

其中关于 β 的一阶、二阶导数分别为

$$\frac{\partial \ell(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{\mathbf{x}}_i (y_i - p_1(\hat{\mathbf{x}}_i; \beta)), \quad (3.30)$$

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T p_1(\hat{\mathbf{x}}_i; \beta) (1 - p_1(\hat{\mathbf{x}}_i; \beta)). \quad (3.31)$$

3.4 线性判别分析

严格说来 LDA 与 Fisher 判别分析稍有不同, 前者假设了各类样本的协方差矩阵相同且满秩.

线性判别分析(Linear Discriminant Analysis, 简称 LDA)是一种经典的线性学习方法, 在二分类问题上因为最早由 [Fisher, 1936] 提出, 亦称“Fisher 判别分析”.

LDA 的思想非常朴素: 给定训练样例集, 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; 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. 图 3.3 给出了一个二维示意图.

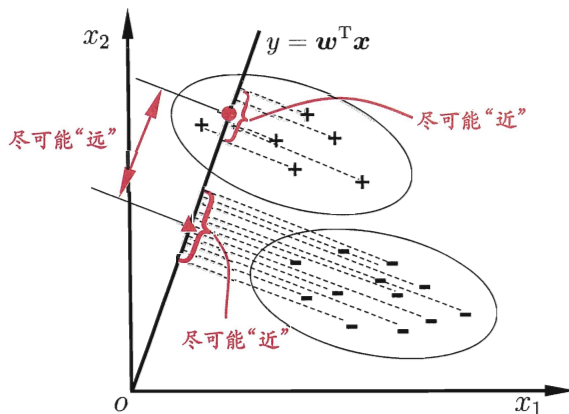


图 3.3 LDA 的二维示意图. “+”、“-”分别代表正例和反例, 椭圆表示数据簇的外轮廓, 虚线表示投影, 红色实心圆和实心三角形分别表示两类样本投影后的中心点.

给定数据集 $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, $y_i \in \{0, 1\}$, 令 X_i 、 μ_i 、 Σ_i 分别表示第 $i \in \{0, 1\}$ 类示例的集合、均值向量、协方差矩阵. 若将数据投影到直线 \mathbf{w} 上, 则两类样本的中心在直线上的投影分别为 $\mathbf{w}^T \mu_0$ 和 $\mathbf{w}^T \mu_1$; 若将所有样本点都投影到直线上, 则两类样本的协方差分别为 $\mathbf{w}^T \Sigma_0 \mathbf{w}$ 和 $\mathbf{w}^T \Sigma_1 \mathbf{w}$. 由于直线是

一维空间, 因此 $w^T \mu_0$ 、 $w^T \mu_1$ 、 $w^T \Sigma_0 w$ 和 $w^T \Sigma_1 w$ 均为实数.

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, 即 $w^T \Sigma_0 w + w^T \Sigma_1 w$ 尽可能小; 而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, 即 $\|w^T \mu_0 - w^T \mu_1\|_2^2$ 尽可能大. 同时考虑二者, 则可得到欲最大化的目标

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \frac{w^T (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w}. \end{aligned} \quad (3.32)$$

定义“类内散度矩阵”(within-class scatter matrix)

$$\begin{aligned} S_w &= \Sigma_0 + \Sigma_1 \\ &= \sum_{x \in X_0} (x - \mu_0)(x - \mu_0)^T + \sum_{x \in X_1} (x - \mu_1)(x - \mu_1)^T \end{aligned} \quad (3.33)$$

以及“类间散度矩阵”(between-class scatter matrix)

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T, \quad (3.34)$$

则式(3.32)可重写为

$$J = \frac{w^T S_b w}{w^T S_w w}. \quad (3.35)$$

这就是 LDA 欲最大化的目标, 即 S_b 与 S_w 的“广义瑞利商”(generalized Rayleigh quotient).

若 w 是一个解, 则对于任意常数 α , αw 也是式(3.35)的解.

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, 因此式(3.35)的解与 w 的长度无关, 只与其方向有关. 不失一般性, 令 $w^T S_w w = 1$, 则式(3.35)等价于

$$\begin{aligned} \min_w \quad & -w^T S_b w \\ \text{s.t.} \quad & w^T S_w w = 1. \end{aligned} \quad (3.36)$$

拉格朗日乘子法参见附录 B.1.

由拉格朗日乘子法, 上式等价于

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

其中 λ 是拉格朗日乘子. 注意到 $\mathbf{S}_b \mathbf{w}$ 的方向恒为 $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$, 不妨令

$$\mathbf{S}_b \mathbf{w} = \lambda(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1), \quad (3.38)$$

代入式(3.37)即得

$$\mathbf{w} = \mathbf{S}_w^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1). \quad (3.39)$$

奇异值分解参见附录
A.3.

考虑到数值解的稳定性, 在实践中通常是对 \mathbf{S}_w 进行奇异值分解, 即 $\mathbf{S}_w = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$, 这里 $\boldsymbol{\Sigma}$ 是一个实对角矩阵, 其对角线上的元素是 \mathbf{S}_w 的奇异值, 然后再由 $\mathbf{S}_w^{-1} = \mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T$ 得到 \mathbf{S}_w^{-1} .

参见习题 7.5.

值得一提的是, LDA 可从贝叶斯决策理论的角度来阐释, 并可证明, 当两类数据同先验、满足高斯分布且协方差相等时, LDA 可达到最优分类.

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 . 我们先定义“全局散度矩阵”

$$\begin{aligned} \mathbf{S}_t &= \mathbf{S}_b + \mathbf{S}_w \\ &= \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T, \end{aligned} \quad (3.40)$$

其中 $\boldsymbol{\mu}$ 是所有示例的均值向量. 将类内散度矩阵 \mathbf{S}_w 重定义为每个类别的散度矩阵之和, 即

$$\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i}, \quad (3.41)$$

其中

$$\mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T. \quad (3.42)$$

由式(3.40)~(3.42)可得

$$\begin{aligned} \mathbf{S}_b &= \mathbf{S}_t - \mathbf{S}_w \\ &= \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T. \end{aligned} \quad (3.43)$$

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

$$\max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}, \quad (3.44)$$

其中 $\mathbf{W} \in \mathbb{R}^{d \times (N-1)}$, $\text{tr}(\cdot)$ 表示矩阵的迹(trace). 式(3.44)可通过如下广义特征值问题求解:

$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}. \quad (3.45)$$

\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. 于是, 可通过这个投影来减小样本点的维数, 且投影过程中使用了类别信息, 因此 LDA 也常被视为一种经典的监督降维技术.

降维参见第 10 章.

3.5 多分类学习

例如上一节最后介绍的 LDA 推广.

现实中常遇到多分类学习任务. 有些二分类学习方法可直接推广到多分类, 但在更多情形下, 我们是基于一些基本策略, 利用二分类学习器来解决多分类问题.

通常称分类学习器为“分类器”(classifier).

关于多个分类器的集成, 参见第 8 章.

OvR 亦称 OvA (One vs. All), 但 OvA 这个说法不严格, 因为不可能把“所有类”作为反类.

亦可根据各分类器的预测置信度等信息进行集成, 参见 8.4 节.

不失一般性, 考虑 N 个类别 C_1, C_2, \dots, C_N , 多分类学习的基本思路是“拆解法”, 即将多分类任务拆为若干个二分类任务求解. 具体来说, 先对问题进行拆分, 然后为拆出的每个二分类任务训练一个分类器; 在测试时, 对这些分类器的预测结果进行集成以获得最终的多分类结果. 这里的关键是如何对多分类任务进行拆分, 以及如何对多个分类器进行集成. 本节主要介绍拆分策略.

最经典的拆分策略有三种: “一对一”(One vs. One, 简称 OvO)、“一对其余”(One vs. Rest, 简称 OvR)和“多对多”(Many vs. Many, 简称 MvM).

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, 例如 OvO 将为区分类别 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, 于是我们将得到 $N(N-1)/2$ 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. 图 3.4 给出了一个示意图.

OvR 则是每次将一个类的样例作为正例、所有其他类的样例作为反例来训练 N 个分类器. 在测试时若仅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, 如图 3.4 所示. 若有多个分类器预测为正类, 则通常考虑各