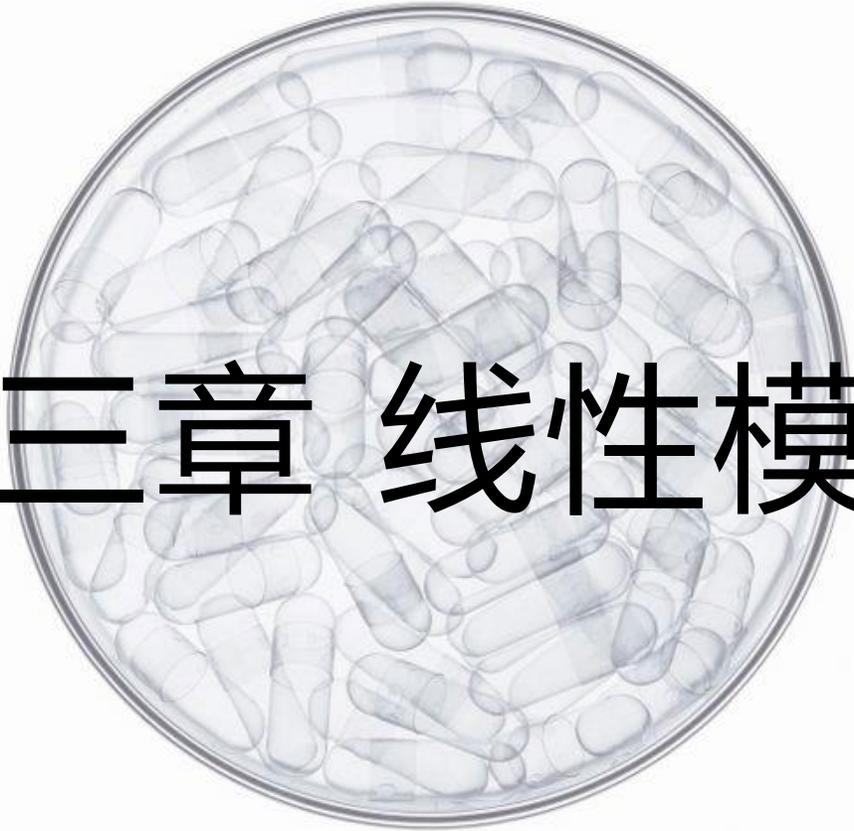


机器学习



第三章 线性模型

汇报人：杜昭武

目录

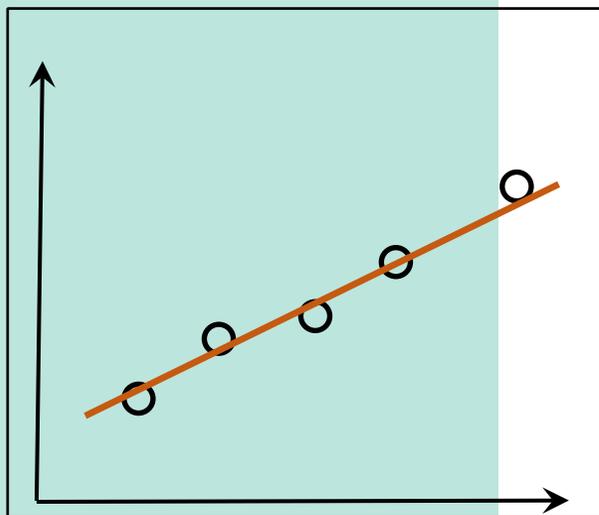
1. 基本形式
2. 线性回归
3. 对数几率回归
4. 线性判别分析
5. 多分类学习
6. 类别不平衡问题

一、基本形式

线性模型 (linear model) 试图学得一个通过属性的线性组合来进行预测的函数, 即向量形式为

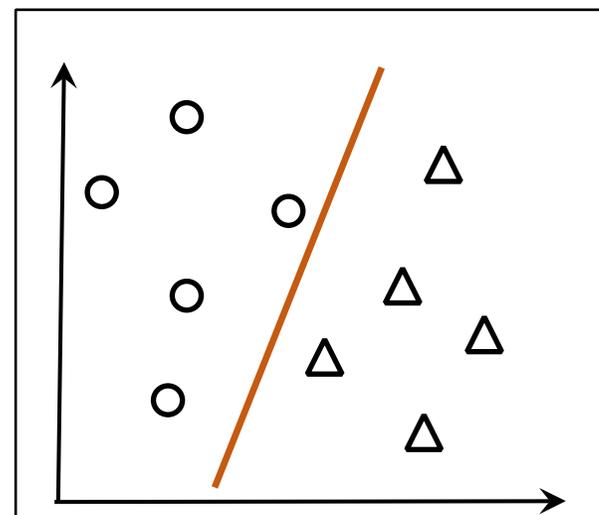
$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$

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



回
归

- 线性回归



分
类

- 对数几率
- 线性判别分析
- 多分类学习

二、线性回归

1. 假设属性数目只有一个：线性回归试图学得

$$f(x_i) = wx_i + b \quad \text{使得} \quad f(x_i) \simeq y_i$$

非数值属性的变换问题：离散属性的处理：若有“序”(order)，则连续化；否则，转化为 k 维向量

- 若属性值之间存在“序”的关系，可通过连续化将其转换为连续值，例如二值属性“身高”的取值“高”“矮”可转化为{1.0, 0.0}
- 若属性之间不存在序的关系。假定有 k 个属性，则通常转换为 k 维向量，例如属性“瓜类”的取值，“西瓜”“黄瓜”“南瓜”可转化为 (0,0,1) (0,1,0) (1,0,0)

令均方误差最小化，有

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

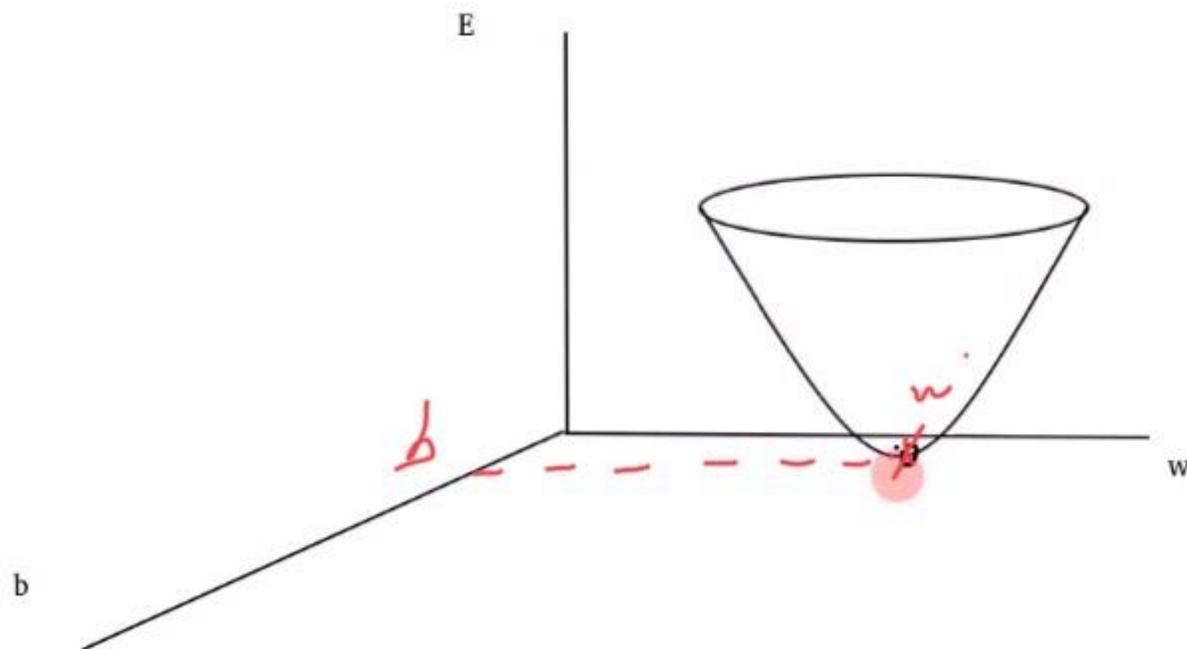
二、线性回归

$$E_{(w,b)} = \sum_{i=1}^m (y_i - wx_i - b)^2 \quad \text{用最小二乘进行估计}$$

这里需要注意的是， w 和 b 是未知数； x 和 y 是常数。选择合适的 w ， b 使得均方误差最小，将式子展开得到

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

$$E(w,b) = (\sum x_i) w^2 + mb^2 + (-2\sum(x_i y_i)) w + (-2\sum y_i) b + (-2\sum x_i) wb + \sum y_i^2$$



二、线性回归

分别对 w 和 b 求导:

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

令导数为 0, 得到闭式(closed-form)解:

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

二、线性回归

2. 一般的情形数据集D, 样本由d个属性描述: 此时线性回归试图学得

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

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

把 w和 b吸收入向量形式, 数据集表示为

$$\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} & 1 \end{pmatrix} * \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \\ b \end{pmatrix} = \begin{pmatrix} f(\vec{x}_1) \\ f(\vec{x}_2) \\ \vdots \\ f(\vec{x}_m) \end{pmatrix}$$

二、线性回归

同样用最小二乘法求解

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

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

附录A

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \quad \text{令其为零可得 } \hat{\mathbf{w}}$$

- 若 满秩或正定, 令上式等于0则可解出
- 若 不满秩, 则可解出多个 $\hat{\mathbf{w}}$

此时需要求助于归纳偏好, 或引入 **正则化 (regularization)** 第6、11章

$$2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) = 0$$

$$\cancel{2}\mathbf{X}^T \mathbf{X}\hat{\mathbf{w}} = \cancel{2}\mathbf{X}^T \mathbf{y}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

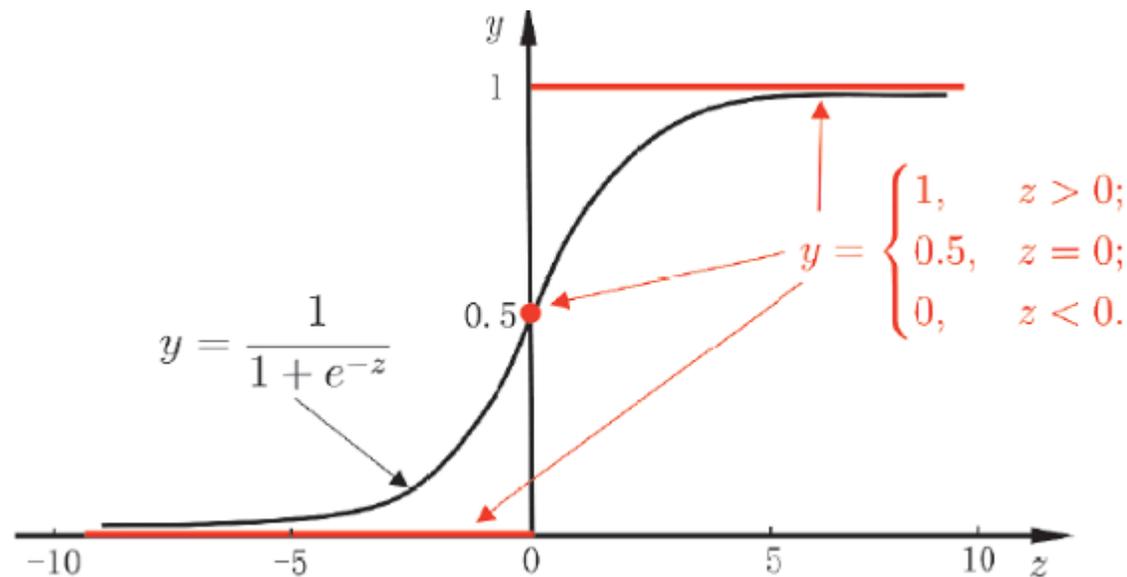
三、对数几率回归

1. 二分类任务:

线性回归模型产生的实值输出 $z = \mathbf{w}^T \mathbf{x} + b$
期望输出 $y \in \{0, 1\}$ } 找 z 和 y 的联系函数

理想的“单位阶跃函数” (unit-step function)

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



- 阶跃函数性质不好，需要找“替代函数” (surrogate function)
- 常用单调可微且任意阶可导的函数为

$$y = \frac{1}{1 + e^{-z}}$$

对数几率函数 (logistic function) 简称“对率函数”

三、对数几率回归

以对率函数为联系函数：

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

即： $\ln \left(\frac{y}{1-y} \right) = \mathbf{w}^T \mathbf{x} + b$

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，例如 Y 为看作“好瓜”的概率， $1-y$ 则是看作“坏瓜”的概率。 Y 与 $1-y$ 的比值称为“几率” (odds)，反映了 x 作为正例的相对可能性，对几率取对数则得到“对数几率” (log odd,亦称为logit)

“对数几率回归”(logistic regression) 简称“对率回归”

- 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

三、对数几率回归

求解思路：

若将 y 看作类后验概率估计 $p(y = 1 | \mathbf{x})$ 则

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

解的

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

于是，使用“极大似然法” (maximum likelihood method)

给定数据集

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

$$P(y_1 | \vec{x}_1) * P(y_2 | \vec{x}_2) * \dots * P(y_m | \vec{x}_m)$$

对率回归模型最大化“对数似然” (log-likelihood) 函数

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

三、对数几率回归

求解思路:

$$\text{令 } \boldsymbol{\beta} = (\boldsymbol{w}; b) \quad \hat{\boldsymbol{x}} = (\boldsymbol{x}; 1) \quad \boldsymbol{w}^T \boldsymbol{x} + b \quad \text{可简写为 } \boldsymbol{\beta}^T \hat{\boldsymbol{x}}$$

$$\text{再令 } p_1(\hat{\boldsymbol{x}}_i; \boldsymbol{\beta}) = p(y = 1 \mid \hat{\boldsymbol{x}}_i; \boldsymbol{\beta}) = \frac{e^{\boldsymbol{w}^T \boldsymbol{x} + b}}{1 + e^{\boldsymbol{w}^T \boldsymbol{x} + b}}$$

$$p_0(\hat{\boldsymbol{x}}_i; \boldsymbol{\beta}) = p(y = 0 \mid \hat{\boldsymbol{x}}_i; \boldsymbol{\beta}) = 1 - p_1(\hat{\boldsymbol{x}}_i; \boldsymbol{\beta}) = \frac{1}{1 + e^{\boldsymbol{w}^T \boldsymbol{x} + b}}$$

则似然项可重写为 $p(y_i \mid \boldsymbol{x}_i; \boldsymbol{w}_i, b) = y_i p_1(\hat{\boldsymbol{x}}_i; \boldsymbol{\beta}) + (1 - y_i) p_0(\hat{\boldsymbol{x}}_i; \boldsymbol{\beta})$

于是, 最大化似然函数 $\ell(\boldsymbol{w}, b) = \sum_{i=1}^m \ln p(y_i \mid \boldsymbol{x}_i; \boldsymbol{w}, b)$

等价于最小化 $\ell(\boldsymbol{\beta}) = \sum_{i=1}^m \left(-y_i \boldsymbol{\beta}^T \hat{\boldsymbol{x}}_i + \ln \left(1 + e^{\boldsymbol{\beta}^T \hat{\boldsymbol{x}}_i} \right) \right)$

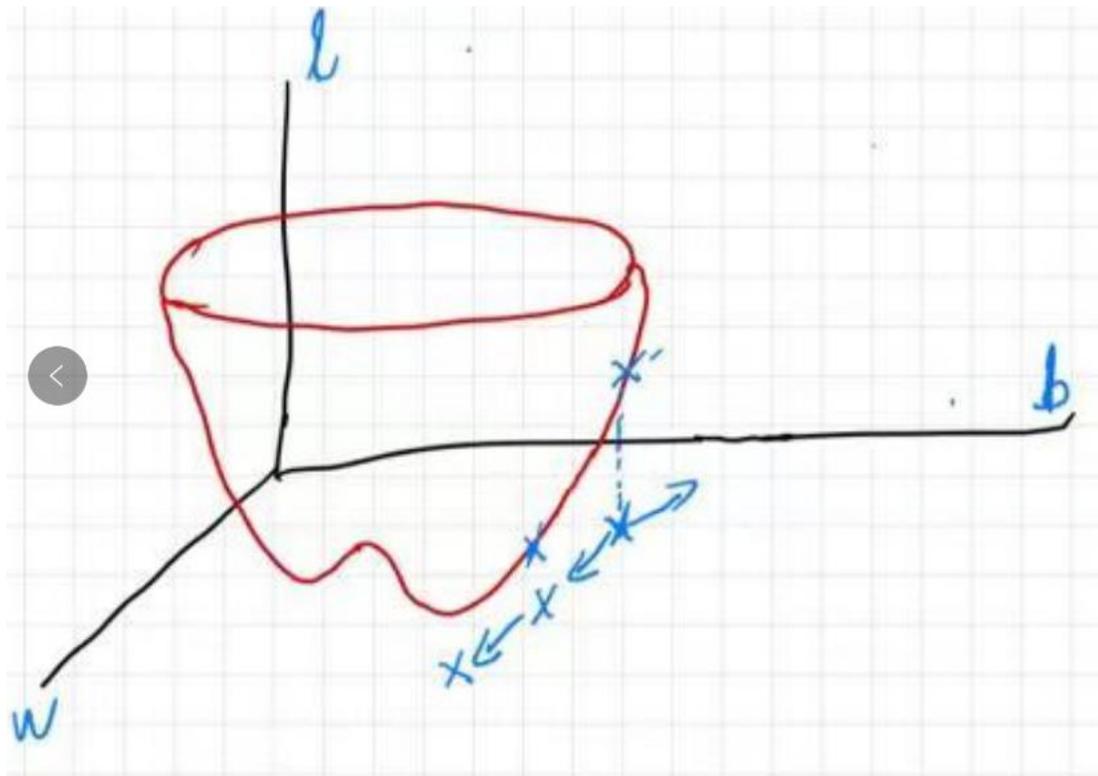
三、对数几率回归

求解思路：

高阶可导连续凸函数，可用经典的数值优化方法 如梯度下降法/牛顿法 [Boyd and Vandenberghe, 2004]

过程：假设只有两个参数的情形下，数值化寻找极小点

1. 取曲面上随机点，计算梯度
2. 朝着梯度下降的方向走一步，
3. 在新的点上计算梯度
4. 最终有一定概率找到极小点



四、线性判别分析

线性判别分析 (LDA) 的思想非常朴素：给定训练集例集，设法将样例投影到一条直线上，使得同类样例的投影点尽可能接近、异类样例的投影点尽可能远离。

由于将样例投影到一条直线（低维空间），因此也被视为一种“监督降维”技术

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

第 i 类样例的集合 X_i

第 i 类样例的均值向量 μ_i

第 i 类样例的协方差矩阵 Σ_i

两类样本的中心在直线上的投影： $w^T \mu_0$ 和 $w^T \mu_1$

两类样本的协方差： $w^T \Sigma_0 w$ 和 $w^T \Sigma_1 w$

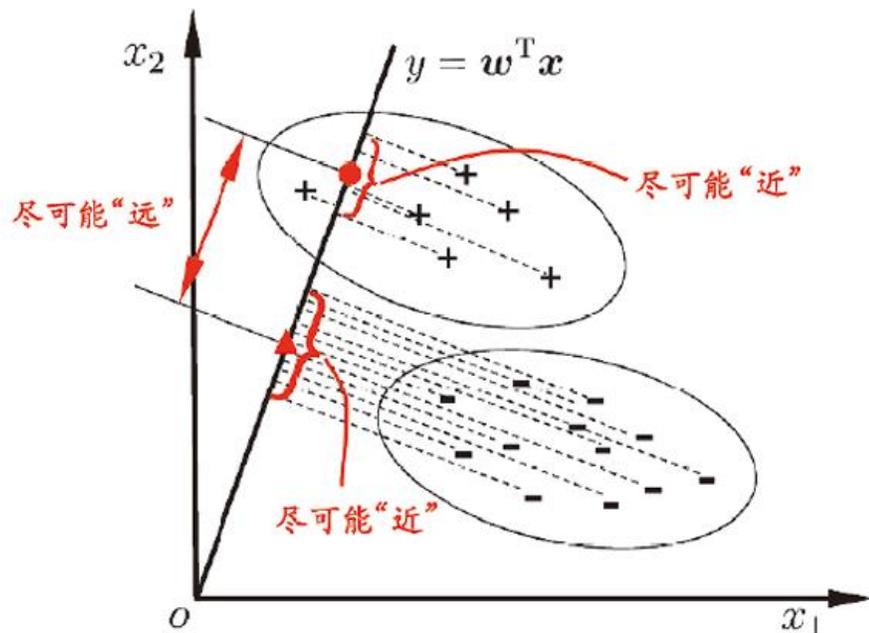
同类样例的投影点尽可能接近 $\rightarrow w^T \Sigma_0 w + w^T \Sigma_1 w$ 尽可能小

异类样例的投影点尽可能远离 $\rightarrow \|w^T \mu_0 - w^T \mu_1\|_2^2$ 尽可能大

于是，最大化

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

$$\Sigma_0 = \sum (\mathbf{x} - \mu_0)(\mathbf{x} - \mu_0)^T$$
$$w^T \Sigma_0 w = \sum w^T (\mathbf{x} - \mu_0)(\mathbf{x} - \mu_0)^T w$$



四、线性判别分析

求解思路：

类内散度矩阵 (within-class scatter matrix):

$$\begin{aligned} \mathbf{S}_w &= \mathbf{\Sigma}_0 + \mathbf{\Sigma}_1 \\ &= \sum_{\mathbf{x} \in X_0} (\mathbf{x} - \boldsymbol{\mu}_0) (\mathbf{x} - \boldsymbol{\mu}_0)^T + \sum_{\mathbf{x} \in X_1} (\mathbf{x} - \boldsymbol{\mu}_1) (\mathbf{x} - \boldsymbol{\mu}_1)^T \end{aligned} \rightarrow$$

类间散度矩阵 (between-class scatter matrix):

$$\mathbf{S}_b = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T$$

LDA的目标：最大化广义瑞利商 (generalized Rayleigh quotient):

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

四、线性判别分析

求解思路：

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

令 $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$ 最大化广义瑞利商等价形式为

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

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

运用拉格朗日乘子法, 有 $\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$

$$\text{由 } \mathbf{S}_b \text{ 定义, 有 } \mathbf{S}_b \mathbf{w} = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \mathbf{w}$$

注意到 $(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \mathbf{w}$ 标量, 令其等于 λ

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

拉格朗日乘子法：

$$\min(-\mathbf{w}^T \mathbf{S}_b \mathbf{w} + \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1))$$

$$\partial[-\mathbf{w}^T \mathbf{S}_b \mathbf{w} + \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1)] / \partial \mathbf{w}$$

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



五、多分类学习

基本思路：“拆解法”，将多分类学习任务拆为若干个二分类任务求解，最经典的三种拆分策略：

“一对一” (OvO)、 “一对其余” (OvR) 和 “多对多” (MvM)

OvO:

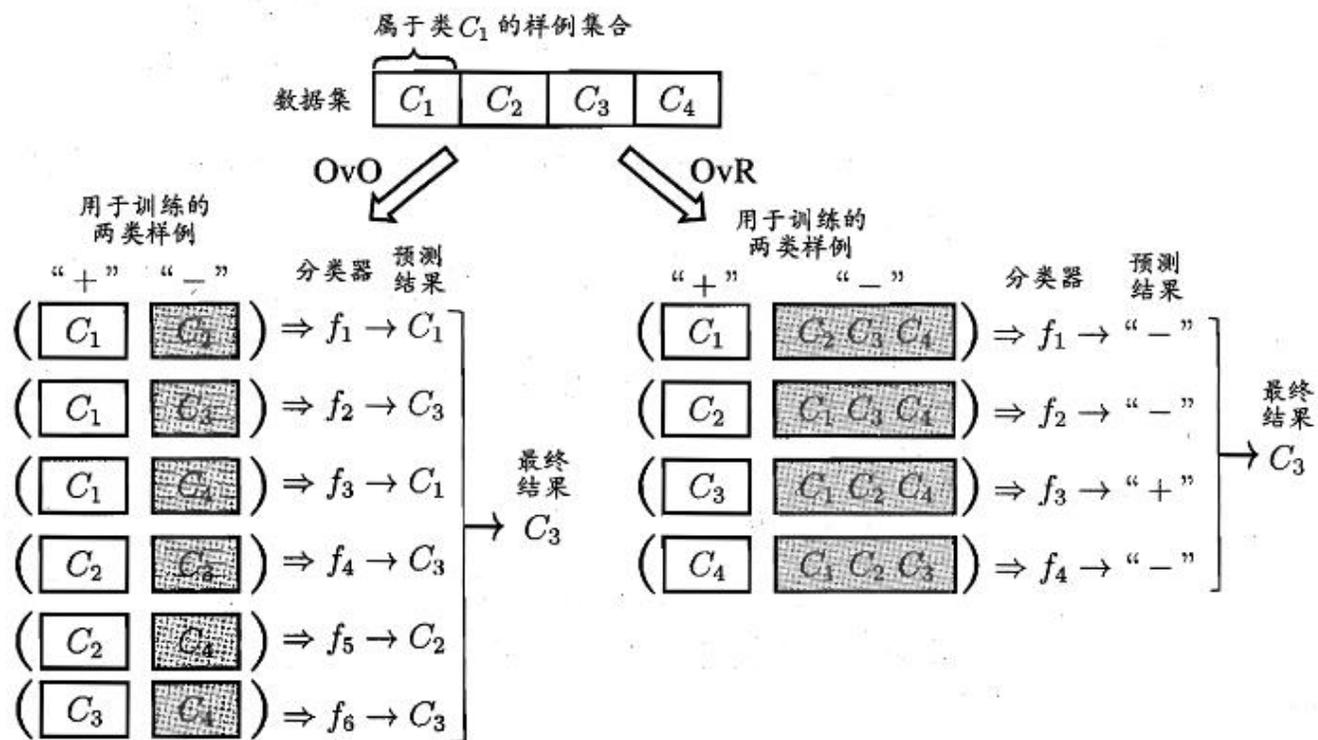
给定数据集

$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}, y_i \in (C_1, C_2, \dots, C_N)$, OvO将这N个类别两两配对，从而产生 $N(N-1)/2$ 个分类任务，最终结果可以通过投票产生：

把被预测得最多的类别作为最终分类结果。

OvR:

则是每次将一个类的样例作为正例，所有其他的样例作为反例来训练N个分类器。在测试时若只有一个分类器预测为正类，则对应的类别标记作为最终分类结果。



五、多分类学习

MVM技术

纠错输出码 (Error Correcting Output Codes, ECOC), ECOC是将编码的思想引入类别拆分, 并尽可能在编码过程中具有容错性。ECOC工作过程主要分为两步:

编码: 对N个类别做M次划分, 每次划分将一部分类别化为正类, 另一部分分为反类, 产生M个训练集——M个分类器。

解码: M个分类器分别预测, 这些预测标记组成一个编码, 将其与每个类别的编码比较, 区别最小的就是最终结果。

类别划分通过“编码矩阵”指定, 即如下二元阵: 将类别指定为正类反类

	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

欧式距离: 在a图中分类器f2将C1和C3类的样例作为正例; 若基于欧式距离, 预测结果是C3。以测试集和第一行为例, $\sqrt{[(1 - (-1))^2 + (-1 - 1)^2 + (1 - (-1))^2]} = \sqrt{12}$

海明距离: 两个合法代码对应位上编码不同的位数为码距, 又称海明码, 例如测试示例: 101010, 分类示例: 100101, 此时海明距离为4

六、类别不均衡问题

定义：类别不平衡，就是指分类任务中不同类别的训练样例数目差别很大的情况

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

用上式对新样本分类，实际上在用预测的 y 值与一个阈值比较。若 $\frac{y}{1-y} > 1$ ，则预测为正例

然而，正反例数目不同时，观测几率 m^+/m^- ，由于我们假设无偏采样，则观测几率就代表真实几率，于是若

$$\frac{y}{1-y} > \frac{m^+}{m^-}，则预测为正例$$

原来的分类器基于上式进行决策，所以需要进行缩放操作，即

$$\frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^+}{m^-}$$

但是很难保证“训练集是真实样本总体的无偏采样”。现有技术有三种做法：

- “欠采样”：直接除去一些反例
- “过采样”：增加一些正例
- “阈值移动”：不增不减，但是预测时采取 $\frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^+}{m^-}$

七、课后习题

3.1 试析在什么情形下式(3.2)中不必考虑偏置项 b . $f(x) = w^T x + b$

偏置项在数值上代表了当自变量取0时, 因变量的取值, 也就是拟合直线的截距, 加偏置项是为了更好的拟合数据。

1. 当学习到的模型恰好经过原点时, 可以不考虑偏置项 b

2. 如果对数据进行了归一化处理, 即对目标变量减去均值向量, 此时就不需要考虑偏置项了

3.2 试证明, 对于参数 w , 对率回归的目标函数(3.18)是非凸的, 但其对数似然函数(3.27)是凸的.

$$y = \frac{1}{1 + e^{-(w^T x + b)}} \quad (3.18)$$

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

3.18 $y = \frac{1}{1 + e^{-(w^T x + b)}}$ 证明是非凸的.
求 = 所有 $y'' < 0$. 则是非凸的
 $\frac{dy}{dw} = \frac{e^{-(w^T x + b)} \cdot x}{[1 + e^{-(w^T x + b)}]^2}$
 $\frac{d^2 y}{dw^2} = \frac{-x^2 \cdot e^{-(w^T x + b)}}{[1 + e^{-(w^T x + b)}]^3} < 0$ y 是非凸的

3.27 $\ell(\beta) = \sum_{i=1}^m (-y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i}))$ $\beta = (w; b)$
 $\frac{\partial \ell}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{x}_i \hat{x}_i^T p(\hat{x}_i; \beta) (1 - p(\hat{x}_i; \beta)) > 0$ y 是凸的.
(3.31式)

七、课后习题

3.6 线性判别分析仅在线性可分数据上能获得理想结果, 试设计一个改进方法, 使其能较好地用于非线性可分数据

对于非线性可分的数据, 要想使用判别分析, 一般思想是将其映射到更高维的空间上, 使它在高维空间上线性可分进一步使用判别分析。

3.7 令码长为 9, 类别数为 4, 试给出海明距离意义下理论最优的 ECOC 二进制码并证明之。

类别数为4, 因此1V3有四种分法, 2V2有六种分法, 3V1同样有四种分法。按照书上的话, 理论上任意两个类别之间的距离越远, 则纠错能力越强。那么可以等同于让各个类别之间的累积距离最大。对于1个2V2分类器, 4个类别的海明距离累积为4; 对于3V1与1V3分类器, 海明距离均为3, 因此认为2V2的效果更好。因此码长为9, 类别数为4的最优EOOC二进制码由6个2V2分类器和3个3v1或1v3分类器构成。

七、课后习题

3.8* ECOC 编码能起到理想纠错作用的重要条件是: 在每一位编码上出错的概率相当且独立. 试析多分类任务经 ECOC 编码后产生的二分类器满足该条件的可能性及由此产生的影响.

重要条件为: ①出错概率相当; ②出错概率相互独立

①书上p66说, 将多个类解分为2个“类别子集”, 所形成的两个类别子集区分难度往往不同, 即其导致的二分类问题难度不同。所以每个编码拆解后类别之间差异越相同, 区分难度越相当, 则满足条件①的可能性越大。但其实很难实现。

②一个好的纠错输出码应该满足条件“各个位上的分类器相互独立”。当类别越多时, 满足条件②的可能性越大。

影响: 一个理论纠错性质很好但二分类问题较难的编码, 与一个理论纠错性质差一些但导致的二分类问题较简单的编码, 难分孰强孰弱。

3.9 使用 OvR 和 MvM 将多分类任务分解为二分类任务求解时, 试述为何无需专门针对类别不平衡性进行处理.

对OvR和MvM来说, 由于每个类进行了相同的处理, 其拆解出来的二分类任务中类别不平衡的影响会相互抵消, 因此通常不需专门处理

谢谢！