

第五讲 - 文本聚类

张建章

阿里巴巴商学院
杭州师范大学

2023-02-22



- 1 文本聚类任务定义
- 2 常用文本聚类方法
- 3 凝聚层次聚类
- 4 HAC 类簇相似度计算方法
- 5 HAC 文档聚类代码示例
- 6 聚类评价指标
- 7 课后实践

文本聚类的含义和用途

文本聚类：将文本集合划分成不同的类别，使得同一类别内的文本具有较高的主题相似度，而不同类别的文本具有较低的主题相似度。最常用的探索性文本挖掘技术之一，通常采用无监督机器学习算法，广泛用于商业文档管理、新闻聚合、客服反馈分析。

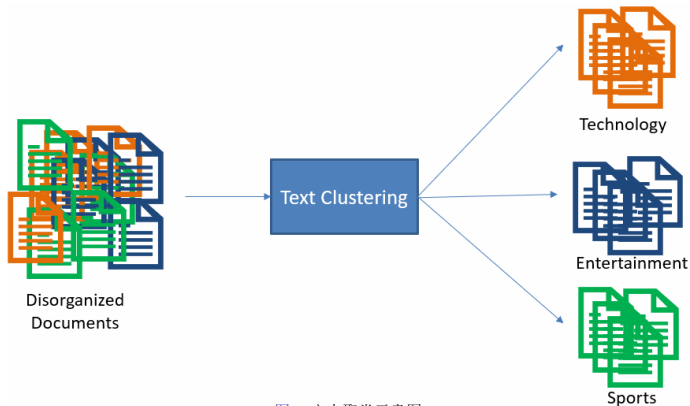


图 1: 文本聚类示意图

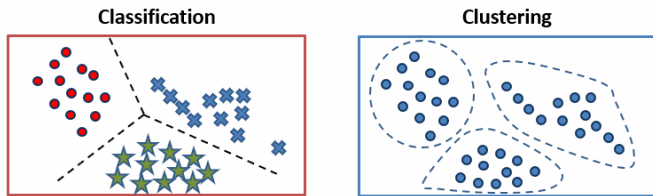
文本聚类与文本分类的区别

两者都是常用的文本挖掘方法，其主要区别如下：

- **学习方式：**文本分类是监督学习方法，需要事先已标注的训练集，来训练模型从而对未知文本分类；文本聚类是无监督学习方法，不需要事先标注数据，根据文本之间的相似度自动将文本分组。

- **应用场景：**文本分类用于有明确分类标准的场景，如垃圾邮件分类、情感分析等；文本聚类适用于探索性数据分析，如新闻聚合、挖掘用户反馈中隐含的主题等。

- **学习算法：**文本分类常用算法有朴素贝叶斯、支持向量机、神经网络等。文本聚类常用算法有 K -均值、层次聚类、DBSCAN 等。



任务形式化定义

给定一个文本集合 $D = \{d_1, d_2, \dots, d_n\}$, d_i 表示一个文档。聚类结果是通过一个分配函数 $f: D \rightarrow C$, 将这些文档分成 K 个类簇, $C = \{C_1, C_2, \dots, C_K\}$, 使每个类簇内的文档相似度较高, 而不同类簇间的文档相似度较低。

文档相似度通常使用余弦相似度:

$$\text{sim}(d_i, d_j) = \cos(\theta) = \frac{\langle d_i, d_j \rangle}{\|d_i\| \|d_j\|} = \frac{d_i \cdot d_j}{\|d_i\| \|d_j\|}$$

最大化类内相似度和最小化类间相似度分别表示如下:

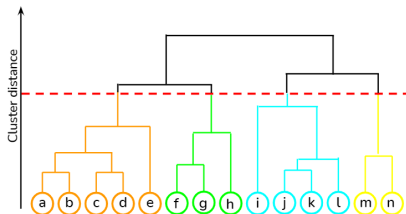
$$\operatorname{argmax} \sum_{k=1}^K \sum_{d_i, d_j \in C_k} \text{sim}(d_i, d_j)$$

$$\operatorname{argmin} \sum_{k=1}^K \sum_{d_i \in C_k, d_j \notin C_k} \text{sim}(d_i, d_j)$$

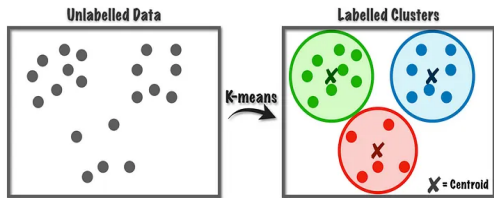
层次聚类和分区聚类

层次聚类：生成树状的聚类结构，可以生成多个聚类层次，展示数据点之间的层次关系，通过截断树状结构来得到不同数量的聚类，代表性算法有自底向上的凝聚层次聚类 (HAC) 和自顶向下的分裂层次聚类。

分区聚类：将数据点分割成一组不重叠的子集或类簇，没有层次结构，通常需要预先确定聚类的数量 K ，代表性算法为 K -均值聚类。



层次聚类



分区聚类

HAC 算法原理

输入：一组向量化表示的文档 (N 个文档)，如 TF-IDF，嵌入向量等；

- ❶ **初始化：**将每个数据点看作一个独立的簇；
- ❷ **计算类簇相似性：**通过余弦相似度计算类簇之间的相似性，计算方式有单链接、全链接、平均链接、质心法；
- ❸ **合并类簇：**将相似度最高的两个类簇合并为一个新的类簇，更新类簇相似性，
- ❹ **重复：**重复步骤❸ $N - 1$ 次，直到全部文档合并成一个类簇。

输出：树形结构 (Dendrogram) 表示文档之间的层次关系，可通过截断树状结构得到所需数量的类簇。

[HAC 原理动图演示 \(点我观看\)](#)

HAC 算法流程

```

SIMPLEHAC( $d_1, \dots, d_N$ )
1  for  $n \leftarrow 1$  to  $N$ 
2  do for  $i \leftarrow 1$  to  $N$ 
3      do  $C[n][i] \leftarrow \text{SIM}(d_n, d_i)$ 
4       $I[n] \leftarrow 1$  (keeps track of active clusters)
5   $A \leftarrow []$  (assembles clustering as a sequence of merges)
6  for  $k \leftarrow 1$  to  $N - 1$ 
7  do  $\langle i, m \rangle \leftarrow \arg \max_{\langle i, m \rangle : i \neq m \wedge I[i]=1 \wedge I[m]=1} C[i][m]$ 
8       $A.\text{APPEND}(\langle i, m \rangle)$  (store merge)
9      for  $j \leftarrow 1$  to  $N$ 
10         do  $C[i][j] \leftarrow \text{SIM}(i, m, j)$ 
11              $C[j][i] \leftarrow \text{SIM}(i, m, j)$ 
12          $I[m] \leftarrow 0$  (deactivate cluster)
13 return  $A$ 

```

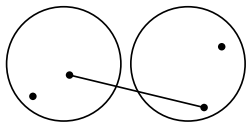

HAC 算法原理解析

- Line 1-3: 计算一个 $N \times N$ 维的相似度矩阵 C ;
- Line 4: 初始化一个 $1 \times N$ 维的 0-1 向量 I , 保存可被继续合并的类簇, 例如, $I[1] = 1$ 表示第二个文档处于未被合并状态, $I[3] = 0$ 表示第三个文档已被合并到某个类簇中, 处于冻结状态;
- Line 5: 初始化一个序列 A , 保存类簇合并过程;
- Line 6-12: 合并最相似的类簇, 并更新矩阵 C 、向量 I 、序列 A 。

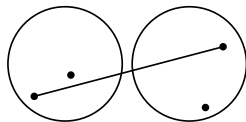
注意: 理解上述算法流程伪代码的关键是将每一个类簇命名为一个数字编号 i , 该编号为类簇中索引值最小的文档的索引值。

4. HAC 类簇相似度计算方法

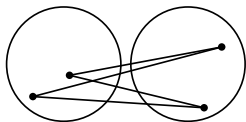
HAC 中的类簇相似度有四种：单链接、全链接、质心法、平均链接：



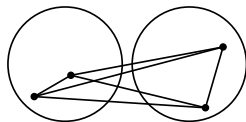
(a) single-link: maximum similarity



(b) complete-link: minimum similarity



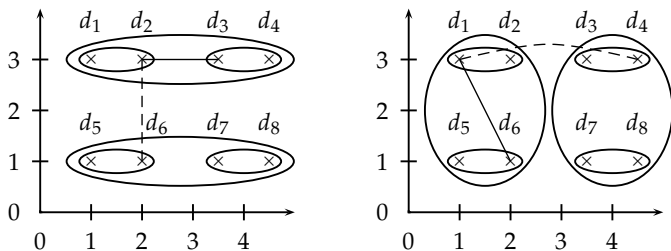
(c) centroid: average inter-similarity



(d) group-average: average of all similarities

类簇相似性示例

使用单链接 (左) 和全链接 (右) 进行 HAC 聚类:



前四次迭代, 左右两图中的聚类过程相同, 可表示为:

$\langle 1, 2 \rangle \rightarrow \langle 3, 4 \rangle \rightarrow \langle 5, 6 \rangle \rightarrow \langle 7, 8 \rangle$

单链接第四次之后的迭代聚类过程, 可表示为:

$\langle 1, 3 \rangle \rightarrow \langle 5, 7 \rangle \rightarrow \langle 1, 5 \rangle$

全链接第四次之后的迭代聚类过程, 可表示为:

$\langle 1, 5 \rangle \rightarrow \langle 3, 7 \rangle \rightarrow \langle 1, 3 \rangle$

平均链接和质心法

单链接和全链接只依赖两个文档的余弦相似度来度量类簇的相似度，无法完整反映一个类簇中的文档的分布情况。平均链接和质心法在计算类簇相似度时会考虑类簇中的全部文档。

质心法计算两个类簇质心的相似度：

$$\left(\frac{1}{N_i} \sum_{d_m \in w_i} \vec{d}_m\right) \cdot \left(\frac{1}{N_j} \sum_{d_n \in w_j} \vec{d}_n\right) = \frac{1}{N_i N_j} \sum_{d_m \in w_i} \sum_{d_n \in w_j} \vec{d}_m \cdot \vec{d}_n$$

平均链接计算两个类簇中全部文档的成对相似度：

$$\frac{1}{(N_i + N_j)(N_i + N_j - 1)} \sum_{d_m \in w_i \cup w_j} \sum_{d_n \in w_i \cup w_j, d_n \neq d_m} \vec{d}_m \cdot \vec{d}_n$$

```
import numpy as np
from sklearn.cluster import AgglomerativeClustering

# 使用TF-IDF向量表示文档集合
tfidf_matrix = np.array([
    [0.1, 0.3, 0.1],
    [0.1, 0.5, 0.2],
    [0.6, 0.1, 0.1],
    [0.7, 0.1, 0.2],
    [0.2, 0.4, 0.6],
    [0.1, 0.2, 0.7]])

# 设定聚类数量
n_clusters = 3
# 创建AgglomerativeClustering对象
clustering = AgglomerativeClustering(n_clusters=n_clusters,
    ↪ affinity='cosine', linkage='average')
# 使用TF-IDF矩阵进行聚类
labels = clustering.fit_predict(tfidf_matrix)
# 打印聚类结果
print("聚类标签:", labels)
```

外部指标和内部指标

外部指标是指在评估聚类效果时，使用与聚类结果独立的外部信息，通常是数据的真实标签，主要用来衡量聚类结果与真实分类之间的一致性。常见的外部指标包括：

纯度 (Purity): 衡量每个聚类中占比最大的真实类别的数量。

兰德指数 (Rand Index): 计算所有成对数据点中，聚类分配与真实标签一致与不一致的比例。

内部指标不依赖于任何外部信息，仅基于数据本身的结构和聚类结果来评估聚类的质量，主要用来衡量聚类的凝聚度和分离度。常见的内部指标包括：

轮廓系数 (Silhouette Coefficient): 评估每个样本到其聚类内部点的平均距离与到最近聚类的平均距离的比率。

戴维斯-布尔丁指数 (Davies-Bouldin Index): 基于聚类内的紧密度和聚类间的分离度，较低的值表示更好的聚类效果。

聚类纯度

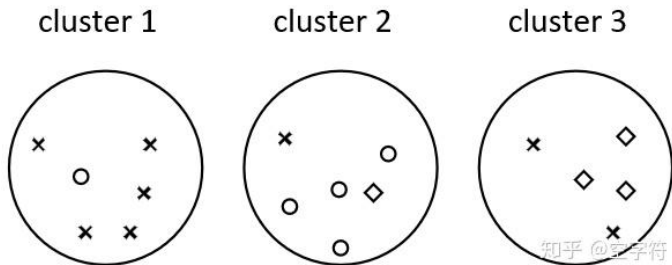
聚类纯度通过检查每个类簇中最多的真实类别成员数来评估聚类结果的纯净度。计算过程为：(1) 对每个聚类，找出数量最多的类别；(2) 将每个聚类的主要类别的数量相加，除以数据集中的总数据点数。

$$\text{纯度} = \frac{1}{N} \sum_{k=1}^K \max_j |C_k \cap T_j|$$

其中， N 是数据集中的总数据点数， K 是聚类的数量， C_k 是第 k 个聚类中的数据点集， T_j 是第 j 个真实类别的数据点集， \max_j 表示选取最大交集的操作。

聚类纯度的值范围从 0 到 1，其中 1 表示完全纯净的聚类，每个聚类完全由单一类别的数据点组成。这种度量对于了解聚类与实际类别的对应关系非常有用，但它不考虑聚类的结构或形状，也不考虑非主导类别的信息。

聚类纯度



$$\text{纯度} = \frac{5 + 4 + 3}{17} = 0.706$$

轮廓系数 I

轮廓系数同时考虑了聚类内的凝聚度和聚类间的分离度。轮廓系数为每个样本点提供了一个度量，用来判断这个点是否被分配到了合适的聚类中。这个系数的值介于-1 到 1 之间，其中接近 1 的值表明样本被很好地分配到了聚类中，接近-1 的值表明样本更适合被分配到另一个聚类中，而接近 0 的值则表明样本在两个聚类的边界上，整体轮廓系数是所有样本点轮廓系数的平均值。单个样本点的轮廓系数计算步骤如下：

1. 计算样本点的凝聚度 a ：凝聚度 a 是指一个样本点与其同一个聚类中所有其他点的平均距离，表示了聚类内部的紧密程度。对于聚类 C_i 中的某个点 x ，其凝聚度计算为：

$$a(x) = \frac{1}{|C_i| - 1} \sum_{x' \in C_i, x' \neq x} d(x, x')$$

其中 $d(x, x')$ 是点 x 和 x' 之间的距离。

轮廓系数 II

2. 计算样本点的分离度 b : 分离度 b 是指一个样本点与最近的其他聚类中所有点的平均距离, 表示不同聚类之间的分隔程度。对于聚类 C_i 中的某个点 x , 找到与其平均距离最小的其他类簇 C_j ($j \neq i$):

$$b(x) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{x' \in C_j} d(x, x')$$

3. 计算轮廓系数 s :

$$s(x) = \frac{b(x) - a(x)}{\max(a(x), b(x))}$$

$s(x)$ 介于 -1 和 1 之间。当 $a(x) < b(x)$ (即分离度大于凝聚度), 轮廓系数是正的, 这表明聚类分配得当; 当 $a(x) > b(x)$, 轮廓系数是负的, 表明 x 可能被分配错了聚类。

1. 使用本课程提供的中文新闻语料库，使用 `scikit-learn` 中聚类算法进行文本聚类，具体要求如下：

- ① 从每个类别的新闻中随机采样 100 条新闻，构造待聚类文档集合；
- ② 分别使用 **HAC** 和 K 均值方法对文档集合进行聚类；
- ③ 分析两种聚类方法聚类结果的差异；

THE END