



燕山大学

Python 机器学习实验指导书

Python machine learning Experiment Instruction Book

实验四：K-Means

教 务 处

2023 年 5 月

实验四 K-Means

一、实验目的

1. 理解并掌握聚类算法模型。
2. 能够基于 K-Means 聚类算法模型实现鸢尾花聚类、基于经纬度的城市聚类。
3. 能够举一反三，基于 K-Means 聚类算法实现果汁饮料聚类分析。

二、实验原理

聚类 (Clustering) 是一种无监督学习 (unsupervised learning)，简单地说就是把相似的对象归到同一簇中。簇内的对象越相似，聚类的效果越好。

聚类算法是无监督学习中最常见的一种，给定一组数据，需要聚类算法去挖掘数据中的隐含信息。聚类算法的应用很广：顾客行为聚类，google 新闻聚类等。

(一) K-Means 算法

K-Means 算法是最为经典的基于划分的聚簇方法，是十大经典数据挖掘算法之一。简单的说 K-Means 就是在没有任何监督信号的情况下将数据分为 K 份的一种方法。

K-Means 算法的思想：对于给定的样本集，按照样本之间的距离大小，将样本集划分为 K 个簇。让簇内的点尽量紧密的连在一起，而让簇间的距离尽量的大。

具体的算法步骤如下：

1. 随机选择 K 个中心点。
2. 把每个数据点分配到离它最近的中心点。
3. 重新计算每类中的点到该类中心点距离的平均值。
4. 分配每个数据到它最近的中心点。
5. 重复步骤 3 和 4，直到所有的观测值不再被分配或是达到最大的迭代次数（R 把 10 次作为默认迭代次数）。

K-Means 算法也可以采用 Python 里的 Scikit-learn 库来实现，在 Scikit-learn 中的 cluster（聚类）程序包可以调用 KMeans() 函数直接实现对 K-Means 算法的运行，无需过多的编程或者对参数的调整。

(二) 鸢尾花数据集 iris

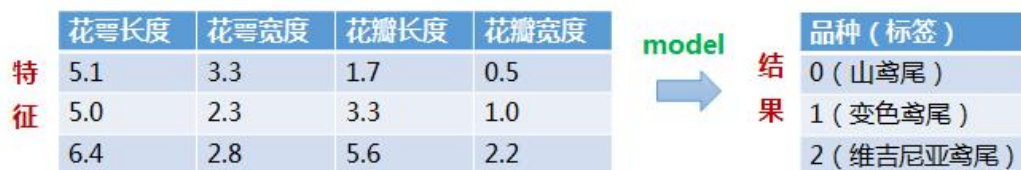
鸢尾花数据集总共包含 150 行数据。每一行数据由 4 个特征值及一个目标值组成。

4 个特征值分别为：萼片长度 (SepalLengthCm)、萼片宽度 (SepalWidthCm)、花瓣长度 (PetalLengthCm)、花瓣宽度 (PetalWidthCm)。

目标值为三种不同类别的鸢尾花,分别为: Iris-setosa(山鸢尾), Iris-versicolor (杂色鸢尾), Iris-virginica (维吉尼亚鸢尾)。

```
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
5.0,3.6,1.4,0.2,Iris-setosa
5.4,3.9,1.7,0.4,Iris-setosa
4.6,3.4,1.4,0.3,Iris-setosa
5.0,3.4,1.5,0.2,Iris-setosa
4.4,2.9,1.4,0.2,Iris-setosa
4.9,3.1,1.5,0.1,Iris-setosa
5.4,3.7,1.5,0.2,Iris-setosa
```

鸢尾花分类数据集



鸢尾花分类

(三) 中国部分二级城市经纬度数据集

数据集总共包含 351 行数据, 每行数据包含一个城市, 共 4 列。其中前两列为城市所在省和市, 最后两列为城市的经纬度。

省级行政区	城市	北纬	东经
北京	北京市	39.904690	116.407170
天津	天津市	39.085100	117.199370
上海	上海市	31.230370	121.473700
重庆	重庆市	29.564710	106.550730
香港特别行政区	九龙	22.327115	114.17495
香港特别行政区	新界	22.341766	114.202408
香港特别行政区	香港岛	22.266416	114.177314
澳门特别行政区	路环岛	22.116226	113.564857
澳门特别行政区	澳门半岛	22.198751	113.549134
澳门特别行政区	氹仔岛	22.156838	113.577669
台湾省	台中市	24.138620	120.679510
台湾省	台北市	25.037798	121.565170
台湾省	台南市	23.172478	120.279363
台湾省	嘉义市	23.481568	120.452538
台湾省	高雄市	22.620856	120.286795
台湾省	基隆市	25.130741	121.746248
台湾省	新北市	25.012366	121.465746
河北省	石家庄市	38.042760	114.514300
河北省	唐山市	39.630480	118.180580
河北省	秦皇岛市	39.935450	119.599640
河北省	邯郸市	36.625560	114.539180

中国部分二级城市经纬度数据集

（四）果汁饮料含量数据集

该数据集是采集自某企业流水线生产的一种果汁饮料含量的数据。

该数据集共有 59 个样本数据。每个样本有 2 个变量，包括 juice（该饮料的果汁含量偏差）、sweet（该饮料的糖分含量偏差），单位均为 mg/ml。

所有特征变量都为与标准含量相比的偏差，该数据集没有目标类别标签变量。

juice	sweet
2.1041	0.8901
-1.0617	-0.4111
0.3521	-1.7488
-0.1962	2.5952
1.4158	1.0928
-0.3951	-1.3293
1.4717	-1.6269
1.1289	1.943
-1.3395	1.4619
-1.784	0.7658
0.2253	-1.6511
-1.7436	-0.8622
1.3344	0.7974
-1.5782	1.5956
1.5828	-1.9999
-1.3934	-1.5497
1.2468	-1.5924
-1.2568	-2.1163
0.9322	1.5899
-1.5856	1.6365

果汁饮料含量数据集

三、实验环境

计算机；网络环境。

四、实验内容及步骤

（一）实验内容

1. 对于给定的例题，基于 K-Means 聚类算法进行鸢尾花分类、基于经纬度的城市聚类等练习。

2. 对于给定的项目，自行编写程序，使用 K-Means 聚类算法实现果汁饮料聚类分析。

（二）实验步骤

1. 【选做】进入指定实验课程，可通过阅览“知识讲解”进行实验相关知识的查缺补漏。

➤ 知识讲解：K-Means

2. 【必做】在熟悉了实验原理的基础上，进行“实验四：K-Means”例题部分的实操练习。

➤ 例题 1：鸢尾花分类（K-Means）（可观看扩展实验 K-Means 讲解视频辅助理解）

➤ 例题 2：基于经纬度的城市聚类

3. 【必做】在步骤 2 例题练习的基础上，对给定的实操项目，自行编写程序，使用 K-Means 聚类算法实现果汁饮料的聚类分析。

➤ 实操项目——不同含量果汁饮料的聚类

4. 【必做】完成实验报告。

请严格基于实验报告的模板撰写实验报告。

本课程所有实验全部结束后再统一打印。

附：实操项目要求

➤ 实操项目——不同含量果汁饮料的聚类

某企业通过采集企业自身流水线生产的一种果汁饮料含量的数据集，来实现 K-Means 算法。通过聚类以判断该果汁饮料在一定标准含量偏差下的生产质量状况，对该饮料进行类别判定。

【数据集】

该数据集共有样本 59 个，变量 2 个，包括 juice（该饮料的果汁含量偏差）、sweet（该饮料的糖分含量偏差），单位均为 mg/ml。

所有特征变量都为与标准含量相比的偏差，该数据集没有目标类别标签变量。

【实验要求】

1. 加载数据集，读取数据，探索数据。

（数据集路径：data/data76878/4_beverage.csv）

2. 样本数据转化（可将 pandasframe 格式的数据转化为数组形式），并进行可视化（绘制散点图），观察数据的分布情况，从而可以得出 k 的几种可能取值。

3. 针对每一种 k 的取值，进行如下操作：

（1）进行 K-Means 算法模型的配置、训练。

（2）输出相关聚类结果，并评估聚类效果。

这里可采用 CH 指标来对聚类有效性进行评估。在最后用每个 k 取值时评估的 CH 值进行对比，可得出 k 取什么值时，聚类效果更优。

注：这里缺乏外部类别信息，故采用内部准则评价指标（CH）来评估。
(`metrics.calinski_harabaz_score()`)

(3) 输出各类簇标签值、各类簇中心，从而判断每类的果汁含量与糖分含量情况。

(4) 聚类结果及其各类簇中心点的可视化（散点图），从而观察各类簇分布情况。
(不同的类表明不同果汁饮料的果汁、糖分含量的偏差情况。)

4. 【扩展】（选做）：设置 k 一定的取值范围，进行聚类并评价不同的聚类结果。

参考思路：设置 k 的取值范围；对不同取值 k 进行训练；计算各对象离各类簇中心的欧氏距离，生成距离表；提取每个对象到其类簇中心的距离，并相加；依次存入距离结果；绘制不同 K 值对应的总距离值折线图。

封面设计： 贾丽

地 址： 中国河北省秦皇岛市河北大街 438 号

邮 编： 066004

电 话： 0335-8057068

传 真： 0335-8057068

网 址： <http://jwc.ysu.edu.cn>