

# 第三章 线性学习机器

2022.3.1



多选题 10分

⚙ 设置

复习：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



# 线性分类器

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

特征  
权值

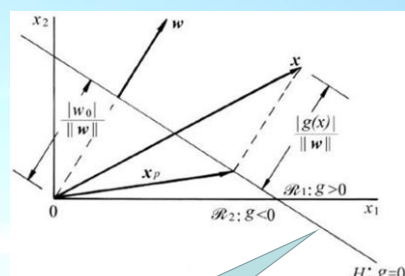
线性判别函数

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

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

3



## 3.1 Fisher线性判别(FLD)

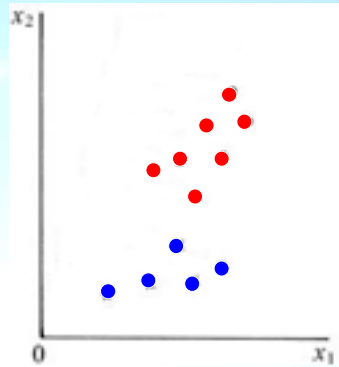
Fisher's Linear Discriminant

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





- 假设数据分布如图所示，如何求解线性分类器？



Xuegong Zhang

5

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

(17 February 1890 – 29 July 1962)

- British statistician and geneticist
  - “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

6

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

(17 February 1890 – 29 July 1962)

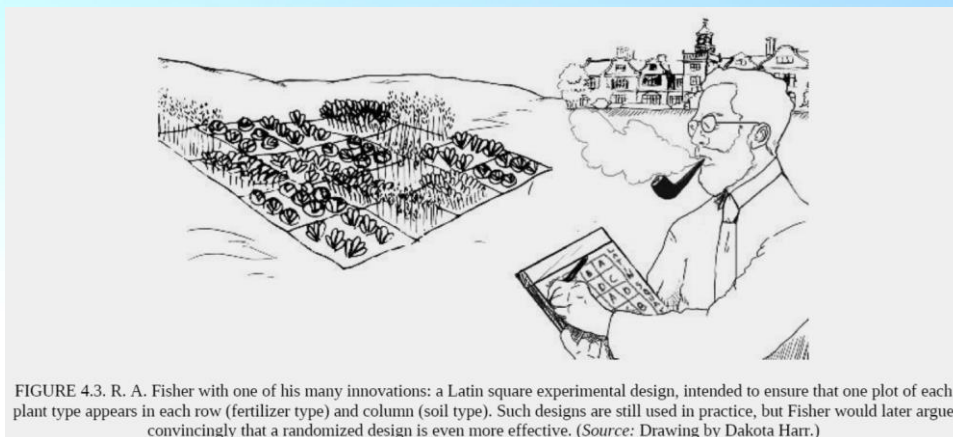


FIGURE 4.3. R. A. Fisher with one of his many innovations: a Latin square experimental design, intended to ensure that one plot of each plant type appears in each row (fertilizer type) and column (soil type). Such designs are still used in practice, but Fisher would later argue convincingly that a randomized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

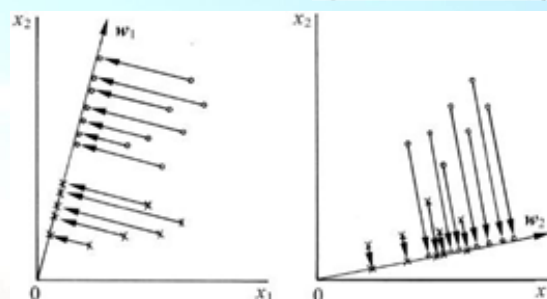
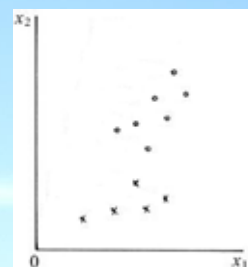
7

## 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

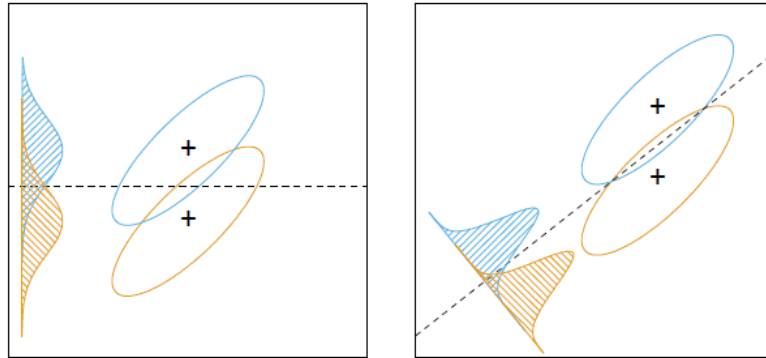
- 样本集：  $\mathcal{X} = \{x_1, \dots, x_N\}$ ,  
 其中，第一类 ( $\omega_1$ ):  $\mathcal{X}_1 = \{x_1^1, \dots, x_{N_1}^1\}$ ,  
 第二类 ( $\omega_2$ ):  $\mathcal{X}_2 = \{x_1^2, \dots, x_{N_2}^2\}$
- 投影：  $\mathcal{X} \rightarrow \mathcal{Y}$ :  $y_i = w^T x_i, i = 1, \dots, N$



Xuegong Zhang

8

## 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

9

## 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  $\mathbf{m}_i = \frac{1}{N_i} \sum_{x_j \in \mathcal{X}_i} \mathbf{x}_j, i = 1, 2$

类内离散度矩阵 within-class scatter matrix

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

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

10





## 考查样本的类内和类间离散度

- 在 $y$ 空间:

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

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

类内离散度

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

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

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

Xuegong Zhang

11



## 考查样本的类内和类间离散度

- 在 $y$ 空间:

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

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

类内离散度

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

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

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

- Fisher准则

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

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

Xuegong Zhang

12



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

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



- Fisher准则:

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

思考：最大化  $J_F(\mathbf{w})$  会遇到什么问题？

Xuegong Zhang

13



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

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

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：

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

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

—— 带有等式约束的优化问题

怎样求解？

Xuegong Zhang

14



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

- 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

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

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

只考虑投影的方向, 得

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

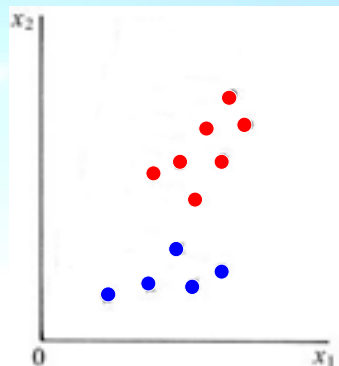


Xuegong Zhang

15



- 不忘初心: 求分类器的目标实现了吗?



Xuegong Zhang

16





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

- 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  $w_0$ ?

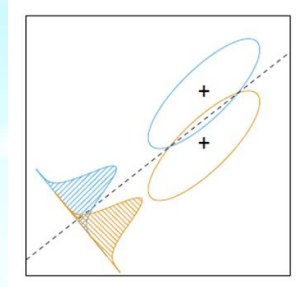
- 根据对数据的不同认识，可以有多种选择方法，比如

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

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

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

- 可以根据对错误率的要求来选择



Xuegang Zhang

17

- Commonly used thresholds:

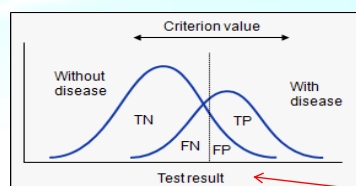
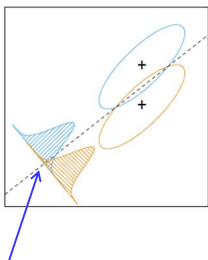
$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2)$$

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

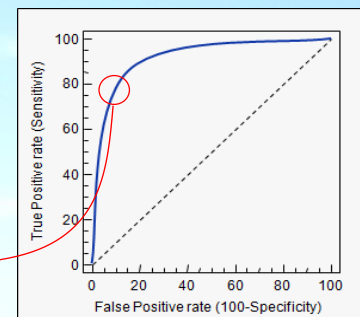
$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2) + \frac{1}{N_1 + N_2 - 2} \ln \frac{P(\omega_1)}{P(\omega_2)}$$



- Choose the threshold with ROC curve



adjusting the  
threshold



Xuegang Zhang

18

- 优秀!
- 不过，这算机器学习吗?



人  
↓  
机器

“看”到东西 → 认出东西、产生想法

观察 → 判断

观测 → 分类决策

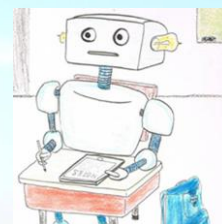
量化观测 → 类别标签

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

模式识别

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。

基于数据的机器学习



Xuegong Zhang

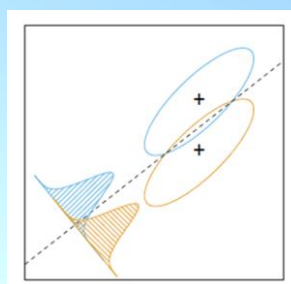
19



objects



sensor



Feature  
extractor

Classifier  
Predictor  
Generator

OUTPUT  
(decision,  
action, ...)

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

20



# 休息10秒钟



Xuegong Zhang

21



## 3.2 感知器 Perceptron

Xuegong Zhang

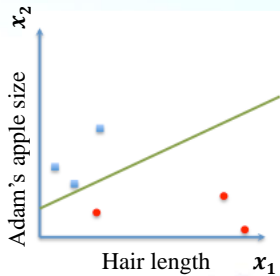


# 一个简单的例子

- 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.



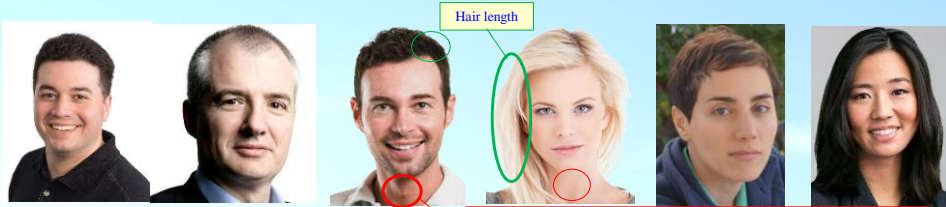
**For the Green line:**  
 $w_1x_1 + w_2x_2 + w_0 = 0$   
**For each blue dots:**  
 $w_1x_1 + w_2x_2 + w_0 > 0$   
**For each red dots:**  
 $w_1x_1 + w_2x_2 + w_0 < 0$

Xuegong Zhang

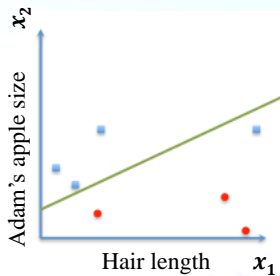


# 一个简单的例子

- 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:**  
 $w_1x_1 + w_2x_2 + w_0 < 0$

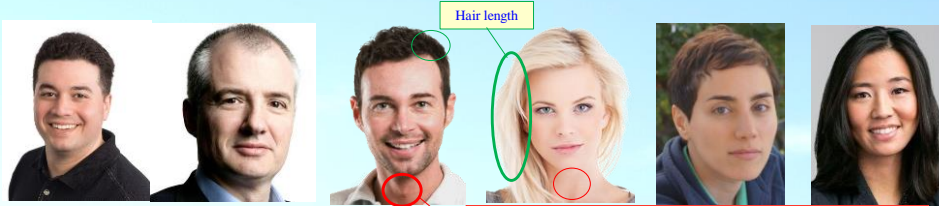
Xuegong Zhang



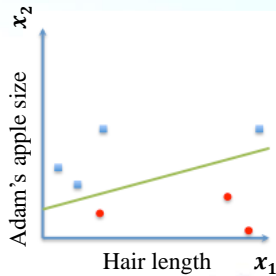


## 一个简单的例子

- 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:

$$w_1x_1 + w_2x_2 + w_0 < 0$$



25

Xuegong Zhang

## 感知器 (Perceptron)

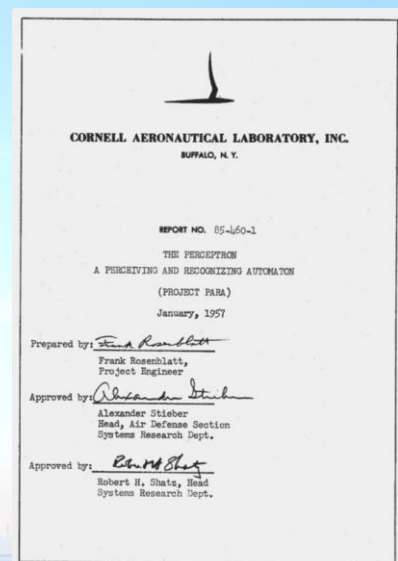


Frank Rosenblatt, *The Perceptron – a perceiving and recognizing automaton*,  
Report 85-460-1, Cornell Aeronautical Laboratory, Jan. 1957

- 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？



Xuegong Zhang





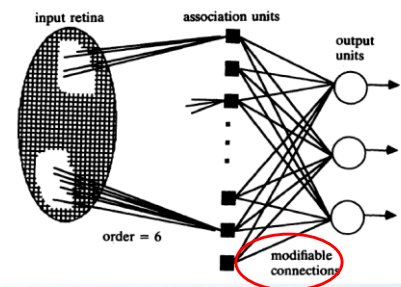
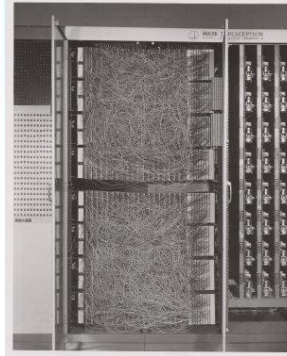
## 感知器

### • 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！

- 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

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



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

Xuegong Zhang

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

27



## 感知器

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

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

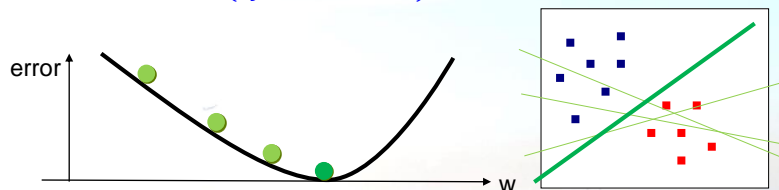
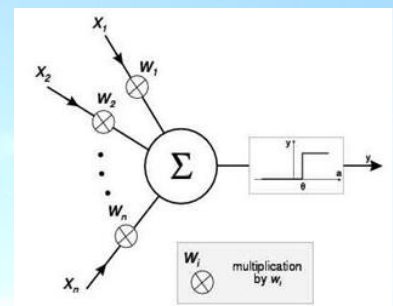
- 目标：最优化目标函数  $J(w)$

$J(w)$  = 训练错误数

- 学习算法

- 梯度下降法

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



Xuegong Zhang

28





# 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

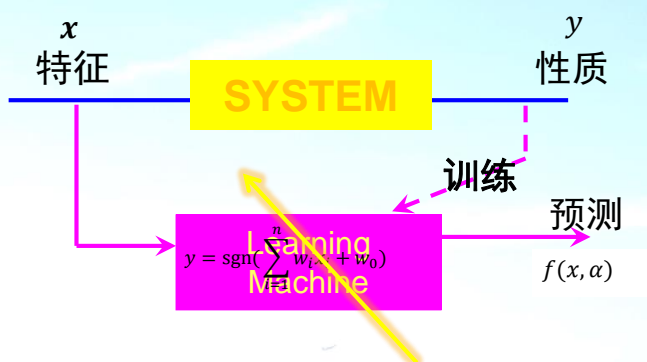


Xuegong Zhang



## 监督学习 Supervised Learning

- 用已知答案的数据去训练 → 监督学习



若  $y$  是离散类  
就是模式识别  
→ 监督模式识别

比如：

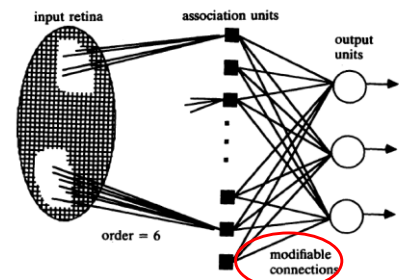
- $x$  身高体重 →  $y$  性别
- $x$  图像特征 →  $y$  性别
- $x$  图像像素 →  $y$  图像物体
- $x$  司机视觉 →  $y$  方向盘角度
- $x$  音频信号 →  $y$  语音内容
- $x$  基因表达 →  $y$  疾病类型
- $x$  肺CT影像 →  $y$  新冠肺炎
- $x$  近期疫情 →  $y$  未来走向
- ...

Xuegong Zhang



## 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

## 感知器的求解



- 线性判别函数的齐次简化

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

$\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$  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：

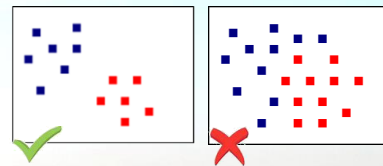
$$\mathbf{y}'_i = \begin{cases} \mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_1, \\ -\mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_2, \end{cases} \quad i = 1, \dots, N$$

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

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

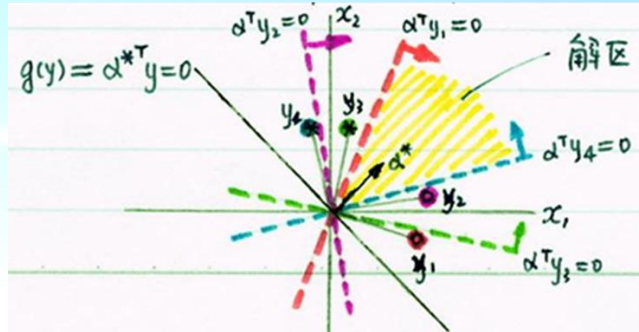
- 线性可分性：

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





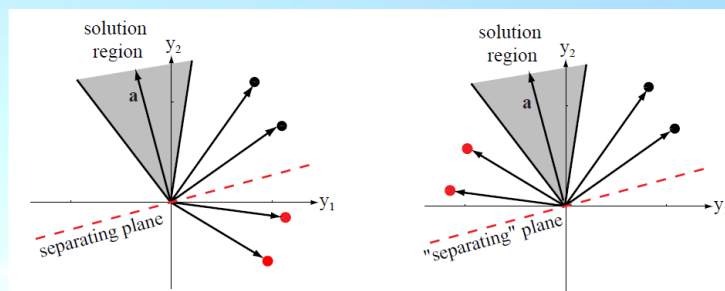
- 解向量  $\alpha^*$ : 满足  $\alpha^T y_i > 0, i = 1, \dots, N$
- 解区: 权值空间中所有解向量组成的区域



- 对每个样本  $y_i$ ,  $\alpha^T y_i = 0$  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

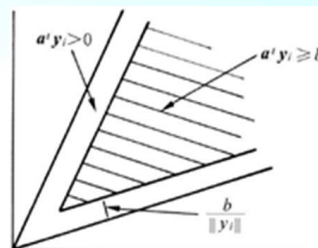
Xuegong Zhang

33



Duda et al, Pattern Classification

- 引入余量  $\alpha^T y_i \geq b > 0$



Xuegong Zhang

34



## 感知器准则函数

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

$\mathcal{Y}^k$ : 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

35



## 感知器准则函数

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

$\mathcal{Y}^k$ : 在第 $k$ 步被 $\alpha$ 错分的样本集合

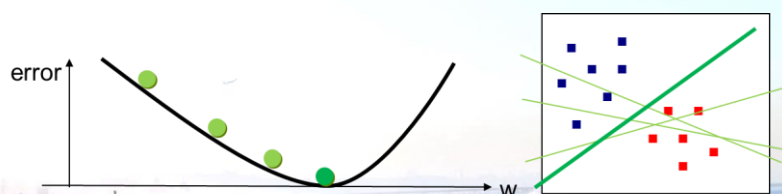
- 感知器算法 (Rosenblatt, 1957):

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

- 用梯度下降法(Gradient descent)迭代求解

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

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



Xuegong Zhang

36



## 感知器学习算法

### 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  $J_p = 0$

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

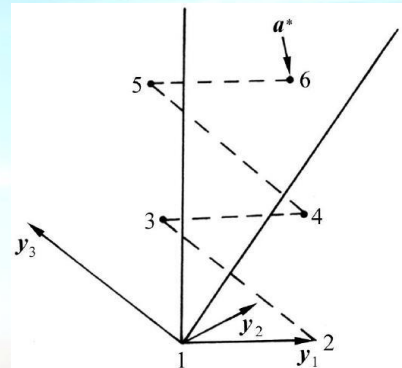
### 收敛性：

- 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



- 变增量法，如绝对修正法

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



Xuegong Zhang

37

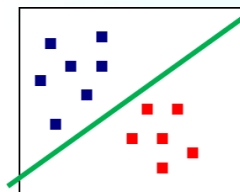


## 讨论

- 感知器有什么问题？



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



Xuegong Zhang

38





讨论

- 感知器有什么问题？

- 样本线性不可分呢？

- 线性可分时多解？

- 多类呢？
- 有问题怎么办？

- 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

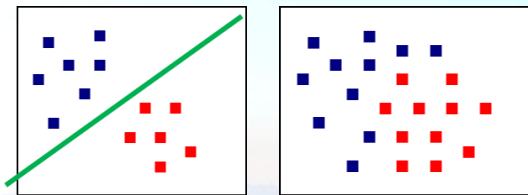
- 寻求“最优分类器”

• 比如支持向量机SVM

- 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

39

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



提交





## 3.3 线性回归 Linear Regression



## 线性分类器

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

↑ 预测    ↑ 分类器    ↑ 样本    ↑ 样本    ← 特征

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





# 线性回归

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

↑

↑

↑

预测

回归

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

↑

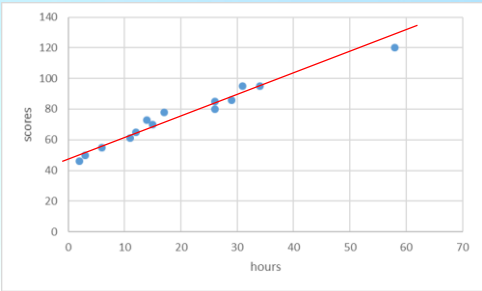
样本

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



# 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

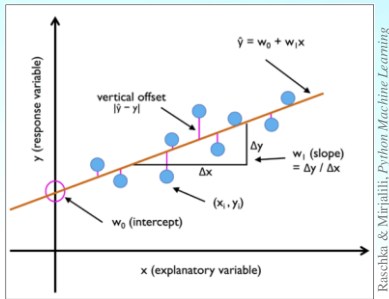
$$y = w_0 + w_1 x$$



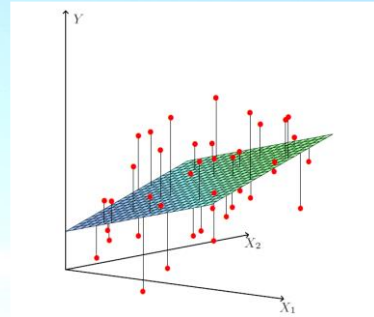
## 多元线性回归

### Simple Linear Regression

$$y = w_0 + w_1 x$$



Raschka & Mirjalili, Python Machine Learning



Hastie, Tibshirani, Friedman, Elements of Statistical Learning

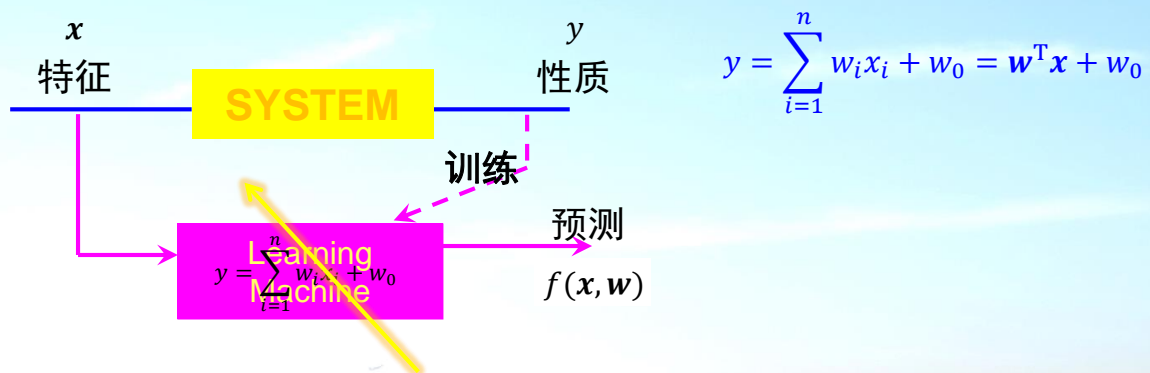
### Multiple Linear Regression

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

Xuegong Zhang

45

## 线性回归版的机器学习



Xuegong Zhang

46



## 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



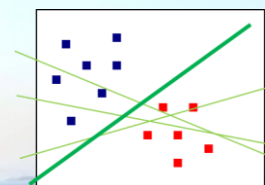
Xuegong Zhang

47

## 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）： \_\_\_\_\_
  - 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



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

#### – 它需要训练/学习材料

→ 训练数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_j \in R^{d+1}, y_j \in \{-1, 1\}$

#### – 我们需要为它树立学习的目标

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

#### – 我们需要告诉它怎样学

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



Xuegong Zhang

49



## 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）  $f(x) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

#### – 它需要训练/学习材料

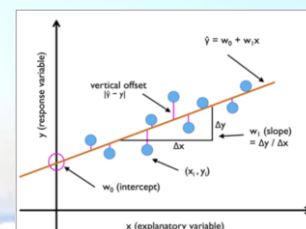
→ 训练数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_j \in R^{d+1}, y_j \in R$

#### – 我们需要为它树立学习的目标

→ 目标函数、学习准则 \_\_\_\_\_

#### – 我们需要告诉它怎样学

→ 学习/训练算法 ?



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）  $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

- 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

51

## 线性回归算法



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

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

解：

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

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

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

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

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

$$\text{dim: } \begin{matrix} (d+1) \times N & N \times (d+1) & (d+1) \times N \end{matrix}$$

Xuegong Zhang





### Linear regression algorithm:

- 1: Construct the matrix  $X$  and the vector  $y$  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

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

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

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

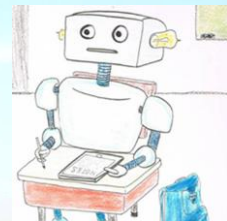
Xuegong Zhang

53

- 优秀！
- 不过，这算机器学习吗？



- 这解析解也算是“机器学习”？
  - OLS由Legendre在1805和Gauss在1809发明，远早于机器学习概念的诞生
  - 如FLD一样，算法从数据中算出“预测规律”
    - ➔ 也算是“机器学习”
  - 线性回归也可迭代求解，如感知器那样迭代“学习”



Xuegong Zhang

54



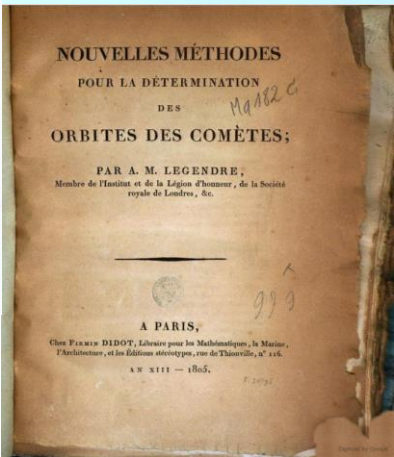
A.M. Legendre. Nouvelles mthodes pour la dtermination des orbites des comtes, Firmin Didot, Paris, 1805. "Sur la Mthode des moindres quarrs" appears as an appendix.

Adrien-Marie Legendre



1820 watercolor caricature of Adrien-Marie Legendre by French artist Julien-Leopold Boilly (see portrait debate), the only existing portrait known<sup>[1]</sup>

**Born** 18 September 1752  
Paris, France  
**Died** 10 January 1833 (aged 80)  
Paris, France



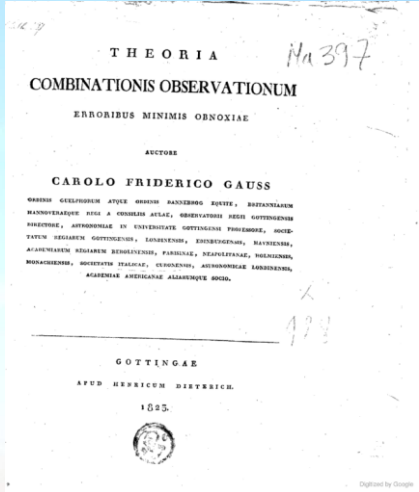
C.F. Gauss. Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientum. (1809)

Carl Friedrich Gauss



Carl Friedrich Gauß (1777–1855), painted by Christian Albrecht Jensen

**Born** Johann Carl Friedrich Gauss  
30 April 1777  
Brunswick, Principality of Brunswick-Wolfenbüttel  
**Died** 23 February 1855 (aged 77)  
Göttingen, Kingdom of Hanover, German Confederation



- Wikipedia

# 算法本身有没有问题？



## 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  $w_{\text{lin}} = X^\dagger y$ .



## 算法本身有没有问题？

- 若 $(X^T X)$ 可逆，则  $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$ ； 如果不可逆呢？

- 可逆（非奇异，非退化，满秩）

$$X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} \quad \end{bmatrix}_{N \times (d+1)}$$

- $(X^T X)$ 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

57



## 如何评价回归效果？

Xuegong Zhang

58



## Evaluation of regression models



$R^2$ : the **goodness-of-fit**, the **coefficient of determination**, ...

R平方/拟合度/决定系数

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

Unexplained variation

Total variation

For OLS,

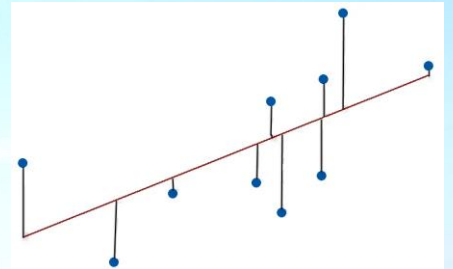
$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

and  $0 \leq R^2 \leq 1$ .

$R^2 = 1$ : Perfect regression.

$R^2 = 0$ : Baseline model. Predictions are the average.

$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2} \in [-\infty, 1]$  for other types of regression as a general measure of goodness-of-fit, and should no longer be called  $R^2$ .



Xuegong Zhang

59

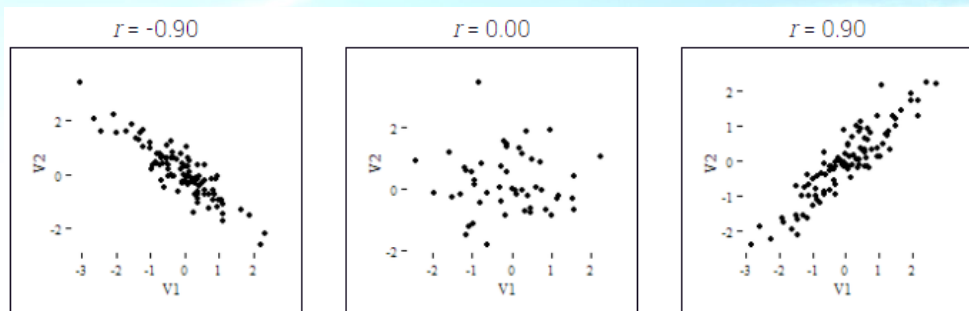
## Pearson Correlation Coefficient



皮尔森相关系数



$$r_{XY} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)} \cdot \sqrt{\text{var}(Y)}} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \cdot \sqrt{\sum (Y_i - \bar{Y})^2}}$$



What are the relation between  $R^2$  and  $r$  ?



Xuegong Zhang

60

## $R^2$ is not enough for evaluating regression

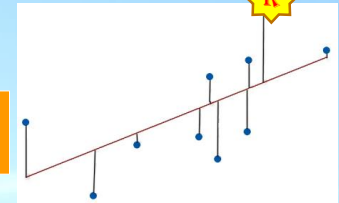


$$y_i = w_0 + \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \quad i = 1, \dots, N$$

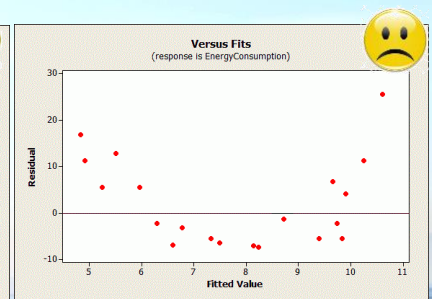
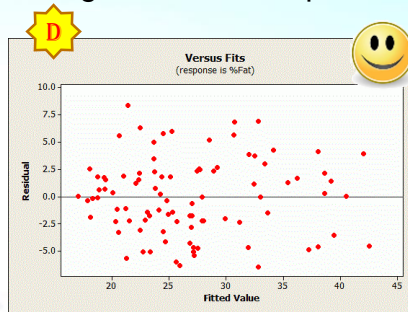
Dependent  
variable

Fitted values  
(deterministic)

Error, residual, noise  
(stochastic)



- $R^2$  : percentage of dependent variable variations that the linear model explains.
- $R^2$  does not indicate if the regression model provides an adequate fit to the data.
- Use **residual plots** to check whether the model is adequate.
- Poor fitting if error is not random



Xuegong Zhang

## Evaluating each coefficient



$$y_i = w_0 + \sum_{j=1}^d w_j x_{ij} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2), i = 1, \dots, N,$$

Does each  $w_j$   
contribute?

Test for statistical significance of regression coefficients

$$\frac{\hat{w}_j - w_j}{s_{\hat{w}_j}} \sim t_{N-d-1}, \quad j = 0, 1, \dots, d.$$

Xuegong Zhang

62





# 休息10秒钟



Xuegong Zhang

63



## 3.4 最小平方误差分类器 Minimum Squared Error (MSE)

Xuegong Zhang

64

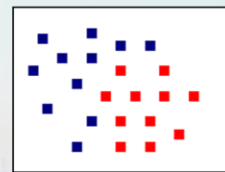
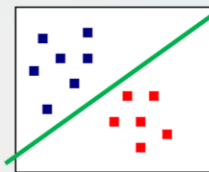


## 回顾

### 讨论

- 感知器有什么问题？
  - 样本线性不可分呢？
  - 线性可分时多解？
  - 多类呢？
- 有问题怎么办？
  - 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
  - 寻求“最优分类器”
    - 比如支持向量机SVM
  - 多类分类方法
  - 用两类分类器完成多类分类

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



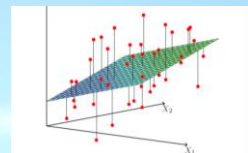
Xuegong Zhang

65

## MSE方法的思想

- 对线性不可分情况，怎样最小化线性分类器的错误？
- 线性回归：求 $w$ 使 $y_i = w^T x_i, i = 1, \dots, N$ 
  - 不可能全部样本正好都满足，于是最小化平方误差，方法是：

$$\min E = \frac{1}{N} \sum_{j=1}^N (f(x_j) - y_j)^2$$



- 线性分类器：求 $\alpha$ 使 $\alpha^T y_i > 0, i = 1, \dots, N$ 
  - 当样本集线性不可分时，使尽可能多的样本满足不等式
- 考虑：为每个样本引入 $b_i$ ，并令

$$\alpha^T y_i = b_i > 0, \quad j = 1, \dots, N$$

不等式组  $\rightarrow$  等式组，可用最小二乘求解

注意：这里切换回感知器时的规范化推广向量形式

Xuegong Zhang

66



## MSE分类器准则

$$\alpha^T y_i > 0, i = 1, \dots, N \iff \alpha^T y_i = b_i > 0, i = 1, \dots, N$$

- 不等式转化成等式

$$\mathbf{Y}\alpha = \mathbf{b}, \mathbf{b} = [b_1, b_2, \dots, b_N]^T$$

- MSE准则  $\alpha^*: \min_{\alpha} J_S(\alpha)$

$$J_S(\alpha) = \|\mathbf{Y}\alpha - \mathbf{b}\|^2 = \sum_{i=1}^N (\alpha^T y_i - b_i)^2$$

- 解法:

- 最小二乘伪逆解  $\alpha^* = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} = \mathbf{Y}^+ \mathbf{b}, \mathbf{Y}^+ = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$
- 梯度下降学习

$$\nabla J_S(\alpha) = 2\mathbf{Y}^T (\mathbf{Y}\alpha - \mathbf{b})$$

$$\alpha(k+1) = \alpha(k) + \rho_k (b_k - \alpha(k)^T y^k) y^k$$

--- Widrow-Hoff算法, **ADALINE**

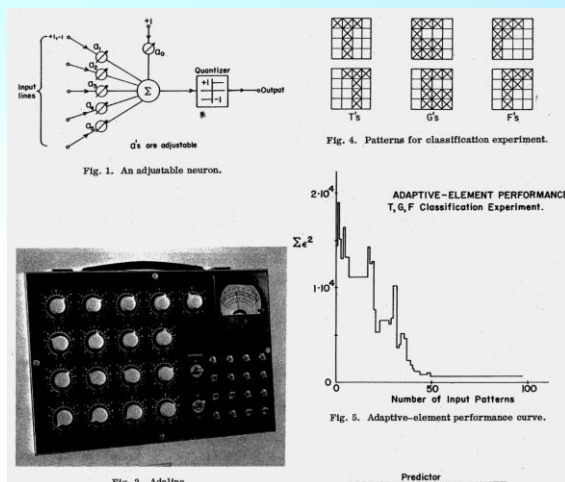
Xuegong Zhang

67

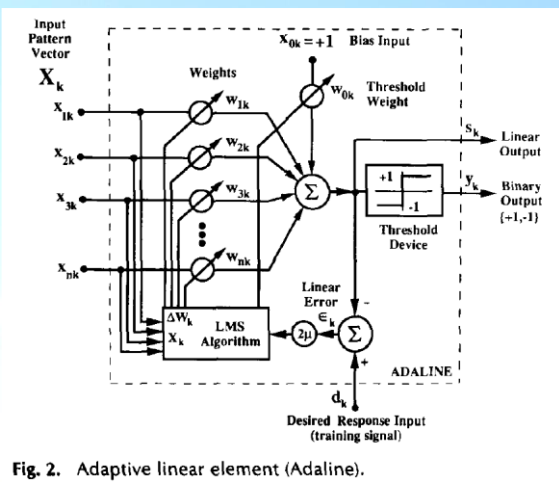
## ADALINE



Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960



Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960



Widrow & Lehr, 30 years of adaptive neural networks: Perceptron, Madaline, and Backpropagation, *Proceedings of the IEEE*, 78(9): 1415-1442, 1990

Xuegong Zhang

68



## 有没有问题？



如何给定  $\mathbf{b} = [b_1, b_2, \dots, b_N]^T$  ?

- 可以证明，如果  $\mathbf{b}$  选为

$$b_i = \begin{cases} N/N_1, & \text{if } y_i \in \omega_1 \\ N/N_2, & \text{if } y_i \in \omega_2 \end{cases},$$



则MSE解等价于  $w_0 = -\hat{m}$  的FLD解。

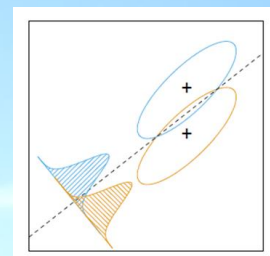
- 如果  $\mathbf{b}$  选为：

$$b_i = 1, i = 1, \dots, N,$$

则当  $N \rightarrow \infty$  时，MSE解以最小均方误差最优逼近贝叶斯判别函数

$$g_0(\mathbf{x}) = P(\omega_1|\mathbf{x}) - P(\omega_2|\mathbf{x})$$

即  $\alpha_{\text{MSE}}$  使  $e^2 = \int [\alpha^T \mathbf{y} - g_0(\mathbf{x})]^2 p(\mathbf{x}) d(\mathbf{x})$  最小。



Xuegong Zhang

69

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

70



# 3.5 罗杰斯特回归

## Logistic Regression

注：单词logistic与logic没有关系，不应译作“逻辑回归”，  
应音译作“罗杰斯特回归”或“罗杰斯蒂回归”

Xuegong Zhang

71

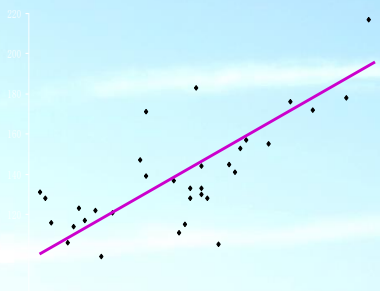


# 简单线性回归

Age	SBP	Age	SBP	Age	SBP
22	131	41	139	52	128
23	128	41	171	54	105
24	116	46	137	56	145
27	106	47	111	57	141
28	114	48	115	58	153
29	123	49	133	59	157
30	117	49	128	63	155
32	122	50	183	67	176
33	99	51	130	71	172
35	121	51	133	77	178
40	147	51	144	81	217

Table 1 Age and systolic blood pressure (SBP) among 33 adult women

$SBP = 81.54 + 1.222 \cdot Age$



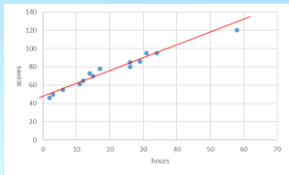
Xuegong Zhang

72



# 回归用来预测二值结果？

Student id	Final score	Study Hours per Week
1	0	3
2	1	34
3	0	17
4	0	6
5	0	12
6	1	15
7	1	26
8	1	29
9	0	14
10	1	58
11	0	2
12	1	31
13	1	26
14	0	11



Simple Learn Regression  
 $y = w_0 + w_1x$

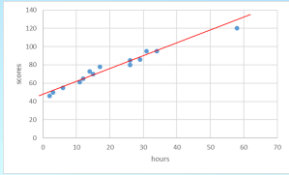
Xuegong Zhang

73

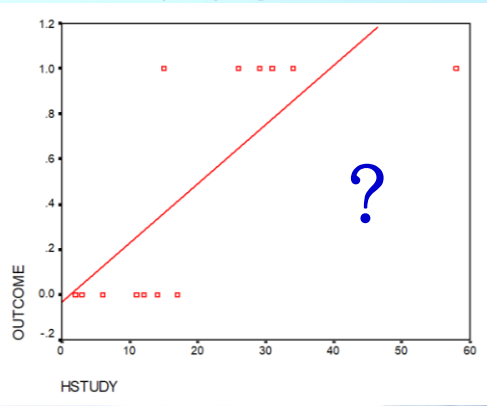


# 回归用来预测二值结果？

Student id	Final score	Study Hours per Week
1	0	3
2	1	34
3	0	17
4	0	6
5	0	12
6	1	15
7	1	26
8	1	29
9	0	14
10	1	58
11	0	2
12	1	31
13	1	26
14	0	11



Simple Learn Regression  
 $y = w_0 + w_1x$



Xuegong Zhang

74

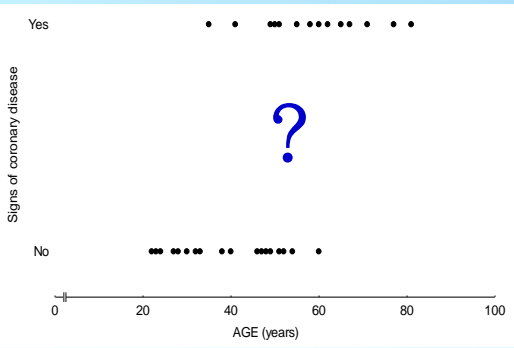




# 用年龄回归冠心病

Age and signs of coronary heart disease (CD)

Age	CD	Age	CD	Age	CD
22	0	40	0	54	0
23	0	41	1	55	1
24	0	46	0	58	1
27	0	47	0	60	1
28	0	48	0	60	0
30	0	49	1	62	1
30	0	49	0	65	1
32	0	50	1	67	1
33	0	51	0	71	1
35	1	51	1	77	1
38	0	52	0	81	1



Xuegong Zhang

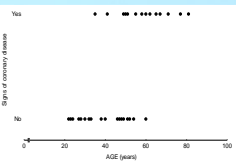
75



## 让我们来计数

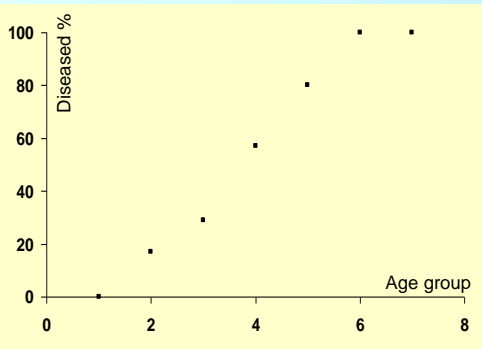
Age and signs of coronary heart disease (CD)

Age	CD	Age	CD	Age	CD
22	0	40	0	54	0
23	0	41	1	55	1
24	0	46	0	58	1
27	0	47	0	60	1
28	0	48	0	60	0
30	0	49	1	62	1
30	0	49	0	65	1
32	0	50	1	67	1
33	0	51	0	71	1
35	1	51	1	77	1
38	0	52	0	81	1



Prevalence (%) of signs of CD according to age group

Age group	# in group	Diseased	
		#	%
20 - 29	5	0	0
30 - 39	6	1	17
40 - 49	7	2	29
50 - 59	7	4	57
60 - 69	5	4	80
70 - 79	2	2	100
80 - 89	1	1	100



Xuegong Zhang

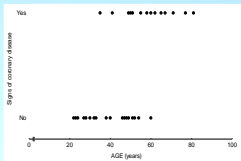
76



让我们来计数

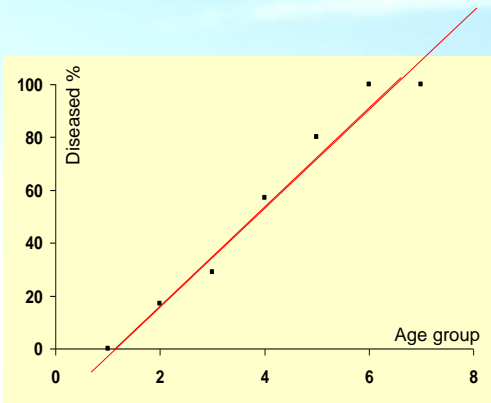
Age and signs of coronary heart disease (CD)

Age	CD	Age	CD	Age	CD
22	0	40	0	54	0
23	0	41	1	55	1
24	0	46	0	58	1
27	0	47	0	60	1
28	0	48	0	60	0
30	0	49	1	62	1
30	0	49	0	65	1
32	0	50	1	67	1
33	0	51	0	71	1
35	1	51	1	77	1
38	0	52	0	81	1



Prevalence (%) of signs of CD according to age group

Age group	# in group	Diseased	
		#	%
20 - 29	5	0	0
30 - 39	6	1	17
40 - 49	7	2	29
50 - 59	7	4	57
60 - 69	5	4	80
70 - 79	2	2	100
80 - 89	1	1	100



Xuegon

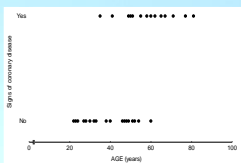
77



让我们来计数

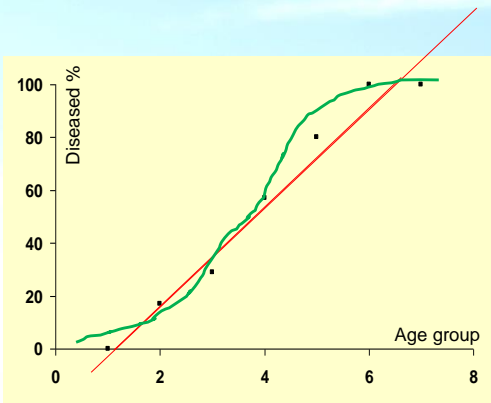
Age and signs of coronary heart disease (CD)

Age	CD	Age	CD	Age	CD
22	0	40	0	54	0
23	0	41	1	55	1
24	0	46	0	58	1
27	0	47	0	60	1
28	0	48	0	60	0
30	0	49	1	62	1
30	0	49	0	65	1
32	0	50	1	67	1
33	0	51	0	71	1
35	1	51	1	77	1
38	0	52	0	81	1



Prevalence (%) of signs of CD according to age group

Age group	# in group	Diseased	
		#	%
20 - 29	5	0	0
30 - 39	6	1	17
40 - 49	7	2	29
50 - 59	7	4	57
60 - 69	5	4	80
70 - 79	2	2	100
80 - 89	1	1	100



Xuegon

78

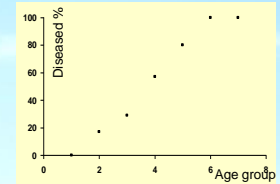
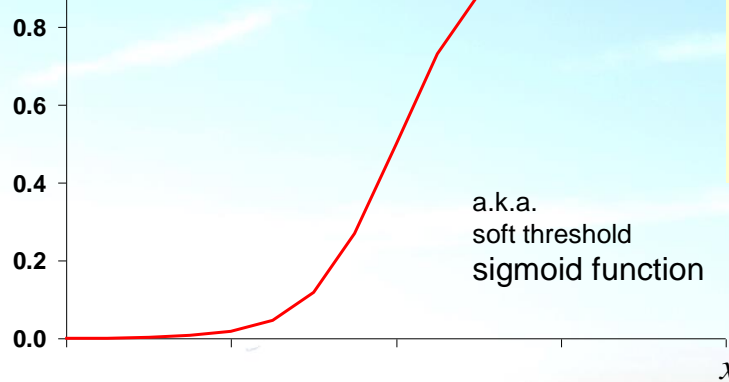


## 罗杰斯特函数Logistic function



患病比例

$$P(y = 1|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$



Xuegong Zhang

79

## 罗杰斯特回归与几率



$$P(y|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

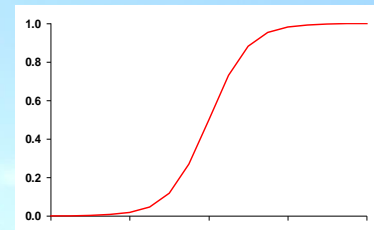
Odds 几率

$$\frac{P(y|x)}{1 - P(y|x)} = e^{\alpha + \beta x}$$

Log odds 对数几率

$$\ln \left[ \frac{P(y|x)}{1 - P(y|x)} \right] = \alpha + \beta x$$

$P(y|x)$ 的罗杰特 (logit) 函数



注意中文中的混淆词：几率、机率、概率、可能性、...

Xuegong Zhang

80



## 多元罗杰斯特回归

$$P(y|\mathbf{x}) = \frac{e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}{1 + e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}$$

$$\ln\left(\frac{P(y|\mathbf{x})}{1 - P(y|\mathbf{x})}\right) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$odds = \frac{P(y|\mathbf{x})}{1 - P(y|\mathbf{x})} = e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

### • $\beta_i$ 的解释

- 在其他因素不变的情况下，因素 $x_i$ 增加一个单位带来的对数几率的增加
- 可以用来从流行病学数据中研究各种因素与患病的关系

Xuegong Zhang

81

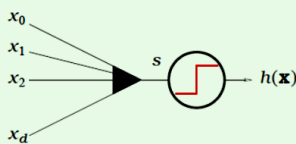


## 三种线性机器

$$s = \sum_{i=0}^d w_i x_i$$

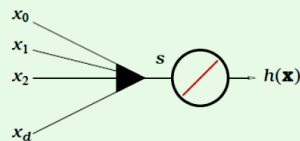
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



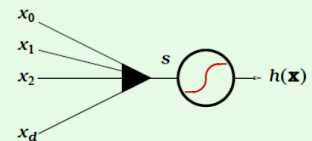
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

$$h(\mathbf{x}) = \theta(s)$$

Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

Xuegong Zhang

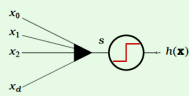
82

## 罗杰斯特回归的概率意义

$$s = \sum_{i=0}^d w_i x_i$$

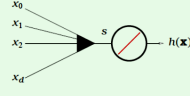
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



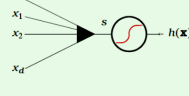
linear regression

$$h(\mathbf{x}) = s$$

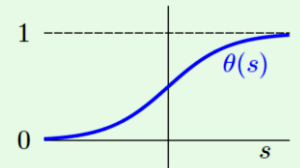


logistic regression

$$h(\mathbf{x}) = \theta(s)$$



$$\theta(s) = \frac{e^s}{1 + e^s}$$



Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

- $s = \mathbf{w}^T \mathbf{x}$  是事件的总信息，是各种因素（特征）的加权求和
- $h(\mathbf{x}) = \theta(s)$  是对事件  $y = 1$  概率的估计

Xuegong Zhang

83

## 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

84





## 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）  $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$
  - 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  $\min J_P(\boldsymbol{\alpha}) = \sum_{\mathbf{y}_j \in \mathcal{Y}^k} (-\boldsymbol{\alpha}^T \mathbf{y}_j)$
  - 我们需要告诉它怎样学
    - 学习/训练算法  $\boldsymbol{\alpha}(k+1) = \boldsymbol{\alpha}(k) - \rho_k \nabla J = \boldsymbol{\alpha}(k) + \rho_k \sum_{\mathbf{y}_j \in \mathcal{Y}^k} \mathbf{y}_j$

Xuegong Zhang

85



## 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）  $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$
  - 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  $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$

Xuegong Zhang

86



## 机器学习的基本要素：罗杰斯特回归版？

- 怎样造一个学习机器？

- 它需要老师

→ 我们设计它（特征和模型）  $h(x) = \theta(w^T x)$

- 它需要训练/学习材料

→ 训练数据  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_j \in R^{d+1}, y_j \in \{-1, 1\}$

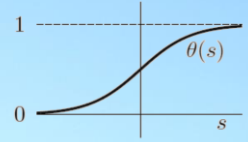
- 我们需要为它树立学习的目标

→ 目标函数、学习准则 ？

- 我们需要告诉它怎样学

→ 学习/训练算法 ？

$$\theta(s) = \frac{e^s}{1 + e^s}$$



Xuegong Zhang

87



## 休息10秒钟



Xuegong Zhang

88



## 机器学习的基本要素：罗杰斯特回归版？

- 怎样造一个学习机器？

- 它需要老师

- 我们设计它（特征和模型）  $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$

- 它需要训练/学习材料

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

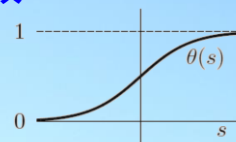
- 我们需要为它树立学习的目标

- 目标函数、学习准则？

- 我们需要告诉它怎样学

- 学习/训练算法？

$$\theta(s) = \frac{e^s}{1 + e^s}$$



Xuegong Zhang

89

## 考虑样本的产生模型



- 设独立同分布(i.i.d.)样本集  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in \mathbb{R}^{d+1}, y_j \in \{-1, 1\}$  依以下概率产生：

$$P(y|\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1 \\ 1 - f(\mathbf{x}) & \text{for } y = -1 \end{cases}$$

生成模型  
Generative model

- 罗杰斯特回归用  $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$  估计  $f(\mathbf{x})$

Xuegong Zhang

90



## 考虑样本的产生模型

- 设独立同分布(i.i.d.)样本集  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_j \in R^{d+1}, y_j \in \{-1, 1\}$  依以下概率产生：

$$P(y|x) = \begin{cases} f(x) & \text{for } y = +1 \\ 1 - f(x) & \text{for } y = -1 \end{cases}$$

生成模型  
Generative model

- 罗杰斯特回归用  $h(x) = \theta(w^T x)$  估计  $f(x)$

- 似然函数 (Likelihood) :



– 对数据中的一个实例  $(x_j, y_j)$ , 如果  $h = f$ , 我们有多大可能对  $x_j$  得到  $y_j$ ?

$$P(y_j|x_j) = \begin{cases} h(x_j) & \text{for } y_j = +1 \\ 1 - h(x_j) & \text{for } y_j = -1 \end{cases}$$

– 换言之, 已经有这个数据实例,  $h$  有多大可能是产生数据的模型?

Xuegong Zhang

91



## 似然函数

- 设独立同分布(i.i.d.)样本集  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_j \in R^{d+1}, y_j \in \{-1, 1\}$  依以下概率产生：

$$P(y|x) = \begin{cases} f(x) & \text{for } y = +1 \\ 1 - f(x) & \text{for } y = -1 \end{cases}$$

生成模型  
Generative model

- 罗杰斯特回归用  $h(x) = \theta(w^T x)$  估计  $f(x)$

- 似然函数 (Likelihood) :

– 对数据中的一个实例  $(x_j, y_j)$ , 如果  $h = f$ , 我们有多大可能对  $x_j$  得到  $y_j$ ?

$$P(y_j|x_j) = \begin{cases} h(x_j) & \text{for } y_j = +1 \\ 1 - h(x_j) & \text{for } y_j = -1 \end{cases}$$

– 换言之, 已经有这个数据实例,  $h$  有多大可能是产生数据的模型?

- 注意到  $\theta(-s) = 1 - \theta(s)$ ,  $x_j$  上的似然函数可写为:



$$P(y_j|x_j) = \theta(y_j w^T x_j)$$

Xuegong Zhang

92



## 似然函数

- 设独立同分布(i.i.d.)样本集  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_j \in R^{d+1}, y_j \in \{-1, 1\}$  依以下概率产生：

$$P(y|x) = \begin{cases} f(x) & \text{for } y = +1 \\ 1 - f(x) & \text{for } y = -1 \end{cases}$$

生成模型  
Generative model

- 罗杰斯特回归用  $h(x) = \theta(w^T x)$  估计  $f(x)$

- 似然函数 (Likelihood) ?

问题：似然函数是谁的函数？



- 对数据中的一个实例  $(x_j, y_j)$ , 如果  $h = f$ , 我们有多大可能对  $x_j$  得到  $y_j$ ?

$$P(y_j|x_j) = \begin{cases} h(x_j) & \text{for } y_j = +1 \\ 1 - h(x_j) & \text{for } y_j = -1 \end{cases}$$

- 换言之, 已经有这个数据实例,  $h$  有多大可能是产生数据的模型?

- 注意到  $\theta(-s) = 1 - \theta(s)$ ,  $x_j$  上的似然函数可写为:

$$P(y_j|x_j) = \theta(y_j w^T x_j)$$

Xuegong Zhang

93

单选题 10分

设置

似然函数是谁的函数？

$$P(y_j|x_j) = \theta(y_j w^T x_j)$$

- ☐ A 是  $y_i$  的函数
- ☒ B 是  $w$  的函数
- ☐ C 是  $x_i$  的函数
- ☐ D 是  $(x_j, y_j)$  的函数

提交

94





## 似然函数

- 设独立同分布(i.i.d.)样本集  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$  依以下概率产生：

$$P(y|\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1 \\ 1 - f(\mathbf{x}) & \text{for } y = -1 \end{cases}$$

生成模型  
Generative model

- 罗杰斯特回归用  $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$  估计  $f(\mathbf{x})$
- 似然函数 (Likelihood) :
  - 对数据中的一个实例  $(\mathbf{x}_j, y_j)$ , 如果  $h = f$ , 我们有多大可能对  $\mathbf{x}_j$  得到  $y_j$ ?

$$P(y_j|\mathbf{x}_j) = \begin{cases} h(\mathbf{x}_j) & \text{for } y_j = +1 \\ 1 - h(\mathbf{x}_j) & \text{for } y_j = -1 \end{cases}$$

- 换言之, 已经有这个数据实例,  $h$  有多大可能是产生数据的模型?
- 注意到  $\theta(-s) = 1 - \theta(s)$ , 模型在  $\mathbf{x}_j$  上的似然函数可写为:

$$l(h|(\mathbf{x}_j, y_j)) = P(y_j|\mathbf{x}_j, h) = \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$

Xuegong Zhang

95

## 罗杰斯特回归的目标：最大化似然函数



- 参数为  $\mathbf{w}$  的罗杰斯特模型在 i.i.d. 数据  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$  上的似然函数是

$$L(\mathbf{w}) = \prod_{j=1}^N P(y_j|\mathbf{x}_j) = \prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$

Xuegong Zhang

96



## 罗杰斯特回归的目标：最大化似然函数

- 参数为 $\mathbf{w}$ 的罗杰斯特模型在i.i.d.数据 $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_j \in \mathbb{R}^{d+1}, y_j \in \{-1, 1\}$ 上的似然函数是

$$L(\mathbf{w}) = \prod_{j=1}^N P(y_j | \mathbf{x}_j) = \prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$

- 最大化似然函数，等价于最小化以下目标函数：

$$\begin{aligned} \min \quad E(\mathbf{w}) &= -\frac{1}{N} \ln(L(\mathbf{w})) = -\frac{1}{N} \ln \left( \prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j) \right) \\ &= \frac{1}{N} \sum_{j=1}^N \ln \left( \frac{1}{\theta(y_j \mathbf{w}^T \mathbf{x}_j)} \right) \\ &= \frac{1}{N} \sum_{j=1}^N \ln (1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j}) \end{aligned}$$

$$\left[ \theta(s) = \frac{1}{1 + e^{-s}} \right] \text{H}$$

Xuegong Zhang

97

## 求解：梯度下降法

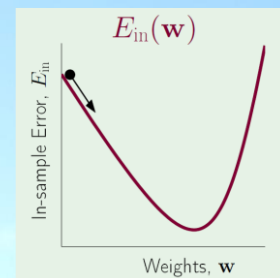


梯度下降的一般原理：

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \eta \hat{\mathbf{v}}$$

其中， $\eta$ ：步长(学习率)

$$\hat{\mathbf{v}} = -\nabla E(\mathbf{w}(k))$$



Xuegong Zhang

Abu-Mostafa, Magdon-Ismail, Lin, Learning from Data, Lecture 9

98



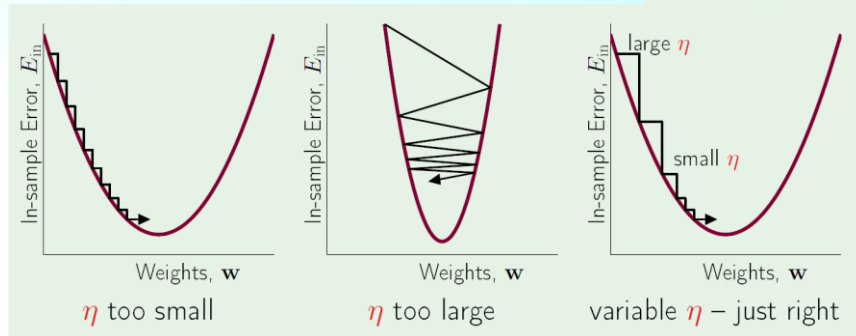
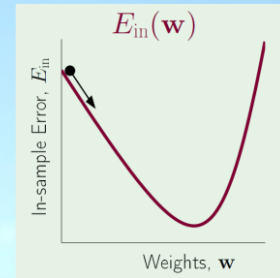
## 求解：梯度下降法

梯度下降的一般原理：

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \eta \hat{\mathbf{v}}$$

其中， $\eta$ ：步长(学习率)

$$\hat{\mathbf{v}} = -\nabla E(\mathbf{w}(k))$$



Xuegong Zhang

Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

99



## 罗杰斯特回归算法

1. Set  $k = 0$ , initialize  $\mathbf{w}(0)$

2. Do

– Compute the gradient  $\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j x_j}{1 + e^{y_j \mathbf{w}(k)^T x_j}}$

– Update the weights  $\mathbf{w}(k+1) = \mathbf{w}(k) - \eta \nabla E$ , set  $k = k + 1$

Until the stopping criterion met

3. Return the final weights  $\mathbf{w}$



Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

Xuegong Zhang

100



## 罗杰斯特回归算法

1. Set  $k = 0$ , initialize  $\mathbf{w}(0)$
2. Do
  - Compute the gradient  $\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j x_j}{1 + e^{y_j \mathbf{w}(k)^T x_j}}$
  - Update the weights  $\mathbf{w}(k+1) = \mathbf{w}(k) - \eta \nabla E$ , set  $k = k + 1$
 Until the stopping criterion met
3. Return the final weights  $\mathbf{w}$

Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

- 初始化：可以全零，更好是小随机数，比如0均值、小方差的正态分布
- 终止条件：梯度低于一定阈值，或迭代次数达到预设上限

Xuegong Zhang

101

## 机器学习的基本要素：罗杰斯特回归版



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

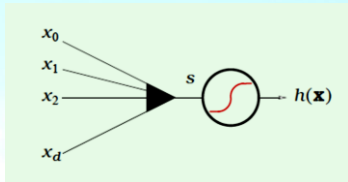
Xuegong Zhang

102

## 有没有问题？



$$h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$$



- 如何分类决策？分类器是什么？
  - $h(\mathbf{x}) \geq 0.5$ 
    - 最小错误率决策
  - 根据ROC曲线
  - 根据两类错误率的相对损失（风险）
    - 最小风险决策

Xuegong Zhang

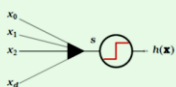
103

## 线性学习机器小结

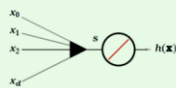
## 模型

$$s = \sum_{i=0}^d w_i x_i$$

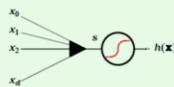
linear classification  
 $h(\mathbf{x}) = \text{sign}(s)$



linear regression  
 $h(\mathbf{x}) = s$



logistic regression  
 $h(\mathbf{x}) = \theta(s)$



## 目标

- For perceptron

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

- For linear regression

$$\min E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N (\mathbf{w}^T \mathbf{x}_j - y_j)^2$$

- For logistic regression

$$\min E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \ln(1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j})$$

## 学习算法

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$$

Xuegong Zhang

104





## 本章知识点

- 机器学习的基本概念
- 线性学习机器的基本思想
  - 模型、目标函数、梯度下降优化
- Fisher线性判别、感知器、线性回归、MSE(ADALINE)、罗杰斯特回归
- 似然函数的概念

Xuegong Zhang

105

## 下周内容

## 非线性分类与经典人工神经网络

欢迎关注B站XGlab视频/直播



Xuegong Zhang

106