



南京農業大學

NANJING AGRICULTURAL UNIVERSITY

数字人文下的文本聚类



目录

CONTENT

- ◆ 文本聚类基本知识
- ◆ 非物质文化遗产的自动聚类





南京農業大學

NANJING AGRICULTURAL UNIVERSITY

文本聚类基本知识

1、文本聚类

◆聚类又称群分析，是数据挖掘的一种重要的思想。文本分类一般和聚类一起探讨，其区别在于，其根据文本内部的特征以及类别数量自动将待分类文本聚合成若干个类别。别在于有没有预先划分好的类别。文本分类任务需要预先定义好具体类别，再接着将待分类文本划分入这些类别中；而文本聚类任务则不需要预先定义具体类别，只需提前规定好划分的类别

2、特征提取

◆所提取的特征在一定程度上决定了整个分类的性能。如何从非结构化、半结构化文本中提取相应特征并转化为结构化信息是特征提取的关键。目前常用的特征提取模型如下：One-Hot、BOW词袋模型、连续词袋模型（CBOW）、Skip-Gram模型和Word2vec模型。

3、K-means

◆在整个聚类算法体系中，K-means算法是典型的基于距离的一种算法。距离是该算法相似性计算和评价的基础，如果两个文本对象在计算过程中距离较近则相似度的值较大。在具体的实现聚类的过程中，由距离靠近的文本对象所组成的簇是判定聚类结果的关键。



南京農業大學

NANJING AGRICULTURAL UNIVERSITY

非物质文化遗产的 自动聚类

语料准备

一个非遗项目字段样例如下所示，分别抽取出项目名称和详细介绍作为本章聚类的一个非遗项目。一个项目代表一条数据，如下表所示，共有2690条数据。

项目名称	苗族古歌	项目序号	1
项目编号	I -1	公布时间	2006（第一批）
类别	民间文学	所属地区	贵州省
类型	新增项目	申报地区或单位	贵州省黄平县
详细信息	申报地区或单位：贵州省黄平县。苗族分布在我国西南数省区。按方言划分，大致可分为湘西方言区、黔东南方言区、川滇黔方言区。黔东南清水江流域一带是全国苗族最大的聚居区，大致包括凯里、剑河、黄平、台江、雷山、丹寨、施秉、黄平、镇远、三穗，以及广西三江和湖南靖县等地。在此广大苗族聚居区普遍流传着一种以创世为主体内容的诗体神话，俗称“古歌”或“古歌古词”...（文本内容较长，不予全文显示）		
相关继承人信息	编号： 01-0004, 姓名： 王明芝, 性别： 女, 出生日期： 1939.06, 民族： 类别： 民间文学, 项目编号： I -1, 项目名称： 苗族古歌. 申报地区或单位： 贵州省黄平县. 编号： 01-0003, 姓名： 龙通珍, 性别： 女, 出生日期： 1936.04, 民族： 类别： 民间文学, 项目编号： I -1, 项目名称： 苗族古歌. 申报地区或单位： 贵州省黄平县		
相关项目信息	I -1, 民间文学, 湖南省花垣县, I -1, 民间文学, 贵州省台江县		

数据预处理

◆经过数据清洗后，对文本进行分词。分词工具采用jieba分词，停用词表为哈工大停用词表。

一条数据的分词结果如下：

苗族古歌 申报 地区 单位 贵州省 黄平县 苗族 分布 我国 西南 数 省区 方言 划分 分为 湘西 方言 区 黔 东 方 言 区 川 滇 黔 方 言 区 黔 东 南 清 水 江 流 域 一 带 全 国 苗 族 聚 居 区 包 括 凯 里 剑 河 黄 平 台 江 雷 山 丹 寨 施 秉 黄 平 镇 远 三 穗 广 西 三 江 湖 南 靖 县 苗 族 聚 居 区 流 传 一 种 创 世 主 体 内 容 诗 体 神 话 俗 称 古 歌 古 歌 古 词 苗 族 古 歌 内 容 包 罗 万 象 宇 宙 诞 生 人 类 物 种 起 源 开 天 辟 地 初 民 时 期 滔 天 洪 水 苗 族 迁 徙 苗 族 古 代 社 会 制 度 日 常 生 产 生 活 无 所 不 包 苗 族 古 代 神 话 总 汇 苗 族 古 歌 古 词 神 话 鼓 社 祭 婚 丧 活 动 亲 友 聚 会 节 日 场 合 演 唱 演 唱 者 多 为 中 老 年 人 巫 师 歌 手 酒 席 演 唱 古 歌 场 合 苗 族 古 歌 古 词 神 话 民 族 心 灵 记 忆 苗 族 古 代 社 会 百 科 全 书 经 典 史 学 民 族 学 哲 学 人 类 学 多 方 面 价 值 古 歌 古 词 神 话 民 间 流 传 唱 诵 文 化 市 场 经 济 冲 击 苗 族 古 歌 濒 临 失 传 台 江 为 例 全 县 13 万 苗 族 同 胞 中 唱 完 整 部 古 歌 寥 寥 无 几 二 百 余 人 能 唱 完 整 古 歌 中 老 年 人 传 承 古 歌 老 人 年 事 已 高 抓 紧 抢 救 保 护 苗 族 古 歌 这 一 民 族 瑰 宝 最 终 世 间 消 失

文本表示

◆将词映射到向量空间中，以下选取word2vec方法，使用的工具是Gensim，语言模型采用CBOW，训练方法采用Negative Sampling，最小词频min_count设置为0，上下文最大距离window设置为5，维度size设置了100维，可根据情况改变维度，进行不同维度的对比。其余参数为默认。最后使用均值文本表示法表示文本向量，即将该文档内所含有词的对应词向量相加求平均。

K-means聚类

- ◆实现K-means算法的方法有很多，本次使用的工具是sklearn中k-means++算法。关于k值范围确定，考虑到非遗分类法主要有六类法、八类法、十类法、十三类法和十六类，其中十三类法和十六类法是更为细致的划分，但是目前基于层级划分的类目都不会过多，为了更加全面对比各分类法中类别划分的科学性，同时考虑计算指标时图表的连续性和可预测性，本次K-means聚类k值范围确定为2到17，其余参数为默认。
- ◆轮廓系数是评价聚类效果的常用指标，将簇中所有非遗样本的轮廓系数求平均值就得到非遗平均轮廓系数，取值范围为 $[-1,1]$ 。程序实现了轮廓系数的计算，并使用python画图工具实现每次聚类效果的轮廓系数可视化。如要人工评估聚类效果，可直接输出estimator.labels_观察。

K-means聚类

◆轮廓系数是评价聚类效果的常用指标，将簇中所有非遗样本的轮廓系数求平均值就得到非遗平均轮廓系数，取值范围为 $[-1,1]$ 。程序实现了轮廓系数的计算，并使用python画图工具实现每次聚类效果的轮廓系数可视化。如要人工评估聚类效果，可直接输出 `estimator.labels_` 观察。