

Clustering with sklearn

姓名：窦金峰

学号：201834862

一、实验任务

Homework3: Clustering with sklearn

- The Tweets dataset is in format of JSON like follows:
 - {"text": "centrepoint winter white gala london", "cluster": 65}
 - {"text": "mourinho seek killer instinct", "cluster": 96}
 - {"text": "roundup golden globe won seduced johansson voice", "cluster": 72}
 - {"text": "travel disruption mount storm cold air sweep south florida", "cluster": 140}

- 测试sklearn中以下聚类算法在tweets数据集上的聚类效果。
- 使用NMI(Normalized Mutual Information)作为评价指标。

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large $n_{samples}$, medium $n_{clusters}$ with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with $n_{samples}$	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with $n_{samples}$	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium $n_{samples}$, small $n_{clusters}$	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large $n_{samples}$ and $n_{clusters}$	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large $n_{samples}$ and $n_{clusters}$	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large $n_{samples}$, medium $n_{clusters}$	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers

<https://scikit-learn.org/stable/modules/clustering.html#>

二、聚类算法

1)、K-means

1、算法原理

K-means 算法是硬聚类算法，是典型的基于原型的目标函数聚类方法的代表，它是数据点到原型的某种距离作为优化的目标函数，利用函数求极值的方法得到迭代运算调整规则。

2、算法步骤

- a) 随机选择 K 个随机的点，成为聚类质心；
- b) 对剩余数据点计算它到每个质心的距离，并把它归到最近的质心的类；
- c) 重新计算已经得到的各个类的质心；

d) 重复 a),b),c)直到新的质心与原质心相等或距离小于指定阈值;

3、参数分析

a) K: 表示聚类个数, 不同的聚类个数会影响聚类效果;

2) 、Affinity Propagation

1、算法原理

AP 算法的基本思想是将全部样本看作网络的节点, 然后通过网络中各条边的消息传递计算出个样本的聚类中心。聚类过程中共有两种消息在各节点间传递, 分别是吸引度和归属度, 直到产生 m 个高质量的 exemplar, 同时将剩余的数据点分配到相应的聚类中。

2、算法步骤

a) 计算初始的相似度矩阵, 将各点之间的吸引度和归属度初始化为 0;

b) 更新各点之间的吸引度, 随之更新各点之间的归属度;

c) 确定当前样本的代表样本点 k ;

d) 重复 b),c), 直到所有的样本的所属不再发生变化为止;

3、参数分析

该聚类算法不需要额外设定参数, 因为不需要事先指定聚类的数量, 而且聚类的结果不会发生变化;

3) 、Spectral Clustering

1、算法原理

谱聚类是一种基于图论的聚类算法, 主要思想是把所有的数据看成空间中的点, 这些点之间可以用边连接起来。距离较远的两个点之间的边权重值较低, 而距离较近的两个点之间的权重值较高, 通过对所有数据点组成的

e) 重复 a),b),c)直到新的质心与原质心相等或距离小于指定阈值;

3、参数分析

b) K: 表示聚类个数, 不同的聚类个数会影响聚类效果;

2) 、Affinity Propagation

1、算法原理

AP 算法的基本思想是将全部样本看作网络的节点, 然后通过网络中各条边的消息传递计算出个样本的聚类中心。聚类过程中共有两种消息在各节点间传递, 分别是吸引度和归属度, 直到产生 m 个高质量的 exemplar, 同时将剩余的数据点分配到相应的聚类中。

2、算法步骤

e) 计算初始的相似度矩阵, 将各点之间的吸引度和归属度初始化为 0;

f) 更新各点之间的吸引度, 随之更新各点之间的归属度;

g) 确定当前样本的代表样本点 k ;

h) 重复 b),c), 直到所有的样本的所属不再发生变化为止;

3、参数分析

该聚类算法不需要额外设定参数, 因为不需要事先指定聚类的数量, 而且聚类的结果不会发生变化;

3) 、Spectral Clustering

1、算法原理

谱聚类是一种基于图论的聚类算法, 主要思想是把所有的数据看成空间中的点, 这些点之间可以用边连接起来。距离较远的两个点之间的边权重值较低, 而距离较近的两个点之间的权重值较高, 通过对所有数据点组成的

图进行切图，让切图后不同的子图间边权重和尽可能的低，而子图内的边权重和尽可能的高，从而达到聚类的目的。

2、算法步骤

- a) 对样本构建相似度矩阵 S ;
- b) 根据相似度矩阵 S 构建邻接矩阵 W 、度矩阵 D
- c) 计算出拉普拉斯矩阵 L
- d) 对拉普拉斯矩阵 L 进行标准化
- e) 计算最小的 K 个特征值所对应的特征向量
- f) 特征向量标准化，并组成特征矩阵 F
- g) 将特征向量按照某种聚类方式聚类

3、参数分析

- a) K : 聚类个数

4)、Agglomerative Clustering

1、算法原理

Agglomerative Clustering 是一种自底而上的层次聚类方法，可以在不同的层次上对数据集进行划分，形成树状的聚类结构。

2、算法步骤

- a) 将每个样本都作为一个簇;
- b) 计算聚类簇之间的距离，找出距离最近的两个簇，将这两个簇合并
- c) 重复 b)，直到聚类簇的数量为 K

3、参数分析

- a) $n_clusters$: 整数，代表聚类个数

b) affinity: 一个字符串或者可调用对象，用于计算距离，可以为：

“euclidean”, “1”, “2”, “manhattan”, “cosine”, 或 ‘precomputed’。但是作为字符串使用的时候，只有“Euclidean”可用！

c) linkage: 一个字符串，用于指定链接算法，有三种方式：

i. ‘ward’: 单链接 single-linkage, 采用；

ii. ‘complete’: 全链接 complete-linkage 算法，采用；

iii. ‘average’: 均连接 average-linkage 算法，采用；

5)、Mean Shift 聚类算法

1、算法原理

Mean shift 算法是基于核密度估计的爬山算法；简单的说，mean shift 就是沿着密度上升的方向寻找同属一个簇的数据点；具体来说就是要定义一个均值漂移，然后不断迭代这个均值漂移的过程。

a) 均值漂移：给定 d 维空间的 n 个数据点集 X ，那么对于空间中的任意

点 x 的 mean shift 向量基本形式可以为： h

这个向量就是漂移向量，其中 h 表示数据集内到 x 的距离小于圆的半

径 h 的数据点。而漂移过程就是利用漂移向量更新球心 x 的位置：

$$x = x + h。$$

2、聚类流程

a) 在未被标记的数据点中随机选择一个点作为 center；

b) 找出离 center 距离在 bandwidth 之内的所有点，记作集合 M ，认为这些点属于簇 c ，同时把集合 M 内的点属于簇 c 的频率加 1，这个参数将用于最后步骤的分类；

- c) 以 center 为中心, 计算漂移向量 shift;
- d) $\text{center} = \text{center} + \text{shift}$, 即将 center 按照 shift 进行移动, 并将得到的新的 center 作为中心;
- e) 重复 b),c),d), 直至收敛, 该过程中遇到的所有的点都应被标记为 c;
- f) 重复 a),b),c),d),e)直至所有的点都被标记;
- g) 将数据点分给被标记频率最高的簇, 即完成分类;

3、参数分析

- a) bandwidth: 浮点数, bandwidth 越小, 分类越多;

6)、DBSCAN

1、算法原理

不同于划分和层次聚类方法, DBSCAN 将密度相连的点的最大集合定义为一个簇, 即由密度可达关系导出的最大密度相连的样本集合就可以作为一个簇。

2、算法步骤:

- a) 检测数据库中尚未检查过的对象 p, 如果 p 未被处理(归为某个簇或者标记为噪声), 则检查其邻域, 若包含的对象数不小于 minPts, 建立新簇 C, 将其中的所有点加入候选集 N;
- b) 对候选集 N 中所有尚未被处理的对象 q, 检查其邻域, 若至少包含 minPts 个对象, 则将这些对象加入 N; 如果 q 未归入任何一个簇, 则将 q 加入 C;
- c) 重复步骤 b), 继续检查 N 中未处理的对象, 当前候选集 N 为空;

d) 重复步骤 a)~c), 直到所有对象都归入了某个簇或标记为噪声。

3、参数分析

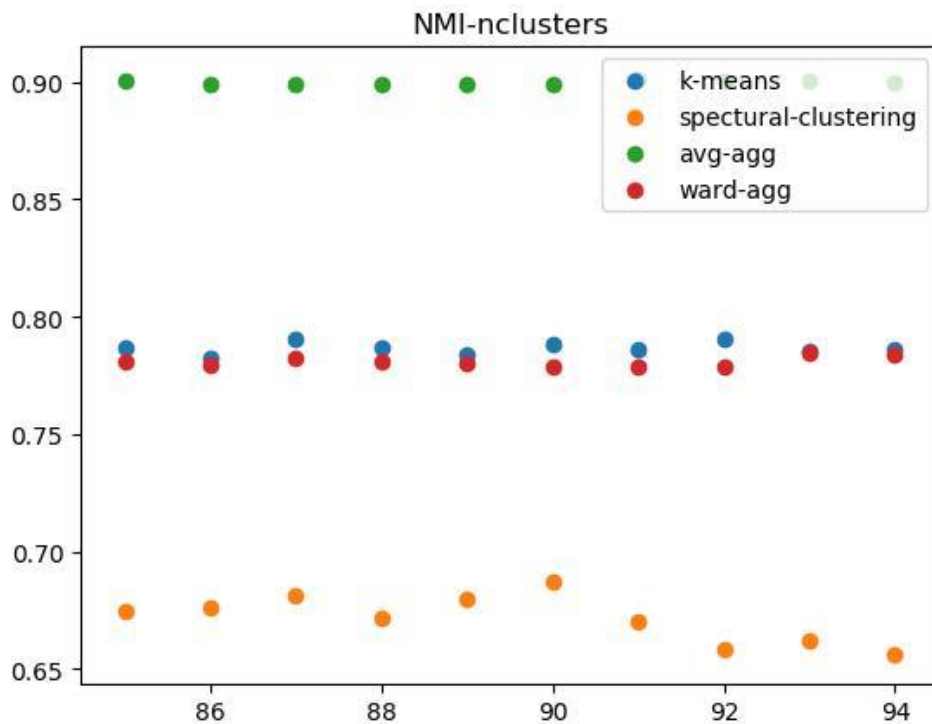
a) Eps: 浮点数, 表示两个互为邻居的 samples 之间的最大距离;

b) Min_samples: 整数, sample 被视为 core point 的最小邻居数;

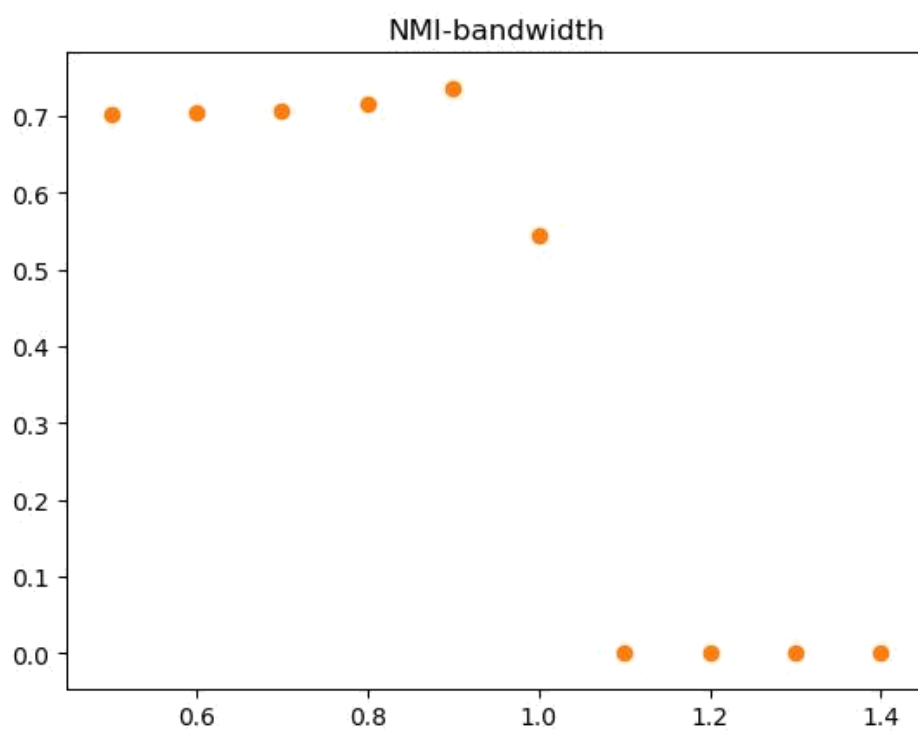
三、聚类效果

1、k-means、spectral-clustering、avg- Agglomerative、ward- Agglomerative 四

种聚类算法的 NMI 随着聚类个数的值分别稳定在 0.8、0.65、0.9、0.78、



2、Mean-shift 聚类算法的聚类效果随着 bandwidth 的变化如下图所示:



3、DBSCAN 聚类算法的聚类效果随着 eps 的变化情况如下图所示：

