



2022年秋季 《机器学习概论》课程

# 第三章：线性模型

主讲：连德富 特任教授 | 博士生导师

邮箱：[liandefu@ustc.edu.cn](mailto:liandefu@ustc.edu.cn)

手机：13739227137

主页：<http://staff.ustc.edu.cn/~liandefu>

# 基本形式

- 线性模型一般形式

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \cdots + w_dx_d + b$$

$\mathbf{x} = (x_1; x_2; \cdots, x_d)$ 是由属性描述的示例，其中 $x_i$ 是 $\mathbf{x}$ 在第 $i$ 个属性上的取值

- 向量形式

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

$\mathbf{w} = (w_1; w_2; \cdots, w_d)$ 是属性的权重

# 线性模型优点

- 形式简单、易于建模
- 可解释性
- 非线性模型的基础：引入层级结构或高维映射
- 一个例子
  - 综合考虑色泽、根蒂和敲声来判断西瓜好不好
  - 其中根蒂的系数最大，表明根蒂最要紧；而敲声的系数比色泽大，说明敲声比色泽更重要

$$f_{\text{好瓜}}(\mathbf{x}) = 0.2 \cdot x_{\text{色泽}} + 0.5 \cdot x_{\text{根蒂}} + 0.3 \cdot x_{\text{敲声}} + 1$$

# 线性回归

- 给定数据集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ,  
其中  $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots, x_{id})$ ,  $y_i \in \mathbb{R}$
- 线性回归目标
  - 学得一个线性模型以尽可能准确地预测实值输出标记
- 离散属性处理
  - 有“序”关系
    - 连续化为连续值
  - 无“序”关系
    - 有k个属性值，则转换为k维向量

# 线性回归

- 单一属性的线性回归目标

$$f(x_i) = w x_i + b \text{ 使得 } f(x_i) \approx y_i$$

- 参数/模型估计：最小二乘法（least square method）

$$\begin{aligned}(w^*, b^*) &= \arg \min_{w, b} \sum_{i=1}^m (f(x_i) - y_i)^2 \\ &= \arg \min_{w, b} \sum_{i=1}^m (w x_i + b - y_i)^2\end{aligned}$$

# 线性回归 - 最小二乘法

- 最小化均方误差

$$E(w, b) = \sum_{i=1}^m (y_i - wx_i - b)^2$$

- 分别对 $w$ 和 $b$ 求导, 可得

$$\frac{\partial E(w, b)}{\partial w} = 2 (w \sum_i x_i^2 - \sum_i (y_i - b)x_i)$$

$$\frac{\partial E(w, b)}{\partial b} = 2 (mb - \sum_i (y_i - wx_i))$$

# 线性回归 - 最小二乘法

- 令导数梯度等于0，得到闭形式解

$$\begin{aligned}\frac{\partial E(w, b)}{\partial w} &= w \sum_i x_i^2 - \sum_i (y_i - b) x_i \\&= w \sum_i x_i^2 - \sum_i (y_i - \bar{y} + w\bar{x}) x_i \\&= w (\sum_i x_i^2 - \bar{x} \sum_i x_i) - \sum_i (y_i - \bar{y}) x_i \\&= w (\sum_i x_i^2 - \bar{x} \sum_i x_i) - (\sum_i y_i x_i - \sum_i \bar{y} x_i) \\&= w (\sum_i x_i^2 - \bar{x} \sum_i x_i) - (\sum_i y_i x_i - \sum_i y_i \bar{x}) \\&= w (\sum_i x_i^2 - \bar{x} \sum_i x_i) - \sum_i y_i (x_i - \bar{x})\end{aligned}$$

$$b = \frac{1}{m} \sum_i (y - wx_i) = \bar{y} - w\bar{x}$$

$$w = \frac{\sum_i y_i (x_i - \bar{x})}{(\sum_i x_i^2 - \frac{1}{m} (\sum_i x_i)^2)}$$

# 多元线性回归

- 给定数据集

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

其中  $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots, x_{id}), y_i \in \mathbb{R}$

- 多元线性回归目标

$$f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i + b \text{ 使得 } f(\mathbf{x}_i) \approx y_i$$



# 多元线性回归

- 把 $\mathbf{w}$ 和 $b$ 吸收入向量形式  $\hat{\mathbf{w}} = (\mathbf{w}; b)$ ，数据集表示为

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} & 1 \\ x_{21} & x_{22} & \cdots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix}$$

$$\mathbf{y} = (y_1; y_2; \cdots; y_m)$$

# 多元线性回归 - 最小二乘法

- 最小二乘法 (least square method)

$$\hat{\mathbf{w}}^* = \arg \min_{\hat{\mathbf{w}}} \|\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}\|_2^2 = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

$E(\hat{\mathbf{w}})$

- 求 $E(\hat{\mathbf{w}})$ 关于变量 $\hat{\mathbf{w}}$ 的导数得到

$$\nabla_{\hat{\mathbf{w}}} E(\hat{\mathbf{w}}) = 2\mathbf{X}^\top (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})$$

# 多元线性回归 - 满秩讨论

- $\mathbf{X}^\top \mathbf{X}$  是满秩矩阵或正定矩阵, 则

$$\nabla_{\hat{\mathbf{w}}} E(\hat{\mathbf{w}}) = 2\mathbf{X}^\top (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) = 0 \quad \rightarrow \quad \boxed{\hat{\mathbf{w}}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$

- 把  $\hat{\mathbf{w}}^*$  代回  $f(\mathbf{x}_i)$ , 线性回归模型为

$$\boxed{f(\hat{\mathbf{x}}_i) = \mathbf{x}_i (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$

- 如果  $\mathbf{X}^\top \mathbf{X}$  不是满秩矩阵
  - 根据归纳偏好选择解
  - 引入正则化

# 一元线性回归

- 重新考虑一个特征的情形

$$\mathbf{X} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{pmatrix} \quad \mathbf{X}^\top \mathbf{X} = \begin{pmatrix} \sum_i x_i^2 & \sum_i x_i \\ \sum_i x_i & m \end{pmatrix} \quad \mathbf{X}^\top \mathbf{y} = \begin{pmatrix} \sum_i x_i y_i \\ \sum_i y_i \end{pmatrix}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} = \frac{1}{m \sum_i x_i^2 - (\sum_i x_i)^2} \begin{pmatrix} m & -\sum_i x_i \\ -\sum_i x_i & \sum_i x_i^2 \end{pmatrix}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \frac{1}{m \sum_i x_i^2 - (\sum_i x_i)^2} \begin{pmatrix} m \sum_i x_i y_i - \sum_i x_i \sum_j y_j \\ \sum_i x_i^2 \sum_j y_j - \sum_i x_i y_i \sum_j x_j \end{pmatrix}$$

# 一元线性回归

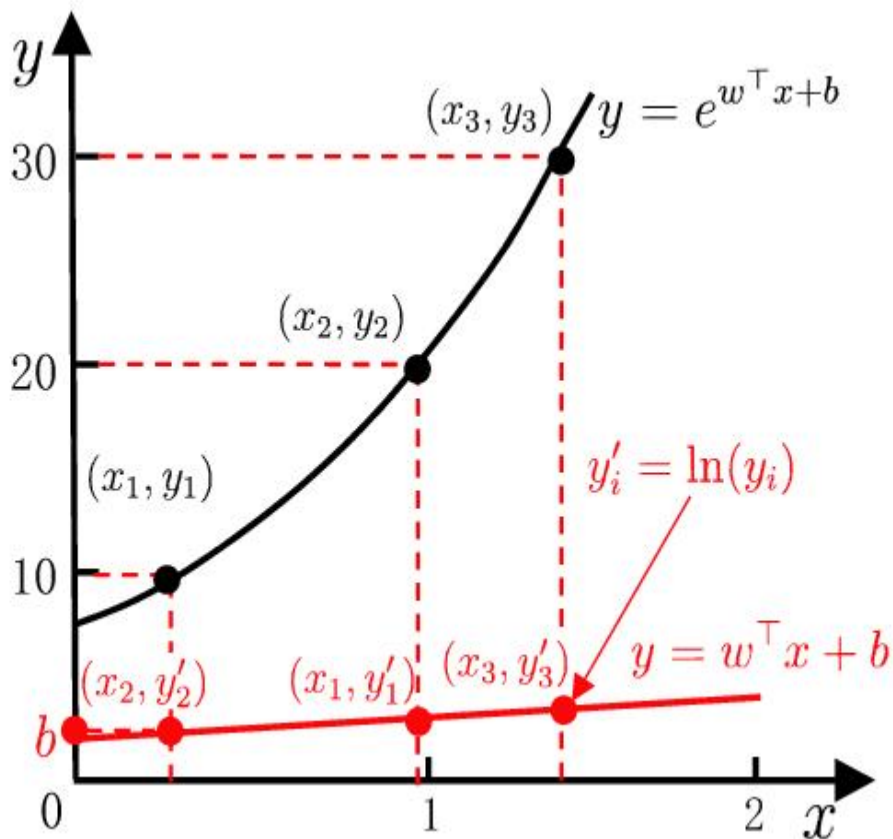
$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \frac{1}{m \sum_i x_i^2 - (\sum_i x_i)^2} \begin{pmatrix} m \sum_i x_i y_i - \sum_i x_i \sum_j y_j \\ \sum_i x_i^2 \sum_j y_j - \sum_i x_i y_i \sum_j x_j \end{pmatrix}$$

$$w = \frac{m \sum_i x_i y_i - \sum_i x_i \sum_j y_j}{m \sum_i x_i^2 - (\sum_i x_i)^2} = \frac{\sum_i x_i y_i - \sum_j y_j \bar{x}}{\sum_i x_i^2 - \frac{1}{m} (\sum_i x_i)^2} = \frac{\sum_i (x_i - \bar{x}) y_i}{\sum_i x_i^2 - \frac{1}{m} (\sum_i x_i)^2}$$

$$b = \frac{\sum_i x_i^2 \sum_j y_j - \sum_i x_i y_i \sum_j x_j}{m \sum_i x_i^2 - (\sum_i x_i)^2} = \frac{\bar{x}^2 \sum_j y_j - \sum_i x_i y_i \bar{x}}{\sum_i x_i^2 - \frac{1}{m} (\sum_i x_i)^2} = \frac{\sum_j y_j (\bar{x}^2 - \bar{x} \bar{x})}{\sum_i x_i^2 - \frac{1}{m} (\sum_i x_i)^2}$$

# 对数线性回归

- 输出标记的对数为线性模型逼近的目标



$$y = \mathbf{w}^T \mathbf{x} + b$$



$$\ln y = \mathbf{w}^T \mathbf{x} + b$$

# 线性回归 - 广义线性模型

- 一般形式

$$y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$$

- $g(\cdot)$ 称为链接函数 (link function)
  - 单调可微
- 对数线性回归  $g(\cdot) = \ln(\cdot)$  就是广义线性模型的特例

# 二分类任务

- 预测值与输出标记

$$z = \mathbf{w}^T \mathbf{x} + b \quad y \in \{0, 1\}$$

- 寻找函数将分类标记与线性回归模型输出联系起来
- 最理想的函数——单位阶跃函数

$$y = \begin{cases} 0, & z < 0 \\ 0.5, & z = 0 \\ 1, & z > 0 \end{cases}$$

- 预测值大于0就判别为正例，小于0判别为负例，预测值为临界值0时可以任意判别



# 二分类任务

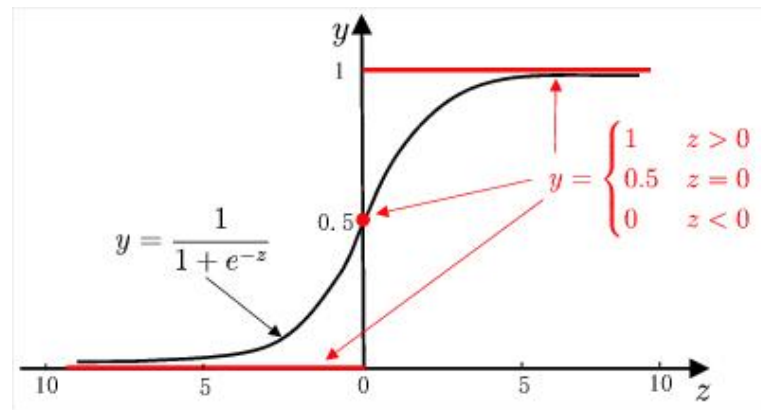
- 单位阶跃函数缺点
  - 不连续，无法用在广义线性模型中
- 替代函数——对数几率函数（logistic function）

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

单调可微、任意阶可导

$$\sigma^{-1}(y) = \ln \frac{y}{1 - y}$$

单位阶跃函数与对数几率函数的比较



# 对数几率回归

- 运用对数几率函数

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \rightarrow \quad \hat{y} = \frac{1}{1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)}}$$

$\hat{y}$  视为样本  $\mathbf{x}$  作为正例的可能性

$$\mathbf{w}^\top \mathbf{x} + b = \ln \frac{\hat{y}}{1 - \hat{y}}$$

$\frac{\hat{y}}{1 - \hat{y}}$  称为几率，反映了  $\mathbf{x}$  作为正例的相对可能性

## 对数几率回归优点

- 无需事先假设数据分布
- 可得到“类别”的近似概率预测
- 可直接应用现有数值优化算法求取最优解

# 对数几率回归 - 极大似然法

如何优化参数 $(\mathbf{w}, b)$ ?

$$\hat{y} = \frac{1}{1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)}} \triangleq P(y = 1 | \mathbf{x}; \mathbf{w}, b)$$

↓

$$P(y = 0 | \mathbf{x}; \mathbf{w}, b) = \frac{e^{-(\mathbf{w}^\top \mathbf{x} + b)}}{1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)}} = \frac{1}{1 + e^{\mathbf{w}^\top \mathbf{x} + b}}$$

给定数据集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ,  $y_i \in \{0, 1\}$

- 极大似然法 最大化 对数似然

$$\begin{aligned}\ell(\mathbf{w}, b) &= \sum_{i=1}^m y_i \log P(y = 1 | \mathbf{x}_i; \mathbf{w}, b) + (1 - y_i) \log P(y = 0 | \mathbf{x}_i; \mathbf{w}, b) \\ &= \sum_{i=1}^m \log P(y = y_i | \mathbf{x}_i; \mathbf{w}, b)\end{aligned}$$

# 对数几率回归 - 极大似然法

- 极大似然法 最小化 负对数似然

$$\ell(\hat{\mathbf{w}}) = - \sum_{i=1}^m \log P(y = y_i | \hat{\mathbf{x}}_i; \hat{\mathbf{w}})$$

$$P(y = 1 | \mathbf{x}; \mathbf{w}, b) = \frac{e^{\mathbf{w}^\top \hat{\mathbf{x}}}}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}}} = \frac{e^{1\mathbf{w}^\top \hat{\mathbf{x}}}}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}}} = \frac{e^{y\mathbf{w}^\top \hat{\mathbf{x}}}}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}}}$$

$$P(y = 0 | \mathbf{x}; \mathbf{w}, b) = \frac{1}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}}} = \frac{e^{0\mathbf{w}^\top \hat{\mathbf{x}}}}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}}} = \frac{e^{y\mathbf{w}^\top \hat{\mathbf{x}}}}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}}}$$

$$\ell(\hat{\mathbf{w}}) = - \sum_{i=1}^m \log \frac{e^{y_i \mathbf{w}^\top \hat{\mathbf{x}}_i}}{1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}_i}} = \sum_{i=1}^m -y_i \mathbf{w}^\top \hat{\mathbf{x}}_i + \log(1 + e^{\mathbf{w}^\top \hat{\mathbf{x}}_i})$$

# 对数几率回归 - 极大似然法

- 考察函数  $f(x) = -ax + \ln(1 + \exp(x))$

一阶导数 
$$f'(x) = -a + \frac{1}{1 + e^{-x}}$$

二阶导数 
$$f''(x) = \left(\frac{1}{1 + e^{-x}}\right)' = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x)(1 - \sigma(x))$$

因为  $\sigma(x) \in (0,1)$ ，所以  $f''(x) > 0$ ，因此  $f(x)$  是凸函数。

因为  $\ell(\hat{\mathbf{w}})$  是  $f(x)$  和  $g(\mathbf{w}) = \hat{\mathbf{w}}^\top \hat{\mathbf{x}}_i$  的复合函数，即  $\ell(\hat{\mathbf{w}}) = f(g(\hat{\mathbf{w}}^\top \hat{\mathbf{x}}_i))$ ，  
所以  $\ell(\hat{\mathbf{w}})$  是关于  $\hat{\mathbf{w}}$  的凸函数

# 对数几率回归 - 极大似然法

- 负对数似然  $\ell(\hat{\mathbf{w}}) = \sum_{i=1}^m (-y_i \hat{\mathbf{w}}^\top \hat{\mathbf{x}}_i + \log(1 + e^{\hat{\mathbf{w}}^\top \hat{\mathbf{x}}_i}))$

一阶导数

$$\begin{aligned}\nabla \ell(\hat{\mathbf{w}}) &= \sum_i f'(z_i) \frac{\partial z_i}{\partial \hat{\mathbf{w}}} = \sum_i \left(-y_i + \frac{1}{1 + e^{-z_i}}\right) \mathbf{x}_i \\ &= -\sum_i (y_i - P(y = 1 | \hat{\mathbf{x}}_i; \hat{\mathbf{w}})) \hat{\mathbf{x}}_i\end{aligned}$$



$p_1$

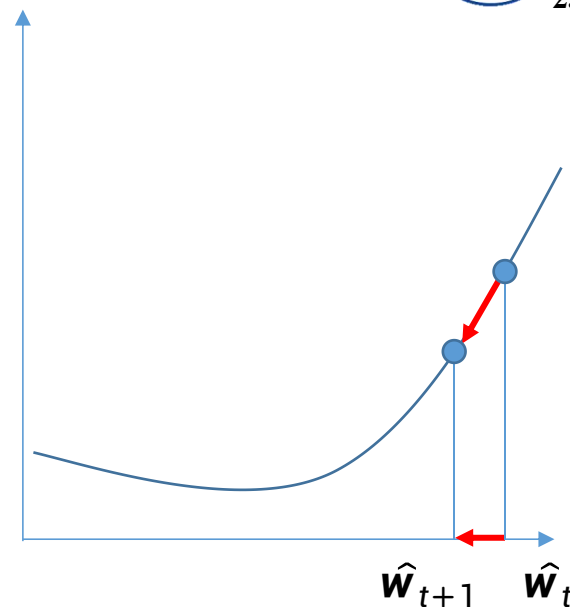
二阶导数

$$\nabla^2 \ell(\hat{\mathbf{w}}) = \frac{\partial \nabla \ell(\hat{\mathbf{w}})}{\partial \hat{\mathbf{w}}^\top} = \sum_i p_1(1 - p_1) \mathbf{x}_i \mathbf{x}_i^\top$$

# 对数几率回归 - 极大似然法

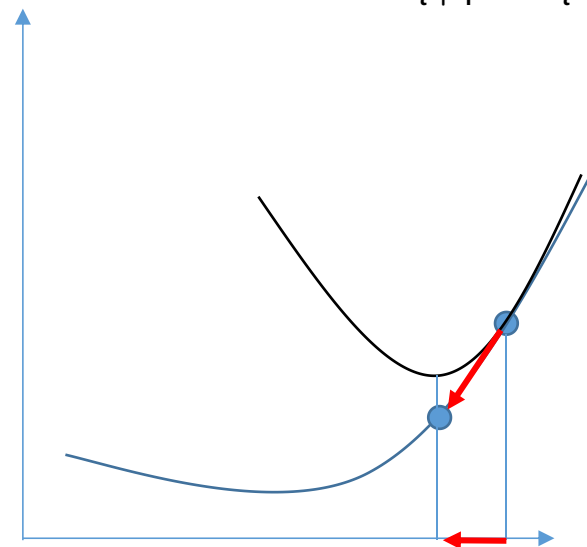
- 梯度下降法

```
while  $\|\nabla \ell(\hat{\mathbf{w}})\| > \delta$  do  
     $\hat{\mathbf{w}}_{t+1} \leftarrow \hat{\mathbf{w}}_t - a \nabla \ell(\hat{\mathbf{w}})$   
end while
```



- 牛顿法

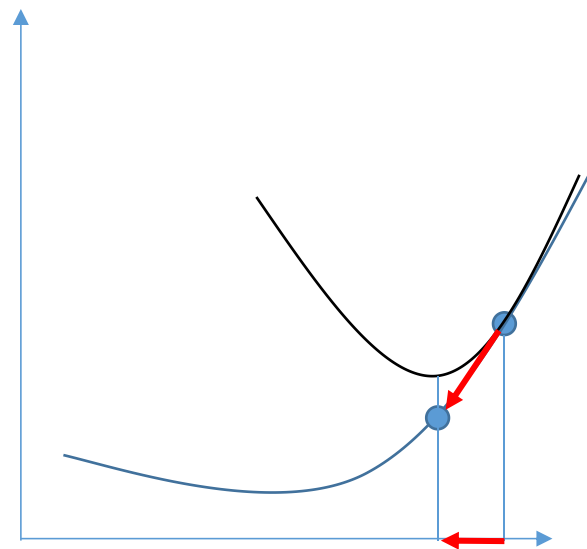
```
while  $\|\nabla \ell(\hat{\mathbf{w}})\| > \delta$  do  
     $\hat{\mathbf{w}}_{t+1} \leftarrow \hat{\mathbf{w}}_t - (\nabla^2 \ell(\hat{\mathbf{w}}))^{-1} \nabla \ell(\hat{\mathbf{w}})$   
end while
```



# 对数几率回归 - 极大似然法

## 牛顿法

```
while  $\|\nabla \ell(\hat{\mathbf{w}})\| > \delta$  do  
     $\hat{\mathbf{w}}_{t+1} \leftarrow \hat{\mathbf{w}}_t - (\nabla^2 \ell(\hat{\mathbf{w}}))^{-1} \nabla \ell(\hat{\mathbf{w}})$   
end while
```



- 考虑 $\ell(\hat{\mathbf{w}})$ 在 $\hat{\mathbf{w}}_t$ 处的二阶泰勒展开

$$\ell(\hat{\mathbf{w}}) \approx \ell(\hat{\mathbf{w}}_t) + \nabla \ell(\hat{\mathbf{w}}_t)^\top (\hat{\mathbf{w}} - \hat{\mathbf{w}}_t) + \frac{1}{2} (\hat{\mathbf{w}} - \hat{\mathbf{w}}_t)^\top \nabla^2 \ell(\hat{\mathbf{w}}) (\hat{\mathbf{w}} - \hat{\mathbf{w}}_t)$$

- 对 $\hat{\mathbf{w}}$ 求导数并令其等于0, 得到 $\hat{\mathbf{w}} = \hat{\mathbf{w}}_t - (\nabla^2 \ell(\hat{\mathbf{w}}))^{-1} \nabla \ell(\hat{\mathbf{w}})$



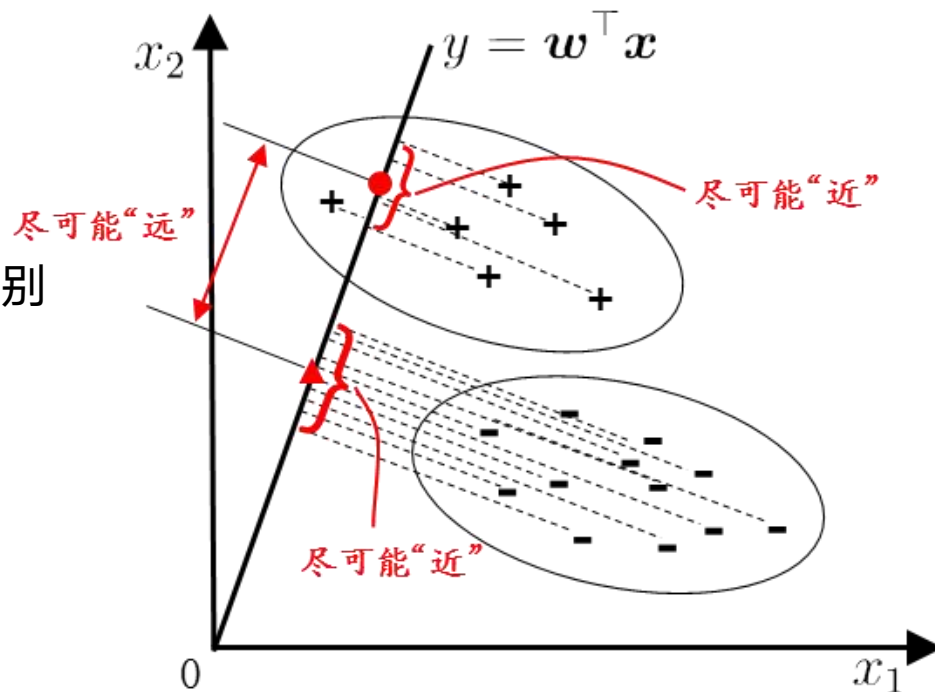
# 二分类任务 – 线性判别分析

- 线性判别分析 (Linear Discriminant Analysis) [Fisher, 1936]

投影到低维空间, 使得

- 欲使同类样例的投影点尽可能接近
- 欲使异类样例的投影点尽可能远离
- 新样本投影后根据投影位置进行判别

LDA也可被视为一种  
监督降维技术



# 二分类任务 – 线性判别分析

- 给定数据集  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m, y_i \in \{0, 1\}$

令

- 第 $c$ 类示例的集合  $X_c$
- 第 $c$ 类示例的均值向量  $\boldsymbol{\mu}_c = \frac{1}{|X_c|} \sum_{\mathbf{x} \in X_c} \mathbf{x}$
- 第 $c$ 类示例的协方差矩阵  $\boldsymbol{\Sigma}_c = \sum_{\mathbf{x} \in X_c} (\mathbf{x} - \boldsymbol{\mu}_c)(\mathbf{x} - \boldsymbol{\mu}_c)^\top$

若将数据投影到方向 $\mathbf{w}$ 确定的直线上

- 两类样本的中心在直线上的投影  $\frac{1}{|X_c|} \sum_{\mathbf{x} \in X_c} \mathbf{w}^\top \mathbf{x} = \mathbf{w}^\top \boldsymbol{\mu}_c$
- 两类样本的协方差  $\sum_{\mathbf{x} \in X_c} (\mathbf{w}^\top \mathbf{x} - \mathbf{w}^\top \boldsymbol{\mu}_c)^2 = \mathbf{w}^\top \boldsymbol{\Sigma}_c \mathbf{w}$

# 二分类任务 – 线性判别分析

- 欲使同类样例的投影点尽可能接近,

可以让同类样例投影点的协方差尽可能小

- 欲使异类样例的投影点尽可能远离,

可以让类中心之间的距离尽可能大

- 同时考虑两者, 则可得到最大化目标

定义类内散度

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T$$

$$J = \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} = \frac{w^T (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} = \frac{w^T S_b w}{w^T S_w w}$$

$$S_w = \Sigma_0 + \Sigma_1$$

定义类间散度

# 二分类任务 – 线性判别分析

$$\max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}} \quad \begin{array}{l} \text{广义瑞利商} \\ \text{(generalized Rayleigh quotient)} \end{array}$$

- 若 $\mathbf{w}$ 是一个解，则对于任意常数 $a$ ,  $a\mathbf{w}$ 也是解
- 不失一般性，令 $\mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1$ ，则等价于

$$\max_{\mathbf{w}} \mathbf{w}^\top \mathbf{S}_b \mathbf{w} \quad \text{s.t. } \mathbf{w}^\top \mathbf{S}_w \mathbf{w} = 1$$

- 引入拉格朗日乘子 $\lambda$ ，并令朗格拉日函数梯度等于0，可以得到

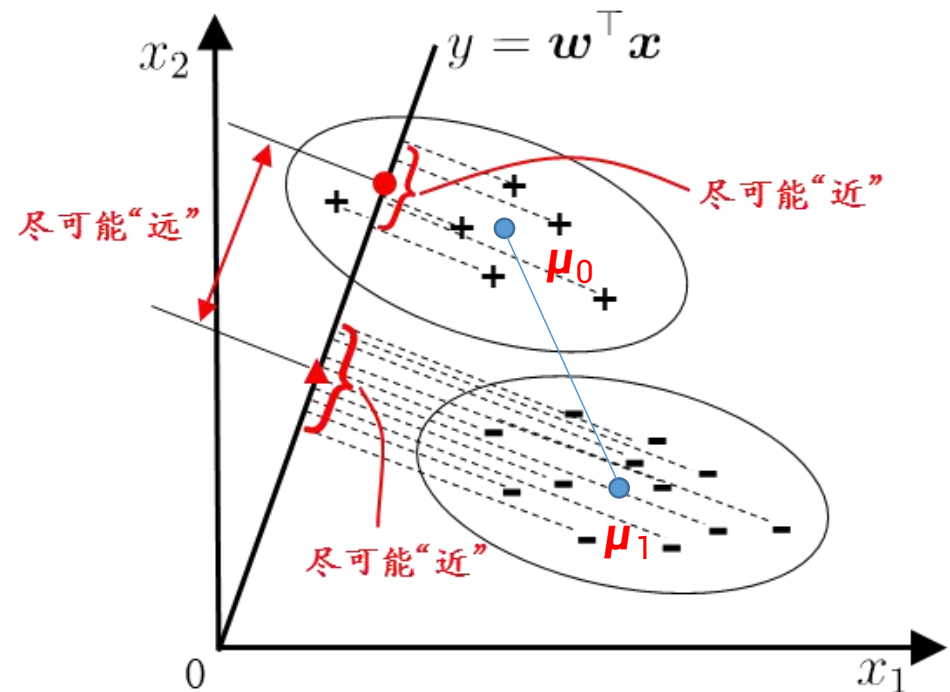
$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

# 二分类任务 – 线性判别分析

最优解  $\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$

•  $\mathbf{S}_b \mathbf{w} = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top \mathbf{w} \propto (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$

• 由此可得  $\mathbf{w} \propto \mathbf{S}_w^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$



# LDA推广 – 多分类任务

- 全局散度矩阵  $\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w = \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top$

- 类内散度矩阵  $\mathbf{S}_w = \sum_c \mathbf{S}_{w_c} \quad \mathbf{S}_{w_c} = \sum_{\mathbf{x} \in X_c} (\mathbf{x} - \boldsymbol{\mu}_c)(\mathbf{x} - \boldsymbol{\mu}_c)^\top$

- 类间散度矩阵  $\mathbf{S}_{w_c} = \sum_{\mathbf{x} \in X_c} \mathbf{x} \mathbf{x}^\top - m_c \boldsymbol{\mu}_c \boldsymbol{\mu}_c^\top \quad \mathbf{S}_t = \sum_i \mathbf{x}_i \mathbf{x}_i^\top - m \boldsymbol{\mu} \boldsymbol{\mu}^\top$

$$\begin{aligned} \mathbf{S}_b &= \mathbf{S}_t - \mathbf{S}_w = \sum_c m_c \boldsymbol{\mu}_c \boldsymbol{\mu}_c^\top - m \boldsymbol{\mu} \boldsymbol{\mu}^\top \\ m \boldsymbol{\mu} &= \sum_c m_c \boldsymbol{\mu}_c \quad m = \sum_c m_c \\ &= \sum_c m_c (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\boldsymbol{\mu}_c - \boldsymbol{\mu})^\top \end{aligned}$$

# LDA推广 – 多分类任务

- 优化目标

$$\max_W \frac{\text{tr}(\mathbf{W}^\top \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^\top \mathbf{S}_w \mathbf{W})} = \frac{\sum_i \mathbf{w}_i^\top \mathbf{S}_b \mathbf{w}_i}{\sum_i \mathbf{w}_i^\top \mathbf{S}_w \mathbf{w}_i}$$

- 若 $\mathbf{W}$ 是一个解，则对于任意常数 $a$ ,  $a\mathbf{W}$ 也是解

- 等价于  $\max_W \text{tr}(\mathbf{W}^\top \mathbf{S}_b \mathbf{W})$  s.t.  $\text{tr}(\mathbf{W}^\top \mathbf{S}_w \mathbf{W}) = 1$

- 引入拉格朗日乘子 $\lambda$ ，并令朗格拉日函数梯度等于0，可以得到广义特征值问题

$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

$\mathbf{W}$ 的闭式解则是  $\mathbf{S}_w^{-1} \mathbf{S}_b$  的 $N-1$ 个最大广义特征值所对应的特征向量组成的矩阵

# 多分类学习

- 多分类学习方法

- 二分类学习方法推广到多类，利用二分类学习器解决多分类问题

- ✓ 对问题进行拆分，为拆出的每个二分类任务训练一个分类器
    - ✓ 对于每个分类器的预测结果进行集成以获得最终的多分类结果

- 拆分策略

- 一对一 (One vs. One, OvO)
  - 一对其余 (One vs. Rest, OvR)
  - 多对多 (Many vs. Many, MvM)

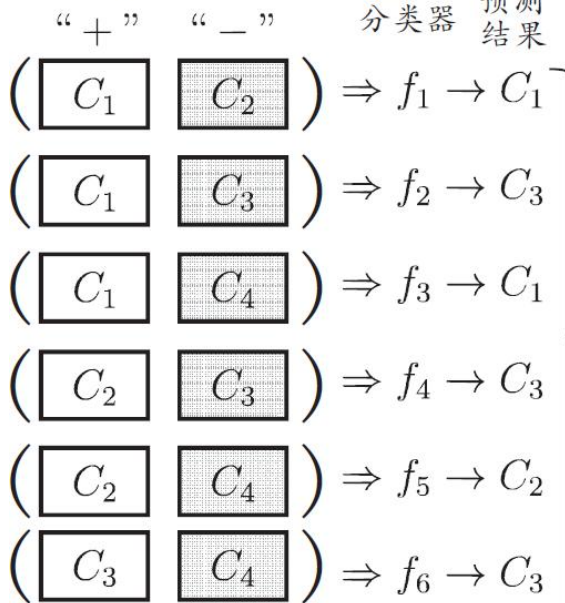


# 多分类学习 - 一对一



OvO

用于训练的  
两类样例



最终  
结果

$\rightarrow C_3$

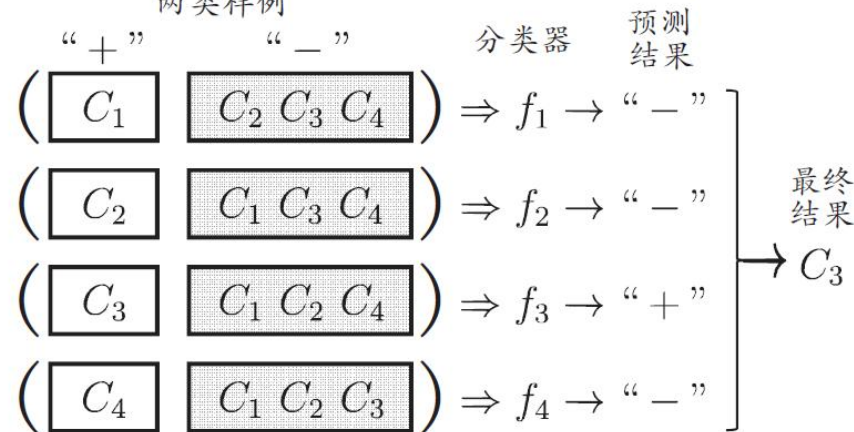
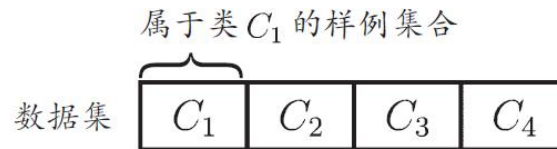
## 拆分阶段

- $N$  个类别两两配对
  - $N(N-1)/2$  个二类任务
- 各个二类任务学习分类器
  - $N(N-1)/2$  个二类分类器

## 测试阶段

- 新样本提交给所有分类器预测
  - $N(N-1)/2$  个分类结果
- 投票产生最终分类结果
  - 被预测最多的类别为最终类别

# 多分类学习 - 一对其余



## 拆分阶段

- 某一类作为正例，其他反例
  - N 个二类任务
- 各个二类任务学习分类器
  - N 个二类分类器

## 测试阶段

- 新样本提交给所有分类器预测
  - N 个分类结果
- 比较各分类器预测置信度
  - 置信度最大类别作为最终类别

# 多分类学习 – 两种策略比较

## 一对一

- 训练 $N(N-1)/2$ 个分类器，存储开销和测试时间大
- 训练只用两个类的样例，训练时间短

## 一对其余

- 训练 $N$ 个分类器，存储开销和测试时间小
- 训练用到全部训练样例，训练时间长

预测性能取决于具体数据分布，多数情况下两者差不多

# 多分类学习 - 多对多

- 多对多 (Many vs Many, MvM)

若干类作为正类，若干类作为反类

- 纠错输出码 (Error Correcting Output Code, ECOC)

编码

对N个类别做M次划分，每次划分将一部分类别划为正类，一部分划为反类，形成二分类训练集

解码

M个分类器分别对测试样本进行预测，预测标记组成一个编码。将距离最小的类别为最终类别

	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	海明距离	欧氏距离
	↓	↓	↓	↓	↓	↓	↓
$C_1 \rightarrow$	-1	+1	-1	+1	+1	3	$2\sqrt{3}$
$C_2 \rightarrow$	+1	-1	-1	+1	-1	4	4
$C_3 \rightarrow$	-1	+1	+1	-1	+1	1	2
$C_4 \rightarrow$	-1	-1	+1	+1	-1	2	$2\sqrt{2}$

测试示例 $\rightarrow$	-1	-1	+1	-1	+1	↑ ↑	
--------------------	----	----	----	----	----	-----	--

# 多分类学习 - 多对多

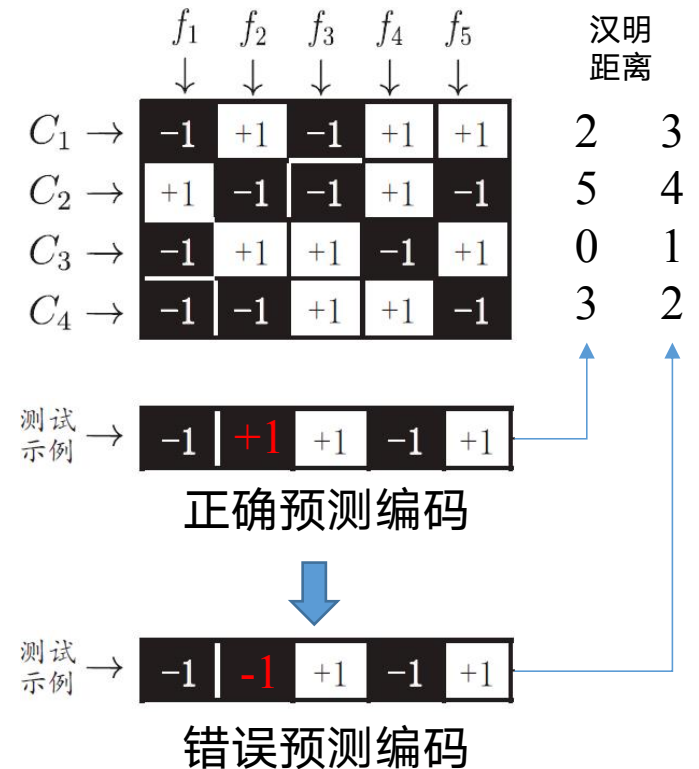
编码

对N个类别做M次划分，每次划分将一部分类别划为正类，一部分划为反类，形成二分类训练集

解码

M个分类器分别对测试样本进行预测，预测标记组成一个编码。将距离最小的类别为最终类别

- 纠错能力
  - 若  $f_2$  预测错误
  - 仍然能产生正确的最终分类



- ✓ ECOC编码越长、对分类器错误纠错能力越强
- ✓ 对同等长度编码，理论上任意两个类别之间的编码距离越远，则纠错能力越强

# 多分类学习 - 多对多

- 纠错输出码(Error Correcting Output Code, ECOC)

	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	海明 距离	欧氏 距离
$C_1 \rightarrow$	-1	+1	-1	+1	+1	3	$2\sqrt{3}$
$C_2 \rightarrow$	+1	-1	-1	+1	-1	4	4
$C_3 \rightarrow$	-1	+1	+1	-1	+1	1	2
$C_4 \rightarrow$	-1	-1	+1	+1	-1	2	$2\sqrt{2}$
测试 示例 $\rightarrow$	-1	-1	+1	-1	+1		

(a) 二元 ECOC 码

[Dietterich and Bakiri,1995]

# 类别不平衡问题

- 类别不平衡 (class imbalance)
  - 不同类别训练样例数相差很大情况 (正类为小类)

## 类别平衡正例预测

分类决策规则:  $\frac{y}{1-y} > 1$ , 则为正例

## 类别不平衡正例预测

分类决策规则:  $\frac{y}{1-y} > \frac{m^+}{m^-}$ , 则正例

$y$  为分类器预测值, 表达了正例可能性, 几率  $\frac{y}{1-y}$  反映相对可能性

- 分类器都基于**类别平衡分类决策规则**决策的, 只能对预测值进行缩放

$$\frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^-}{m^+} \quad \Rightarrow \quad y' = \frac{m^- y}{m^+ (1-y) + m^- y}$$

# 类别不平衡问题

“训练集是真实样本总体的无偏采样” 假设往往不成立，未必能基于训练集观测几率来推断真实几率

## 解决办法

- 欠采样 (undersampling)
  - 去除一些反例使正反例数目接近 (EasyEnsemble [Liu et al.,2009])
- 过采样 (oversampling)
  - 增加一些正例使正反例数目接近 (SMOTE [Chawla et al.2002])
- 阈值移动 (threshold-moving)



# 作业

- 3.2
- 3.7
- 在LDA多分类情形下，试计算类间散度矩阵 $S_b$ 的秩并证明
- 给出公式3.45的推导过程
- 证明 $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ 是投影矩阵，并对线性回归模型从投影角度解释