

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 (2)

2020.2.18



Xuegong Zhang



2.4 最小平方误差分类器 Minimum Squared Error

Xuegong Zhang

2

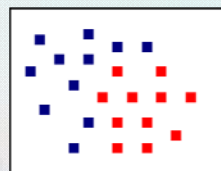
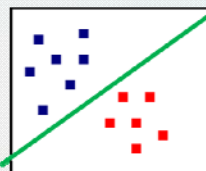


回顾

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

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



Xuegong Zhang

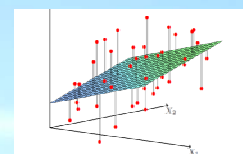
3



MSE方法的思想

- 对线性分类器，怎样最小化其错误？
- 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：

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



Xuegong Zhang

4



MSE方法的思想

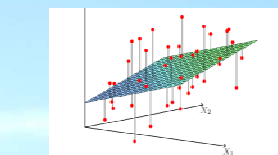
- 对线性分类器，怎样最小化其错误？
- 线性回归：求 \mathbf{w} 使 $y_i = \mathbf{w}^T \mathbf{x}_i, i = 1, \dots, N$
 - 不可能全部样本正好都满足，于是最小化平方误差，方法是：

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

- 线性分类器：求 $\boldsymbol{\alpha}$ 使 $\boldsymbol{\alpha}^T \mathbf{y}_i > 0, i = 1, \dots, N$
 - 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

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

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增广向量形式

Xuegong Zhang

5



MSE分类器准则

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

- 不等式转化成等式

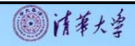
$$\mathbf{Y}\boldsymbol{\alpha} = \mathbf{b}, \quad \mathbf{b} = [b_1, b_2, \dots, b_N]^T$$

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

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

Xuegong Zhang

6



MSE分类器准则

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



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

- 不等式转化成等式

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

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

- 解法:

- 最小二乘伪逆解 $\alpha^* = (Y^T Y)^{-1} Y^T b = Y^+ b, Y^+ = (Y^T Y)^{-1} Y^T$

- 梯度下降学习

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

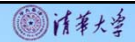
$$\alpha(k+1) = \alpha(k) + \rho_k(b_k - \alpha(k)^T y^k) y^k$$

--- Widrow-Hoff算法, **ADALINE**

Xuegong Zhang

7

ADALINE



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

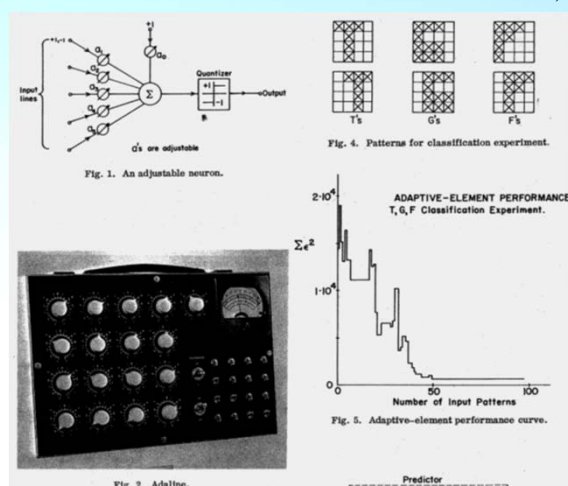


Fig. 2. Adaline.

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

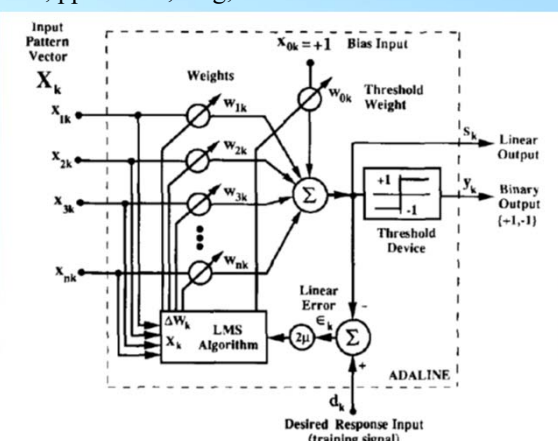


Fig. 2. Adaptive linear element (Adaline).

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

8

Xuegong Zhang



有没有问题？

请在雨课堂中弹幕提问



Xuegong Zhang

9



有没有问题？

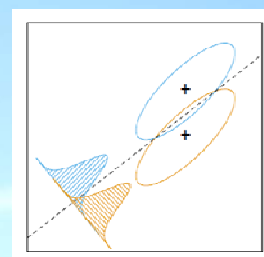
如何给定 $\mathbf{b} = [b_1, b_2, \dots, b_N]^T$?

- 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。



Xuegong Zhang

10

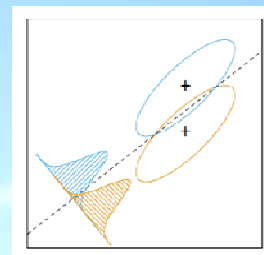


有没有问题？

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

Xuegong Zhang

11

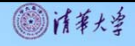


休息1分钟



Xuegong Zhang

12



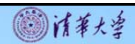
2.5 罗杰斯特回归 Logistic Regression

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

13

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

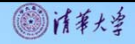
$$SBP = 81.54 + 1.222 \cdot Age$$



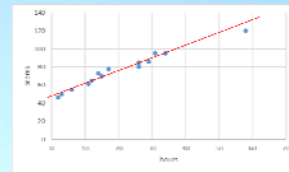
Xuegong Zhang

14

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

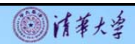


Simple Linear Regression
 $y = w_0 + w_1 x$

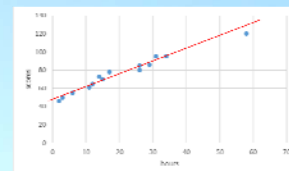
Xuegong Zhang

15

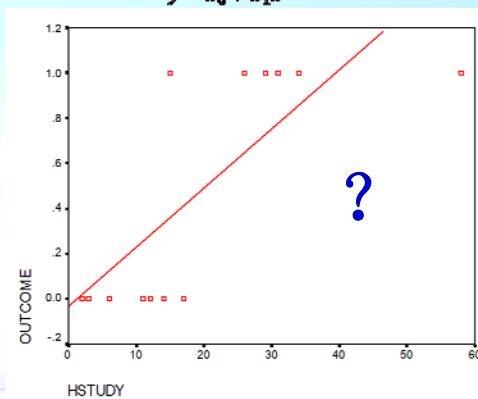
回归用来预测二值结果？



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 Linear Regression
 $y = w_0 + w_1 x$



Xuegong Zhang

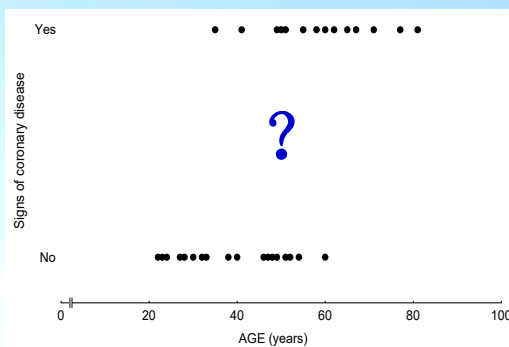
16

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

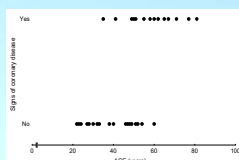
17

让我们来计数



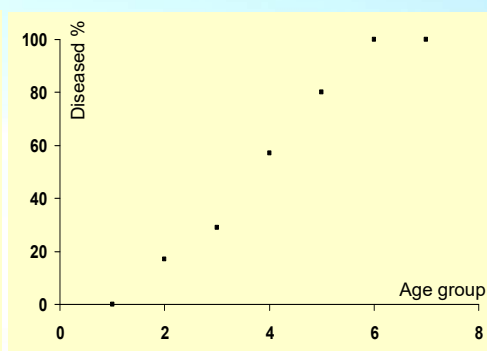
Age and signs of coronary heart disease (CD)

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

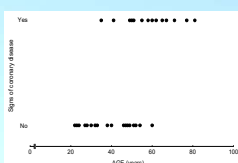
18



让我们来计数

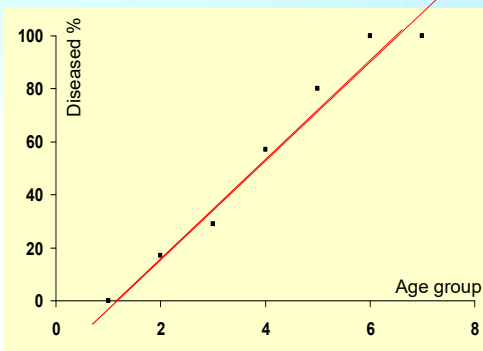
Age and signs of coronary heart disease (CD)

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



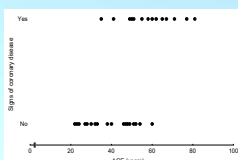
19



让我们来计数

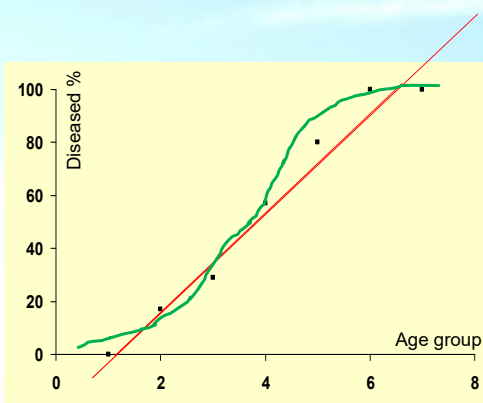
Age and signs of coronary heart disease (CD)

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



20



罗杰斯特函数 Logistic function

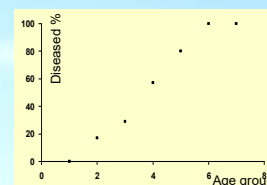


患病比例

1.0
0.8
0.6
0.4
0.2
0.0

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

a.k.a.
soft threshold
sigmoid function



Note: "logistic" is not related with "logic".

So I prefer the translation 罗杰斯特 or 罗杰斯蒂.

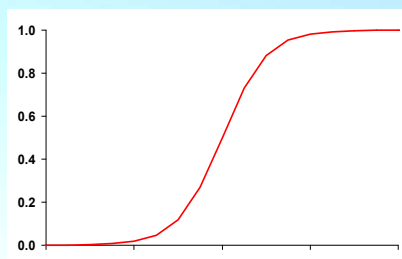
x

Xuegong Zhang

21



罗杰斯特回归



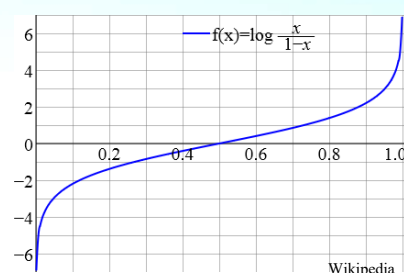
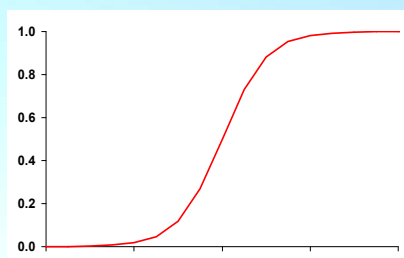
$$P(y|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

Xuegong Zhang

22



罗杰斯特回归与几率



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

反函数

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

$P(y|x)$ 的罗杰特 (logit) 函数

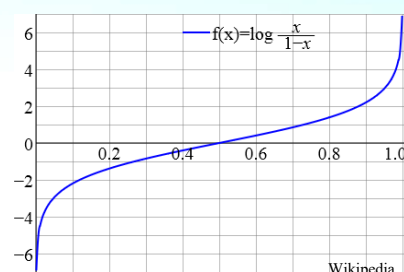
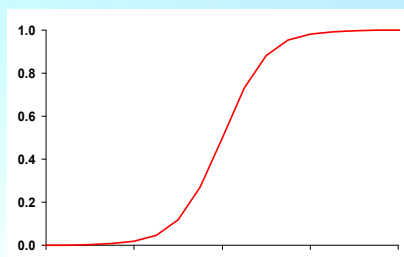
---- 对数几率 (log odds)

Xuegong Zhang

23



罗杰斯特回归与几率



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

反函数

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

$P(y|x)$ 的罗杰特 (logit) 函数

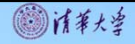
---- 对数几率 (log odds)

表示发生或不发生的概率

注意中文中的混淆词：几率、机率、概率、可能性、...

Xuegong Zhang

24



多元罗杰斯特回归

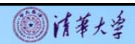
$$P(y|\mathbf{x}) = \frac{e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}{1 + e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}$$

$$\ln \left(\frac{P}{1-P} \right) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i$$

$$odds = \frac{P}{1-P} = e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

Xuegong Zhang

25



多元罗杰斯特回归

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

$$\ln \left(\frac{P}{1-P} \right) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i$$

$$odds = \frac{P}{1-P} = e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

• β_i 的解释

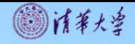
- 在其他因素不变的情况下，因素 x_i 增加一个单位带来的对数几率的增加
- 可以用来从流行病学数据中研究各种因素与患病的关系

Xuegong Zhang

26

用线性加权和激活函数来看
三种线性机器

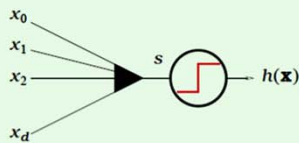
三种线性机器



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

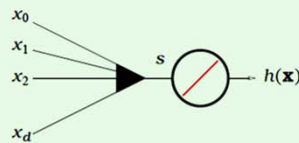
linear classification

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



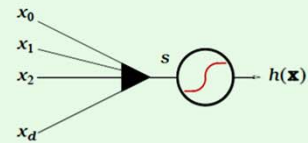
linear regression

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



logistic regression

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

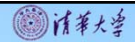


Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

Xuegong Zhang

27

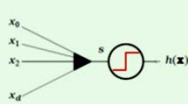
罗杰斯特回归的概率意义



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

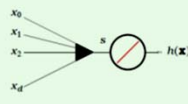
linear classification

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



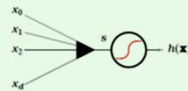
linear regression

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



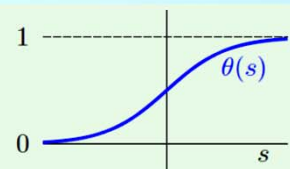
logistic regression

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



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

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



- $s = \mathbf{w}^T \mathbf{x}$ 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

28

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

29



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 $\min J_P(\boldsymbol{\alpha}) = \sum_{\mathbf{y}_j \in Y^k} (-\boldsymbol{\alpha}^T \mathbf{y}_j)$
 - 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

30



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

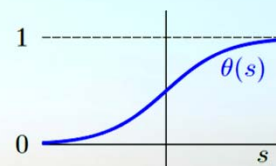
31



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 ?
 - 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

32



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 $f(\mathbf{x})$

生成模型的概念

Xuegong Zhang

33



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?

Xuegong Zhang

34



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

35



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

36

单选题 10分

设置

似然函数是谁的函数？

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

- ☐ A 是 y_i 的函数
- ☒ B 是 \mathbf{w} 的函数
- ☐ C 是 \mathbf{x}_i 的函数
- ☐ D 是 (\mathbf{x}_j, y_j) 的函数

提交

37

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

38

罗杰斯特回归的目标：似然函数最大化



- 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

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

Xuegong Zhang

39

罗杰斯特回归的目标：似然函数最大化



- 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：
取对数，防止小于1的概率相乘导致取值太小。

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

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

Xuegong Zhang

40



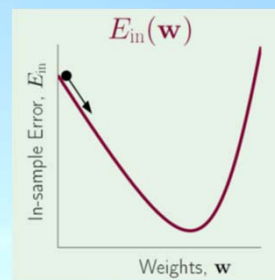
求解：梯度下降法

梯度下降的一般原理：

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

其中， η ：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

41



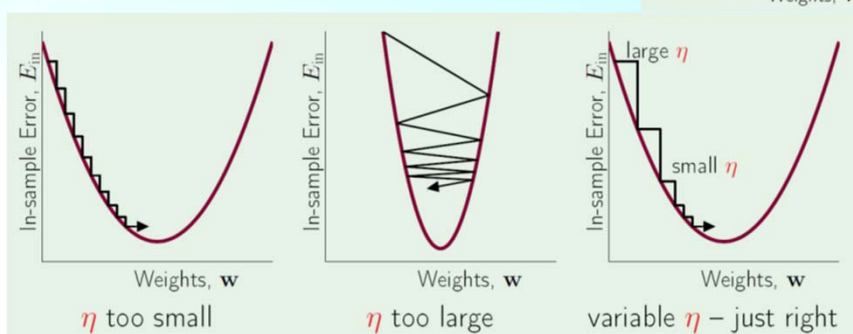
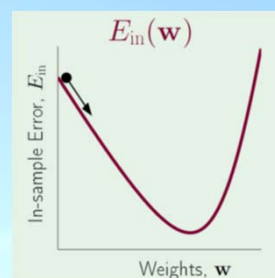
求解：梯度下降法

梯度下降的一般原理：

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

其中， η ：步长（学习率）

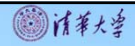
$$\hat{\mathbf{v}} = -\nabla E(\mathbf{w}(k))$$



Xuegong Zhang

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

42



罗杰斯特回归算法

1. Set $k = 0$, initialize $\mathbf{w}(0)$
2. Do
 - Compute the gradient $\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j \mathbf{x}_j}{1 + e^{y_j \mathbf{w}(k)^T \mathbf{x}_j}}$
 - 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

43



罗杰斯特回归算法

1. Set $k = 0$, initialize $\mathbf{w}(0)$
2. Do
 - Compute the gradient $\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j \mathbf{x}_j}{1 + e^{y_j \mathbf{w}(k)^T \mathbf{x}_j}}$
 - 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

44



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

45



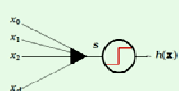
线性学习机器小结

模型

$$s = \sum_{i=1}^d w_i x_i$$

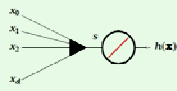
linear classification

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



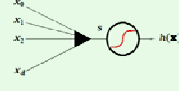
linear regression

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



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

46



休息1分钟



Xuegong Zhang

47



2.6 多类分类器

Multicategory Linear Classifiers

Xuegong Zhang

48



多类分类器

- FLD、感知器、MSE、罗杰斯特回归都是两类分类器

- 如何实现多类分类？

- 用多个两类分类器实现多类分类
- 构造多类分类器



请用弹幕回答

Xugong Zhang

49



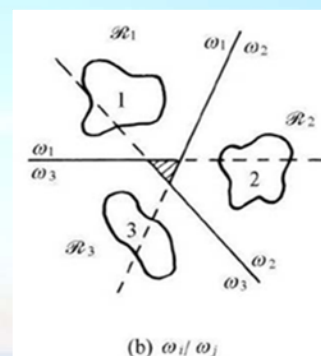
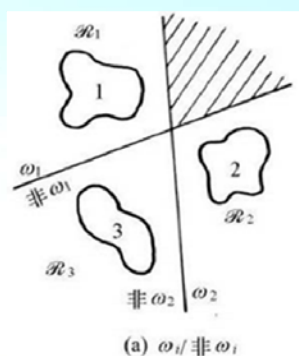
多类分类器

- FLD、感知器、MSE、罗杰斯特回归都是两类分类器

- 如何实现多类分类？

- 用多个两类分类器实现多类分类
- 构造多类分类器

- 用两类分类器实现多类分类的两种策略
 - 一对多
(One-vs-rest, one-vs-all, one-over-all)
 - 一对一 (Pairwise classification)



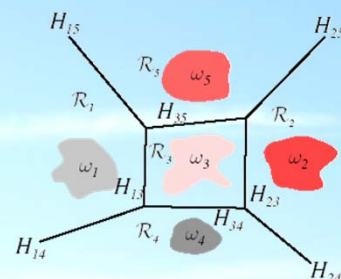
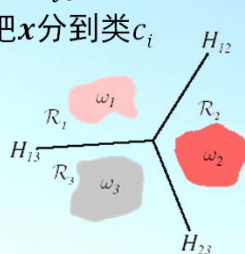
Xugong Zhang

50



多类线性机器（多类线性判别）

- 给定 C 类，定义 C 个判别函数 $g_i(\mathbf{y}) = \mathbf{a}_i^T \mathbf{y}$
- 决策：若 $g_i(\mathbf{y}) > g_j(\mathbf{y}), \forall j \neq i$ ，则把 x 分到类 c_i



from Duda et al.

Xuegong Zhang

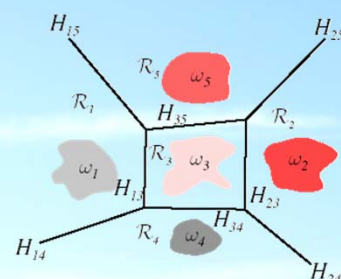
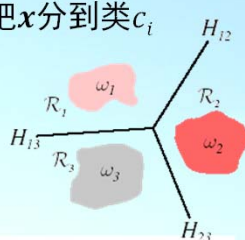
from Duda et al.

51



多类线性机器（多类线性判别）

- 给定 C 类，定义 C 个判别函数 $g_i(\mathbf{y}) = \mathbf{a}_i^T \mathbf{y}$
- 决策：若 $g_i(\mathbf{y}) > g_j(\mathbf{y}), \forall j \neq i$ ，则把 x 分到类 c_i



from Duda et al.

- 学习算法：

If $\mathbf{y}^k \in \mathcal{Y}_i$ but $\alpha_i(k)^T \mathbf{y}^k \leq \alpha_j(k)^T \mathbf{y}^k, j \neq i$

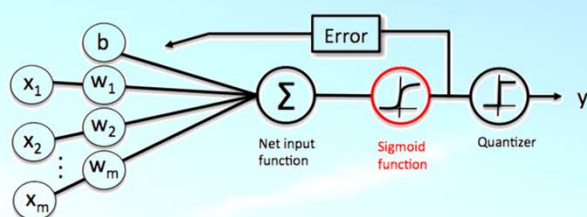
Then

$$\begin{cases} \alpha_i(k+1) = \alpha_i(k) + \rho_k \mathbf{y}^k \\ \alpha_j(k+1) = \alpha_j(k) - \rho_k \mathbf{y}^k \\ \alpha_l(k+1) = \alpha_l(k) \quad l \neq i, j \end{cases}$$

Xuegong Zhang

52

多类罗杰斯特回归与软最大 (Softmax)



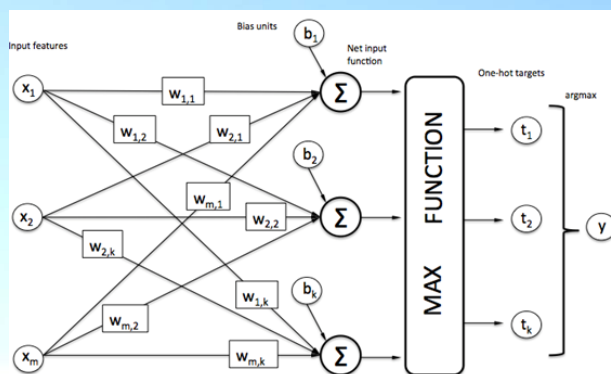
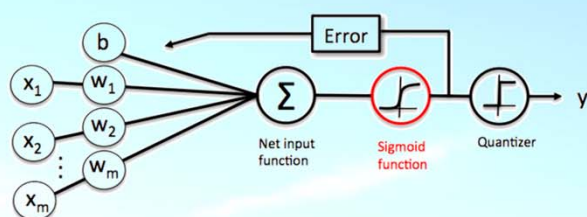
- 两类：
 - 线性加权求和 \rightarrow Logistic函数 \rightarrow 类别概率与阈值比较 \rightarrow 分类

Xuegang Zhang

Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

53

多类罗杰斯特回归与软最大 (Softmax)



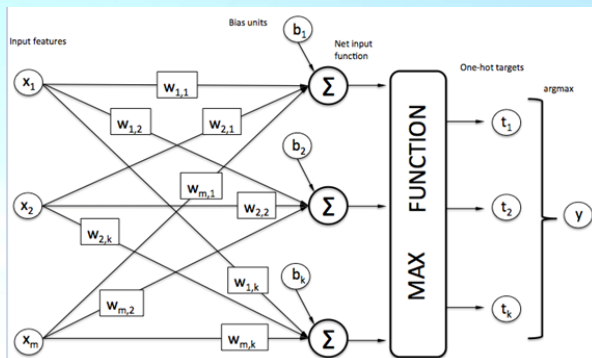
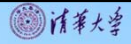
- 两类：
 - 线性加权求和 \rightarrow Logistic函数 \rightarrow 类别概率与阈值比较 \rightarrow 分类
- 多类：
 - 线性加权求和 \rightarrow 各类的Logistic函数 \rightarrow 各类概率比较 \rightarrow 分类

Xuegang Zhang

Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

54

多类罗杰斯特回归与软最大 (Softmax)



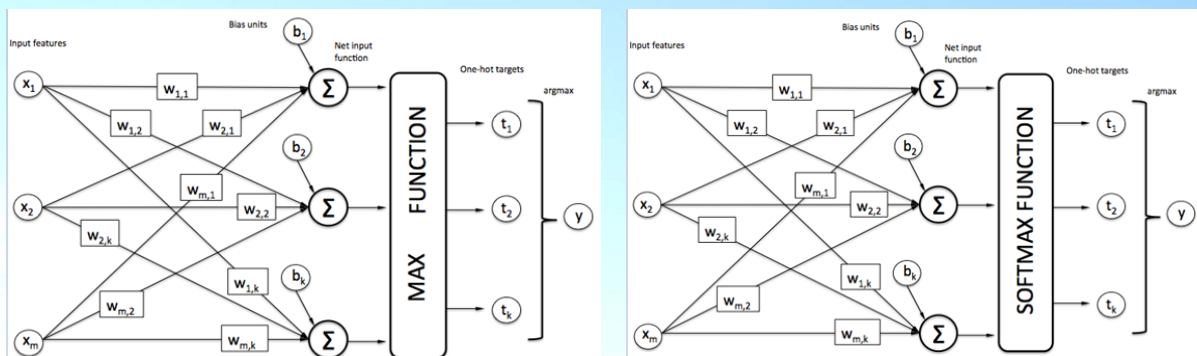
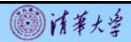
- 问题：未考虑同一样本只能属于一类

Xuegang Zhang

Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

55

多类罗杰斯特回归与软最大 (Softmax)



- 问题：未考虑同一样本只能属于一类
- 解决：多类logistic函数 \rightarrow Softmax (软最大)

Xuegang Zhang

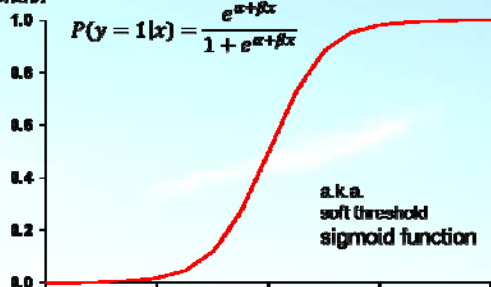
Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

56

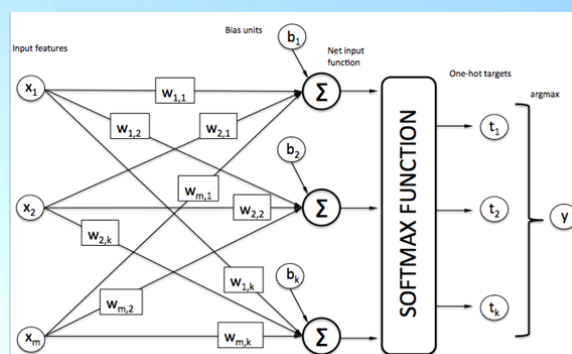


软最大Softmax

患病比例



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



$$P(y=j|\mathbf{x}) = \frac{e^{w_j \cdot \mathbf{x}}}{\sum_{k=1}^K e^{w_k \cdot \mathbf{x}}}, \quad j=1, \dots, K$$

Softmax函数（归一化指数函数）

Xuegong Zhang

注：机器学习领域中有很多random names 😊

57



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、罗杰斯特回归、Softmax
- 似然函数的概念

Xuegong Zhang

58

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



59

提交