

Homework3 实验报告

姓名:	程倩楠	学号:	201834858
实验课程:	Data Mining	实验名称:	Clustering with sklearn
指导老师:	尹建华	完成时间:	2018. 12. 08

一、实验目的

1. 理解各种聚类算法的原理。
2. 学会调用 sklearn 模块中的各种聚类函数和评估函数。

二、实验环境

python3.7

三、实验内容

1. 测试 sklearn 中以下聚类算法在 Tweets 数据集上的聚类效果。

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> 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 <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	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

2. 使用 NMI(Normalized Mutual Information)作为评价指标, 对各种聚类算法进行比较。

四、实验过程

(一) 数据处理

定义一个数据处理函数, 读取 json 文件, 返回 sklearn 聚类函数及评估函数可直接调用的数据形式, 包括以下几个步骤:

1. 读取 json 数据文件，得到文本列表['text1','text2','text3',...]和标签数组[label1 label2 label3 ...]。
2. 调用 sklearn 模块中的 TfidfVectorizer 函数从文本 list 中提取特征，获取 tweets 文本的表示向量，得到一个[n_samples,n_features]维度的稀疏矩阵，每一行对应一个文本表示向量。

TfidfVectorizer 函数实参

```
26 #文本向量化
27 #Extracting features from the dataset, using a sparse vectorizer
28 #Vectorizer results are normalized
29 vectorizer=TfidfVectorizer(max_df=0.5, #过滤掉df>(0.5*doc_num)的单词
30                             max_features=3000, #构建词汇表仅考虑max_features
31                             min_df=2, #过滤掉df<2的单词
32                             stop_words='english', #过滤掉英文停用词
33                             use_idf=True) #启动idf重新计算权重
34 X=vectorizer.fit_transform(data) #稀疏矩阵
```

(二) 聚类算法调用与评估

定义函数，分别调用 sklearn 中的八种聚类算法对 Tweets 数据进行聚类并返回聚类标签。利用 sklearn 自带的 NMI 评估函数对每个聚类算法聚类效果进行单独评价，得到评估分数。

1. K-Means

```
13 #定义函数实现KMeans聚类算法
14 # @param X 数据
15 # @param k 簇的数目
16 # @param minibatch 布尔型(True:MiniBatchKMeans False:KMeans)
17 # @return y_pred,km
18 def KMeansAlgorithm(X, k, minibatch):
19     if minibatch:
20         km=MiniBatchKMeans(n_clusters=k, #形成的簇数以及生成的中心数
21                             init='k-means++', #用智能的方式选择初始聚类中心以加速收敛
22                             n_init=3, #随机初始化的次数
23                             batch_size=100, #Size of the mini batches
24                             )
25     else:
26         km=KMeans(n_clusters=k,
27                   init='k-means++',
28                   max_iter=300, #一次单独运行的最大迭代次数
29                   n_init=10, #使用不同的聚类中心进行初始化的次数
30                   )
31     km.fit(X)
32     y_pred=km.labels_
33     return y_pred,km
```

2. Affinity Propagation

```

11 #定义函数实现AffinityPropagation聚类算法
12 # @param X 样本特征数据
13 # @return y_pred
14 def AffinityPropagationAlgorithm(X):
15     '''
16     1. 参数damping: float, optional, default: 0.5
17       Damping factor (between 0.5 and 1) is the extent to which the current
18       value is maintained relative to incoming values (weighted 1 - damping).
19       This in order to avoid numerical oscillations when updating these values
20       (messages).
21     2. 参数preference: array-like, shape (n_samples,) or float, optional
22       The number of exemplars, ie of clusters, is influenced by the input
23       preferences value.
24       If the preferences are not passed as arguments, they will be set to the
25       median of the input similarities.
26     '''
27     ap=AffinityPropagation(damping=0.5, preference=None)
28     ap.fit(X)
29     y_pred=ap.labels_
30     return y_pred

```

3. Mean-shift

```

11 #定义函数执行AffinityPropagation聚类算法
12 # @param X 样本特征数据
13 # @return y_pred
14 def MeanShiftAlgorithm(X):
15     '''
16     参数bandwidth: float, optional
17     Bandwidth used in the RBF kernel.
18     '''
19     ms=MeanShift(bandwidth=0.9)
20     ms.fit(X)
21     y_pred=ms.labels_
22     return y_pred

```

4. Spectral Clustering

```

12 #定义函数执行SpectralClustering聚类算法
13 # @param X 样本特征数据
14 # @param k 簇的数目
15 # @return y_pred
16 def SpectralClusteringAlgorithm(X,k):
17     #参数n_clusters: integer, optional
18     #The dimension of the projection subspace.
19     sc=SpectralClustering(n_clusters=k)
20     sc.fit(X)
21     y_pred=sc.labels_
22     return y_pred

```

5. Ward Hierarchical Clustering

```

12 #定义函数执行WardHierarchicalClustering聚类算法
13 # @param X 样本特征数据
14 # @param k 簇的数目

```

```

15 # @return y_pred
16 def WardHierarchicalClusteringAlgorithm(X,k):
17     #参数linkage: optional (default="ward")
18     #ward minimizes the variance of the clusters being merged.
19     #参数n_clusters: int, default=2
20     #The number of clusters to find.
21     whc=AgglomerativeClustering(linkage='ward',
22                                 n_clusters=k)
23     whc.fit(X)
24     y_pred=whc.labels_
25     return y_pred

```

6. Agglomerative Clustering

```

12 #定义函数执行AgglomerativeClustering聚类算法
13 # @param X 样本特征数据
14 # @param k 簇的数目
15 # @return y_pred
16 def AgglomerativeClusteringAlgorithm(X,k):
17     '''
18     1. 参数linkage: {"ward", "complete", "average", "single"}, optional
19         average uses the average of the distances of each observation of the
20         two sets.
21     2. 参数n_clusters: int, default=2
22         The number of clusters to find.
23     '''
24     ac=AgglomerativeClustering(linkage='average',
25                                n_clusters=k)
26     ac.fit(X)
27     y_pred=ac.labels_
28     return y_pred

```

7. DBSCAN

```

11 #定义函数执行DBSCAN聚类算法
12 # @param X 样本特征数据
13 # @return y_pred
14 def DBSCANAlgorithm(X):
15     '''
16     1. 参数eps: float, optional
17         The maximum distance between two samples for them to be considered as in
18         the same neighborhood.
19     2. 参数min_samples: int, optional
20         The number of samples (or total weight) in a neighborhood for a point to
21         be considered as a core point.
22         This includes the point itself.
23     '''
24     ds=DBSCAN(eps=0.5,
25               min_samples=1,
26               metric='cosine')
27     ds.fit(X)
28     y_pred=ds.labels_
29     return y_pred

```

8. Gaussian Mixture


```

12 #定义函数执行GaussianMixture聚类算法
13 # @param X 样本特征数据
14 # @param k 簇的数目
15 # @return y_pred
16 def GaussianMixtureAlgorithm(X,k):
17     #参数n_components: int, defaults to 1
18     #The number of mixture components.
19     gm=GaussianMixture(n_components=k)
20     gm.fit(X)
21     y_pred=gm.predict(X)
22     return y_pred

```

(三) 比较 8 种聚类算法

主文件 main.py 循环调用 8 种聚类算法，得到在 Tweets 数据集上的聚类标签，对 NMI 评估分数进行直观的比较。主要包括以下过程：

1. 调用数据处理函数，得到 sklearn 函数可以直接调用的数据格式。
2. 定义聚类算法列表，循环调用 8 种聚类算法对 Tweets 文本进行聚类，得到聚类标签，将 NMI 评估分数保存到评估列表。
3. 绘制条形图，直观比较 8 种聚类算法的效果。
4. 排序并输出聚类算法及其对应的评估分数。

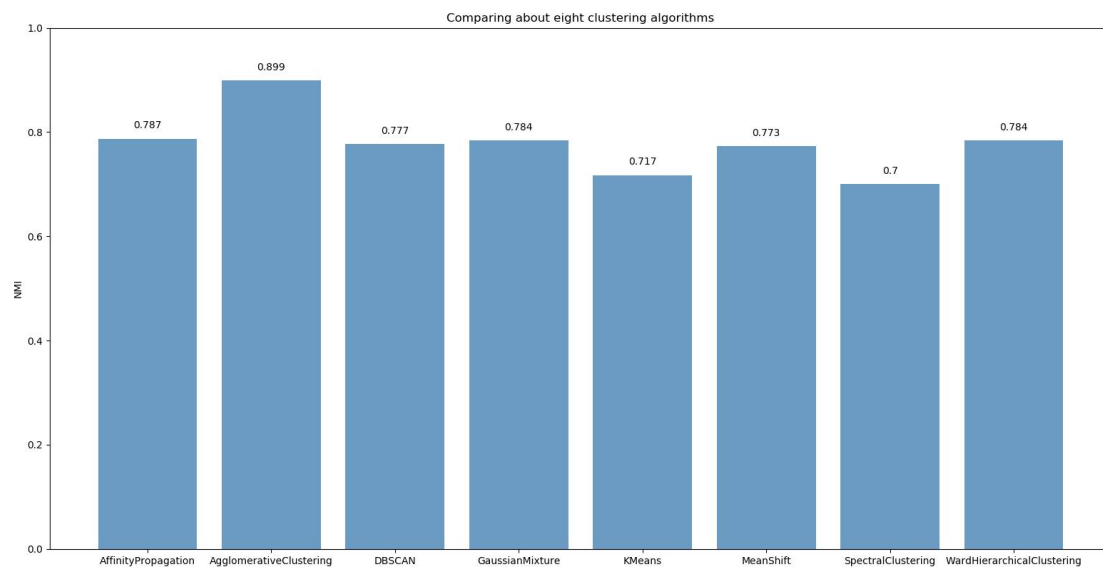
五、实验结果

(一) 单独调用和评估聚类算法

聚类算法	NMI
K-Means	0.708
Affinity Propagation	0.787
Mean-shift	0.773
Spectral Clustering	0.708
Ward Hierarchical Clustering	0.784
Agglomerative Clustering	0.899
DBSCAN	0.777
Gaussian Mixture	0.789

(二) 循环调用和比较聚类算法

1. 绘制条形图，直观比较 8 种聚类算法的效果。其中，横坐标代表 8 种聚类算法，纵坐标代表 NMI 评估分数。



2. 依据 NMI 评估分数对算法进行排序

```
Clustering algorithm and corresponding NMI score:  
1. AgglomerativeClustering: 0.899  
2. AffinityPropagation: 0.787  
3. GaussianMixture: 0.784  
4. WardHierarchicalClustering: 0.784  
5. DBSCAN: 0.777  
6. MeanShift: 0.773  
7. KMeans: 0.717  
8. SpectralClustering: 0.7
```