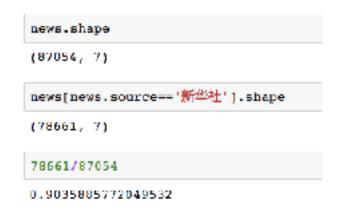
# Project-1 最终汇报文档

## 准备工作

- 读取csv数据为dataframe, 此数据集维度是(89611, 7)。
- 我们的目标是根据content字段的文本判断source字段是否新华社,所以首先我查看content是否有缺失值,content有缺失值的样本直接删除,因为没有文本无法预测而且有缺失值的样本较少删除对整个数据集影响力较小,删去缺失样本后数据集维度是(87054,7)。



- 观察数据发现新华社文章约占全部样本的90.36%,据此推测模型的准确率应该相当高,比较容易过拟合。
- 定义一个函数首先去掉文本中的特殊字符、中文标点符号,分词后过滤掉停用词最后输出以空格分开的文本,具体效果如下:

```
def split_text(text):return ''.join([w for w in list(jieba.cut(re.sub('\s|[%s]' % (punctuation),'',text))) if w not in split_text(news.iloc[1].content)

Building prefix dict from the default dictionary ...
Loading model from cache /var/fulders/vr/sknbl6y310vyq7zp6cb3w76c0000qm/T/jieba.cache
Loading model cost 0.893 seconds.

Prefix dict has been built successfully.

'機定 835 唯一 Windows10 桌面 平台 认证 ARM 处理器 高速 強调 不会 只 考虑 性能 屏蔽掉 小 核心 相反 正 联手 微软 找到 一种 适合 桌面 平台 美頭 性能 功耗 完美 方案 报道 微软 已经 拿到 一些 源局 Windows10 更好 理解 big little 架构 资料 显示 微定 835Win10 电脑 预计 均 二合一 形态 产品 当然 高速 碳龙 未来 它许 见到 三星 Exynos 取发料 华为 麒麟 小米 澎湃 进入 Windows10 桌面 平台'
```

- 最后把这个函数应用到所有文本构造出语料库

# 构造特征和模型训练

- 首先用sklearn的Countvectorizer设置min\_df=0.015(忽略文档中词频不足词汇数 1.5%)建立词频矩阵,然后再用TfidfTransformer将结果转化成tfidf矩阵shape (87054, 883)。
- 预测值为source = '新华社' 标记为 1, 其它标记为0。
- 用train\_test\_split函数吧样本分为70%的训练集和30%的测试集,把测试集代入朴素贝叶斯分类器(MultinomialNB)
- 最后用cross-validate计算accuracy, precision, recall和f1score这几个指标,与测试集在这几个指标上的表现做对比。

#### Cross-validation result

# {'fit\_time': array([0.66403699, 0.431885 , 0.42695594]), 'score\_time': array([0.15592504, 0.15189481, 0.21946597]), 'test\_precision': array([0.95943065, 0.96208067, 0.96164526]), 'train\_precision': array([0.96233212, 0.96128883, 0.96252095]), 'test\_recall': array([0.90745081, 0.91132065, 0.90608307]), 'train\_recall': array([0.91001008, 0.90821138, 0.90775059]), 'test\_accuracy': array([0.8817506 , 0.88746123, 0.88252571]), 'train\_accuracy': array([0.8865449 , 0.88405868, 0.88475361]), 'test\_fl': array([0.93271709, 0.93601299, 0.93303772]), 'train\_fl': array([0.93544004, 0.93399664, 0.9343338 ])}

#### **Testset result**

```
show_test_reslt(y_test,y_predict)
accuracy: 0.8870467511582494
precison: 0.9620043806714049
recall: 0.9110960585919309
```

f1 score: 0.9358584101582883

这些指标在cross-validation中表现稳定,而且测试集和训练集的数据差别不大,模型 拟合效果较好。

# 检测抄袭

-首先用训练好的模型预测全量样本,把预测为正类但是真实值为负类的样本index拿出来存在一个list里、这里我们获得2813个疑似样本。

#### - 聚类分析

- 首先用Normalizer将tfidf矩阵将每个属性规则化为单位向量(因为sklearn中的 kmeans默认使用欧氏距离,而单位向量组成的矩阵计算欧氏距离与余弦距离有线 性关系,这样相当于使用余弦距离度量相似度)
- 然后创建一个kmeans聚类器,在此我分类25个簇。(最佳中心点k的选择可以使用 MSE,Silhouette等指标选择elbow确定合适的值,但是因为电脑配置问题跑不出 来,暂且自己定一个)
- 最后根据聚类分析的结果构建id-class对应字典和class-id对应字典,这样就能找出疑似抄袭的文本在哪个簇,其中最相似的文本有哪些

# 聚类的同一类中相似文本效果如下:

19

新华老领州,义乌(郑江),2017年5月20日\\*\*第九届中国国际旅游商品博定会开幕\n3月20日。宋白海南的参展悉展示黎族的终工艺。\n当日,第九届中国 国际旅游商品·博克会在浙江义乌国际博克中心开幕,为第4天的职会以-旅游新品,银行天下-为主题。没标准原位2000余个,昭引境内外1300余家企业参展。集 中具尽版新装备用息、户种件利用息、旅游器所、智慧旅游产品、旅游纪念品等五大类产品。\\*\*新华社发(冀献明理)

新华赴景片,义乌(浙江),2017年5月26日、第九届中国国际旅游商品博览会开幕、65月26日,参属客商进购非洲绿起石作品。1。当日,萧九届中国国际旅游商品博览会在浙江公岛国际博览中心开幕。为明4天的复令区"旅游新品"和行天下"为主题,设标者属位2100余个,吸引境内外1300余家企业参属,集中展示旅游器各用品、户外体采用品、旅游服所、智慧旅游产品、旅游纪念品等五大体产品、1点新华社发(吕廷相)

新华社银片、义乌(浙江),2017年5月26日\b第九届中国国际旅游商品境定会开幕\n5月26日,观众在拍摄可移动强星。\n当日,第九届中国国际旅游商品 博见会在浙江父乌国时情况中心开幕。为第4天的居会以"旅游附品·银行天下"为主题。设标准居位2000余个,吸引境为介1903步家企业参报,集中展示旅游表 衡用品、产外体制用品、旅游刷片、智慧旅游广场、旅游和水品等五大类产品。\n数学社发(调献明数)

新华杜凯片,义乌(尼江),2017年5月26日NE某几届中国国际旅游商品牌定会开幕ND5月26日,参展商调证可对话选督机提入。NE当日,第九届中国国际旅游商品牌资会在浙江义岛国际博义中心开幕。为明4天的民会以"旅游新品"智行天下"为主题。这标准展位2000余个,吸引填内外1300余家企业参展。集中展示旅游装备用品、广外外闲用品、旅游期所、智慧旅游产品、旅游机场品等五大关产品、Na新年社发(清晰明摄)

- 根据以上的结果,我们只需要遍历疑似抄袭文档和它所在簇中的新闻文本,计算 edit distance就定位可以找出抄袭文档(或者大篇幅引用)和原文。

### 大致效果如下图:

#### 序便抄袭:

中国5月份56座城市新建商品住宅价格区比上涨。4月份为58座上涨,5月份15个一级和热点二级城市房地产市场基本稳定。5月份房地产调路改美处理路板 思现。

. 納計局: 15个一线和路点二线城市房价同比涨端全部回答 国家统计局城市司高级统计师刘建伟解读5月世房价数据

\_据城市另位平均活础继续回路

国家武计局与目录等了2017年3月份70个大牛城市住宅贸务价格统计数据。对此,国家统计局城市可高级统计师对建伟进行了解读。 一、15个一线和热点二级城市新越商品住宅价格同比涨偏全部回序、9个城市环比下降或持平 5月份,因地制宜,因城路装的房地产调控政策效果继续显现,15个一线和融点二线城市房地产市场基本稳定,从同比看,15个城市新建商品住宅价格涨幅

连续6个月和A个月到降。5月份比4月份分别到降0.8和0.5个百分点。 三、70个人中域中中的介环比下降及阻偏回搭域中个数均有所能加

70个城市中新建商总位宅价格区比下降的城市有9个。比上月增加1个;连编四落的城市有26个,比上月增加3个。二手住宅价格区比下降的城市有 7个、比上月增加2个:涨幅回落的城市有30个、比上月增加8个。

医家虎针周19日度布数据。5月份。15个一线和热点工线城市斯提腾品任宅价格同比张瑞全即回答。其中9个城市环化下降或寿平。这9个价格环比下降或持平的城市均、北京、上海、南京、杭州、白色、福州、郑州、汉圳、成都。 --5月份。运地制宜、园城前颇的房地产周拉政策及果继续登现。15个一线和热点工线城市房地产市场基本稳定。--医家虎针局城市司高贵虎针师刘建体说。从内比看。15个城市新建商品住宅价格还强均比上月回席,回席城里在0.5至6.4个百分点之间。从环比看,9个城市新建商品住宅价格下降城特中:5个城市压 幅在0.51以内。

国家统计局当天还发布了5月份70个大中城市住宅销售曾格统计数据。刘建伟介绍。5月份,73个大中城市中新建商品住宅和二手住宅给韩同比张幅七上月 因落的城市分别有29和18个。其中,一二线城市同比张辅回客之其明显。据则算,一线城市新建商品住宅和二手任宅价格同比张幅均连续8个月包落,5月份比4 月份分别回路2.2和1.7个百分点;二线城市新建商品住宅和二手任宅价格同比张幅分别连续6个月和4个月回路,5月份比4月份分别回路0.8和0.5个百分点。 此外,70个大平城市中房市外北下南及金融的路域市个政均有市域则。第12层示,5月67、70个大平城市中房间的战士市任务外北下海的城市有70个,比上月增加7个。比上月增加7个。比上月增加7个。比上月增加7个。

aditdistince: 305

- 但是最后发现、疑似抄袭文本中有些有的词语或者主题雷同、实际上不是抄袭或者 无关系的新闻,相似度也很高,但是找不到改动。