

2.2线性判别分析LDA

线性判别分析的基本思想是将高维的模式样本投影到最佳鉴别矢量空间，以达到抽取分类信息和压缩特征空间维数的效果，投影后保证模式样本在新的子空间有最大的类间距离和最小的类内距离，即模式在该空间中有最佳的可分离性。

1.LDA假设以及符号说明：

假设对于一个 R^n 空间有 m 个样本分别为 x_1, x_2, \dots, x_m 即 每个 x 是一个 n 行的矩阵，其中 n_i 表示属于 i 类的样本个数，假设有一个有 c 个类，则 $n_1 + n_2 + \dots + n_i + \dots + n_c = m$ 。

S_b 类间离散度矩阵

S_w 类内离散度矩阵

n_i 属于 i 类的样本个数

x_i 第 i 个样本

u 所有样本的均值

u_i 类 i 的样本均值

2.公式推导，算法形式化描述

根据符号说明可得类 i 的样本均值为：

$$u_i = \frac{1}{n_i} \sum_{x \in \text{class } i} x \tag{1}$$

同理我们也可以得到总体样本均值：

$$u = \frac{1}{m} \sum_{i=1}^m x_i \tag{2}$$

根据类间离散度矩阵和类内离散度矩阵定义，可以得到如下式子：

$$S_b = \sum_{i=1}^c n_i (u_i - u)(u_i - u)^T \tag{3}$$

$$S_w = \sum_{i=1}^c \sum_{x_k \in \text{class } i} (u_i - x_k)(u_i - x_k)^T \tag{4}$$

LDA做为一个分类的算法，我们当然希望它所分的类之间耦合度低，类内的聚合度高，即类内离散度矩阵的中的数值要小，而类间离散度矩阵中的数值要大，这样的分类的效果才好。