

2020.2.18 - 2020.6.2 9:50-12:15@6C102雨课堂+腾讯会议  
《模式识别与机器学习》



## 第二章 线性学习机器

2020.2.18



Xuegong Zhang

多选题 10分

设置

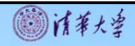
以下陈述正确的是：

- ☒ A 模式识别既指对样本分类的任务，也指实现分类的一些方法
- ☒ B 机器学习是一大类方法的总称，包括模式识别方法
- ☒ C 机器学习除了能够分类外，还可以用于其他任务
- ☐ D 机器学习是把任务分解成规则、由计算机自动执行规则进行决策
- ☒ E 在很多情况下，模式识别与机器学习基本是同义词

提交

2

## 一个简单的例子



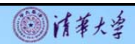
- 如何教小孩从照片识别性别？



Xuegong Zhang

3

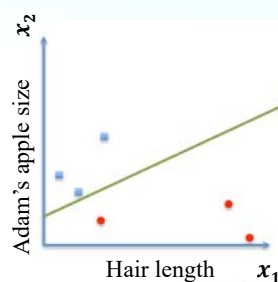
## 一个简单的例子



- 如何教小孩从照片识别性别？



The Adam's apple, or laryngeal prominence, is a feature of the human neck, and is the lump or protrusion that is formed by the angle of the thyroid cartilage surrounding the larynx seen especially in males.



For the Green line:

$$w_1x_1 + w_2x_2 + w_0 = 0$$

For each blue dots:

$$w_1x_1 + w_2x_2 + w_0 > 0$$

For each red dots:

$$w_1x_1 + w_2x_2 + w_0 < 0$$

Xuegong Zhang

4



## 线性分类器

$$g(\mathbf{x}) = \sum_{i=1}^d w_i x_i + w_0$$

特征  
权值

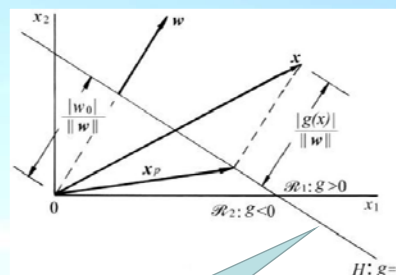
线性判别函数

$$y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$$

决策规则

$$y = \begin{cases} +1 & \Rightarrow \text{class A or } \mathbf{x} \in \omega_1 \\ -1 & \Rightarrow \text{class B or } \mathbf{x} \in \omega_2 \end{cases}$$

分类标签



决策面 (线)

Xuegong Zhang

5

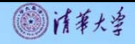


## 2.1 Fisher线性判别(FLD)

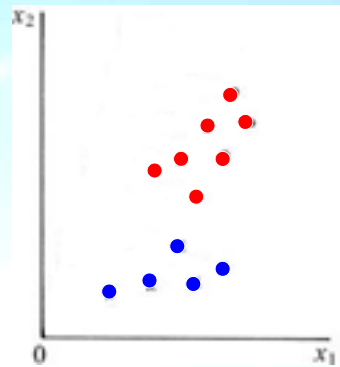
Fisher's Linear Discriminant

a.k.a. Fisher's Linear Discriminant Analysis  
(FLD or LDA)

Xuegong Zhang



- 假设数据分布如图所示，如何求解线性分类器？



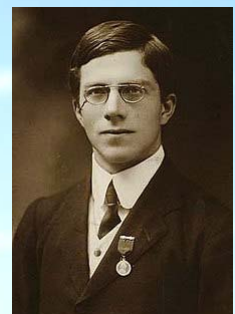
Xuegong Zhang

7

## Sir Ronald Aylmer Fisher (R.A. Fisher)

(17 February 1890 – 29 July 1962)

- British statistician and geneticist **20世纪最伟大的统计学家之一**
  - “a genius who almost single-handedly created the foundations for modern statistical science”
  - “the single most important figure in 20th century statistics”
  - “the greatest of Darwin's successors”.
- Some of the stuff he invented or popularized
  - ANOVA (analysis of variance)
  - Maximum likelihood
  - Fisher's z-distribution (F distribution)
  - Fisher's method for data fusion (meta-analysis)
  - The 0.05 cutoff of p-value, the notion of null hypothesis
  - Fisher's exact test
  - [Fisher's Discriminant Analysis \(in 1936\)](#)
  - .....
  - *The Genetical Theory of Natural Selection* (1930)
  - *The Design of Experiments* (1935)



From Wikipedia

Xuegong Zhang

8

## 不能推断因果关系!

### Sir Ronald Aylmer Fisher (R.A. Fisher)

(17 February 1890 – 29 July 1962)

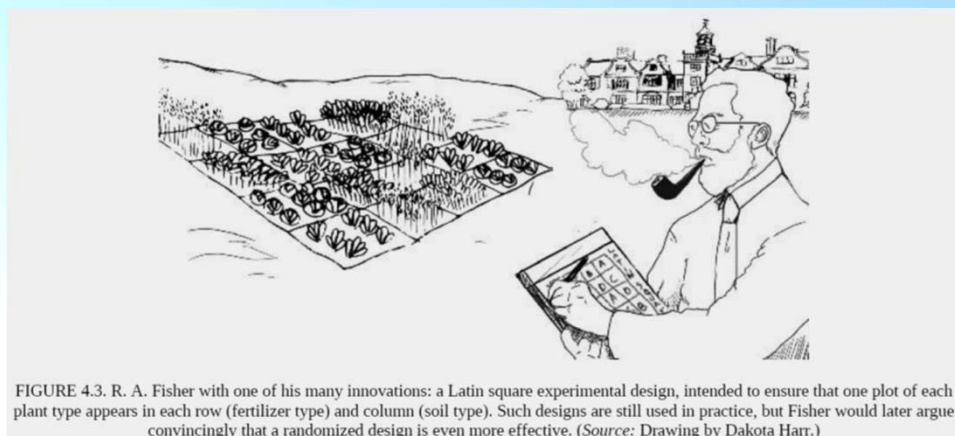


FIGURE 4.3. R. A. Fisher with one of his many innovations: a Latin square experimental design, intended to ensure that one plot of each plant type appears in each row (fertilizer type) and column (soil type). Such designs are still used in practice, but Fisher would later argue convincingly that a randomized design is even more effective. (Source: Drawing by Dakota Harr.)

Judea Pearl, *The Book of Why*, 2018

Suggestion: Find stories about Fisher's role in denying the harm of smoking on health, and learn about the topics of causal inference from the story.

Xuegong Zhang

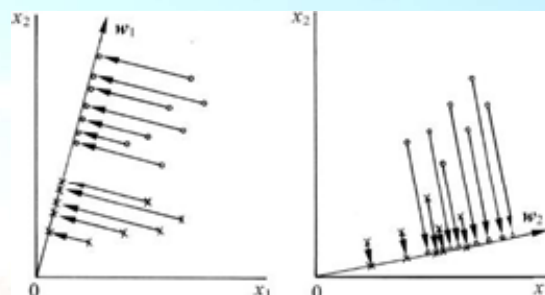
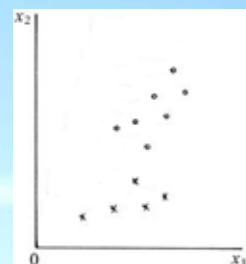
9

### Fisher判别准则



- 求解最佳投影方向，把样本投影到一维上再分类
  - 类内：样本越紧密越好
  - 类间：两类离越远越好

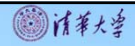
- 样本集：  $\mathcal{X} = \{x_1, \dots, x_N\}$ ,  
 其中，第一类 ( $\omega_1$ ):  $\mathcal{X}_1 = \{x_1^1, \dots, x_{N_1}^1\}$ ,  
 第二类 ( $\omega_2$ ):  $\mathcal{X}_2 = \{x_1^2, \dots, x_{N_2}^2\}$
- 投影：  $\mathcal{X} \rightarrow \mathcal{Y}$ :  $y_i = w^T x_i$ ,  $i = 1, \dots, N$



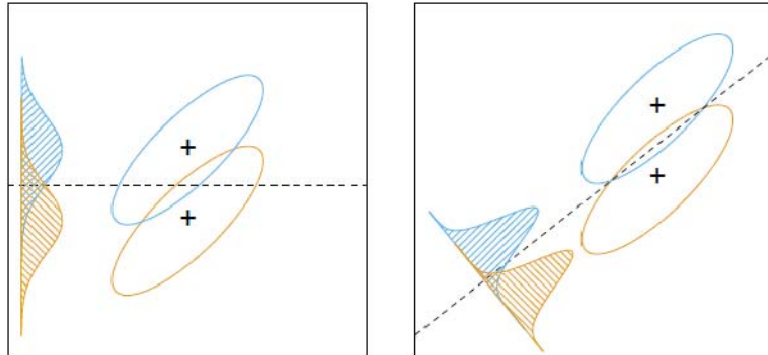
Xuegong Zhang

10





## Fisher判别准则



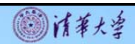
**FIGURE 4.9.** Although the line joining the centroids defines the direction of greatest centroid spread, the projected data overlap because of the covariance (left panel). The discriminant direction minimizes this overlap for Gaussian data (right panel).

Xuegong Zhang

T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd Edition, Springer

11

这里的类内和类间有什么特殊的关系？



## 考查样本的类内和类间离散度

求解最佳投影方向，把样本投影到一维上再分类

- 类内：样本越紧密越好
- 类间：两类离越远越好

### • 在 $\mathcal{X}$ 空间：

类均值向量  $\mathbf{m}_i = \frac{1}{N_i} \sum_{x_j \in \mathcal{X}_i} x_j, i = 1, 2$

类内离散度矩阵 within-class scatter matrix

$$\mathbf{S}_i = \sum_{x_j \in \mathcal{X}_i} (x_j - \mathbf{m}_i)(x_j - \mathbf{m}_i)^T, \quad i = 1, 2$$

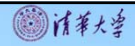
总类内离散度矩阵  $\mathbf{S}_w = \mathbf{S}_1 + \mathbf{S}_2$

类间离散度矩阵 between-class scatter matrix

$$\mathbf{S}_b = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

Xuegong Zhang

12



## 考查样本的类内和类间离散度

- 在 $y$ 空间:

类均值  $\tilde{m}_i = \frac{1}{N_i} \sum_{y_j \in y_i} y_j, \quad i = 1, 2$

类内离散度

$$\tilde{S}_i = \sum_{y_j \in y_i} (y_j - \tilde{m}_i)(y_j - \tilde{m}_i)^T, \quad i = 1, 2$$

总类内离散度  $\tilde{S}_w = \tilde{S}_1 + \tilde{S}_2$

类间离散度矩阵  $\tilde{S}_b = (\tilde{m}_1 - \tilde{m}_2)^2$

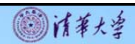
求解最佳投影方向，把样本投影到一维上再分类

- 类内：样本越紧密越好

- 类间：两类离越远越好

Xuegong Zhang

13



## 考查样本的类内和类间离散度

- 在 $y$ 空间:

类均值  $\tilde{m}_i = \frac{1}{N_i} \sum_{y_j \in y_i} y_j, \quad i = 1, 2$

类内离散度

$$\tilde{S}_i = \sum_{y_j \in y_i} (y_j - \tilde{m}_i)(y_j - \tilde{m}_i)^T, \quad i = 1, 2$$

总类内离散度  $\tilde{S}_w = \tilde{S}_1 + \tilde{S}_2$

类间离散度矩阵  $\tilde{S}_b = (\tilde{m}_1 - \tilde{m}_2)^2$

求解最佳投影方向，把样本投影到一维上再分类

- 类内：样本越紧密越好

- 类间：两类离越远越好

- Fisher准则

$$\max_{\mathbf{w}} J_F(\mathbf{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{S}_1 + \tilde{S}_2}$$

$y_i = \mathbf{w}^T \mathbf{x}_i$  这里是 doubleU

Xuegong Zhang

14



- Fisher准则:

$$\max_{\mathbf{w}} J_F(\mathbf{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{S}_1 + \tilde{S}_2}$$

求解最佳投影方向，把样本投影到一维上再分类

- 类内：样本越紧密越好

- 类间：两类离越远越好

- 代入  $y = \mathbf{w}^T \mathbf{x}$ ，得最优投影方向

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} J_F(\mathbf{w})$$

- Fisher准则:

$$J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

存在多解，导致求解困难

思考：最大化  $J_F(\mathbf{w})$  会遇到什么问题？

雨课堂随机点名，  
视频会议中语音回答

Xuegong Zhang

15



$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} J_F(\mathbf{w})$$

$$J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

求解:

- 问题：改变  $\mathbf{w}$  的幅度， $J_F(\mathbf{w})$  不会改变  $\rightarrow$  无唯一解
- 不妨令分母  $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \neq 0$ ，最大化分子  $\mathbf{w}^T \mathbf{S}_b \mathbf{w}$ ，即：

$$\max \mathbf{w}^T \mathbf{S}_b \mathbf{w}$$

$$s. t. \mathbf{w}^T \mathbf{S}_w \mathbf{w} = c$$

拉格朗日乘子法

--- 带有等式约束的优化问题

怎样求解？

雨课堂弹幕抢答

Xuegong Zhang

16





$$\begin{aligned} \max \quad & \mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \end{aligned}$$

- 拉格朗日乘子法求最优:
- 定义拉格朗日函数

$$L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_b \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - c)$$

令  $\frac{\partial L}{\partial \mathbf{w}} = 0$ , 可得

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w}^* = \lambda \mathbf{w}^*$$

只考虑投影的方向, 得

$$\mathbf{w}^* \propto \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

即:  $\mathbf{w}^*$  为  $\mathbf{S}_w^{-1} \mathbf{S}_b$  矩阵的本征向量(eigenvector)

eigen: 自我

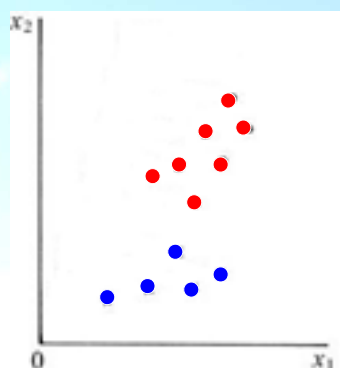


Xuegong Zhang

17



- 不忘初心: 求分类器的目标实现了吗?



请在雨课堂中弹幕回答

Xuegong Zhang

18



$$\mathbf{w}^* \propto \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

- 有了投影方向，还需要确定决策的分界点

$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right) = \text{sgn}(\mathbf{w}^T \mathbf{x} + w_0), \quad y = \begin{cases} +1 & \Rightarrow \mathbf{x} \in \omega_1 \\ -1 & \Rightarrow \mathbf{x} \in \omega_2 \end{cases}$$

- 如何选择  $w_0$ ?

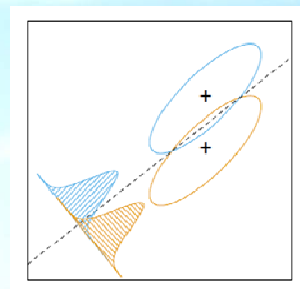
- 根据对数据的不同认识，可以有多种选择方法，比如

$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2)$$

$$w_0 = -\tilde{m}$$

$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2) + \frac{1}{N_1 + N_2 - 2} \ln \frac{P(\omega_1)}{P(\omega_2)}$$

- 可以根据对错误率的要求来选择（见下周课）



Xuegong Zhang

19



- 完美！
- 不过，这算机器学习吗？

**不完全算：**  
学习的成分不够强大，而是完全通过解析的方法



请在雨课堂中弹幕回答



Xuegong Zhang

20



- 完美!
- 不过，这算机器学习吗?



“看”到东西 → 认出东西、产生想法

观察 → 判断

观测 → 分类决策

量化观测 → 类别标签

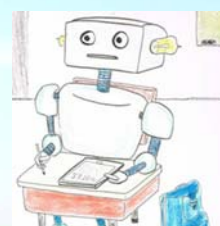
机器

$$x \in R^d \rightarrow y \in \{\omega_1, \omega_2, \dots\}$$

如果机器通过实例学会识别，而不是在程序里写好怎样识别，那就是机器学习。

模式识别

基于数据的机器学习



Xuegong Zhang

21



两种视角看FLD:

- 解析求解线性分类器
- 提取一维最优特征 + 在一维上用阈值分类器

Xuegong Zhang

22



# 休息1分钟



Xuegong Zhang

23

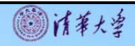


## 2.1 感知器 Perceptron



Xuegong Zhang

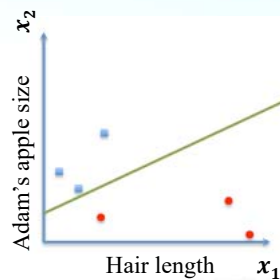
## 一个简单的例子



- 如何教小孩从照片识别性别？



The Adam's apple, or laryngeal prominence, is a feature of the human neck, and is the lump or protrusion that is formed by the angle of the thyroid cartilage surrounding the larynx seen especially in males.



For the Green line:

$$w_1x_1 + w_2x_2 + w_0 = 0$$

For each blue dots:

$$w_1x_1 + w_2x_2 + w_0 > 0$$

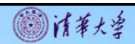
For each red dots:

$$w_1x_1 + w_2x_2 + w_0 < 0$$

Xuegong Zhang

25

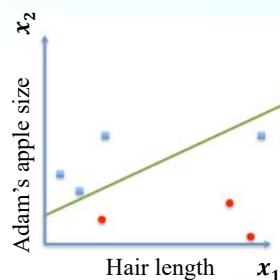
## 一个简单的例子



- 如何教小孩从照片识别性别？



The Adam's apple, or laryngeal prominence, is a feature of the human neck, and is the lump or protrusion that is formed by the angle of the thyroid cartilage surrounding the larynx seen especially in males.



For the Green line:

$$w_1x_1 + w_2x_2 + w_0 = 0$$

For each blue dots:

$$w_1x_1 + w_2x_2 + w_0 > 0$$

For each red dots:

$$w_1x_1 + w_2x_2 + w_0 < 0$$

Xuegong Zhang



26



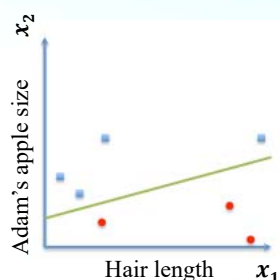


## 一个简单的例子

- 如何教小孩从照片识别性别？



The Adam's apple, or laryngeal prominence, is a feature of the human neck, and is the lump or protrusion that is formed by the angle of the thyroid cartilage surrounding the larynx seen especially in males.



For the Green line:

$$w_1x_1 + w_2x_2 + w_0 = 0$$

For each blue dots:

$$w_1x_1 + w_2x_2 + w_0 > 0$$

For each red dots:

$$w_1x_1 + w_2x_2 + w_0 < 0$$



Xuegong Zhang

27

## 感知器 (Perceptron)



Frank Rosenblatt, *The Perceptron – a perceiving and recognizing automaton*,  
Report 85-460-1, Cornell Aeronautical Laboratory, Jan. 1957

- 感知器：第一台学习机器

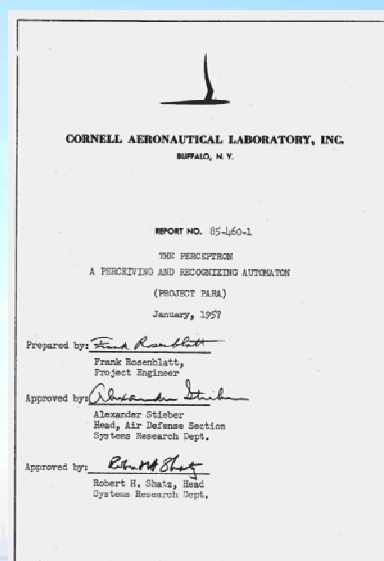
$$y = \text{sgn}\left(\sum_{i=1}^d w_i x_i + w_0\right)$$

- 为什么把它叫做学习机器？



请在雨课堂中弹幕回答

因为在当时就是一台机器！



Xuegong Zhang





## 感知器

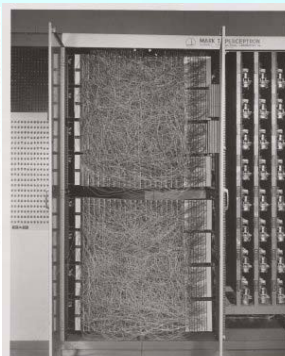
- 为什么把它叫做学习机器？

- ① 因为它是一台机器
- ② 因为它会学习！

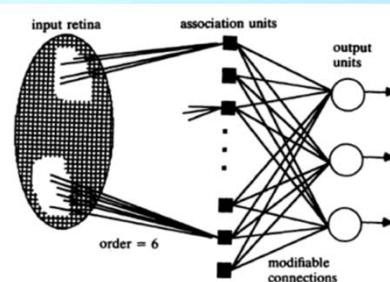
$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right)$$



Xuegong Zhang



<https://en.wikipedia.org/wiki/Perceptron>



M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996

29



## 感知器

- 为什么把它叫做学习机器？

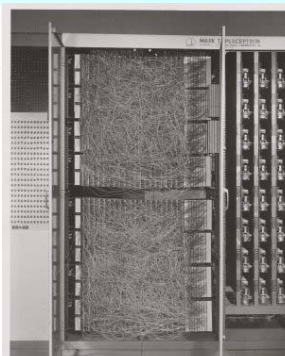
- ① 因为它是一台机器
- ② 因为它会学习！

- 它不是编好程序的冯诺依曼计算机，是一台根据训练数据自我调整的学习机器

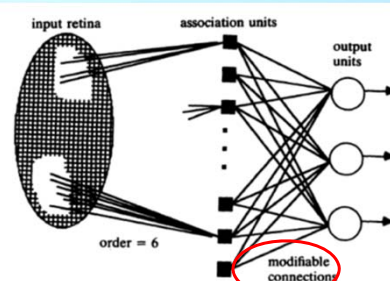
$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right)$$



Xuegong Zhang



<https://en.wikipedia.org/wiki/Perceptron>



M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996

30

## 感知器



- 用数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  训练线性机器

$$y = \text{sgn}\left(\sum_{i=1}^d \textcircled{w_i} x_i + \textcircled{w_0}\right)$$

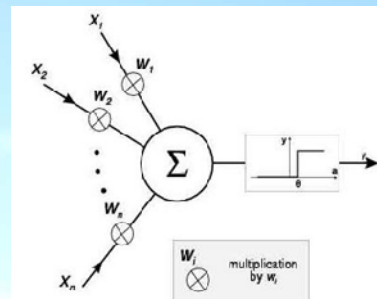
- 目标：最优化目标函数  $J(w)$

$J(w) = \text{训练错误数}$

- 学习算法

- 梯度下降法

do  $w(t+1) = w(t) - \eta \nabla J$   
until  $(\nabla J < \text{threshold})$



Xuegong Zhang

31

## 感知器



- 用数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  训练线性机器

$$y = \text{sgn}\left(\sum_{i=1}^d \textcircled{w_i} x_i + \textcircled{w_0}\right)$$

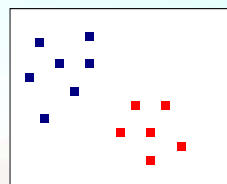
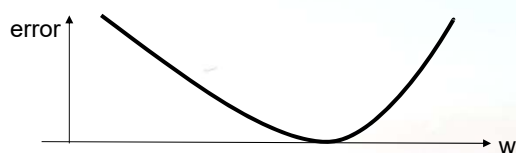
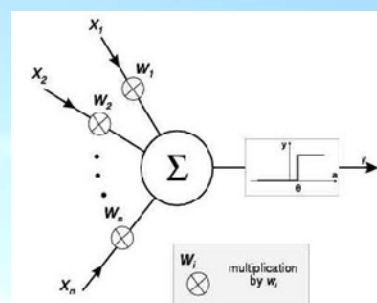
- 目标：最优化目标函数  $J(w)$

$J(w) = \text{训练错误数}$

- 学习算法

- 梯度下降法

do  $w(t+1) = w(t) - \eta \nabla J$   
until  $(\nabla J < \text{threshold})$



Xuegong Zhang

32



## 感知器

- 用数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  训练线性机器

$$y = \text{sgn}\left(\sum_{i=1}^d \textcircled{w_i} x_i + \textcircled{w_0}\right)$$

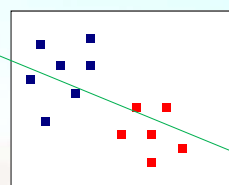
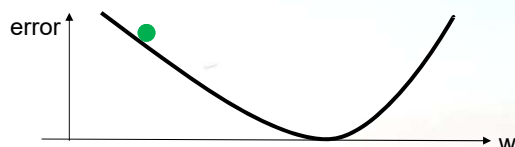
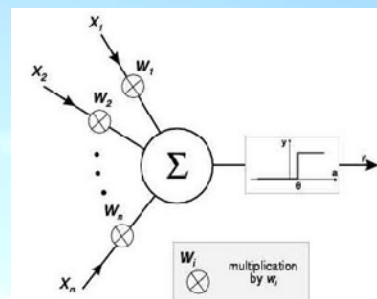
- 目标：最优化目标函数  $J(w)$

$J(w)$  = 训练错误数

- 学习算法

- 梯度下降法

do  $w(t+1) = w(t) - \eta \nabla J$   
until  $(\nabla J < \text{threshold})$



Xuegong Zhang

33



## 感知器

- 用数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  训练线性机器

$$y = \text{sgn}\left(\sum_{i=1}^d \textcircled{w_i} x_i + \textcircled{w_0}\right)$$

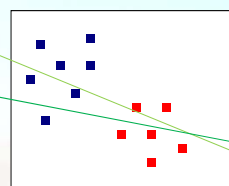
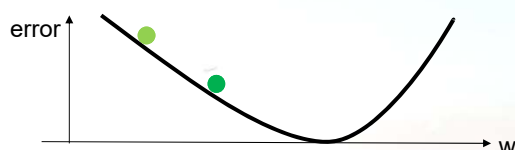
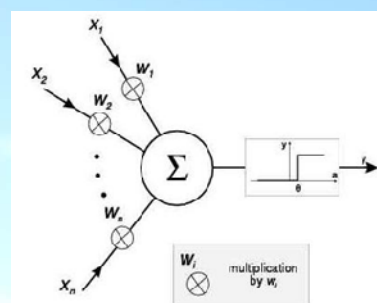
- 目标：最优化目标函数  $J(w)$

$J(w)$  = 训练错误数

- 学习算法

- 梯度下降法

do  $w(t+1) = w(t) - \eta \nabla J$   
until  $(\nabla J < \text{threshold})$



Xuegong Zhang

34

## 感知器



- 用数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  训练线性机器

$$y = \text{sgn}\left(\sum_{i=1}^d \underbrace{w_i}_{\text{red circle}} x_i + \underbrace{w_0}_{\text{red circle}}\right)$$

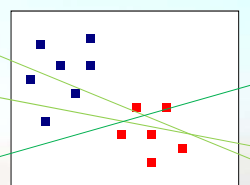
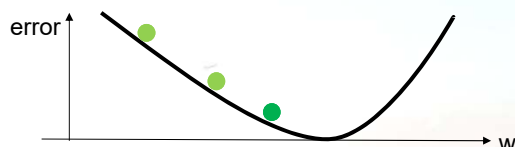
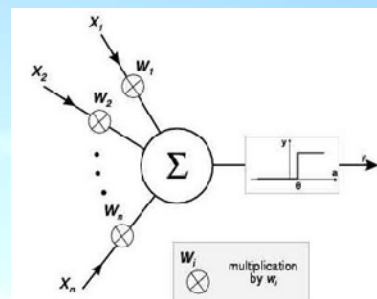
- 目标：最优化目标函数  $J(w)$

$J(w)$  = 训练错误数

- 学习算法

- 梯度下降法

do  $w(t+1) = w(t) - \eta \nabla J$   
until  $(\nabla J < \text{threshold})$



Xuegong Zhang

35

## 感知器



- 用数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  训练线性机器

$$y = \text{sgn}\left(\sum_{i=1}^d \underbrace{w_i}_{\text{red circle}} x_i + \underbrace{w_0}_{\text{red circle}}\right)$$

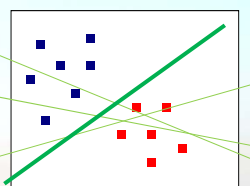
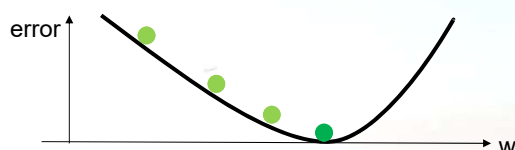
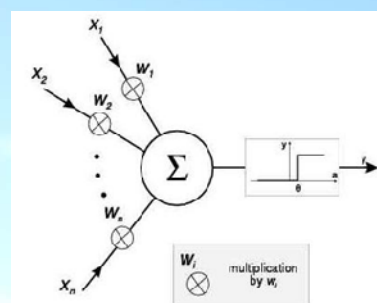
- 目标：最优化目标函数  $J(w)$

$J(w)$  = 训练错误数

- 学习算法

- 梯度下降法

do  $w(t+1) = w(t) - \eta \nabla J$   
until  $(\nabla J < \text{threshold})$



Xuegong Zhang

36



## 机器学习的基本要素

- 怎样造一个学习机器？
  - 它需要老师
    - 我们设计它（特征和模型）
  - 它需要训练/学习材料
    - 训练数据
  - 我们需要为它树立学习的目标
    - 目标函数、学习准则
  - 我们需要告诉它怎样学
    - 学习/训练算法

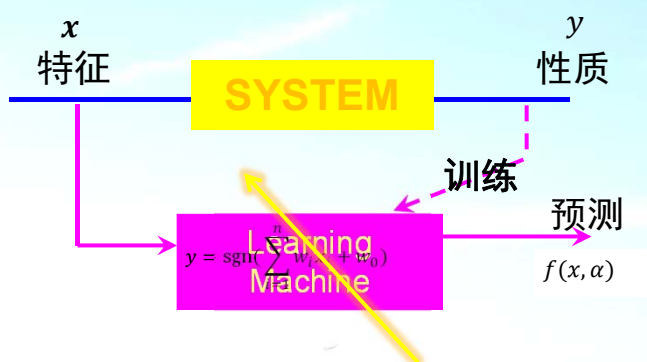


Xuegong Zhang



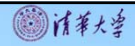
## 监督学习 Supervised Learning

- 用已知答案的数据去训练 → 监督学习



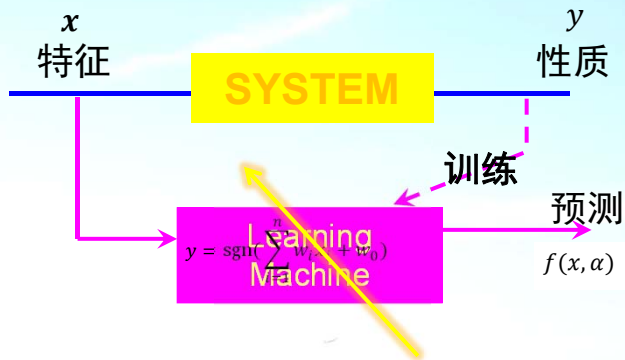
38





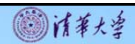
## 监督学习 Supervised Learning

- 用已知答案的数据去训练 → 监督学习



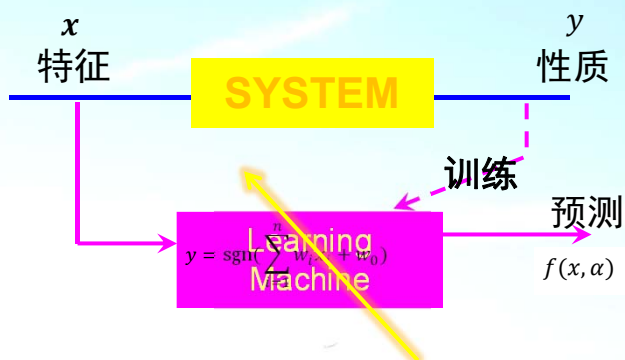
比如：

- $x$  身高体重 →  $y$  性别
- $x$  图像特征 →  $y$  性别
- $x$  图像像素 →  $y$  图像物体
- $x$  司机视觉 →  $y$  方向盘角度
- $x$  音频信号 →  $y$  语音内容
- $x$  基因表达 →  $y$  疾病类型
- $x$  肺CT影像 →  $y$  新冠病毒
- $x$  近期疫情 →  $y$  未来走向
- ...



## 监督学习 Supervised Learning

- 用已知答案的数据去训练 → 监督学习



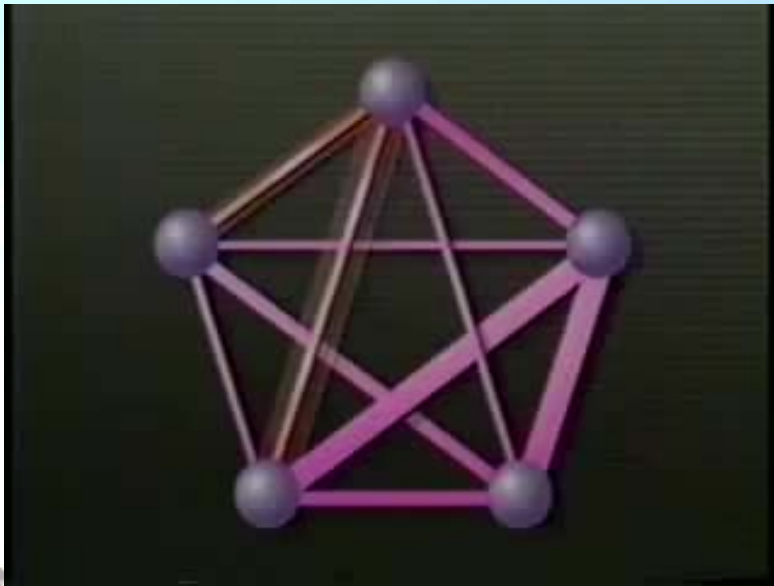
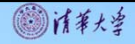
若  $y$  是离散类  
就是模式识别  
→ 监督模式识别

比如：

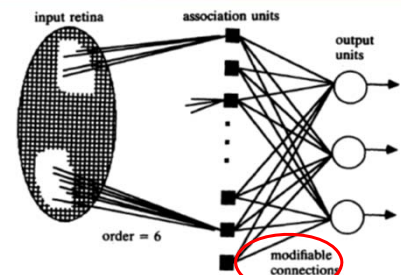
- $x$  身高体重 →  $y$  性别
- $x$  图像特征 →  $y$  性别
- $x$  图像像素 →  $y$  图像物体
- $x$  司机视觉 →  $y$  方向盘角度
- $x$  音频信号 →  $y$  语音内容
- $x$  基因表达 →  $y$  疾病类型
- $x$  肺CT影像 →  $y$  新冠病毒
- $x$  近期疫情 →  $y$  未来走向
- ...



## 感知器是怎么学习的？



$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right)$$



M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996

单选题 1分

设置

休息4分钟，回到座位后请答题

- ☒ A 已回座位
- ☐ B 还没有



Xuegong Zhang

提交

42



## 用什么学习算法?



- 线性判别函数的齐次简化

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \boldsymbol{\alpha}^T \mathbf{y}$$

$\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$  增广的特征向量

$\boldsymbol{\alpha} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$  增广的权向量

Xuegong Zhang

43

## 用什么学习算法?



- 线性判别函数的齐次简化

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \boldsymbol{\alpha}^T \mathbf{y}$$

$\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$  增广的特征向量

$\boldsymbol{\alpha} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$  增广的权向量

为进一步简化推导，把样本向量再进行如下规范化：

$$\mathbf{y}'_i = \begin{cases} \mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_1, \\ -\mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_2, \end{cases} \quad i = 1, \dots, N$$

$\mathbf{y}'_i$ : 规范化增广样本向量，仍记作 $\mathbf{y}_i$

于是，对正确分类样本 $i$ :  $\boldsymbol{\alpha}^T \mathbf{y}_i > 0$

Xuegong Zhang

44



## 用什么学习算法?

- 线性判别函数的齐次简化

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \boldsymbol{\alpha}^T \mathbf{y}$$

$\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$  增广的特征向量

$\boldsymbol{\alpha} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$  增广的权向量

为进一步简化推导, 把样本向量再进行如下规范化:

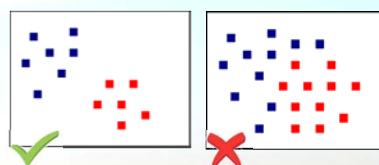
$$\mathbf{y}'_i = \begin{cases} \mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_1, \\ -\mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_2, \end{cases} \quad i = 1, \dots, N$$

$\mathbf{y}'_i$ : 规范化增广样本向量, 仍记作  $\mathbf{y}_i$

于是, 对正确分类样本  $i$ :  $\boldsymbol{\alpha}^T \mathbf{y}_i > 0$

- 线性可分性:

$$\exists \boldsymbol{\alpha}, \quad \boldsymbol{\alpha}^T \mathbf{y}_i > 0, i = 1, \dots, N$$

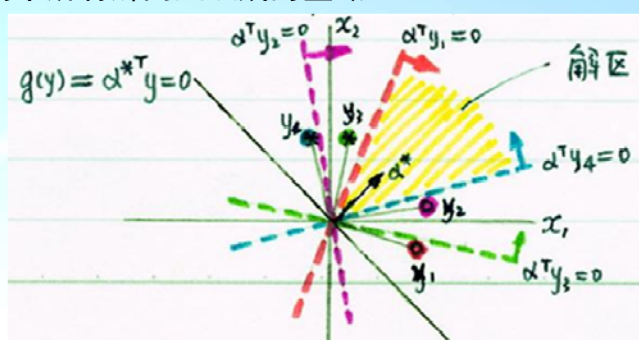


Xuegong Zhang

45



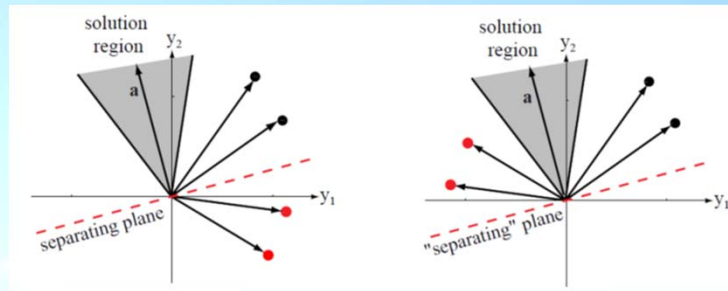
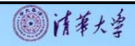
- 解向量  $\boldsymbol{\alpha}^*$ : 满足  $\boldsymbol{\alpha}^T \mathbf{y}_i > 0, i = 1, \dots, N$
- 解区: 权值空间中所有解向量组成的区域



- 对每个样本  $\mathbf{y}_i$ ,  $\boldsymbol{\alpha}^T \mathbf{y}_i = 0$  为权空间中的一个超平面, 解区只可能在超平面的正侧
- 所有样本对应的超平面的正侧的交集就是解区

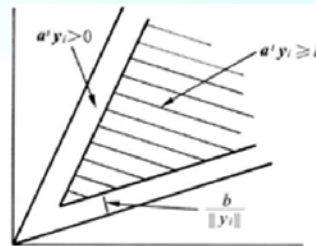
Xuegong Zhang

46



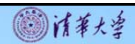
Duda et al, Pattern Classification

- 引入余量  $\alpha^T y_i \geq b > 0$



Xuegong Zhang

47



## 感知器准则函数

$$J_P(\alpha) = \sum_{y_j \in y^k} (-\alpha^T y_j)$$

$y^k$ : 在第  $k$  步被  $\alpha$  错分的样本集合

- 感知器算法 (Rosenblatt, 1957):

$$J_P(\alpha^*) = \min J_P(\alpha) = 0$$

Xuegong Zhang

48



## 感知器准则函数

$$J_P(\alpha) = \sum_{y_j \in Y^k} (-\alpha^T y_j)$$

$Y^k$ : 在第 $k$ 步被 $\alpha$  错分的样本集合

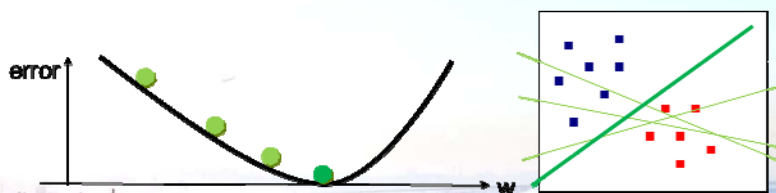
- 感知器算法 (Rosenblatt, 1957):

$$J_P(\alpha^*) = \min J_P(\alpha) = 0$$

- 用梯度下降法 (Gradient descent) 迭代求解

$$\alpha(k+1) = \alpha(k) - \rho_k \nabla J$$

$$\nabla J = \frac{\partial J_P(\alpha)}{\partial \alpha} = \sum_{y_j \in Y^k} (-y_j), \quad \therefore \alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$



Xuegong Zhang

49

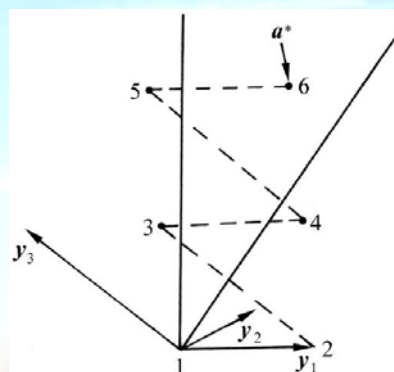


单样本修正:

- 固定增量法:

- ① 初值任意
- ② 对样本 $y_j$ , 若 $\alpha(k)^T y_j \leq 0$  (或 $b$ ), 则 $\alpha(k+1) = \alpha(k) + y_j$
- ③ 对所有样本重复 (2), 直到 $J_P = 0$

$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$



Xuegong Zhang

50



单样本修正：

- 固定增量法：

- ① 初值任意
- ② 对样本  $y_j$ ，若  $\alpha(k)^T y_j \leq 0$  (或  $b$ )，则  $\alpha(k+1) = \alpha(k) + y_j$
- ③ 对所有样本重复 (2)，直到  $J_P = 0$

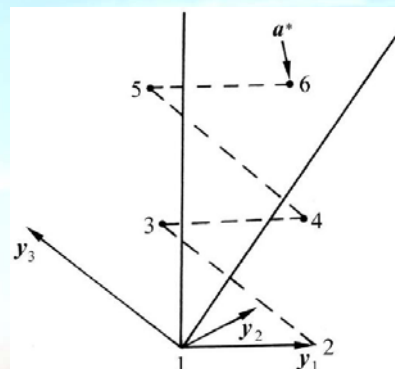
$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$

- 变增量法，如绝对修正法

$$\rho_k = \frac{|\alpha(k)^T y_j|}{\|y_j\|^2}$$

收敛性：

- 对线性可分样本集，  
经过有限次修正后一定可以找到一个解



Xuegang Zhang

51

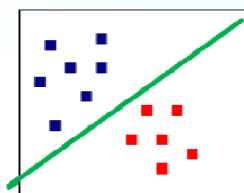


## 讨论

- 感知器有什么问题？

雨课堂随机点名，  
视频会议中语音回答

$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$



Xuegang Zhang

52





## 讨论

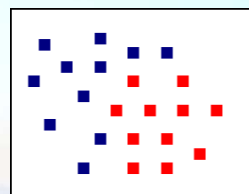
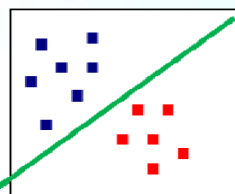
- 感知器有什么问题？
- 有问题怎么办？

– 样本线性不可分呢？

– 线性可分时多解？

– 多类呢？

$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$



Xuegong Zhang

53



## 讨论

- 感知器有什么问题？
- 有问题怎么办？

– 样本线性不可分呢？

– 容忍错误，使错误尽量小

- 比如强制收敛、MSE

– 寻求非线性方法

- 比如神经网络、SVM

– 线性可分时多解？

– 寻求“最优分类器”

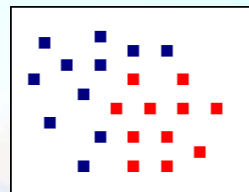
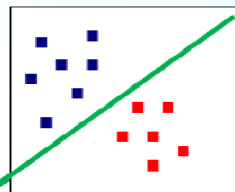
- 比如支持向量机SVM

– 多类呢？

– 多类分类方法

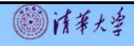
– 用两类分类器完成多类分类

$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$



Xuegong Zhang

54



## 2.3 线性回归 Linear Regression



## 线性分类器

$$\mathbf{y} = f(\mathbf{x})$$

↑  
预测  
↑  
分类器  
↑  
样本

$$\mathcal{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] = \begin{bmatrix} x_{11} & \cdots & x_{m1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \cdots & x_{mn} \end{bmatrix}$$

← 特征  
↑  
样本

$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right) = \text{sgn}(\mathbf{w}^T \mathbf{x} + w_0), \quad y = \begin{cases} +1 & \Rightarrow \mathbf{x} \in \omega_1 \\ -1 & \Rightarrow \mathbf{x} \in \omega_2 \end{cases}$$



## 线性回归

$$\mathbf{y} = f(\mathbf{x}) \quad \mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_m] = \begin{bmatrix} x_{11} & \cdots & x_{m1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \cdots & x_{mn} \end{bmatrix}$$

↑ 预测
↑ 回归
↑ 样本
↑ 样本
← 特征

$$y = \sum_{i=1}^n w_i x_i + w_0 = \mathbf{w}^T \mathbf{x} + w_0$$

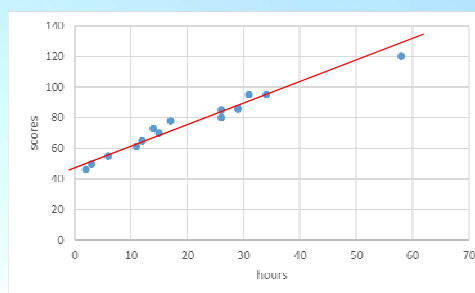
Xuegong Zhang

57



## 简单线性回归

Student id	Final score	Study Hours per Week
1	50	3
2	95	34
3	78	17
4	55	6
5	65	12
6	70	15
7	80	26
8	86	29
9	73	14
10	120	58
11	46	2
12	95	31
13	85	26
14	61	11



Simple Linear Regression

$$y = w_0 + w_1 x$$

Xuegong Zhang

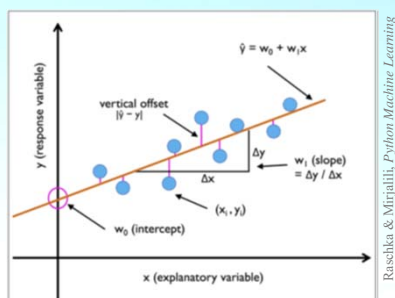
58



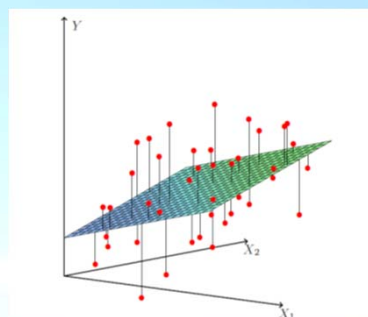
## 多元线性回归

### Simple Linear Regression

$$y = w_0 + w_1 x$$



Raschka & Mirjalili, Python Machine Learning



Hastie, Tibshirani, Friedman, Elements of Statistical Learning

### Multiple Linear Regression

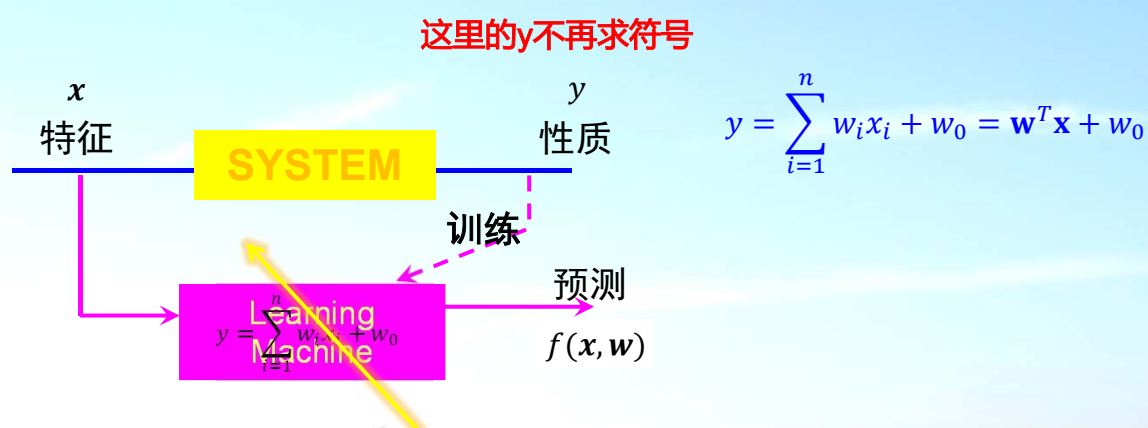
$$y = w_0 + w_1 x + \dots + w_d x_d = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$$

Xuegong Zhang

59



## 线性回归版的机器学习



60



## 机器学习的基本要素

- 怎样造一个学习机器？
  - 它需要老师
    - 我们设计它（特征和模型）
  - 它需要训练/学习材料
    - 训练数据
  - 我们需要为它树立学习的目标
    - 目标函数、学习准则
  - 我们需要告诉它怎样学
    - 学习/训练算法



Xuegong Zhang

61

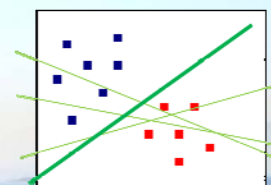


## 机器学习的基本要素：感知器版

- 怎样造一个学习机器？
  - 它需要老师
    - 我们设计它（特征和模型）： 线性判别函数
  - 它需要训练/学习材料
    - 训练数据： 知道分类标签的数据
  - 我们需要为它树立学习的目标
    - 目标函数、学习准则： 错误样本的累计
  - 我们需要告诉它怎样学
    - 学习/训练算法： 梯度下降



雨课堂随机点名，  
视频会议中语音回答



Xuegong Zhang



## 机器学习的基本要素：感知器版

- 怎样造一个学习机器？

- 它需要老师

计算结果 $y$ 的符号

→ 我们设计它（特征和模型）  $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$

- 它需要训练/学习材料

→ 训练数据  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$

- 我们需要为它树立学习的目标

→ 目标函数、学习准则  $\min J_P(\boldsymbol{\alpha}) = \sum_{\mathbf{y}_j \in Y^k} (-\boldsymbol{\alpha}^T \mathbf{y}_j)$

- 我们需要告诉它怎样学

→ 学习/训练算法  $\boldsymbol{\alpha}(k+1) = \boldsymbol{\alpha}(k) - \rho_k \nabla J = \boldsymbol{\alpha}(k) + \rho_k \sum_{\mathbf{y}_j \in Y^k} \mathbf{y}_j$

Xuegong Zhang

63



## 机器学习的基本要素：线性回归版

- 怎样造一个学习机器？

- 它需要老师

直接用 $y$ ，连续版本

→ 我们设计它（特征和模型）  $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

- 它需要训练/学习材料

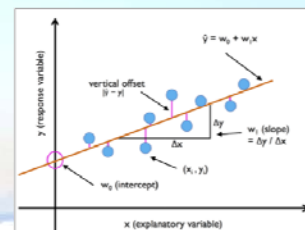
→ 训练数据  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in R^{d+1}, y_j \in R$

- 我们需要为它树立学习的目标

→ 目标函数、学习准则 均方误差

- 我们需要告诉它怎样学

→ 学习/训练算法 ? 最小二乘法



Xuegong Zhang





## 机器学习的基本要素：线性回归版

- 怎样造一个学习机器？
  - 它需要老师
    - 我们设计它（特征和模型）并训练它  $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$
  - 它需要训练/学习材料
    - 训练数据  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in R^{d+1}, y_j \in R$
  - 我们需要为它树立学习的目标
    - 目标函数、学习准则  $\min E = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2$
  - 我们需要告诉它怎样学
    - 学习/训练算法 ?

Xuegong Zhang

65



## 线性回归算法

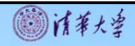
$$\min_{\mathbf{w}} E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2 = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\text{其中 } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}.$$

Xuegong Zhang

66

## 线性回归算法



$$\min_{\mathbf{w}} E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N (f(x_j) - y_j)^2 = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\text{其中 } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}.$$

解:

$$\text{令 } \nabla E(\mathbf{w}) = \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \frac{2}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0,$$

$$\text{有 } \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}.$$

若  $(\mathbf{X}^T \mathbf{X})$  可逆, 则  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .

其中,  $\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  也称作伪逆。

$$\left[ \begin{array}{c} \left[ \begin{array}{c} \phantom{\vdots} \end{array} \right] \left[ \begin{array}{c} \phantom{\vdots} \end{array} \right]^{-1} \left[ \begin{array}{c} \phantom{\vdots} \end{array} \right] \end{array} \right]$$

dim:  $(d+1) \times N \quad N \times (d+1) \quad (d+1) \times N$

Xuegong Zhang

67

### Linear regression algorithm:

- 1: Construct the matrix  $\mathbf{X}$  and the vector  $\mathbf{y}$  from the data set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ , where each  $\mathbf{x}$  includes the  $x_0 = 1$  bias coordinate, as follows

$$\mathbf{X} = \underbrace{\begin{bmatrix} -\mathbf{x}_1^T \\ -\mathbf{x}_2^T \\ \vdots \\ -\mathbf{x}_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

- 2: Compute the pseudo-inverse  $\mathbf{X}^\dagger$  of the matrix  $\mathbf{X}$ . If  $\mathbf{X}^T \mathbf{X}$  is invertible,

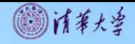
$$\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T.$$

- 3: Return  $\mathbf{w}_{\text{lin}} = \mathbf{X}^\dagger \mathbf{y}$ .

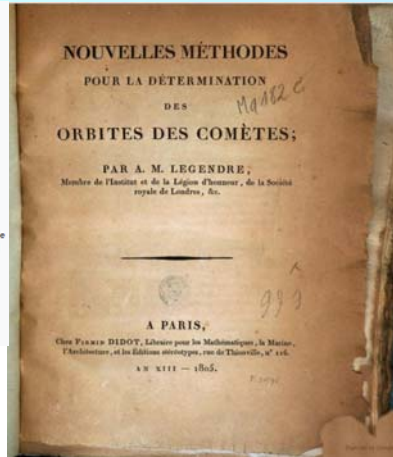
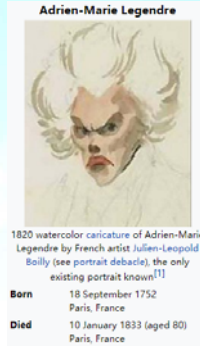
---- Ordinary least squares (OLS) algorithm 最小二乘法

Xuegong Zhang

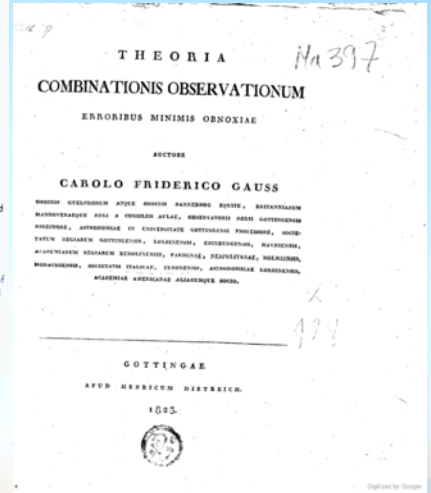
68



A.M. Legendre. Nouvelles méthodes pour la détermination des orbites des comètes, Firmin Didot, Paris, 1805. "Sur la Méthode des moindres carrés" appears as an appendix.



C.F. Gauss. Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientum. (1809)



- Wikipedia

69

Xuegong Zhang

- 完美！
- 不过，这机器学习吗？



请在雨课堂中弹幕回答

- 这解析解也算是“机器学习”？



Xuegong Zhang

70

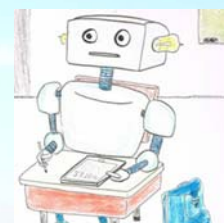


- 完美!
- 不过，这算机器学习吗?



- 这解析解也算是“机器学习”?

- 最小二乘法最早由Legendre在1805和Gauss在1809发明，那时远远没有机器学习的概念
- 如FLD一样，既然算法能从数据中算出可用于预测的规律，通过解析公式算出来人们也倾向于算它是“机器学习”
- 线性回归也可以迭代求解，也可以做成感知器那样学习机器（见下节课）



Xuegong Zhang

71

## 算法本身有没有问题?



请在雨课堂中弹幕提问



### Linear regression algorithm:

- 1: Construct the matrix  $X$  and the vector  $y$  from the data set  $(x_1, y_1), \dots, (x_N, y_N)$ , where each  $x$  includes the  $x_0 = 1$  bias coordinate, as follows

$$X = \underbrace{\begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad y = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

- 2: Compute the pseudo-inverse  $X^\dagger$  of the matrix  $X$ . If  $X^T X$  is invertible,

$$X^\dagger = (X^T X)^{-1} X^T.$$

- 3: Return  $w_{\text{lin}} = X^\dagger y$ .

Xuegong Zhang

72



## 算法本身有没有问题？

- 若  $(X^T X)$  可逆，则  $w^* = (X^T X)^{-1} X^T y$ ；那如果不可逆呢？
- 可逆（非奇异，非退化，满秩）

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} = \begin{bmatrix} & \\ & \\ & \end{bmatrix}_{N \times (d+1)}$$

- $(X^T X)$  可逆： $X$  列满秩：特征间线性独立
  - 当  $N \gg d + 1$  时通常成立
- 当特征不是线性独立时，仍然可以计算伪逆，但解不唯一
- 解决方案：
  - 通过特征选择或变换去除冗余
  - 通过引入其他准则对解加以限制(如SVD或正则化)

Xuegong Zhang

73



## 作业

- 第一章 **该题不计分**
  1. 调研与分析一些系统中的模式识别和机器学习问题  
(注意是分析问题而不是调研具体方法)
- 第二章
  1. 线性判别函数中基本几何关系
  2. Fisher线性判别的证明
  3. 【选做】感知器算法的收敛性证明
  4. 【选做】线性回归与皮尔森相关系数
- **截止日期：2月24日**

Xuegong Zhang

74



