

Linear classification 03

Chen Gong

31 October 2019

本小节为线性分类的第三小节，主要推导了线性判别分析算法，也就是 Fisher 算法。Fisher 算法的主要思想是：**类内小，类间大**。这有点类似于，软件过程里的松耦合，高内聚的思想。这个思想转换成数学语言也就是，同一类数据之间的方差要小，不同类数据之间的均值的差距要大。那么，我们对数据的描述如下所示：

$$X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{32} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}_{N \times P} \quad (1)$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1} \quad (2)$$

那么，我们的数据集可以记为 $\{(x_i, y_i)\}_{i=1}^N$ ，其中， $x_i \in \mathbb{R}^p$ ， $y_i \in \{+1, -1\}$ ，且 $\{y = +1\}$ 为 C_1 类，且 $\{y = -1\}$ 为 C_2 类。那么， X_{c_1} 被定义为 $\{x_i | y_i = +1\}$ ， X_{c_2} 被定义为 $\{x_i | y_i = -1\}$ 。所以，很显然可以得到 $|X_{c_1}| = N_1$ ， $|X_{c_2}| = N_2$ ，并且 $N_1 + N_2 = N$ 。

1 Fisher 线性判别分析

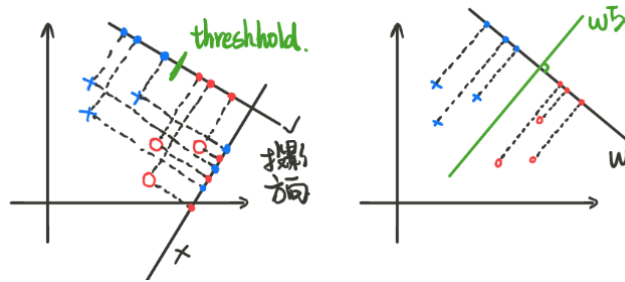


图 1: Fisher 线性判别分析模型图

在左图中，我们设置了两个投影方向，很显然上面那个投影方向可以更好的将两个点之间分开。我们可以在投影方向上找一个点作为两个类别的分界点，也就是阈值 (Threshold)。首先，我们先引入

一个有关投影算法的小知识。

1.1 投影算法

首先，我们需要设定一个投影向量 w ，为了保险起见，对这个投影向量 w 作出约束，令 $\|w\| = 1$ 。那么，在空间中的一个数据点，也就是一个向量，在投影向量上的投影长度可以表述为：

$$x_i \cdot w = |x_i| |w| \cos \theta = |x_i| \cos \theta = \Delta \quad (3)$$

1.2 Fisher 判别分析的损失函数表达式

在这个部分，主要是要得出 Fisher 判别分析的损失函数表达式求法。对于，投影的平均值和方差，我们可以分别表述为：

$$\bar{z} = \frac{1}{N} \sum_{i=1}^N z_i = \frac{1}{N} \sum_{i=1}^N w^T x_i \quad (4)$$

$$S_z = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})(z_i - \bar{z})^T \quad (5)$$

那么对于第一类分类点 X_{c_1} 和第二类分类点 X_{c_2} 可以表述为：

$$C_1 : \quad \bar{z}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i \quad S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (z_i - \bar{z}_1)(z_i - \bar{z}_1)^T \quad (6)$$

$$C_2 : \quad \bar{z}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i \quad S_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} (z_i - \bar{z}_2)(z_i - \bar{z}_2)^T \quad (7)$$

那么类间的距离我们可以定义为： $(\bar{z}_1 - \bar{z}_2)^2$ ，类内的距离被我们定义为 $S_1 + S_2$ 。那么我们的目标函数 Target Function $\mathcal{J}(w)$ ，可以被定义为：

$$\mathcal{J}(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{S_1 + S_2} \quad (8)$$

因为，我们的目的是使方差越小越好，均值之间的差越大越好。

1.3 损失函数表达式的化简

1.3.1 $(\bar{z}_1 - \bar{z}_2)^2$

分子的化简过程如下所示：

$$\begin{aligned} (\bar{z}_1 - \bar{z}_2)^2 &= \left(\frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i \right)^2 \\ &= \left(w^T \left(\frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} x_i \right) \right)^2 \\ &= (w^T (\bar{X}_{c_1} - \bar{X}_{c_2}))^2 \\ &= w^T (\bar{X}_{c_1} - \bar{X}_{c_2})(\bar{X}_{c_1} - \bar{X}_{c_2})^T w \end{aligned} \quad (9)$$

1.3.2 $S_1 + S_2$

分母的化简过程如下所示：

$$\begin{aligned}
 S_1 &= \frac{1}{N_1} \sum_{i=1}^N (z_i - \bar{z}_1)(z_i - \bar{z}_1)^T \\
 &= \frac{1}{N_1} \sum_{i=1}^N \left(w^T x_i - \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i \right) \left(w^T x_i - \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i \right)^T \\
 &= w^T \frac{1}{N_1} \sum_{i=1}^N \left(x_i - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i \right) \left(x_i - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i \right)^T w \\
 &= w^T S_{c_1} w
 \end{aligned} \tag{10}$$

同理可得，

$$S_1 = w^T S_{c_2} w \tag{11}$$

所以，

$$S_1 + S_2 = w^T (S_{c_1} + S_{c_2}) w \tag{12}$$

1.3.3 $\mathcal{J}(w)$ 的最简表达式

$$\mathcal{J}(w) = \frac{w^T (\bar{X}_{c_1} - \bar{X}_{c_2})(\bar{X}_{c_1} - \bar{X}_{c_2})^T w}{w^T (S_{c_1} + S_{c_2}) w} \tag{13}$$

令 S_b 为 between-class 类间方差， S_w 为 within-class，也就是类内方差。那么有

$$S_b = (\bar{X}_{c_1} - \bar{X}_{c_2})(\bar{X}_{c_1} - \bar{X}_{c_2})^T \quad S_w = (S_{c_1} + S_{c_2}) \tag{14}$$

于是，我们可以得到进一步化简的表达式：

$$\mathcal{J}(w) = \frac{w^T S_b w}{w^T S_w w} \tag{15}$$

1.4 损失函数 $\mathcal{J}(w)$ 的梯度

为了方便求导，我们令 $\mathcal{J}(w) = (w^T S_b w)(w^T S_w w)^{-1}$ 。

$$\begin{aligned}
 \frac{\partial \mathcal{J}(w)}{\partial w} &= 2S_b w (w^T S_w w)^{-1} + (-1)(w^T S_b w)(w^T S_w w)^{-2} (2)S_w w = 0 \\
 S_b w (w^T S_w w)^{-1} &= (w^T S_b w)(w^T S_w w)^{-2} S_w w
 \end{aligned} \tag{16}$$

显然， w 的维度是 $p \times 1$ ， w^T 的维度是 $1 \times p$ ， S_w 的维度是 $p \times p$ ，所以， $w^T S_w w$ 是一个实数，同理可得， $w^T S_b w$ 是一个实数所以，可以得到

$$\begin{aligned}
 S_b w &= (w^T S_b w)(w^T S_w w)^{-1} S_w w \\
 S_b w &= \frac{(w^T S_b w)}{(w^T S_w w)} S_w w
 \end{aligned} \tag{17}$$

我们主要是求得梯度的方向，大小不是很重要了。所以，我们可得

$$w = \frac{(w^T S_b w)}{(w^T S_w w)} S_b^{-1} S_w w \propto S_b^{-1} S_w w \tag{18}$$

$$S_w w = (\bar{X}_{c_1} - \bar{X}_{c_2})(\bar{X}_{c_1} - \bar{X}_{c_2})^T w \quad (19)$$

而 $(\bar{X}_{c_1} - \bar{X}_{c_2})^T w$ 是一个实数，所以汇总可得

$$S_b^{-1} S_w w \propto S_w^{-1} (\bar{X}_{c_1} - \bar{X}_{c_2}) \quad (20)$$

那么，我们就可以求得梯度的方向为 $S_w^{-1} (\bar{X}_{c_1} - \bar{X}_{c_2})$ 。如果， S_w^{-1} 是一个各向同性的对角矩阵，那么 $S^{-1} \propto I$ 。所以， $w \propto (\bar{X}_{c_1} - \bar{X}_{c_2})$ 。既然，求得了梯度的方向，其实梯度的大小就不重要的。