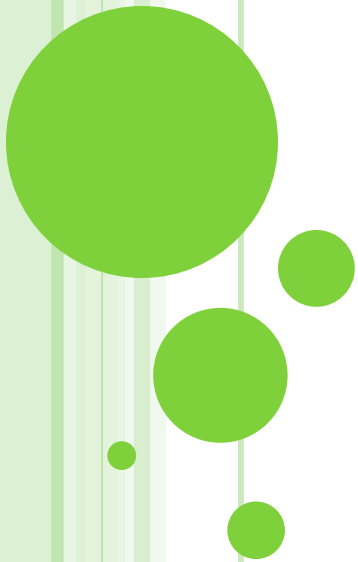


第三章： LINEAR MODELS



目录

- **Linear Regression**
 - 最小二乘法
- **Binary Classification**
 - 对数几率回归
 - 线性判别分析
- **Multi-Class Classification**
 - 一对一
 - 一对其余
 - 多对多
- 类别不平衡问题



基本形式

■ Linear regression model 一般形式

给定由 d 个属性描述的示例 $\boldsymbol{x} = (x_1; x_2; \dots; x_d)$,

线性模型试图学到一个通过属性的线性组合来进行预测的函数

$$f(\boldsymbol{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$

其中 x_i 是 \boldsymbol{x} 在第 i 个属性上的取值

■ 一般写成向量形式:

$$f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + b$$

其中 $\boldsymbol{w} = (w_1; w_2; \dots; w_d)$. \boldsymbol{w} 和 b 学得之后, 模型就得以确定.

线性模型优点

○ 形式简单、易于建模

蕴含着机器学习的重要思想

○ 非线性模型的基础

许多功能强大的非线性模型可在线性模型基础上获得，例如

- 引入层级结构或高维映射

○ 可解释性

w直观的表达了各个属性在预测中的重要性，例如：

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

- 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



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

- 线性回归 (linear regression) 目的

- 学得一个线性模型以尽可能准确地预测实值输出标记
- 考虑最简单的情况：输入属性只有一个(忽略下标)

$$D = \{(x_i, y_i)\}_{i=1}^m, \text{ 其中 } x_i \in \mathbb{R}.$$



线性回归

$D = \{(x_i, y_i)\}_{i=1}^m$, 其中 $x_i \in \mathbb{R}$.

离散属性:

若属性间存在:

有“序”关系

可通过连续化转化为连续值, 例二值属性:

“身高”的取值“高”“矮”可转化为 $\{1.0, 0.0\}$

若无“序”关系

有 k 个属性值, 则转换为 k 维向量, 如属性的“瓜类”的取值:

西瓜、南瓜、黄瓜, 可以转化为:

$(0,0,1), (0,1,0), (1,0,0)$



线性回归

- 单一属性的线性回归目标 $f(x) = wx_i + b$

$$\text{使得 } f(x_i) \simeq y_i$$

如何确定 w, b ???

关键：如何衡量 $f(x)$ 与 y 的差别

均方误差是回归任务中最常用的性能度量，

所以最小化均方误差

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

线性回归

- 参数/模型估计：最小二乘法（least square method）

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

用均方误差最小化来进行模型求解的方法称为“最小二乘法”

线性回归中，最小二乘法就是：

试图找到一条直线，使得所有样本到直线上的欧式距离之和最小。

求解 w 和 b 使 $E_{(w, b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$ 最小化的过程

----称为线性回归模型的最小二乘参数估计（parameter estimation）

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

$$\frac{\partial E_{(w,b)}}{\partial w} = 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right)$$

$$\frac{\partial E_{(w,b)}}{\partial b} = 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right)$$



然后令式(3.5)和(3.6)为零可得到 w 和 b 最优解的闭式(closed-form)解

线性回归 - 最小二乘法

- 得到闭式（closed-form）解

$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2}$$

$$b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i)$$

其中 $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$ 均值



多元线性回归

更一般的，给定数据集 D ，样本由 d 个属性描述

- 给定数据集

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

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

- 多元线性回归目标

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b \quad \text{使得} \quad f(\mathbf{x}_i) \simeq y_i$$

这称为“多元线性回归” (multivariate linear regression)



多元线性回归

类似的，可以用最小二乘法估计 w, b

- 把 w 和 b 吸收入向量形式 $\hat{w} = (w; b)$

数据集 D 表示为 $m \times (d+1)$ 大小的矩阵 X

每行为一个示例

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

该行前 d 个元素对于 d 个属性值

最后一个元素恒置为1

再把标记也写成向量形式:

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

多元线性回归 - 最小二乘法

类似的,

□ 最小二乘法 (least square method)

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

$$\hat{\mathbf{w}}^* = \arg \min_{\hat{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}^T) (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

令 $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$, 对 $\hat{\mathbf{w}}$ 求导得到

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})$$

令上式为零可得 $\hat{\mathbf{w}}$ 最优解的闭式解

下面做简单的讨论



多元线性回归 - 最小二乘法

□ 最小二乘法 (least square method)

矩阵求导推导过程:

$$E_{\hat{w}} = (\mathbf{y} - \mathbf{X}\hat{w})^T (\mathbf{y} - \mathbf{X}\hat{w})$$
$$\frac{\partial E_{\hat{w}}}{\partial \hat{w}} = 2\mathbf{X}^T (\mathbf{X}\hat{w} - \mathbf{y})$$

将 $E_{\hat{w}} = (\mathbf{y} - \mathbf{X}\hat{w})^T (\mathbf{y} - \mathbf{X}\hat{w})$ 展开可得

$$E_{\hat{w}} = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{w} - \hat{w}^T \mathbf{X}^T \mathbf{y} + \hat{w}^T \mathbf{X}^T \mathbf{X} \hat{w}$$

对 \hat{w} 求导可得

$$\frac{\partial E_{\hat{w}}}{\partial \hat{w}} = \frac{\partial \mathbf{y}^T \mathbf{y}}{\partial \hat{w}} - \frac{\partial \mathbf{y}^T \mathbf{X} \hat{w}}{\partial \hat{w}} - \frac{\partial \hat{w}^T \mathbf{X}^T \mathbf{y}}{\partial \hat{w}} + \frac{\partial \hat{w}^T \mathbf{X}^T \mathbf{X} \hat{w}}{\partial \hat{w}}$$

由矩阵微分公式 $\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}$, $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$ 可得

$$\frac{\partial E_{\hat{w}}}{\partial \hat{w}} = 0 - \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \hat{w}$$

$$\frac{\partial E_{\hat{w}}}{\partial \hat{w}} = 2\mathbf{X}^T (\mathbf{X}\hat{w} - \mathbf{y})$$

多元线性回归 - 满秩讨论

□ $\mathbf{X}^T \mathbf{X}$ 是满秩矩阵或正定矩阵, 令 $\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) = \mathbf{0}$
可逆 则 $\hat{\mathbf{w}}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

其中 $(\mathbf{X}^T \mathbf{X})^{-1}$ 是 $\mathbf{X}^T \mathbf{X}$ 的逆矩阵, 令 $\hat{\mathbf{x}}_i = (\mathbf{x}_i, 1)$

学得线性回归模型为 $f(\hat{\mathbf{x}}_i) = \hat{\mathbf{x}}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

多元线性回归 – 满秩讨论

然而,现实任务中 $\mathbf{X}^T\mathbf{X}$ 往往不是满秩矩阵. 例如:

在许多任务中我们会遇到大量的变量, 其数目甚至超过样例数, 导致 \mathbf{X} 的列数多于行数, $\mathbf{X}^T\mathbf{X}$ 显然不满秩.

此时: 可解出多个 \mathbf{w} , 它们都能使均方误差最小化.

选择哪一个解作为输出, 将由学习算法的归纳偏好决定,

常见的做法是: 引入正则化(regularization)项.

□ $\mathbf{X}^T\mathbf{X}$ 不是满秩矩阵

- 根据归纳偏好选择解 (参见1.4节)
- 引入正则化 (参加6.4节, 11.4节)



对数线性回归

线性模型虽简单，却有丰富的变化。例如对于样例 (\mathbf{x}, y) ， $y \in \mathbb{R}$ ，希望线性模型的预测值逼近真实标记 y 时候，就得到了线性回归模型；为了便于观察，我们把线性回归模型简写为：

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

可否逼近 **y** 的衍生物？

例如： 输出标记在指数尺度上变化，

那就可将输出标记的对数作为线性模型逼近的目标，即

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

这就是“对数线性回归” (log-linear regression)，

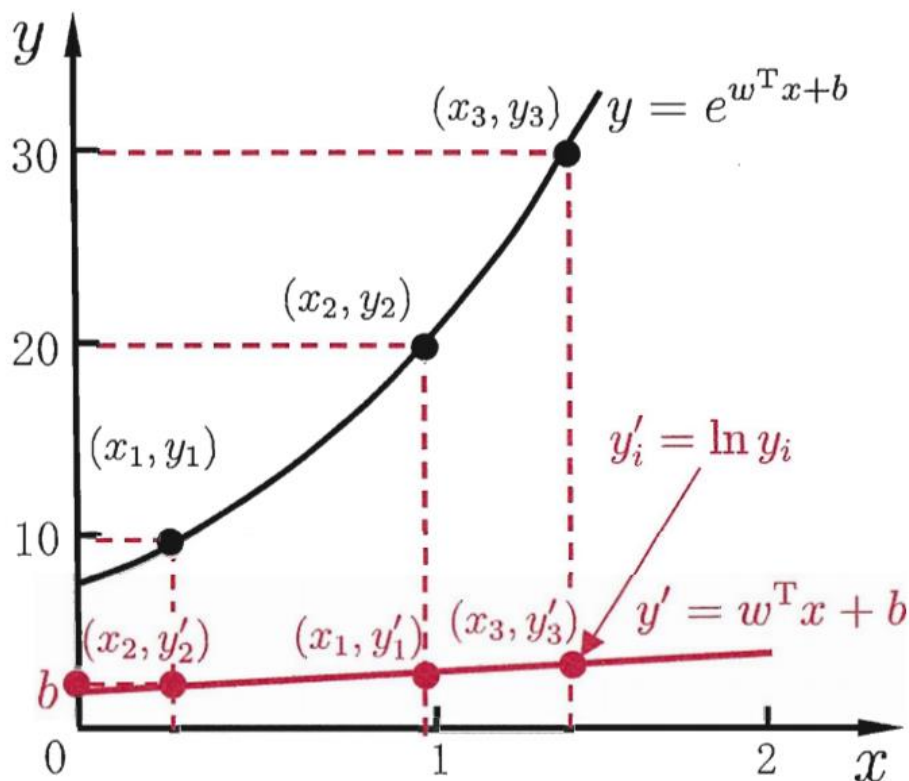
它实际上是在试图让 $e^{\mathbf{w}^T \mathbf{x} + b}$ 逼近 **y**

在形式上仍然是线性回归，但实质上是在求取输入空间到输出空间的非线性函数映射

对数线性回归

如图3.1所示，这里的**对数函数**起到了**将线性模型的预测值与真实标记**联系起来的作

- **输出标记的对数**为线性模型逼近的**目标**



$$\ln y = w^T x + b$$



$$y = w^T x + b$$

图 3.1 对数线性回归示意图

线性回归 - 广义线性模型

更一般的，考虑单调可微函数 $g(\cdot)$:

- 一般形式，令 $y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$

这样得到的模型称为“广义线性模型” (Generalized linear model)

□ $g(\cdot)$ 称为联系函数 (link function)

- 对数线性回归是 $g(\cdot) = \ln(\cdot)$ 时广义线性模型的特例



对数几率回归

讨论了如何使用线性模型进行回归学习，但若要做的是分类任务，

广义线性模型，
$$y = g^{-1}(w^T x + b)$$

只需要找一个单调可微函数，

将分类任务的真实标记 y 与 线性回归模型的预测值联系起来



二分类任务

考虑二分类任务,其输出标记 $y \in \{0,1\}$,

而线性回归模型产生的预测值

$$z = \boldsymbol{w}^T \boldsymbol{x} + b \text{ 是实值}$$

于是, 我们需将实值 z 转换为0/1值



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

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

预测值大于零就判为正例，
小于零就判为反例，
预测值为临界值零则可任意判别

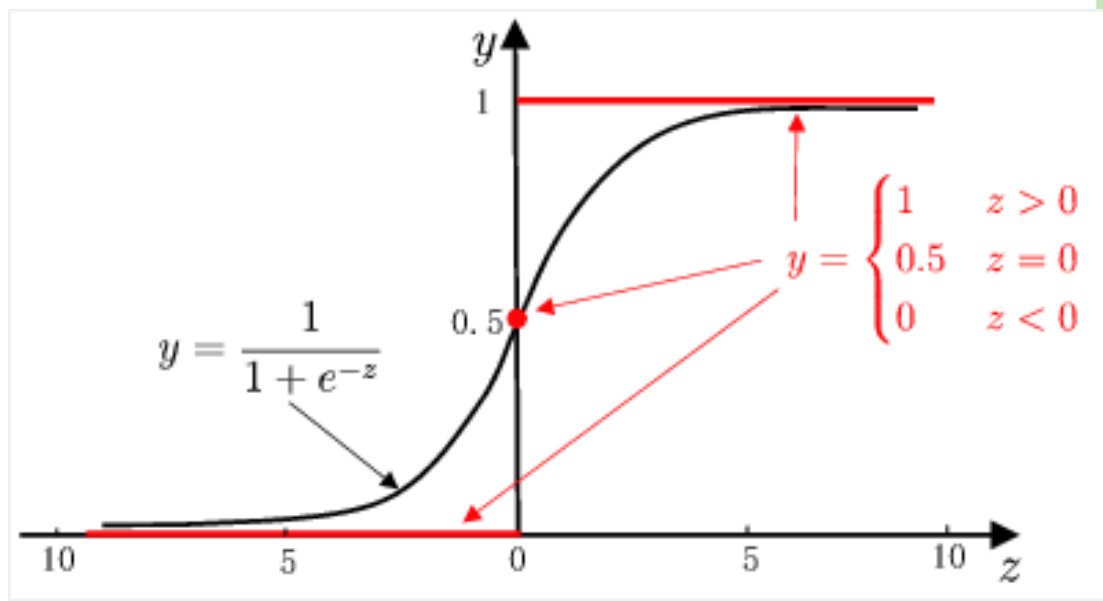


二分类任务

单位阶跃函数与对数几率函数的比较

- 如图

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



单位阶跃函数缺点

- 不连续

但从图 3.2可看出: 单位阶跃函数不连续,因此不能直接用作式(3.15)中的 $g^{-}(\cdot)$

于是我们希望: 找到能在一定程度上近似单位阶跃函数的“替代函数”(surrogate function), 并希望它单调可微.

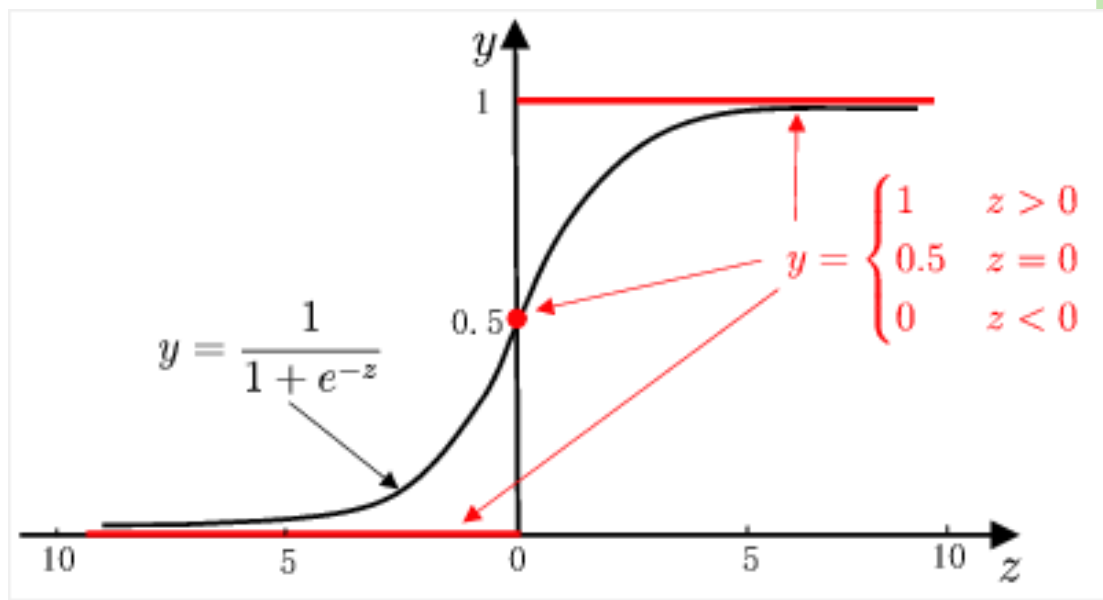
对数几率函数(logistic function)正是这样一个常用的替代函数:

二分类任务

单位阶跃函数与对数几率函数的比较

- 如图

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



单位阶跃函数缺点

- 不连续

替代函数——对数几率函数 (logistic function)

- 单调可微、任意阶可导
- 它将 z 值转化为一个接近0或1的 y 值，并且输出在 $z=0$ 附近变化很陡

对数几率回归

广义线性模型

- 运用对数几率函数代入 $\longrightarrow y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$

$$y = \frac{1}{1 + e^{-z}} \quad \text{变为} \quad y = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

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

则有: $\ln \frac{y}{1 - y} = \mathbf{w}^T \mathbf{x} + b$



对数几率回归

虽然名字是回归，
但是是一种分类算法

■ 对数几率 (log odds)

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

若将 y 视为样本 x 作为正例的可能性,则 $1-y$ 是其反例可能性,

两者的比值 $\frac{y}{1-y}$

称为“几率” (odds), 反映了 x 作为正例的相对可能性.

取对数, 则得到对数几率:

$$\ln \frac{y}{1-y}$$

■ 对数几率回归优点


- 无需事先假设数据分布, 直接对分类可行性建模, 避免假设分布不准确带来的问题
- 不仅预测类别, 可得到“类别”的近似概率预测, 对利用概率辅助决策的任务有效
- 对率函数是任意阶可导的凸函数, 有很好的的数学性质, 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

$$y = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}} \cdot \text{如何确定式(3.18)中的 } \mathbf{w} \text{ 和 } b.$$

将 y 视为后验概率 $p(y=1 | \mathbf{x})$ 则, $\ln \frac{y}{1-y} = \mathbf{w}^T \mathbf{x} + b$

○对数几率


$$\ln \frac{p(y=1 | \mathbf{x})}{p(y=0 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + b$$

显然有

$$p(y=1 | \mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$

$$p(y=0 | \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$



对数几率回归 - 极大似然法

可通过“极大似然法” (maximum likelihood method) 来估计 W, b

↑
使联合概率最大化

极大似然法 (maximum likelihood)

- 给定数据集

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^m$$

- 最大化样本属于其真实标记的概率

- 最大化对数似然函数

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \mathbf{w}_i, b)$$

即令每个样本属于其真实标记的概率越大越好

对数几率回归 - 极大似然法

转化为最小化负对数似然函数求解

- 令 $\beta = (\mathbf{w}; b)$ $\hat{\mathbf{x}} = (\mathbf{x}; 1)$, 则 $\mathbf{w}^T \mathbf{x} + b$ 可简写为 $\beta^T \hat{\mathbf{x}}$

- 再令

$$p_1(\hat{\mathbf{x}}_i; \beta) = p(y = 1 \mid \hat{\mathbf{x}}_i; \beta)$$

$$p_0(\hat{\mathbf{x}}_i; \beta) = p(y = 0 \mid \hat{\mathbf{x}}_i; \beta) = 1 - p_1(\hat{\mathbf{x}}_i; \beta)$$

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i \mid \mathbf{x}_i; \mathbf{w}, b), \quad (3.25)$$

则3.25似然项可重写为

$$p(y_i \mid \mathbf{x}_i; \mathbf{w}_i, b) = y_i p_1(\hat{\mathbf{x}}_i; \beta) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \beta)$$

- 上式代入3.25, 故等价形式为要最小化

$$\ell(\beta) = \sum_{i=1}^m \left(-y_i \beta^T \hat{\mathbf{x}}_i + \ln \left(1 + e^{\beta^T \hat{\mathbf{x}}_i} \right) \right)$$

对数几率回归 - 极大似然法

转化为最小化负对数似然函数求解

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \mathbf{w}, b), \quad (3.25)$$

$$p(y_i | \mathbf{x}_i; \mathbf{w}_i, b) = y_i p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta})$$

- 上式代入3.25, 可得

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^m \ln (y_i p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta}))$$

其中 $p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}, p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) = \frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}},$ 代入上式可得

$$\begin{aligned} \ell(\boldsymbol{\beta}) &= \sum_{i=1}^m \ln \left(\frac{y_i e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} + 1 - y_i}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \right) \\ &= \sum_{i=1}^m \left(\ln(y_i e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} + 1 - y_i) - \ln(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}) \right) \end{aligned}$$

对数几率回归 - 极大似然法

转化为最小化负对数似然函数求解

$$\begin{aligned}\ell(\beta) &= \sum_{i=1}^m \ln \left(\frac{y_i e^{\beta^T \hat{\mathbf{x}}_i} + 1 - y_i}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right) \\ &= \sum_{i=1}^m \left(\ln(y_i e^{\beta^T \hat{\mathbf{x}}_i} + 1 - y_i) - \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i}) \right)\end{aligned}$$

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \mathbf{w}, b), \quad (3.25)$$

由于 $y_i=0$ 或 1 , 则

$$\ell(\beta) = \begin{cases} \sum_{i=1}^m (-\ln(1 + e^{\beta^T \hat{\mathbf{x}}_i})), & y_i = 0 \\ \sum_{i=1}^m (\beta^T \hat{\mathbf{x}}_i - \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i})), & y_i = 1 \end{cases}$$

两式综合可得

$$\ell(\beta) = \sum_{i=1}^m \left(y_i \beta^T \hat{\mathbf{x}}_i - \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i}) \right)$$

此式为极大似然估计的似然函数, 所以最大化似然函数等价于最小化似然函数的相反数,
故等价形式为要最小化

$$\ell(\beta) = \sum_{i=1}^m \left(-y_i \beta^T \hat{\mathbf{x}}_i + \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i}) \right)$$

对数几率回归

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^m \left(-y_i \boldsymbol{\beta}^T \hat{\mathbf{x}}_i + \ln \left(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} \right) \right) . \quad (3.27)$$

式(3.27)是关于 $\boldsymbol{\beta}$ 的高阶可导连续凸函数,

经典数值优化算法**梯度下降法**, **牛顿法**均可以求解

□ 求解得
$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta})$$

□ 牛顿法第**t+1**轮迭代解的更新公式

$$\boldsymbol{\beta}^{t+1} = \boldsymbol{\beta}^t - \left(\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \right)^{-1} \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}$$



对数几率回归

□ 求解得

$$\beta^* = \arg \min_{\beta} \ell(\beta)$$

□ 牛顿法第t+1轮迭代解的更新公式

$$\beta^{t+1} = \beta^t - \left(\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta}$$

其中关于 β 的一阶、二阶导数分别为

$$\frac{\partial \ell(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{\mathbf{x}}_i (y_i - p_1(\hat{\mathbf{x}}_i; \beta))$$

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T p_1(\hat{\mathbf{x}}_i; \beta) (1 - p_1(\hat{\mathbf{x}}_i; \beta))$$

高阶可导连续凸函数，梯度下降法/牛顿法 [Boyd and Vandenberghe, 2004]

二分类任务 - 线性判别分析

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

LDA也可被视为一种
监督降维技术

LDA思想非常朴素:

给定训练样本集, 设法将样本投影到
一条直线, 使得:

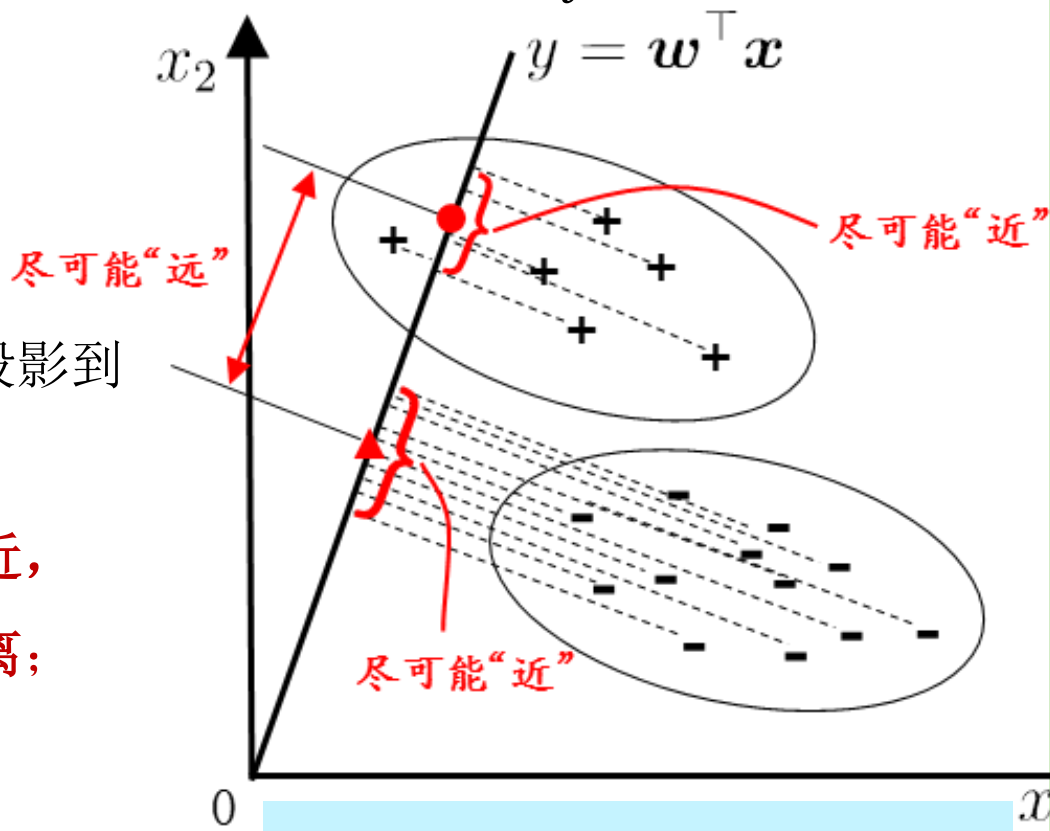
同类样例的投影点尽可能的接近,

异类样例的投影点尽可能的远离;

在对新样本进行分类时,

将其投影到同样的这条直线上,

再根据投影点的位置来确定新样本的
类别, 如图



“+”、“-”分别代表正例和反例,
椭圆表示数据簇外轮廓,虚线表示投影,
红色实心圆和实心三角形分别表示:
两类样本投影后的中心点.

二分类任务 - 线性判别分析

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

LDA的思想

- 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

一些变量

- 第 i 类示例的集合 X_i
- 第 i 类示例的均值向量 μ_i
- 第 i 类示例的协方差矩阵 Σ_i
- 两类样本的中心在直线上的投影: $w^T \mu_0$ 和 $w^T \mu_1$
- 若将所有样本都投影到直线上, 两类样本的协方差:

因为直线为一维空间, 所有均为实数

$w^T \Sigma_0 w$ $w^T \Sigma_1 w$ $w^T \Sigma w$

二分类任务 - 线性判别分析

要让同类样本投影点尽可能接近，可以让同类样本投影点协方差尽可能小

即 $w^T \Sigma_0 w + w^T \Sigma_1 w$ 尽可能小

要让异类样本投影点尽可能远，可以让类中心距离尽可能的大

即 $\|w^T \mu_0 - w^T \mu_1\|_2^2$ 尽可能大

同时考虑，则得到

○ 最大化目标

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \frac{w^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \end{aligned}$$



二分类任务 - 线性判别分析

最大化目标

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \frac{w^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \end{aligned}$$

要让同类样本投影点尽可能接近，可以

让同类样本投影点协方差尽可能小

即 $w^T \Sigma_0 w + w^T \Sigma_1 w$ 尽可能小

要让异类样本投影点尽可能远，可以类中心距

离尽可能的大

即 $\|w^T \mu_0 - w^T \mu_1\|_2^2$ 尽可能大

推导：

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T (\Sigma_0 + \Sigma_1) w} \\ &= \frac{\|(w^T \mu_0 - w^T \mu_1)^T\|_2^2}{w^T (\Sigma_0 + \Sigma_1) w} \\ &= \frac{\|(\mu_0 - \mu_1)^T w\|_2^2}{w^T (\Sigma_0 + \Sigma_1) w} \\ &= \frac{[(\mu_0 - \mu_1)^T w]^T (\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \\ &= \frac{w^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \end{aligned}$$

二分类任务 - 线性判别分析

定义

○ **类内**散度矩阵 (within-class scatter matrix)

$$\begin{aligned} S_w &= \Sigma_0 + \Sigma_1 \\ &= \sum_{x \in X_0} (x - \mu_0)(x - \mu_0)^T + \sum_{x \in X_1} (x - \mu_1)(x - \mu_1)^T \end{aligned}$$

○ **类间**散度矩阵 (between-class scatter matrix)

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

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \frac{w^T (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \end{aligned}$$

可重写为

$$J = \frac{w^T S_b w}{w^T S_w w}$$

这就是 LDA 欲最大化的目标, 即 S_b 与 S_w 的“广义瑞利商”

二分类任务 - 线性判别分析

- 广义瑞利商 (generalized Rayleigh quotient)

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \quad (3.35)$$

如何确定 w 呢？注意到式(3.35)的分子和分母都是关于 w 的二次项，式(3.35)的解与 w 的长度无关，只与其方向有关

不失一般性

- 令 $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$ ，最大化广义瑞利商**等价形式**为

$$\begin{aligned} \min_{\mathbf{w}} \quad & -\mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1 \end{aligned}$$



二分类任务 - 线性判别分析

运用拉格朗日乘子法

$$\min_w -w^T S_b w$$

$$\text{s.t. } w^T S_w w = 1$$

等价于



$$S_b w = \lambda S_w w \quad (3.37)$$

推导



拉格朗日函数为

$$L(w, \lambda) = -w^T S_b w + \lambda(w^T S_w w - 1)$$

对 w 求偏导可得

$$\begin{aligned} \frac{\partial L(w, \lambda)}{\partial w} &= -\frac{\partial(w^T S_b w)}{\partial w} + \lambda \frac{\partial(w^T S_w w - 1)}{\partial w} \\ &= -(S_b + S_b^T)w + \lambda(S_w + S_w^T)w \end{aligned}$$

由于 $S_b = S_b^T, S_w = S_w^T$, 所以

$$\frac{\partial L(w, \lambda)}{\partial w} = -2S_b w + 2\lambda S_w w$$

二分类任务 - 线性判别分析

运用拉格朗日乘子法

$$\min_w -\mathbf{w}^T \mathbf{S}_b \mathbf{w}$$

$$\text{s.t. } \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$$

等价于



推导

$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w} \quad (3.37)$$

$$\frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = -2\mathbf{S}_b \mathbf{w} + 2\lambda \mathbf{S}_w \mathbf{w}$$

令导数等于 0 即可得：

$$-2\mathbf{S}_b \mathbf{w} + 2\lambda \mathbf{S}_w \mathbf{w} = 0$$

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

$$(\mu_0 - \mu_1)(\mu_0 - \mu_1)^T \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

若令 $(\mu_0 - \mu_1)^T \mathbf{w} = \gamma$ ，则

$$\gamma(\mu_0 - \mu_1) = \lambda \mathbf{S}_w \mathbf{w}$$

$$\mathbf{w} = \frac{\gamma}{\lambda} \mathbf{S}_w^{-1} (\mu_0 - \mu_1)$$

由于最终要求解的 \mathbf{w} 不关心其大小，只关心其方向，所以 γ/λ 这个常数项可以任意取值

“不妨令 $\mathbf{S}_b \mathbf{w} = \lambda(\mu_0 - \mu_1)$ ” 就等价于令 $\frac{\gamma}{\lambda} = 1$

二分类任务 - 线性判别分析

○ 运用拉格朗日乘子法

$$\begin{array}{ll} \min_w & -\mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s.t.} & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1 \end{array} \quad \begin{array}{c} \text{等价于} \\ \longleftrightarrow \end{array} \quad \mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w} \quad (3.37)$$

λ 是拉格朗日乘子. 注意到 $\mathbf{S}_b \mathbf{w}$ 的方向恒为 $\mu_0 - \mu_1$, 不妨令

$$\mathbf{S}_b \mathbf{w} = \lambda(\mu_0 - \mu_1),$$

代入式(3.37)即得

$$\mathbf{w} = \mathbf{S}_w^{-1}(\mu_0 - \mu_1)$$

二分类任务 - 线性判别分析

○ 同向向量

$$\mathbf{S}_b \mathbf{w} = \lambda (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$$

○ 结果

同向向量

$$\mathbf{w} = \mathbf{S}_w^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$$

○ 求解

- 考虑数值解的稳定性，通常进行 奇异值分解

$$\mathbf{S}_w = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$$

$\boldsymbol{\Sigma}$ 是一个实对角矩阵，其对角线上的元素是 \mathbf{S}_w 的奇异值

再由 $\mathbf{S}_w^{-1} = \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^T$ 得到 \mathbf{S}_w^{-1}

○ LDA可从贝叶斯决策论解释

- 两类数据同先验、满足高斯分布且协方差相等时，LDA达到最优分类

LDA推广 - 多分类任务

可以将 LDA 推广到多分类任务中. 假定存在 N 个类, 且第 i 类示例数为 m_i .

○ 定义 全局散度矩阵

$$\begin{aligned} S_t &= S_b + S_w \\ &= \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T \end{aligned}$$

μ 是所有示例的均值向量. 将类内散度矩阵 S_w 重定义为每个类别的散度矩阵之和, 即

○ 类内散度矩阵

$$S_w = \sum_{i=1}^N S_{w_i}$$

其中

$$S_{w_i} = \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T$$

LDA推广 - 多分类任务

- 定义 全局散度矩阵

$$\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w$$

$$= \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

- 类内散度矩阵

$$\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i}$$

$$\mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$$

- 求解得

$$\mathbf{S}_b = \mathbf{S}_t - \mathbf{S}_w$$

$$= \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

显然, 多分类 LDA 可以有多种实现方法: 使用 \mathbf{S}_b , \mathbf{S}_w , \mathbf{S}_t 三者中的任何两个即可. 常见的一种实现是采用优化目标

LDA推广 - 多分类任务

- 定义 全局散度矩阵 $S_t = S_b + S_w$

$$= \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T$$

$$S_w = \sum_{i=1}^m S_{w_i}$$

- 类内散度矩阵

$$S_{w_i} = \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T$$

求解得

$$S_b = S_t - S_w$$

$$= \sum_{i=1}^N m_i (\mu_i - \mu)(\mu_i - \mu)^T$$

推导: $S_b = S_t - S_w$

$$= \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T - \sum_{i=1}^N \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T$$

$$= \sum_{i=1}^N \left(\sum_{x \in X_i} ((x - \mu)(x - \mu)^T - (x - \mu_i)(x - \mu_i)^T) \right)$$

$$= \sum_{i=1}^N \left(\sum_{x \in X_i} ((x - \mu)(x^T - \mu^T) - (x - \mu_i)(x^T - \mu_i^T)) \right)$$

$$= \sum_{i=1}^N \left(\sum_{x \in X_i} (xx^T - x\mu^T - \mu x^T + \mu\mu^T - xx^T + x\mu_i^T + \mu_i x^T - \mu_i\mu_i^T) \right)$$

$$= \sum_{i=1}^N \left(\sum_{x \in X_i} (-x\mu^T - \mu x^T + \mu\mu^T + x\mu_i^T + \mu_i x^T - \mu_i\mu_i^T) \right)$$

LDA推广 - 多分类任务

推导:

$$S_b = S_t - S_w$$

$$\begin{aligned} &= \sum_{i=1}^m (x_i - \mu)(x_i - \mu)^T - \sum_{i=1}^N \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T \\ &= \sum_{i=1}^N \left(\sum_{x \in X_i} ((x - \mu)(x - \mu)^T - (x - \mu_i)(x - \mu_i)^T) \right) \\ &= \sum_{i=1}^N \left(\sum_{x \in X_i} ((x - \mu)(x^T - \mu^T) - (x - \mu_i)(x^T - \mu_i^T)) \right) \\ &= \sum_{i=1}^N \left(\sum_{x \in X_i} (xx^T - x\mu^T - \mu x^T + \mu\mu^T - xx^T + x\mu_i^T + \mu_i x^T - \mu_i\mu_i^T) \right) \\ &= \sum_{i=1}^N \left(\sum_{x \in X_i} (-x\mu^T - \mu x^T + \mu\mu^T + x\mu_i^T + \mu_i x^T - \mu_i\mu_i^T) \right) \end{aligned}$$

$$\begin{aligned} &= \sum_{i=1}^N \left(- \sum_{x \in X_i} x\mu^T - \sum_{x \in X_i} \mu x^T + \sum_{x \in X_i} \mu\mu^T + \sum_{x \in X_i} x\mu_i^T + \sum_{x \in X_i} \mu_i x^T - \sum_{x \in X_i} \mu_i\mu_i^T \right) \\ &= \sum_{i=1}^N (-m_i\mu_i\mu^T - m_i\mu\mu_i^T + m_i\mu\mu^T + m_i\mu_i\mu_i^T + m_i\mu_i\mu_i^T - m_i\mu_i\mu_i^T) \\ &= \sum_{i=1}^N (-m_i\mu_i\mu^T - m_i\mu\mu_i^T + m_i\mu\mu^T + m_i\mu_i\mu_i^T) \\ &= \sum_{i=1}^N m_i (-\mu_i\mu^T - \mu\mu_i^T + \mu\mu^T + \mu_i\mu_i^T) \\ &= \sum_{i=1}^N m_i (\mu_i - \mu)(\mu_i - \mu)^T \end{aligned}$$

LDA推广 - 多分类任务

■ 优化目标
$$\max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}$$

其中 $\mathbf{W} \in \mathbb{R}^{d \times (N-1)}$, $\text{tr}(\cdot)$ 表示矩阵的迹(trace)

上式可通过如下
广义特征值问题求解

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

\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

- 若将 \mathbf{W} 视为一个投影矩阵, 则多分类LDA将样本投影到 $N-1$ 维空间, $N-1$ 通常远小于数据原有的属性数, 通过投影减少了维数, 并且使用了类别信息, 因此LDA也被视为一种监督降维技术

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（常用）
- **基本思路：拆解法，即将多分类任务拆为若干个二分类任务求解**
 - 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

关键：如何拆分？如何集成？

本节主要介绍拆分策略



多分类学习

- 经典拆分策略
 - 一对一 (One vs. One, OvO)
 - 一对其余 (One vs. Rest, OvR)
 - 多对多 (Many vs. Many, MvM)

不失一般性, 考虑 N 个类别 C_1, C_2, \dots, C_N

给定数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, $y_i \in \{C_1, C_2, \dots, C_N\}$

OvO 将这 N 个类别两两配对, 从而产生 $N(N-1)/2$ 个二分类任务:

例: 为区分类别 C_i 和 C_j 训练一个分类器,

该分类器把 D 中的 C_i 类样例作为正例, C_j 类样例作为反例;

在测试阶段, 新样本将同时提交给所有分类器将得到 $N(N-1)/2$ 个分类结果,
最终结果可通过投票产生:

即把被预测得最多的类别作为最终分类结果



多分类学习 - 一对一

○ 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

○ 任务拆分

- 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

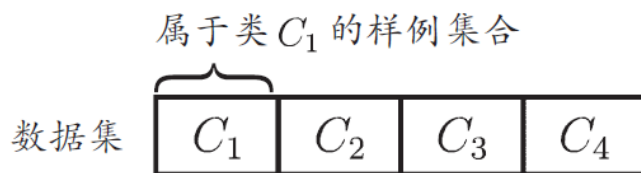
○ 测试阶段

- 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

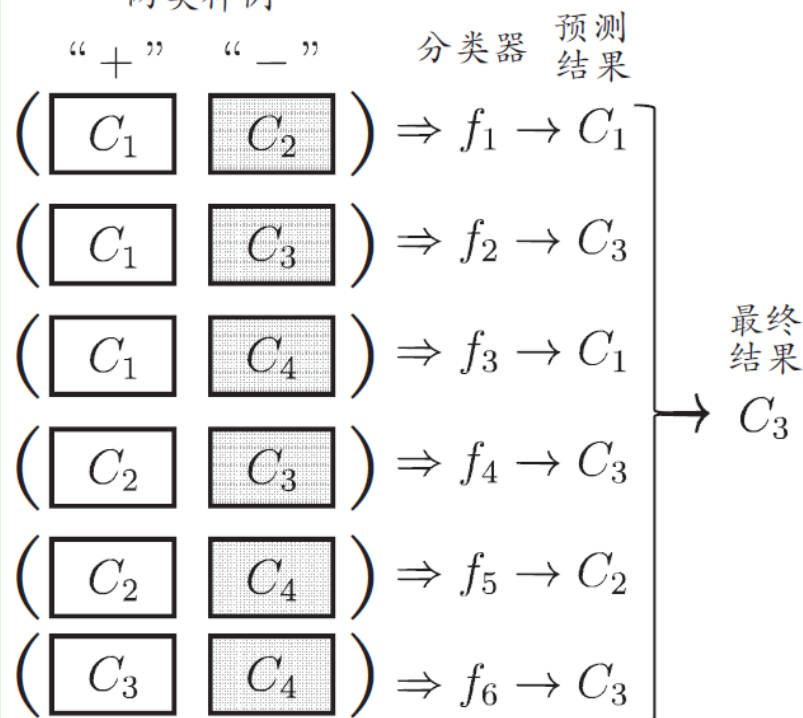


多分类学习 - 两种策略比较

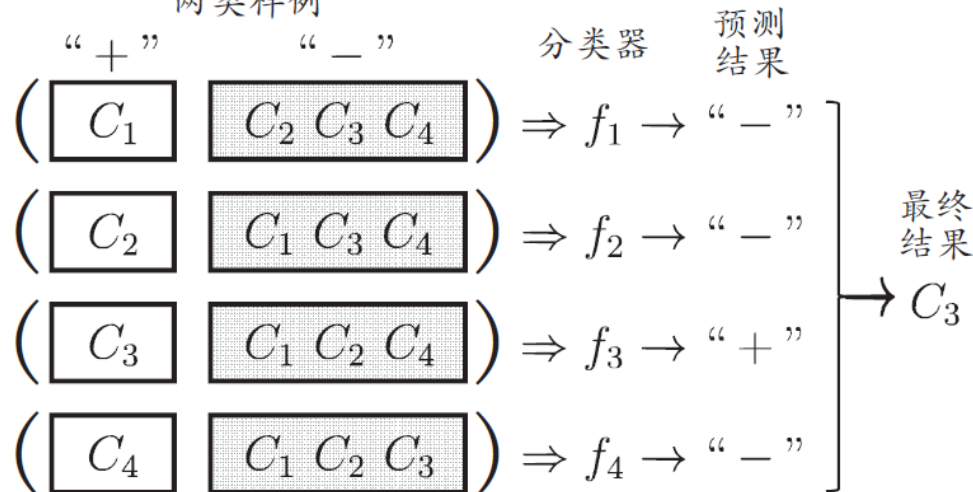
- 一对一 (One vs. One, OvO)
- 一对其余 (One vs. Rest, OvR)
- 多对多 (Many vs. Many, MvM)



用于训练的
两类样例



用于训练的
两类样例



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

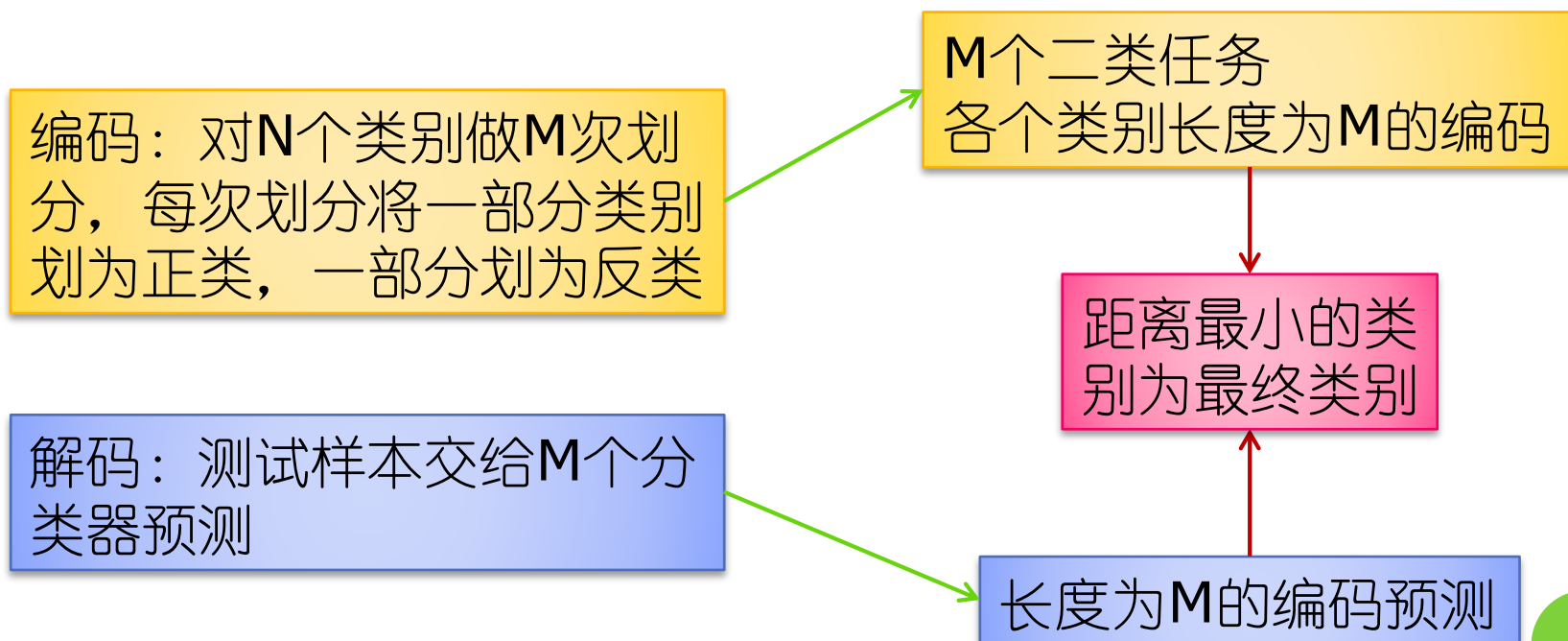
一对其余

- 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)



多分类学习 - 多对多

	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]

	f_1	f_2	f_3	f_4	f_5	f_6	f_7	海明距离	欧氏距离
	↓	↓	↓	↓	↓	↓	↓	↓	↓
$C_1 \rightarrow$	-1	-1	+1	+1	-1	+1	+1	4	4
$C_2 \rightarrow$	-1	0	0	0	+1	-1	0	2	2
$C_3 \rightarrow$	+1	+1	-1	-1	-1	+1	-1	5	$2\sqrt{5}$
$C_4 \rightarrow$	-1	+1	0	+1	-1	0	+1	3	$\sqrt{10}$
测试示例 \rightarrow	-1	+1	+1	-1	+1	-1	+1	↑	↑

(b) 三元 ECOC 码

[Allwein et al. 2000]

- ECOC编码对分类器错误有一定容忍和修正能力，编码越长、纠错能力越强
- 对同等长度的编码，理论上来说，任意两个类别之间的编码距离越远，则纠错能力越强

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 $\frac{y}{1-y} > 1$  $\frac{y}{1-y} > \frac{m^+}{m^-}$ 正负类比例

再缩放

- 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)



优化提要

- 各任务下（回归、分类）各个模型优化的目标
 - 最小二乘法：最小化均方误差
 - 对数几率回归：最大化样本分布似然
 - 线性判别分析：投影空间内最小（大）化类内（间）散度
- 参数的优化方法
 - 最小二乘法：线性代数
 - 对数几率回归：凸优化梯度下降、牛顿法
 - 线性判别分析：矩阵论、广义瑞利商



总结

- 线性回归
 - 最小二乘法（最小化均方误差）
- 二分类任务
 - 对数几率回归
 - 单位阶跃函数、对数几率函数、极大似然法
 - 线性判别分析
 - 最大化广义瑞利商
- 多分类学习
 - 一对一
 - 一对其余
 - 多对多
 - 纠错输出码
- 类别不平衡问题
 - 基本策略：再缩放

