

一. 使用方法理解

我使用的是PCA方法

PCA是一种线性降维算法，也是数据预处理方法。

通过计算数据矩阵的协方差矩阵，然后得到协方差矩阵的特征值特征向量，选择特征值最大的k个特征值所对应的特征向量组成的矩阵，再将数据矩阵转换到新的空间当中，实现数据降维。

下图为协方差矩阵计算公式

$$Cov(X, Y, Z) = \begin{bmatrix} Cov(x, x) & Cov(x, y) & Cov(x, z) \\ Cov(y, x) & Cov(y, y) & Cov(y, z) \\ Cov(z, x) & Cov(z, y) & Cov(z, z) \end{bmatrix}$$

本实验中主要使用了weka中四种方法：

1. Instances对象创建方法Instances(java.io.Reader reader):通过使用FileReader读取arff文件，并分配对每个实例分配一个权重
2. PrincipalComponents中的setMaximumAttributeNames (int value):设置使用PCA时留下最高的属性个数
3. PrincipalComponents中的buildEvaluator(Instances data):初始化PCA对象，并且对数据进行分析

二. 数据集处理的思路

先通过创建FileReader对象，并让其读取对应的arff文件。

通过此Reader对象，创建Instances对象

最后通过通过PrincipalComponents中的setMaximumAttributeNames先设置留下属性的个数，再使用buildEvaluator进行分析

三. 实验结果

