

---

---

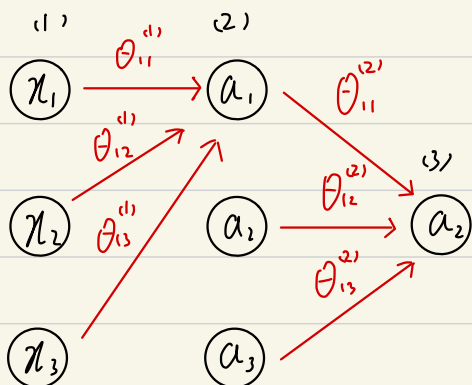
---

---

---



神经网络:



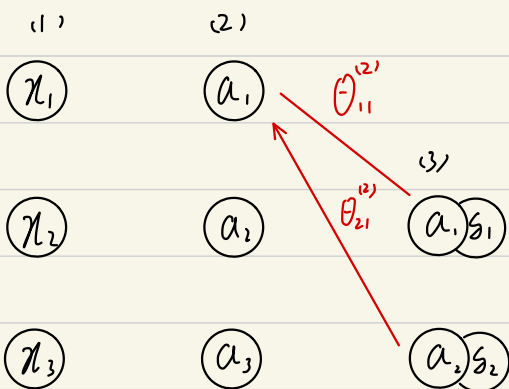
正向传播 (1) → (2)

single:  $z^{(2)} = x_1^{(1)} \theta_{11}^{(2)} + x_2^{(1)} \theta_{12}^{(2)} + x_3^{(1)} \theta_{13}^{(2)} + b_1$

$$= [\theta_{11}^{(2)} \theta_{12}^{(2)} \theta_{13}^{(2)}] \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{bmatrix} + b_1$$

all: 
$$\begin{bmatrix} \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \\ \theta_{21}^{(2)} & \theta_{22}^{(2)} & \theta_{23}^{(2)} \\ \theta_{31}^{(2)} & \theta_{32}^{(2)} & \theta_{33}^{(2)} \end{bmatrix} \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$\theta x + b = z, \quad g(z) = a$$



反向传播: (3) → (2)

$$s^{(3)} = (a^{(3)} - y) g'(z^{(3)})$$

single:  $s_1^{(2)} = s_1^{(3)} \theta_{11}^{(2)} g'(z_1^{(2)}) + s_2^{(3)} \theta_{21}^{(2)} g'(z_2^{(2)})$

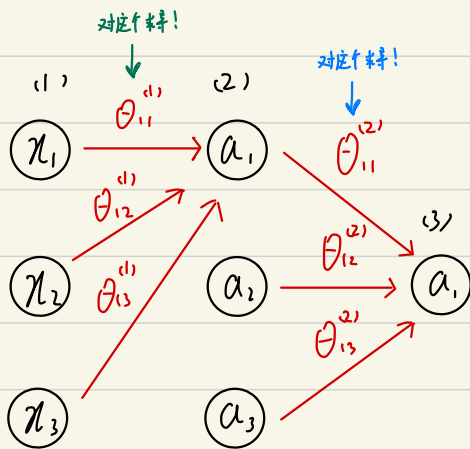
$$= [\theta_{11}^{(2)} \theta_{21}^{(2)}] \begin{bmatrix} s_1^{(3)} \\ s_2^{(3)} \end{bmatrix} \times g'(z^{(2)})$$

偏导:  $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} s_i^{(l+1)} \quad \frac{\partial}{\partial b_i^{(l)}} J(\theta) = s_i^{(l+1)}$

$$\frac{\partial}{\partial \theta^{(l)}} J(\theta) = \begin{bmatrix} x_1 s_1^{(l+1)} & x_2 s_1^{(l+1)} & x_3 s_1^{(l+1)} \\ x_1 s_2^{(l+1)} & x_2 s_2^{(l+1)} & x_3 s_2^{(l+1)} \\ x_1 s_3^{(l+1)} & x_2 s_3^{(l+1)} & x_3 s_3^{(l+1)} \end{bmatrix} = \begin{bmatrix} s_1^{(l+1)} \\ s_2^{(l+1)} \\ s_3^{(l+1)} \end{bmatrix} [x_1 \ x_2 \ x_3]$$

all: 
$$\begin{bmatrix} \theta_{11}^{(2)} & \theta_{21}^{(2)} \\ \theta_{12}^{(2)} & \theta_{22}^{(2)} \\ \theta_{13}^{(2)} & \theta_{23}^{(2)} \end{bmatrix} \begin{bmatrix} s_1^{(3)} \\ s_2^{(3)} \end{bmatrix} \times \begin{bmatrix} g'(z_1^{(2)}) \\ g'(z_2^{(2)}) \\ g'(z_3^{(2)}) \end{bmatrix} = \begin{bmatrix} s_1^{(2)} \\ s_2^{(2)} \\ s_3^{(2)} \end{bmatrix}$$

$$\theta^{(2)T} s^{(3)} \times g'(z^{(2)}) = s^{(2)}$$



$$① J(\theta) \approx \frac{1}{2} (a^{(3)} - y)^2$$

$$\frac{dJ(\theta)}{da^{(3)}} = a^{(3)} - y$$

$$② a^{(3)} = g(z^{(3)})$$

$$\frac{da^{(3)}}{dz^{(3)}} = g'(z^{(3)}) = a^{(3)}(1 - a^{(3)})$$

$$③ z^{(3)} = a_1^{(2)} \theta_{11}^{(2)} + a_2^{(2)} \theta_{12}^{(2)} + a_3^{(2)} \theta_{13}^{(2)} + b$$

$$\frac{dz^{(3)}}{d\theta_{11}^{(2)}} = a_1^{(2)} + 0 + 0 + 0$$

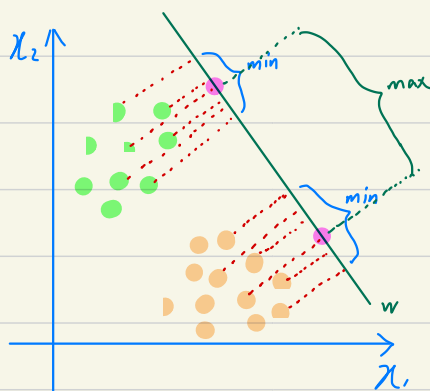
$$④ \frac{dJ(\theta)}{d\theta_{11}^{(2)}} = \frac{dJ(\theta)}{da^{(3)}} \cdot \frac{da^{(3)}}{dz^{(3)}} \cdot \frac{dz^{(3)}}{d\theta_{11}^{(2)}}$$

$$= \underbrace{(a^{(3)} - y) \cdot g'(z^{(3)})}_{\delta^{(3)}} \cdot a_1^{(2)}$$

$$\begin{aligned} \frac{dJ(\theta)}{d\theta_{11}^{(2)}} &= \frac{dJ(\theta)}{da^{(3)}} \cdot \frac{da^{(3)}}{dz^{(3)}} \cdot \frac{dz^{(3)}}{da_1^{(2)}} \\ &\quad \cdot \frac{da_1^{(2)}}{dz_1^{(2)}} \cdot \frac{dz_1^{(2)}}{d\theta_{11}^{(2)}} \\ &= \underbrace{(a^{(3)} - y) \cdot g'(z^{(3)}) \cdot \theta_{11}^{(2)}}_{\delta_1^{(3)}} \cdot g'(z_1^{(2)}) \cdot x_1 \\ &\quad \underbrace{\delta_1^{(2)}}_{\delta_1^{(2)}} \cdot x_1 \end{aligned}$$

# 线性判别分析:

思想: 如图所示, 要找出一条线, 使样本点映射后易于区分(类内小, 类间大)



$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}_{n \times p} \quad (x_i \text{ 均为列向量})$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

设  $w$  为  $(p \times 1)$  列向量

推导: 映射后的值:  $z_i = w^T x_i$ , 均值:  $\bar{z} = \frac{1}{N} \sum_{i=1}^N z_i = \frac{1}{N} \sum_{i=1}^N w^T x_i$

映射后的方差:  $S = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$

$$\text{则类 1: } \begin{cases} \bar{z}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i \\ S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (z_i - \bar{z}_1)(z_i - \bar{z}_1)^T \end{cases}, \text{类 2: } \begin{cases} \bar{z}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i \\ S_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} (z_i - \bar{z}_2)(z_i - \bar{z}_2)^T \end{cases}$$

类间:  $(\bar{z}_1 - \bar{z}_2)^2$ , 类内:  $S_1 + S_2$ , 目标: 类内小, 类间大  $\rightarrow$  loss 函数

$$\text{loss 函数: } J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{S_1 + S_2}, \quad \text{目标: } \operatorname{argmax}(J(w))$$

$$J(w): \text{分子} = \left( \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i \right)^2 = \left[ w^T \left( \frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} x_i \right) \right]^2 = \left[ w^T (\bar{x}_{c1} - \bar{x}_{c2}) \right]^2$$

$$\text{分母以 } S_1 \text{ 为例: } S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} \left( w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j \right) \left( w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j \right)^T \quad \text{把 } w^T \text{ 提出来}$$

$$= \frac{1}{N_1} \sum_{i=1}^{N_1} w^T (x_i - \bar{x}_{c1}) (x_i - \bar{x}_{c1})^T w = w^T \left[ \frac{1}{N_1} \sum_{i=1}^{N_1} (x_i - \bar{x}_{c1}) (x_i - \bar{x}_{c1})^T \right] w = w^T S_{c1} w$$

$$\therefore \text{分母} = w^T S_{c1} w + w^T S_{c2} w = w^T (S_{c1} + S_{c2}) w$$

$$\therefore J(w) = \frac{w^T (\bar{x}_{c1} - \bar{x}_{c2}) (\bar{x}_{c1} - \bar{x}_{c2})^T w}{w^T (S_{c1} + S_{c2}) w}$$

$$\therefore J(w) = \frac{W^T (\bar{x}_{c1} - \bar{x}_{c2}) (\bar{x}_{c1} - \bar{x}_{c2})^T W}{W^T (S_{c1} + S_{c2}) W}, \quad \begin{aligned} \frac{1}{2} S_b &= (\bar{x}_{c1} - \bar{x}_{c2}) (\bar{x}_{c1} - \bar{x}_{c2})^T && \text{类间方差} \\ S_w &= (S_{c1} + S_{c2}) && \text{类内方差} \end{aligned}$$

$$P1) J(w) = W^T S_b W \cdot (W^T S_w W)^{-1}$$

$$\frac{\partial J(w)}{\partial w} = 2 S_b W (W^T S_w W)^{-1} + W^T S_b W \cdot (-1) \cdot (W^T S_w W)^{-2} \cdot 2 S_w \cdot W = 0, \text{ 可得}$$

$$S_b W \cdot \underbrace{(W^T S_w W)}_{\in \mathbb{R}} = \underbrace{W^T S_b W}_{\in \mathbb{R}} S_w \cdot W \quad \begin{aligned} S_w: (p \times p) \quad S_b: (p \times p) \quad W: (n \times p) \quad W^T: (p \times n) \end{aligned}$$

求的是方向向量

$$S_w W = \frac{W^T S_w W}{W^T S_b W} S_b W \rightarrow W = \frac{W^T S_w W}{W^T S_b W} S_w^{-1} S_b W \propto S_w^{-1} S_b W$$

$$S_b W = (\bar{x}_{c1} - \bar{x}_{c2}) \underbrace{(\bar{x}_{c1} - \bar{x}_{c2})^T}_{1 \times p} \cdot \underbrace{W}_{p \times 1} \propto (\bar{x}_{c1} - \bar{x}_{c2})$$

$$\therefore W \propto S_w^{-1} (\bar{x}_{c1} - \bar{x}_{c2})$$

$$\text{经验: } b = -\frac{1}{2} W^T (\bar{x}_{c1} - \bar{x}_{c2})$$

$$\text{模型 } y = W^T x + b$$

# 线性回归

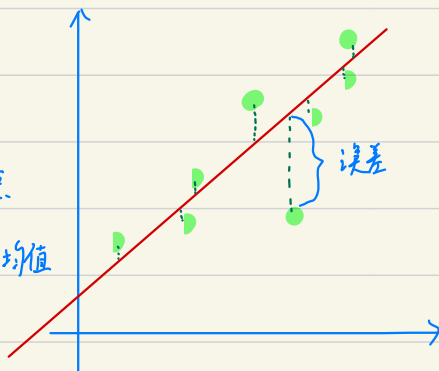
使用场景：目标值与特征线性

思想：使用一个多元一次方程去拟合所有数据点

度量尺度：均方误差(MSE)，即数据误差平方和的平均值

MSE越小，拟合效果越佳

$$MSE: \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$



$$\text{令: } X = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_m \\ 1 & x_1 & x_2 & \dots & x_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1 & x_2 & \dots & x_m \end{bmatrix}_{(n \times m)} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}_{(m \times 1)} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{(n \times 1)} \quad \hat{y} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_m \\ 1 & x_1 & x_2 & \dots & x_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1 & x_2 & \dots & x_m \end{bmatrix}_{(n \times 1)} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix} = X\theta$$

$$\hat{y} = X\theta, \quad \text{损失函数 } J(\theta) = \frac{1}{m} \sum_{i=1}^m (X\theta - y_i)^2, \quad \text{目标: } \min J(\theta)$$

1. 最小二乘法:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (X\theta - y_i)^2 = \frac{1}{m} (X\theta - y)^T (X\theta - y) = \frac{1}{m} [(X\theta)^T - y^T] (X\theta - y)$$

$$= \frac{1}{m} (\theta^T X^T - y^T) (X\theta - y) = \frac{1}{m} (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y)$$

$\therefore y^T X\theta$  与  $\theta^T X^T y$  互为转置，但均为标量

$$\therefore J(\theta) = \frac{1}{m} (\theta^T X^T X\theta - 2\theta^T X^T y + y^T y)$$

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{1}{m} \cdot \frac{\partial}{\partial \theta} (\theta^T X^T X\theta - 2\theta^T X^T y + y^T y) = \frac{1}{m} \cdot (-2X^T y + 2X^T X\theta)$$

$$\text{令 } \frac{\partial J(\theta)}{\partial \theta} = 0, \quad \frac{1}{m} (-2X^T y + 2X^T X\theta) = 0 \quad \text{可得}$$

$$X^T y = X^T X\theta \Rightarrow \theta = (X^T X)^{-1} X^T y$$

## 2. 梯度下降法

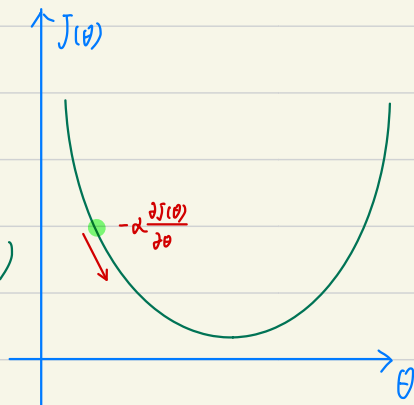
左图为  $J(\theta)$  图象, 要使  $J(\theta)$  最小, 就要令其向 **箭头**

方向移动  $\Rightarrow$  减去 **学习率**  $\times$  **偏导**

推导:

$$\text{由 (1) 可得 } \frac{\partial J(\theta)}{\partial \theta} = \frac{1}{m} (2X^T X \theta - 2X^T y) = \frac{2}{m} X^T (X \theta - y)$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \cdot \frac{2}{m} X^T (X \theta - y)$$



## 归一化/标准化

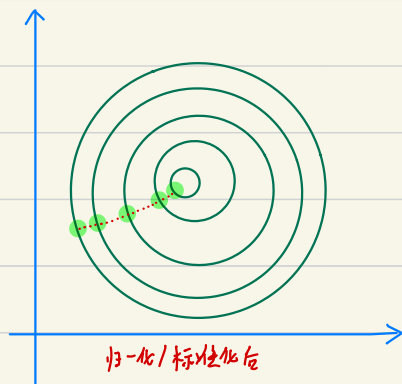
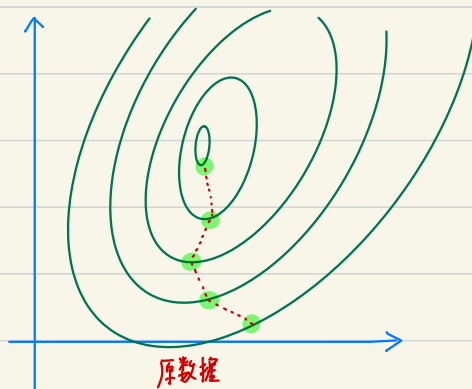
使用梯度下降法时, 归一化/标准化可加快梯度下降的求解速度, 如右图

归一化:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

标准化:

$$x' = \frac{(x - \mu)}{\sigma} \quad \mu \text{ 为期望 } \sigma \text{ 为标准差}$$



## 正则化

常见的向量范数:

$L_0$  范数:  $\|x\|_0$  表示向量  $x$  中非零元素的个数

$L_1$  范数:  $\|x\|_1 = \sum_{i=1}^n |x_i|$  表示非零元素的绝对值之和

$L_2$  范数:  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$

$L_p$  范数:  $\|x\|_p = \sqrt[p]{\sum_{i=1}^n x_i^p}$

过拟合:



红: 恰好,  $y = w_0 + w_1 x_1$

绿: 过拟合,  $y = w_0 + w_1 x_1 + w_2 x_2 + \dots$

如何令  $y = w_0 + w_1 x_1 \rightarrow y = w_0 + w_1 x_1 + w_2 x_2 + \dots$  呢?

减少向量  $w$  的个数 or 令某些参数减小  $\rightarrow$  正则化

推理: 损失函数  $J(w, x, y)$

减小某些  $w$  的值, 但又不知道该减少哪一个  $w$ , 因此可以设 限制条件

$$\textcircled{1} \begin{cases} |w_1| + |w_2| + |w_3| + \dots + |w_n| = \|w\|_1 \leq C \\ \arg \min J(w, x, y) \end{cases} \quad \textcircled{2} \begin{cases} \sqrt{w_1^2 + w_2^2 + w_3^2 + \dots + w_n^2} = \|w\|_2 \leq C \\ \arg \min J(w, x, y) \end{cases}$$

先研究①: 引入拉格朗日函数, 可得

$$L(w) = \underset{\min}{J(w, x, y)} + \lambda (\|w\|_1 - C)$$

要令  $L(w)$  最小, 则令  $\frac{\partial L(w)}{\partial \lambda} = 0$ , 设其解为  $\lambda^*$



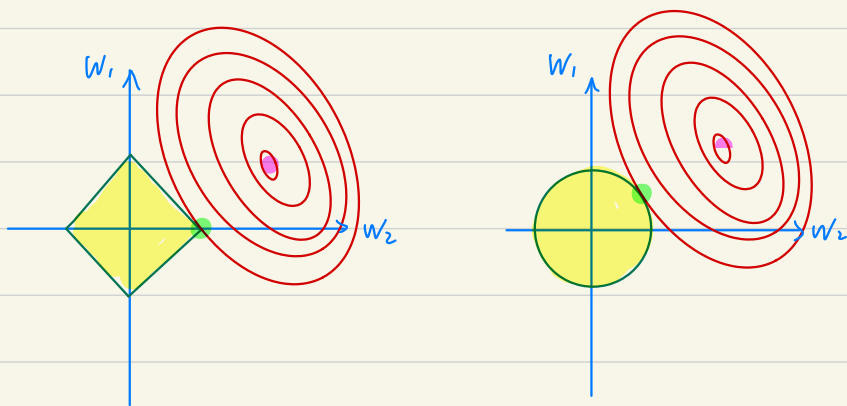
$$\begin{aligned}
 \text{则 } \min L(w, \lambda) &= \min J(w, x, y) + \lambda^* (\|w\|_1 + C) \\
 &= \min J(w, x, y) + \lambda^* \|w\|_1 + \lambda^* C \\
 &\propto \min J(w, x, y) + \lambda^* \|w\|_1
 \end{aligned}$$

$L_1$  正则化

$$\begin{aligned}
 \therefore \text{损失函数: } &J(w, x, y) + \lambda \|w\|_1 \\
 &= J(w, x, y) + \lambda \sum_{i=1}^n |w_i|
 \end{aligned}$$

$L_2$  正则化同理

几何解释:



- :  $w_1$  与  $w_2$  的取值范围
- : 损失函数最小的点
- :  $w_1$  与  $w_2$  取值范围内与 ● 最近的点, 也是  $w_1$  与  $w_2$  取值点

## 主成分分析(PCA)

作用: 降维降噪

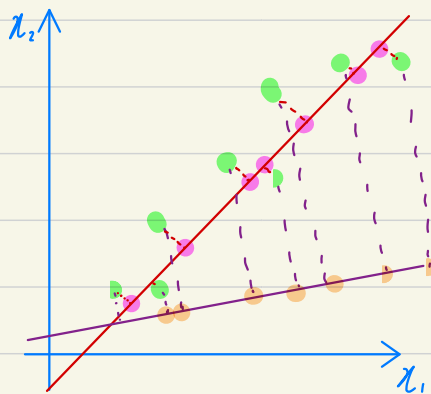
目的: 找到一条线, 使映射后数据之间分隔大

如: 粉色的分隔比 橙色的大  $\rightarrow$

设:

$$X = \begin{bmatrix} x_1^1 & x_1^2 & \dots & x_1^n \\ x_2^1 & x_2^2 & \dots & x_2^n \\ \vdots & \vdots & & \vdots \\ x_m^1 & x_m^2 & \dots & x_m^n \end{bmatrix} \quad W = [w_1 \ w_2 \ \dots \ w_k]$$

$(m \times n)$   $(n \times k)$   $w_i (n \times 1)$



推导: 分隔大  $\rightarrow$  方差大

$$S_j = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2 \quad (\text{单个特征的方差})$$

① 均值归零: 所有数据减去其对应特征的均值

则  $\bar{x} = [0 \ 0 \ 0 \ \dots \ 0]$ ,  $S(w) = \left\{ \begin{array}{l} \frac{1}{m} \sum_{i=1}^m (x_i^1)^2 \\ \frac{1}{m} \sum_{i=1}^m (x_i^2)^2 \\ \vdots \\ \frac{1}{m} \sum_{i=1}^m (x_i^n)^2 \end{array} \right\}$  求和

目标:  $\arg \max S(w)$ , 求偏导, 可得

$$\frac{\partial S(w)}{\partial w} = \begin{bmatrix} \frac{\partial S(w_1)}{\partial w_1} \\ \frac{\partial S(w_2)}{\partial w_2} \\ \vdots \\ \frac{\partial S(w_n)}{\partial w_n} \end{bmatrix} = \frac{2}{m} \begin{bmatrix} \sum_{i=1}^m (x_i^1 w_1 + x_i^2 w_2 + \dots + x_i^n w_n) x_i^1 \\ \sum_{i=1}^m (x_i^1 w_1 + x_i^2 w_2 + \dots + x_i^n w_n) x_i^2 \\ \vdots \\ \sum_{i=1}^m (x_i^1 w_1 + x_i^2 w_2 + \dots + x_i^n w_n) x_i^n \end{bmatrix} = \frac{2}{m} \begin{bmatrix} \sum_{i=1}^m (x_i^1 w) x_i^1 \\ \sum_{i=1}^m (x_i^2 w) x_i^2 \\ \vdots \\ \sum_{i=1}^m (x_i^n w) x_i^n \end{bmatrix}$$

$$= \frac{2}{m} \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \end{bmatrix} \begin{bmatrix} x^{(1)} w \\ x^{(2)} w \\ \vdots \\ x^{(n)} w \end{bmatrix} = \frac{2}{m} \cdot X^T (Xw)$$

使用梯度上升法,  $w_{\text{new}} = w_{\text{old}} + \alpha \frac{\partial S(w)}{\partial w}$ , 取得最大值。

## 基础

条件概率:  $P(B|A) = \frac{P(AB)}{P(A)}$   $P(AB)$  为 A, B 事件一同发生的概率  
 $P(B|A)$  为在 A 发生的前提下 B 发生的概率

贝叶斯公式:  $P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})}$  ( $\bar{A}$  表示非 A)

似然函数: 样本集  $D = \{x_1, x_2, \dots, x_n\}$

$$l(\theta) = P(D|\theta) = P(x_1, x_2, \dots, x_n|\theta) = P(x_1|\theta) \cdot P(x_2|\theta) \cdots P(x_n|\theta) = \prod_{i=1}^n P(x_i|\theta)$$

极大似然估计:  $\arg\max l(\theta) = \arg\max \ln l(\theta)$   
 $\ln l(\theta) = \ln \prod_{i=1}^n P(x_i|\theta) = \sum_{i=1}^n \ln P(x_i|\theta)$

一维高斯分布:  $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (-\infty < x < \infty)$

n 维高斯分布:  $f(x) = \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \cdot e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$

logistic [2] 1]:

$$\text{Sigmoid 概率函数推导: } p(c_i|x) = \frac{P(x|c_i)P(c_i)}{P(x|c_1)P(c_1) + P(x|c_0)P(c_0)} = \frac{1}{1 + \frac{P(x|c_0)P(c_0)}{P(x|c_1)P(c_1)}}$$

$$= \frac{1}{1 + \exp(-\ln \frac{P(x|c_1)P(c_1)}{P(x|c_0)P(c_0)})} \quad \text{令 } a = \ln \frac{P(x|c_1)P(c_1)}{P(x|c_0)P(c_0)} \quad \text{得} \quad \frac{1}{1 + \exp(-a)}$$

$a$  与线性函数  $\theta x + b$  有何关系? 设服从同协方差高斯分布

$$P(x|c_k) = \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\right] \quad \text{其中 } \mu \text{ 为数学期望, } \Sigma \text{ 为协方差矩阵}$$

$$a = \ln \frac{P(x|c_1)P(c_1)}{P(x|c_0)P(c_0)} = \ln \frac{P(x|c_1)P(c_1)}{P(x|c_0)P(c_0)} + \ln \frac{P(c_1)}{P(c_0)}$$

$$= \ln \frac{\frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \cdot \exp\left[-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)\right]}{\frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \cdot \exp\left[-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2)\right]} + \ln \frac{P(c_1)}{P(c_0)}$$

$$= \cancel{\ln \exp} \left[ -\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + \frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2) \right] + \ln \frac{P(c_1)}{P(c_0)}$$

$$= \left[ -\frac{1}{2} \cancel{x^T \Sigma^{-1} x} + \underline{(\Sigma^{-1} \mu_1)^T x} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \cancel{x^T \Sigma^{-1} x} - \underline{(\Sigma^{-1} \mu_2)^T x} + \frac{1}{2} \Sigma^{-1} \mu_2 \right] + \ln \frac{P(c_1)}{P(c_0)}$$

$$= \underline{[\Sigma^{-1}(\mu_1 - \mu_2)]^T x} + \underline{\left[ -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(c_1)}{P(c_0)} \right]}$$

$\theta$

$b$

$$= \theta^T x + b$$

构造损失函数：利用极大似然估计

$$l(w) = \prod_{i=1}^m P(c_i | x^{(i)}; w) = \prod_{i=1}^m (h_w(x^{(i)})^{y^{(i)}} \cdot [1 - h_w(x^{(i)})]^{1-y^{(i)}})$$

$$L(w) = -\ln l(w) = -\sum_{i=1}^m [y^{(i)} \ln h_w(x^{(i)}) + (1-y^{(i)}) \ln [1 - h_w(x^{(i)})]]$$

$$J(w) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \ln h_w(x^{(i)}) + (1-y^{(i)}) \ln [1 - h_w(x^{(i)})]] \quad \text{又解=交叉熵}$$

目标：  $\operatorname{argmin} J(w)$ ，使用梯度下降法 设  $g(x) = \frac{1}{1 + \exp(-x)}$ ， $g'(x) = g(x)(1-g(x))$

$$\frac{\partial J(w)}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \frac{1}{h_w(x^{(i)})} \cdot \frac{\partial h_w(x^{(i)})}{\partial w_j} - (1-y^{(i)}) \frac{1}{1-h_w(x^{(i)})} \cdot \frac{\partial h_w(x^{(i)})}{\partial w_j} \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \frac{1}{h_w(x^{(i)})} - (1-y^{(i)}) \frac{1}{1-h_w(x^{(i)})} \right] \frac{\partial g(x^{(i)} w)}{\partial w_j} \rightarrow g(x^{(i)} w) (1-g(x^{(i)} w)) \cdot x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} (1-g(x^{(i)} w)) - (1-y^{(i)}) (g(x^{(i)} w))] x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - \cancel{y^{(i)} g(x^{(i)} w)} - g(x^{(i)} w) + \cancel{y^{(i)} g(x^{(i)} w)}] x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_w(x^{(i)})) x_j^{(i)} = \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\text{偏导矩阵：} \frac{\partial J(w)}{\partial w} = \frac{1}{m} X^T (h_w(X) - y)$$

$$\text{梯度下降：} \quad w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial J(w)}{\partial w}$$

多分类: softmax 回归 (one vs one)

$$P(c_k|x) = \frac{P(x|c_k)P(c_k)}{P(x|c_1)P(c_1) + P(x|c_2)P(c_2) + \dots + P(x|c_K)P(c_K)}$$

$$= \frac{P(x|c_k)P(c_k)}{\sum_{i=1}^K P(x|c_i)P(c_i)} = \frac{e^{\ln(P(x|c_k)P(c_k))}}{\sum_{i=1}^K e^{\ln[P(x|c_i)P(c_i)]}} = \frac{e^{a_k}}{\sum_{i=1}^K e^{a_i}}$$

其中  $a_k = \ln[P(x|c_k)P(c_k)]$ , 设服从高斯分布

$$P(x|c_k) = \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\right] \quad \text{其中 } \mu \text{ 为数学期望, } \Sigma \text{ 为协方差矩阵}$$

$$a = \ln P(x|c_k)P(c_k) = \ln P(x|c_k) + \ln P(c_k)$$

$$= \ln \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \cdot \exp\left[-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\right] + \ln P(c_k)$$

$$= \left[-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\right] + \ln P(c_k) + \ln\left(\frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{|\Sigma|^{1/2}}\right)$$

$$= \underbrace{(\Sigma^{-1} \mu_k)^T x}_{a} - \underbrace{\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k}_{b} + \ln P(c_k) \quad \text{常数}$$

$$a = \theta^T x + b,$$

同样, 利用极大似然估计

$$\text{似然函数: } L(w) = \prod_{i=1}^m \prod_{k=1}^K P(c_k|x; w) = \prod_{i=1}^m \prod_{k=1}^K \left[ \frac{\exp(x w_k)}{\sum_{l=1}^K \exp(x w_l)} \right]^{y_k^{(i)}}, \quad k \in \{1, 2, \dots, K\}$$

$$J(w) = -\ln L(w) = -\ln \prod_{i=1}^m \prod_{k=1}^K \left[ \frac{\exp(x w_k)}{\sum_{l=1}^K \exp(x w_l)} \right]^{y_k^{(i)}} = \sum_{i=1}^m \sum_{k=1}^K (\ln \exp(x w_k) - \ln \sum_{l=1}^K \exp(x w_l)) \cdot y_k^{(i)}$$

$$\text{对于每个数据来说, } J_m(w) = \sum_{k=1}^K (x w_k - \ln \sum_{i=1}^K e^{x w_i}) \cdot y_m$$

$$\frac{\partial J_m(w)}{\partial w} = \underbrace{x_m^T}_{(n \times 1)} \cdot \left( y_m - \frac{e^{xw}}{\sum_{i=1}^K e^{x w_i}} \right) \quad (1 \times K)$$

再用梯度上升法即可