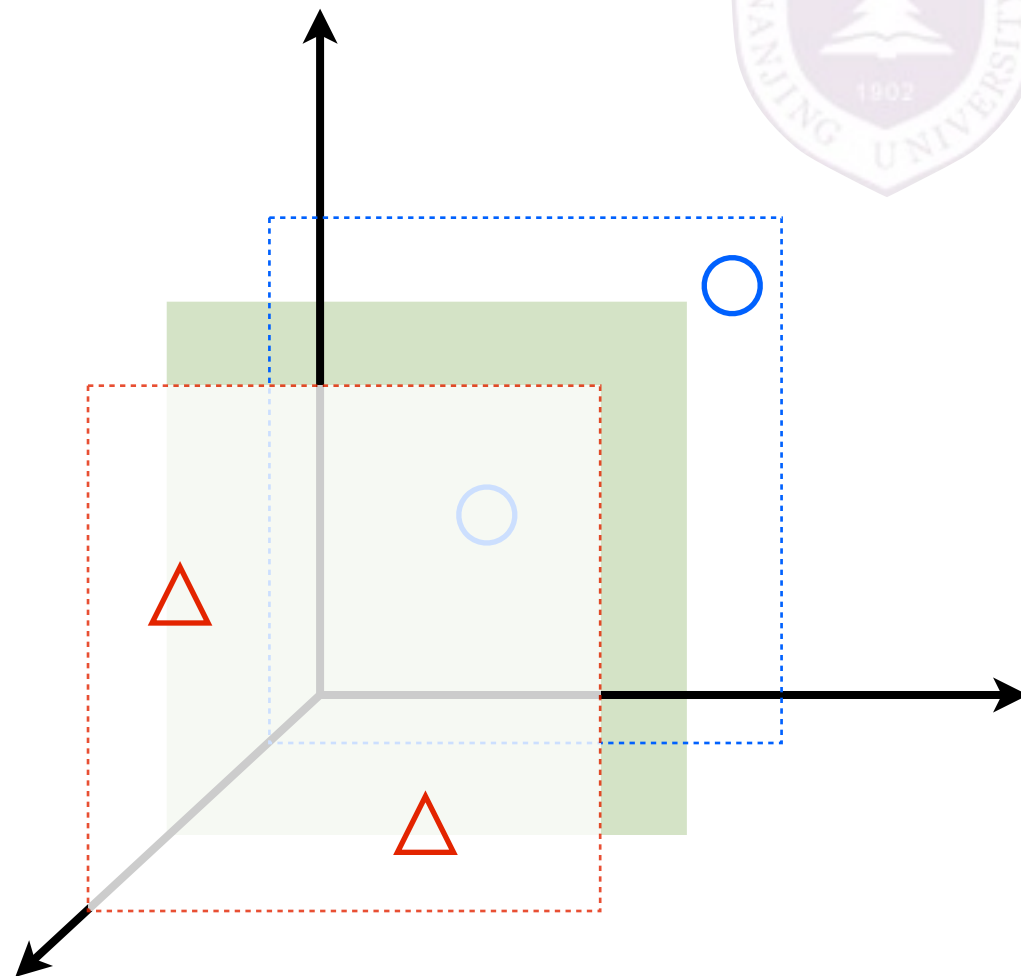
XOR in 2D



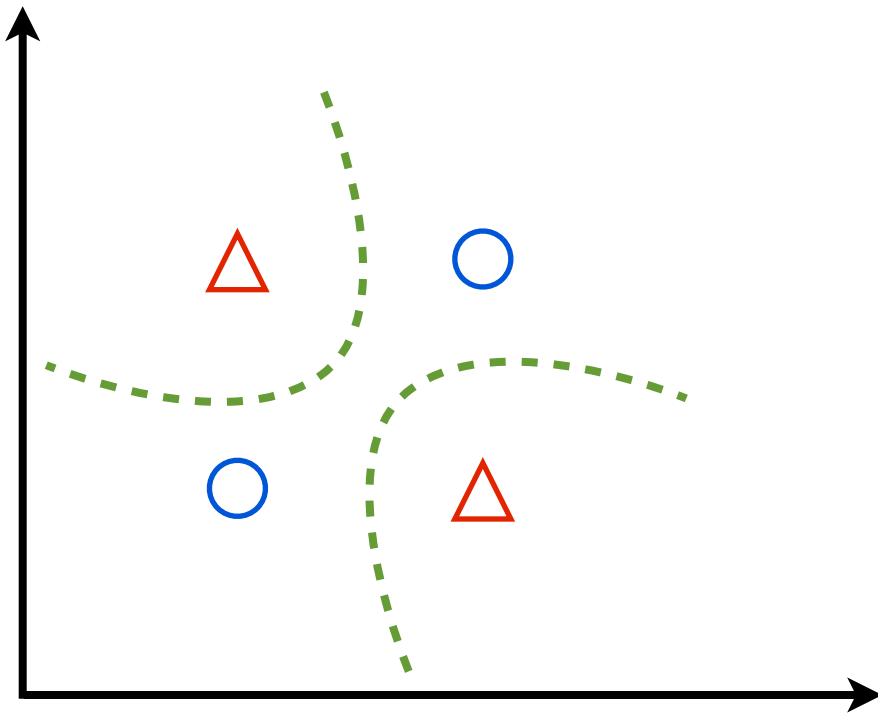
Linearity v.s. dimensionality



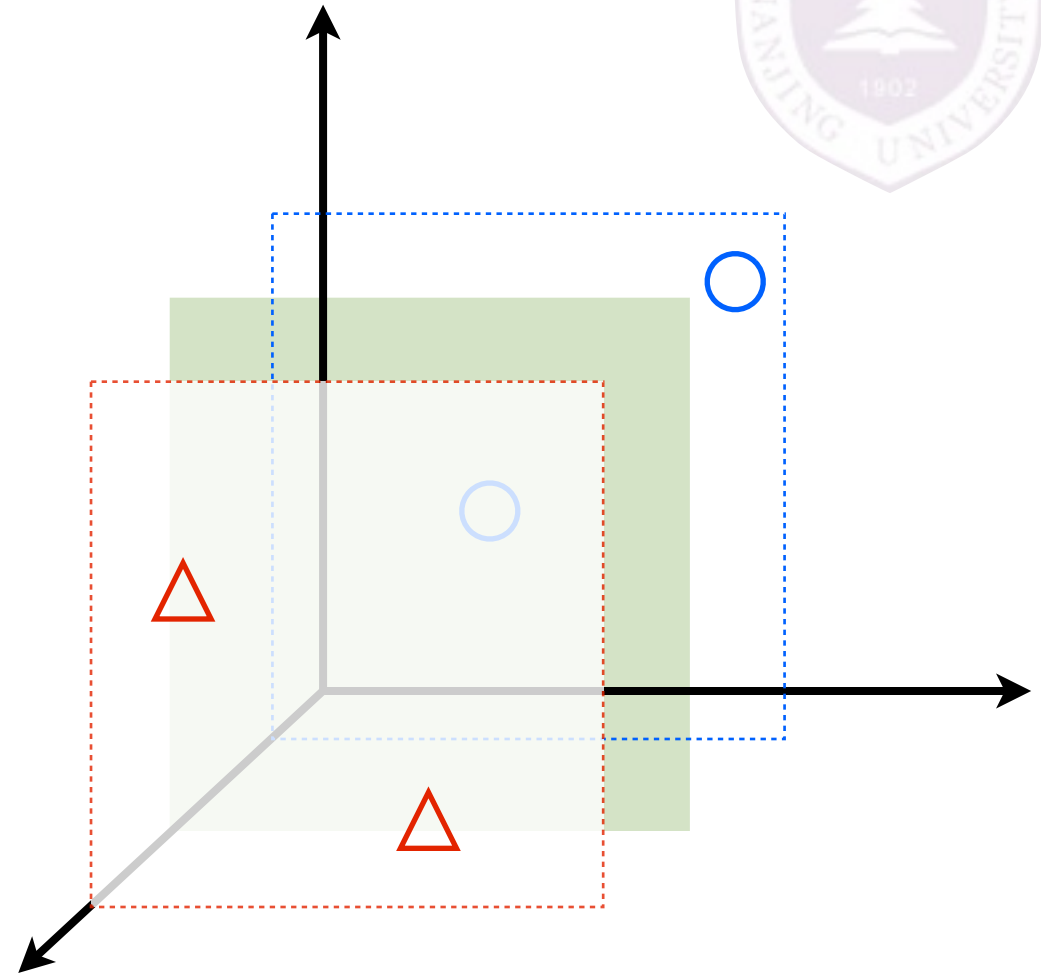
XOR in 2D



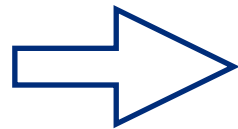
Linearity v.s. dimensionality



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

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



Kernelization

inner product by kernel distance

$$K(\mathbf{x}_1, \mathbf{x}_2) = \langle \phi(\mathbf{x}_1), \phi(\mathbf{x}_2) \rangle$$

polynomial $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^n$

Gaussian radial basis $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{\delta^2}}$

e.g. $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ $\mathbf{x}' = \begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix}$ $\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}$

explicit inner product in higher dimension space:

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = x_1^2x'^2_1 + x_2^2x'^2_2 + 2x_1x_2x'_1x'_2$$

kernel function of the inner product in original space:

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^2 = (x_1x'_1 + x_2x'_2)^2$$

$$= x_1^2x'^2_1 + x_2^2x'^2_2 + 2x_1x_2x'_1x'_2$$

equal

this is easier to calculate

Kernelization



inner product by kernel distance

$$K(\mathbf{x}_1, \mathbf{x}_2) = \langle \phi(\mathbf{x}_1), \phi(\mathbf{x}_2) \rangle$$

polynomial $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^n$

Gaussian radial basis $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{\delta^2}}$

Kernelization



inner product by kernel distance

$$K(\mathbf{x}_1, \mathbf{x}_2) = \langle \phi(\mathbf{x}_1), \phi(\mathbf{x}_2) \rangle$$

polynomial $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^n$

Gaussian radial basis $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{\delta^2}}$

linear model in mapped feature space

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^\top \phi(\mathbf{x}) = \sum_i \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle \\ &= \sum_i \alpha_i K(\mathbf{x}_i, \mathbf{x}) \end{aligned}$$

Kernelization



inner product by kernel distance

$$K(\mathbf{x}_1, \mathbf{x}_2) = \langle \phi(\mathbf{x}_1), \phi(\mathbf{x}_2) \rangle$$

polynomial $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2)^n$

Gaussian radial basis $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{\delta^2}}$

linear model in mapped feature space

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^\top \phi(\mathbf{x}) = \sum_i \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle \\ &= \sum_i \alpha_i K(\mathbf{x}_i, \mathbf{x}) \end{aligned}$$

kernel ridge regression:

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) = Y(K + \lambda I)^{-1} \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}) \\ \dots \\ K(\mathbf{x}_m, \mathbf{x}) \end{pmatrix}$$

Change of basis

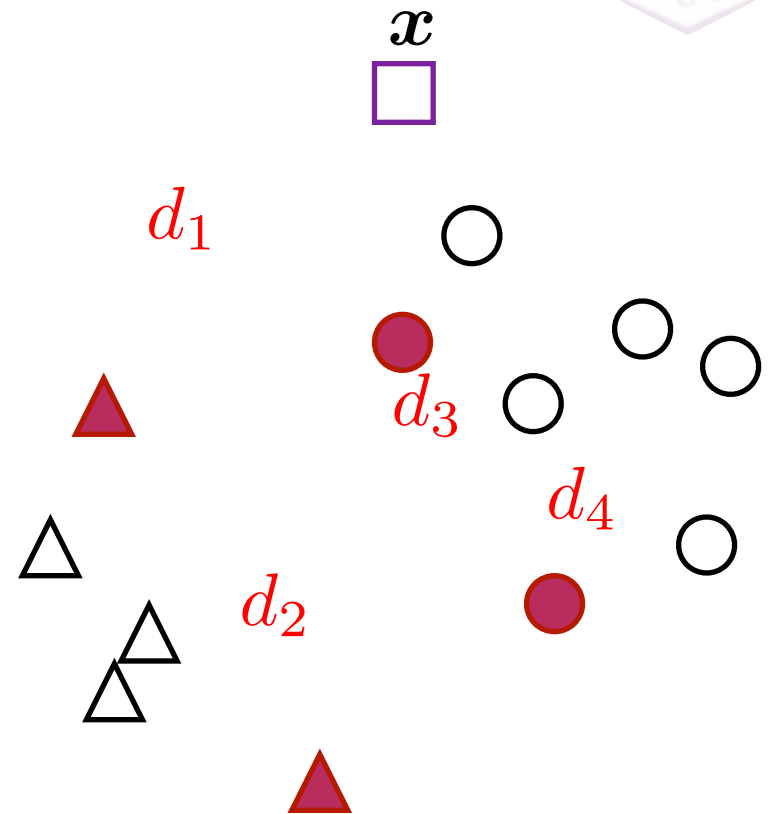
$$d = e^{-\frac{\|\mathbf{x} - \mathbf{x}_2\|_2^2}{\delta^2}}$$

new features:

$$\mathbf{x} \rightarrow \mathbf{z} = (d_1, d_2, d_3, d_4)$$

new decision function:

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \rightarrow f(\mathbf{z}) = \boldsymbol{\alpha}^\top \mathbf{d}$$



Change of basis



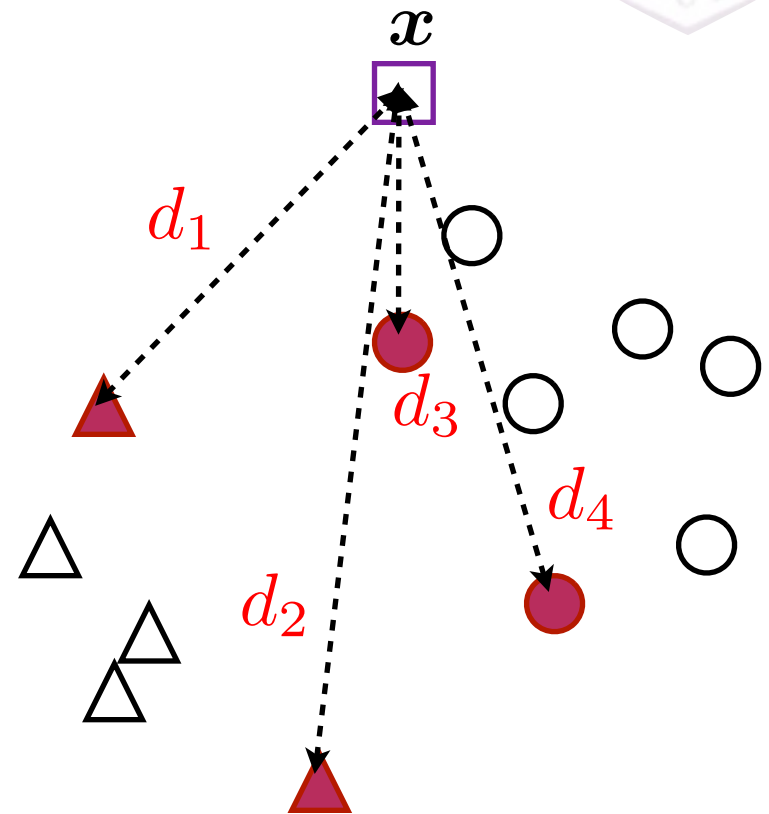
$$d = e^{-\frac{\|\mathbf{x} - \mathbf{x}_2\|_2^2}{\delta^2}}$$

new features:

$$\mathbf{x} \rightarrow \mathbf{z} = (d_1, d_2, d_3, d_4)$$

new decision function:

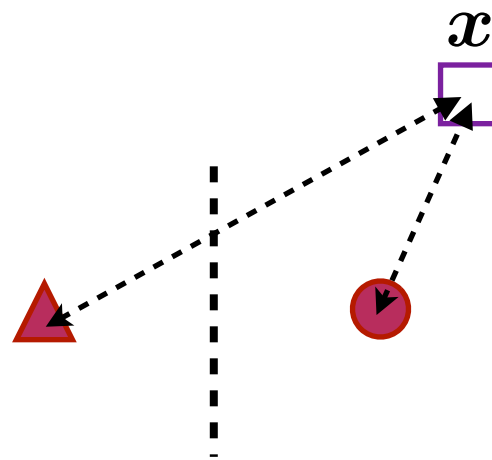
$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \rightarrow f(\mathbf{z}) = \boldsymbol{\alpha}^\top \mathbf{d}$$



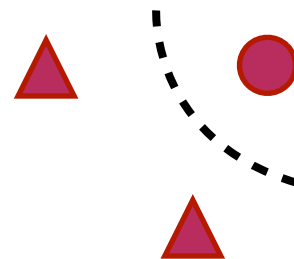
Change of basis



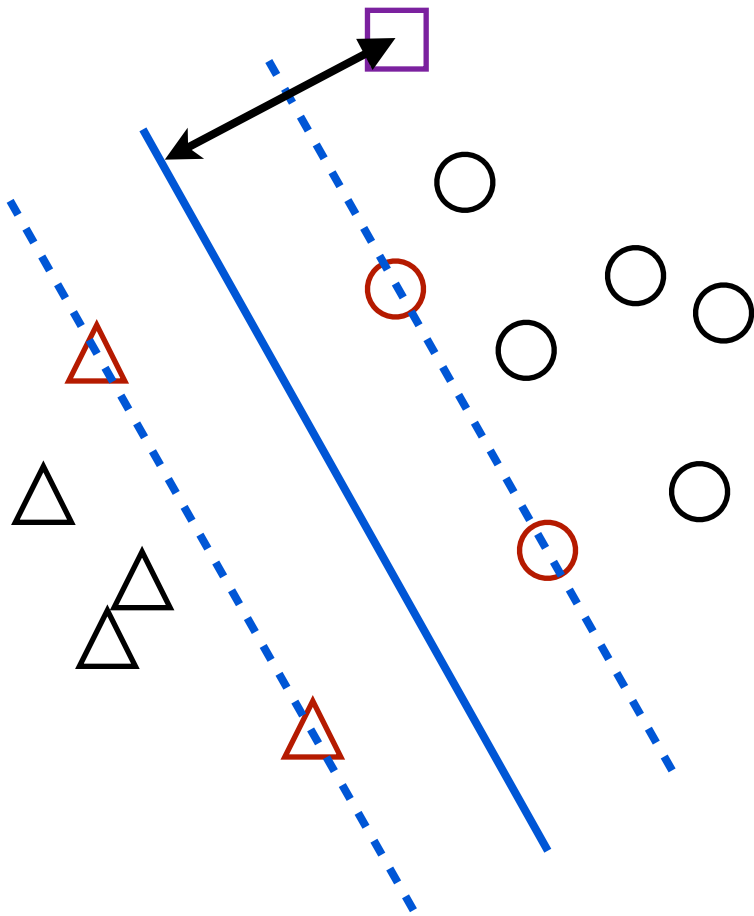
$$f(\mathbf{z}) = \alpha_1 \cdot d_1 + \alpha_2 \cdot d_2 \quad (\alpha \in \mathbb{R})$$



$$f(\mathbf{z}) = \alpha_1 \cdot d_1 + \alpha_2 \cdot d_2 + \alpha_3 \cdot d_3$$



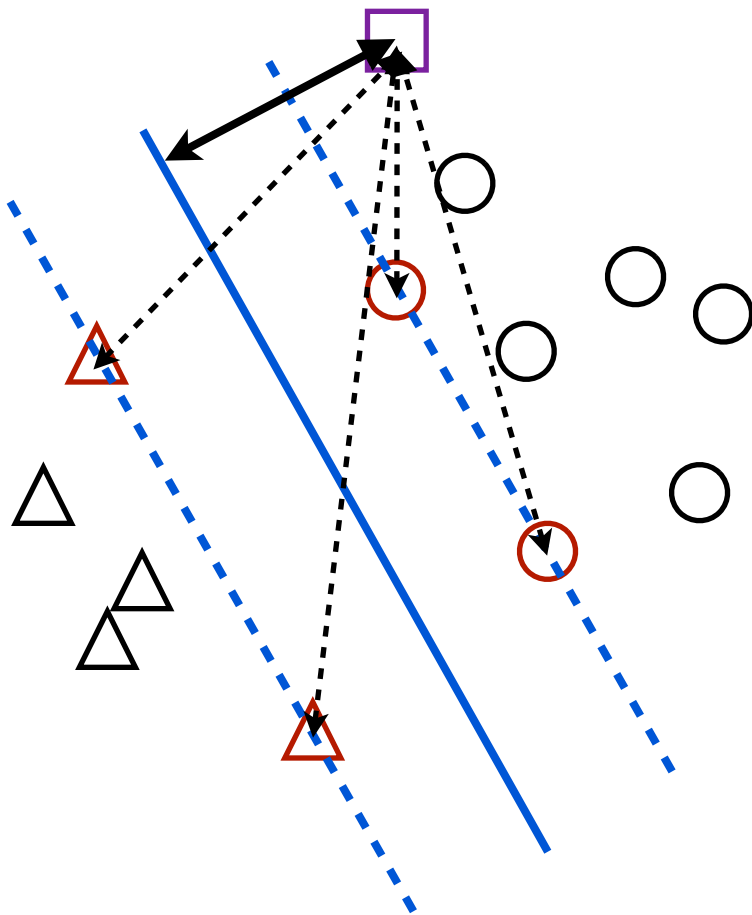
Representer theorem



the optimal function is
in the form of

$$f^*(\cdot) = \sum_i \alpha_i K(\mathbf{x}_i, \cdot)$$

Representer theorem



support vectors

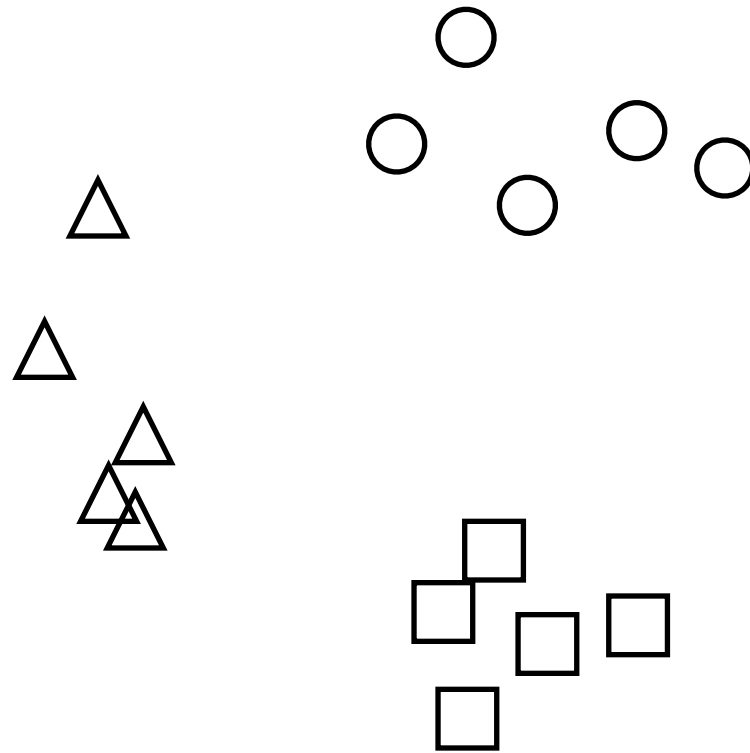
the optimal function is
in the form of

$$f^*(\cdot) = \sum_i \alpha_i K(\mathbf{x}_i, \cdot)$$

Multi-class classification



one-vs-rest

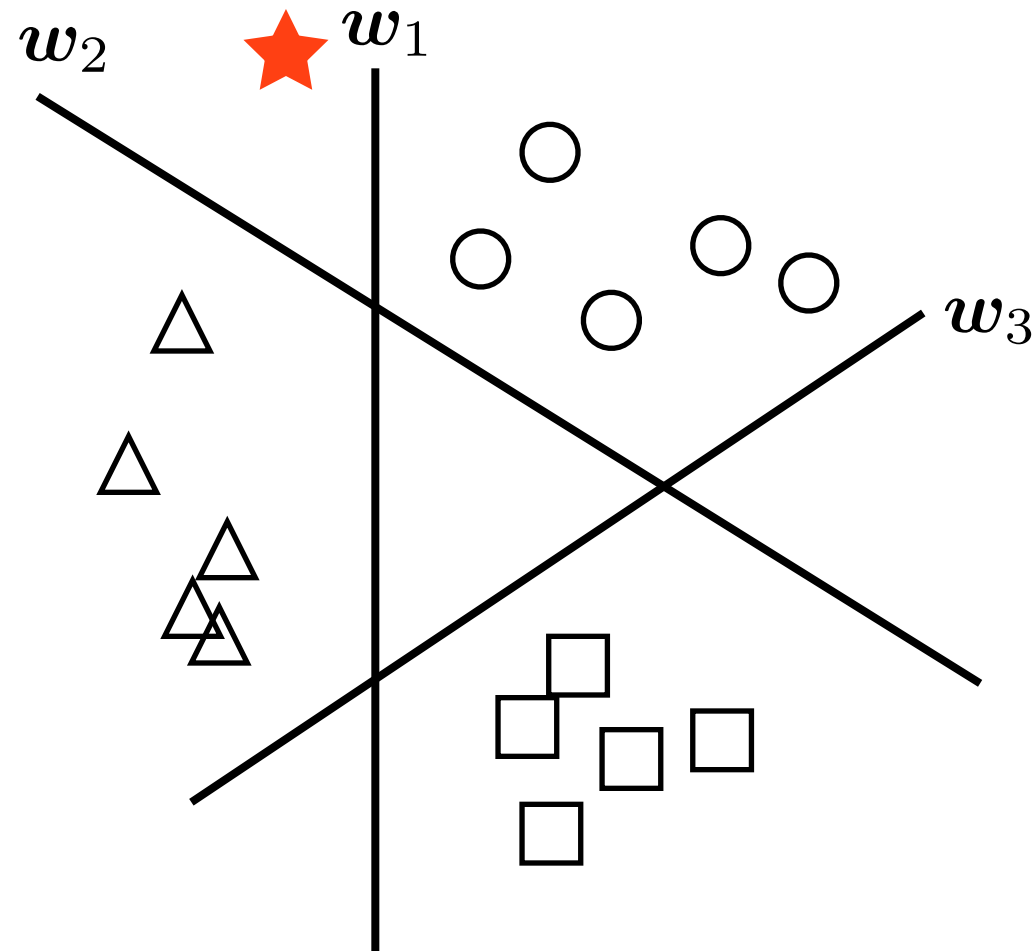


for C classes, need to train C binary classifiers

Multi-class classification



one-vs-rest

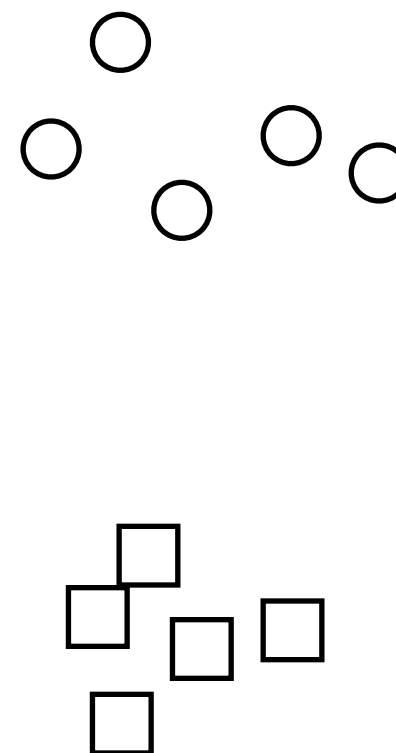
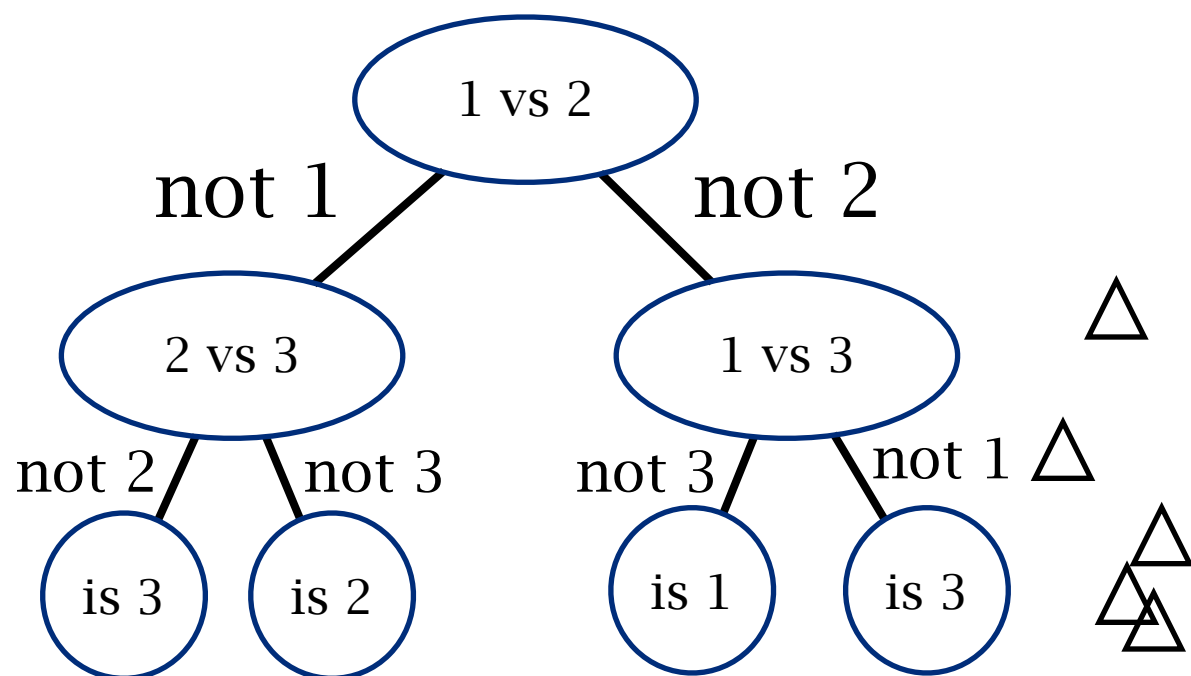


for C classes, need to train C binary classifiers

Multi-class classification



one-vs-one

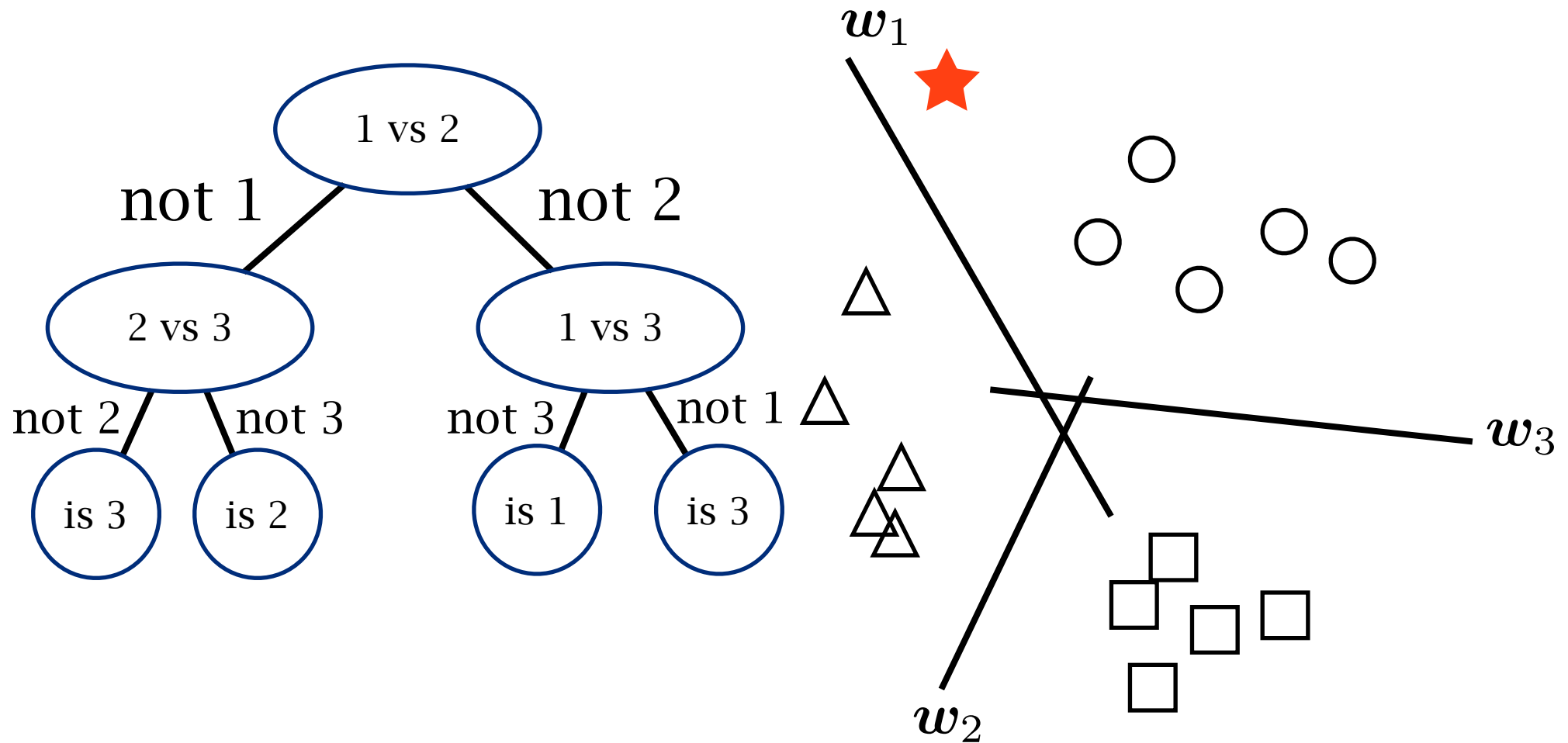


for C classes, need to train $C(C-1)$ binary classifiers

Multi-class classification



one-vs-one



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

习题



L1-norm作为正则化项(regularization)时为何会获得更稀疏(sparse)的解?

Logistic regression是用于回归还是分类?

在低维空间线性不可分的样本是否可以在高维空间线性可分?