

2020 年全国泰迪杯 数据挖掘挑战赛

基于 Python “智慧政务” 中的文本挖掘的应用

摘要

本文基于 Python 研究“智慧政务”中的文本挖掘及应用，通过所给附件内的留言情况建立分类模型，提取热点问题，并结合留言详情制定一套关于答复意见的质量评价方案。

针对问题一：题目要求建立关于留言内容的一级标签分类模型，那么首先读取数据（附件 2），对文本数据进行预处理（数据增强、中文文本分词、停用词过滤等），再将文本向量化，通过 TF-IDF 算法¹构造文本特征向量集，再利用 LinearSVC 类²构造并训练一级标签分类模型，最后使用 sklearn.metrics³中的 F1(精准率与召回率的平衡)分数评估模型，为了使分类模型更为准确，需不断重复上述步骤，不断进行数据增强，以提高结果模型的精准率。

针对问题二：读取附件 3 中的数据，事先使用基于规则的方法（正则表达式处理），从“留言主题”中提取出命名实体（特定地点和特定人群），保存到文档，形成自定义词库。然后把文本数据转换为列表形式，通过 Gensim⁴形成语料库，使用 LsiModel 模型⁵算法，创建 TF-IDF 模型，并传入语料库来训练模型，得出 TfIdf 值，进而计算出相似度，然后对热度问题进行排序，最后通过其它一些基本方法将结果输出，并保存到相应文件中。

针对问题三：根据附件 4 中的“留言详情”和“答复意见”，从多个角度出发，对文本数据进行对比分析，制定一套关于答复意见的质量评价方案。

关键字：分类 TF-IDF SVM 模型 LsiModel 模型 文本相似度 命名实体识别

¹ TF-IDF (term frequency - inverse document frequency, 词频-逆向文件频率) 是一种用于信息检索 (information retrieval) 与文本挖掘 (text mining) 的常用加权技术。

² SVM 可以用于分类、回归、异常检测。SVM 库中包括 SVC、LinearSVC 接口

³ sklearn.metrics 分类指标

⁴ Gensim 是一款开源的第三方 Python 工具包, 用于从原始的非结构化的文本中, 无监督地学习到文本隐层的主题向量表达。它支持包括 TF-IDF, LSA, LDA, 和 word2vec 在内的多种主题模型算法, 支持流式训练, 并提供了诸如相似度计算, 信息检索等一些常用任务的 API 接口

⁵ LSI 一种简单实用的主题模型。LSI 是基于奇异值分解 (SVD) 的方法来得到文本的主题的

Text mining and application in " smart government affairs" based on python

Abstract

This article is based on Python research on the text mining and application of "smart government affairs". A classification model is established based on the message situation in the attached attachment, hot issues are analyzed, and a set of quality evaluation schemes for answering opinions is formulated based on the message details.

Aiming at the problem of the first, The title requires the establishment of a first-level label classification model on message content, then first read the data (Annex 2) , the text data pre-processing (data enhancement, Chinese text segmentation, stop word filtering, etc.) , then the text is vectorized, the text feature vector set is constructed by TF-IDF Algorithm, the first level label classification model is constructed and trained by LINEAR SVC class, and finally sklearn. In order to make the classification model more accurate, it is necessary to repeat the above steps and enhance the data to improve the accuracy of the result model.

Aiming at the problem of the second, read the data in Annex 3, use a rule-based approach (regular expression) beforehand, extract named entities (specific places and specific people) from the "message topics", save them to the document, and form a custom thesaurus. Then the text data is transformed into a list form, the Corpus is formed by Gensim, the TF-IDF model is created by using LsiModel algorithm, and the TF-IDF model is imported into the Corpus to train the model, the TFIDF value is obtained, then the similarity is calculated, and then the heat problem is sorted, finally a some other basic methods will be output, and saved to the corresponding file.

Aiming at the problem of the third, According to the "message details" and "response comments" in Annex 4, the text data comparative analysis, the development of a set of quality evaluation plans for reply opinions was formulated.

Keywords: Classification TF-IDF SVM model LsiModel model Similarity
Named Entity Recongition

目 录

1 问题重述.....	4
1.1 问题背景.....	4
1.2 问题描述.....	4
2 问题分析.....	4
2.1 问题一的分析.....	5
2.1.1 分类模型的定义.....	5
2.1.2 问题分析.....	5
2.2 问题二的分析.....	5
2.2.1 热点问题的定义.....	5
2.2.2 问题分析.....	5
2.2.3 文本相似性热度统计.....	5
2.3 问题三的分析.....	6
3 符号说明.....	6
4 模型建立与求解.....	6
4.1 问题一：群众留言分类.....	6
4.1.1 模型思想.....	6
4.1.2 模型建立.....	8
4.2 问题二：热点问题挖掘.....	11
4.2.1 流程图.....	11
4.2.2 算法描述.....	11
4.2.3 解题方法及过程.....	12
4.2.4 方法所存在的不足.....	15
4.3 问题三：制定评价方案.....	15
5 结果分析.....	16
5.1 问题一结果分析.....	16
5.2 问题二结果分析.....	16
参考文献.....	18

1 问题重述

1.1 问题背景

近年来，随着微信、微博、市长信箱、阳光热线等网络问政平台逐步成为政府了解民意、汇聚民智、凝聚名气的重要渠道，但伴随着各类社情民意相关的文本数据量不断攀升，收集到的社情民意繁多且杂乱，还可能存在恶意投递空白问题，也可能有人心急，多次投递重复的民意，不仅增加了工作人员的工作负担，还可能因为问题解决、反馈的不及时，给民众带来一定的困扰，这些都给以往主要依靠人工来进行留言划分和热点整理的相关部门的工作带来极大的挑战。

为解决这一问题，我们可事先通过计算机自动处理文本，排除空的重复的民意。同时，在如今大数据、云计算、人工智能等技术的发展下，建立基于自然语言处理技术的智慧政务系统已经是社会治理创新发展的新趋势。为此，我们可以利用自然语言处理和文本挖掘的方法，对群众留言进行分类，阶段性的提取热点问题。这不仅能够减轻工作人员的工作负担，而且能够及时解决民众的问题，并作出反馈，有效提升政府的管理水平和施政效率，更好的服务大众。

因此，“智慧政务”的文本挖掘应用对收集社情民意系统有重大意义。

1.2 问题描述

1) 为解决目前好多电子政务系统依靠人工根据经验处理，存在工作量大、效率低，且差错率高等问题，参照附件 1，使用附件 2 的留言详情和一级标签这两列进行分析，建立一级标签分类模型。该模型的准确率较高。

2) 为了及时发现热点问题，以帮助相关部门进行针对性地处理，提升服务效率。根据附件 3 将某一时段内反映特定地点或特定人群问题的留言进行归类，定义合理的热度评价指标，并给出评价结果，按表 1 的格式给出排名前 5 的热点问题，给出相应热点问题对应的留言信息。

3) 根据附件 4 给出的数据，从答复的相关性、完整性、可解释性等角度建立一套答复意见的质量的评价方案并实现。

2 问题分析

2.1 问题一的分析

2.1.1 分类模型的定义

根据留言详情所反映的文本数据特征，按照一定的划分体系进行自动分类。便于后续将留言分派给相应的部门处理。

2.1.2 问题分析

（一）主要任务

根据附件 2 所给数据，训练出关于留言内容的一级标签分类模型。

（二）需解决的问题

数据预处理：解决样本不均衡问题、清洗数据降低计算复杂度和系统开销；

获得分类文本的特征集向量：筛选出最具文本特征的词语并向量化；

建立训练分类器，评估分类模型。

2.2 问题二的分析

2.2.1 热点问题的定义

某一时段内群众集中反映的某一问题可称为热点问题。关键点在于一段时间集中爆发的问题，即多人反映同一问题。

2.2.2 问题分析

（一）主要任务

根据附件 3 将某一时段内反映特定地点或特定人群问题的留言进行归类，定义合理的热度评价指标，并给出评价结果，按格式给出排名前 5 的热点问题及其具体留言信息。

（二）需解决的问题

问题识别：如何从众多留言中识别出相似的留言；

问题归类：把特定地点或人群的数据归并，即把相似的留言归为同一类问题，结果对应表 2；

热度评价：热度评价指标的定义和计算方法，对指标排名之后得出对应表 1。

（三）可能存在的难点

难点一：地点、人群的识别（表达多样化）

难点二：相似的计算复杂（特征多、两两之间计算相似计算量大）

2.2.3 文本相似性热度统计

根据需求可从不同维度进行统计：

1.分组不分句热度统计

根据某列首先进行分组，然后再对描述类列进行相似性统计。

2.分组分句热度统计

根据某列首先进行分组，然后对描述类列按照标点符号进行拆分，然后再对这些句进行热度统计。

3.整句及分句热度统计

对描述类列/按标点符号进行分句，进行热度统计。

4.热词统计

对描述类进行热词统计。

由于给出的文本数据相对较多，所以该题采用整句热度统计这层次对留言详情按标点符号进行分句，进行热度统计。

2.3 问题三的分析

为了能够规范对留言做出的答复意见，为了群众的问题得到解决，我们制定了答复意见评价方案。

3 符号说明

表 1 符号说明

符号	符号含义
TF	词频
stopword	停用词
IDF	逆文档频率
TP (True Positives)	正类判定为正类
TN (True Negatives)	负类判定为负类
FP (False Positives)	负类判定为正类，即“存伪”
FN (False Negatives)	正类判定为负类，即“去真”

4 模型建立与求解

4.1 问题一：群众留言分类

4.1.1 模型思想

将样本数据中的留言详情转为包含数字的词频向量，再转为 TF-IDF 向量，训练基于支持向量机的机器学习模型分类器。

(一) TF-IDF (Term Frequency-Inverse Document Frequency, 词频-逆文件频

率)算法:

TF-IDF (term frequency - inverse document frequency) 是一种用于信息检索与数据挖掘的常用加权技术。TF 是词频(Term Frequency), IDF 是逆文本频率指数(Inverse Document Frequency)。TF-IDF 是在单词计数的基础上,降低了常用高频词的权重,增加罕见词的权重。一个词语在一篇文章中出现次数越多,同时所有文档中出现次数越少(罕见),越能够代表该文章。

1) 计算词频

词频 (term frequency, TF)指的是某一个给定的词语在该文件中出现的次数。为了防止偏向长的文件。将其归一化(词频除以文章总词数), 公式如下:

$$TF_w = \frac{\text{在某一类中词条 } w \text{ 出现的次数}}{\text{该类中所有的词条数目}}$$

2) 计算逆向文件频率

逆向文件频率 (inverse document frequency, IDF)指的是如果包含词条 t 的文档越少, IDF 越大, 则说明词条具有很好的类别区分能力。由总文件数目除以包含某一特定词语之文件的数目, 再将得到的商取对数就可以得到这一词语的 IDF。公式如下:

$$IDF = \log \left(\frac{\text{语料库的文档总数}}{\text{包含词条 } w \text{ 的文档数} + 1} \right)$$

注: 分母之所以要加 1, 是为了避免分母为 0 的情况。

3) 计算词频-逆文件频率

词频-逆文件频率 (Term Frequency-Inverse Document Frequency, TF-IDF)指的是用某一特定文件内的高词语频率, 以及该词语在整个文件集合中的低文件频率, 算出高权重的 TF-IDF。公式如下:

$$TF - IDF = \text{词频 (TF)} \times \text{逆文档频率 (IDF)}$$

(二) 支持向量机 (Support Vector Machine, SVM) 一对一方法

SVM 是基于统计学习理论, 以结构风险最小化原则为理论基础, 通过选择适当函数子集及该子集中的判别函数学习机的实际风险降至最小, 保证了通过有限训练样本得到的小误差分类器对独立测试集的测试误差仍然小, 得到一个具有最优分类能力和推广泛化能力的学习机。

一对一方法（one-against-one）^[1]在 k-类训练样本中构造所有可能的二类分类器，每类仅仅在 k-类中的二类训练样本上训练，结果共构造 $N=k(k-1)/2$ 个分类器。通过投票法组合这些二类分类器，得票最多的类就是新点所属的类。

4.1.2 模型建立

（一）流程图

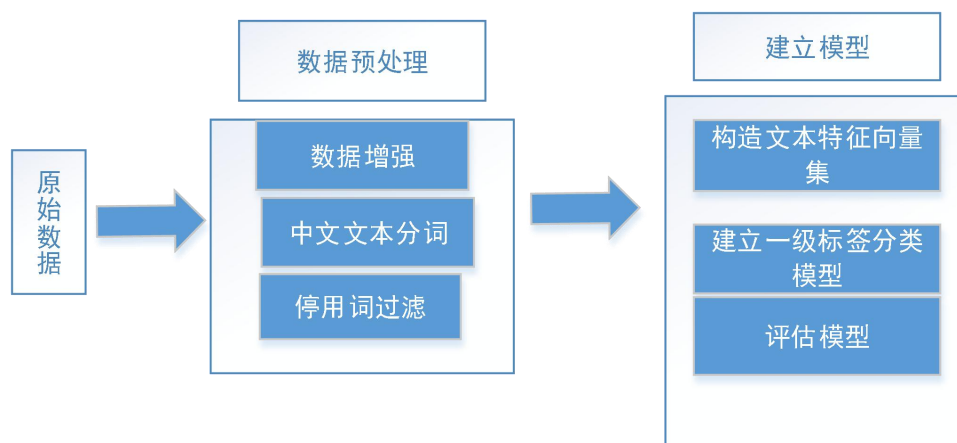


图 1 题目一的流程图

（二）数据预处理

本文把对文本数据的预处理分为四部分：

1) 分析样本数据

读取附件 2 数据进行分析，得出数据总量有 9210，无空行空列。查看各一级标签的数据量发现数量不一致，分布不均匀。最多的一级标签数量有 2009，而最少才 613。

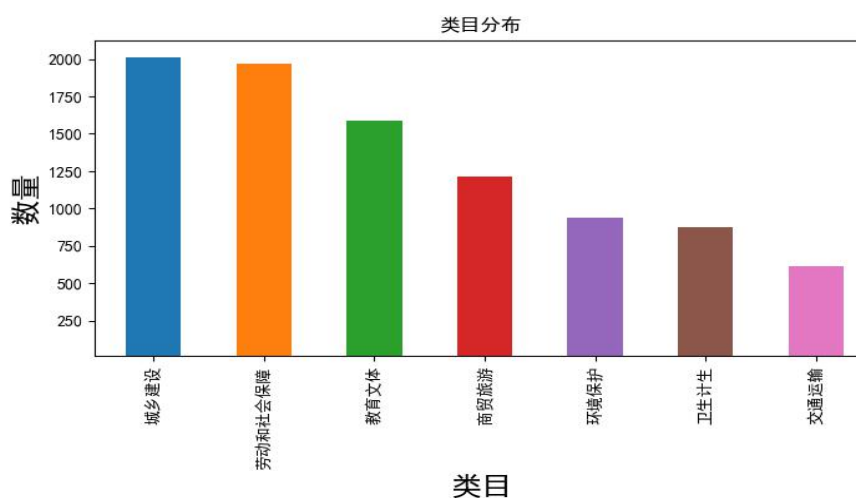


图 2 样本数据各一级标签的分布情况

2) 对样本数据进行数据增强

针对 1) 发现的样本分布不均匀问题,根据一级标签对文本进行分类后,将不同分类的文本通过 EDA 算法进行同义词替换-新增-交换-删除-生成同义句^[2] 进行不同程度的数据增强。这里采用的是 Synonyms 同义词库,生成相对应标签文件,如环境保护.xlsx、交通运输.xlsx 等。增强后的数据总量有 3 万多条,大部分数据量在 4000 左右,分布较为均匀。将生成相对应标签文件合并存储在 output.xls 中。

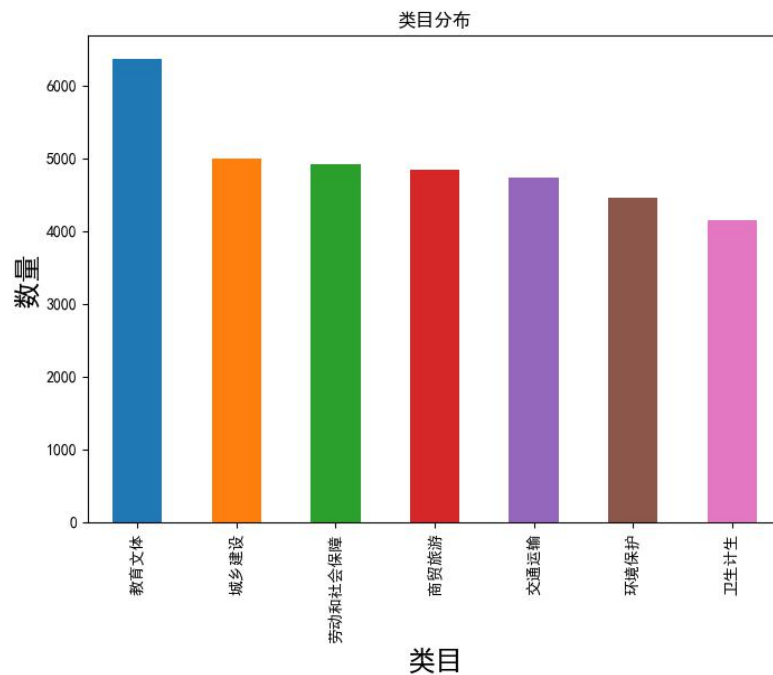


图 3 增强后各一级标签数据的分布情况

3) 对留言详情进行数据清洗

对 output.xls 中数据,运用正则表达式和过滤停用词 (stop word),清洗掉对系统分析预测文本的内容没有任何帮助而会增加计算的复杂度和增加系统开销的符号和词语。

4) 中文分词

将清洗后的数据,运用 python 的中文分词包 jieba 进行分词。部分数据如图:

一级 标签	留言详情	cat_id	clean_review	cut_review
0 交通 运输	座落在月亮岛街道的时代倾城居民小区一二三期之...	0	座落在月亮岛街道的时代倾城居民小区一二三期之间的公共路段两边的人行道常年停满了社会制 度车严重...	座落在 月亮 岛 街道 时代 倾城 居民小 区 二三期 之间 公共 路段 两边 人行道 常年 ...
1 交通 运输	地处月亮岛街道时代倾城小区一二三期之间的公 共...	0	地处月亮岛街道时代倾城小区一二三期之间的公 共道路两边的人行道常年停满了社会车辆严重影 响行人通...	地处 月亮 岛 街道 时代 倾城 小区 二 三期 之间 公共 道路 两边 人行道 常年 停满 ...
2 交通 运输	规划严重影响不知月亮岛街道时代倾城小区一之 间的的道路下的人行道常年停满了社会车辆行人 通行社区...	0	规划严重影响不知月亮岛街道时代倾城小区一之 间的的道路下的人行道常年停满了社会车辆行人 通行社区...	规划 影响 不知 月亮 岛 街道 时代 倾 城 小区 之间 道路 人行道 常年 停满 社会 车 ...
3 交通 运输	地处月亮岛时代倾城小区二三期之间的公共两边人行道 常年了社会车辆严...	0	地处月亮岛时代倾城小区二三期之间的公共两边 人行道常年了社会车辆严重影响行人通行社区和 交警也一...	地处 月亮 岛 时代 倾城 小区 二三期 之间 公共 两边 人行道 常年 社会 车辆 影响 ...
4 交通 运输	地处月亮岛街道时期倾城小区一二三期之间的公 共...	0	地处月亮岛街道时期倾城小区一二三期之间的公 共路段两旁的人行道常年停满了社会车受到影响 行人通行...	地处 月亮 岛 街道 时期 倾城 小区 二 三期 之间 公共 路段 两旁 人行道 常年 停满 ...

图 4 数据预处理后部分结果

（三）基于 TF-IDF 算法的文本特征提取并向量化

将预处理好的数据通过 `sklearn.feature_extraction.text.TfidfVectorizer` 方法提取文本的 TF-IDF 的特征值并将其向量化。这里我们使用了参数 `ngram_range=(1,2)`,这表示我们除了抽取评论中的每个词语外,还要抽取每个词相邻的词并组成一个“词语对”,如:词 1, 词 2, 词 3, 词 4, (词 1, 词 2), (词 2, 词 3), (词 3, 词 4)。这样就扩展了我们特征集的数量,有了丰富的特征集才有可能提高我们分类文本的准确度。

（四）构造并训练基于 SVM 一对一算法的分类器

得到生成各个标签的 TF-IDF 向量后,这里用的是 `scikit-learn` 支持向量机算法库中的 `LinearSVC` 类实现构造并训练分类器^[3]。

（五）对模型进行评价

调用 `sklearn.metrics` 中的 `classification_report` 函数和 `accuracy_score` 方法^[3],显示每个类的精确度、召回率、F1 值、分类准确率分数等信息以此来评估分类模型。

分类准确率分数: 所有分类正确的百分比。

精确率: (Precision) 真正正确的占有预测为正的比例。公式如下:

$$\text{Precision} = \frac{TP}{TP + FP}$$

召回率 (Recall): 真正正确的占有实际为正的比例。公式如下:

$$\text{Recall} = \frac{TP}{TP + FN}$$

F1 值: 精确度和召回率的调和平均值。公式如下:

$$\frac{2}{F_1} = \frac{1}{Precision} + \frac{1}{Recall}$$

4.2 问题二：热点问题挖掘

4.2.1 流程图

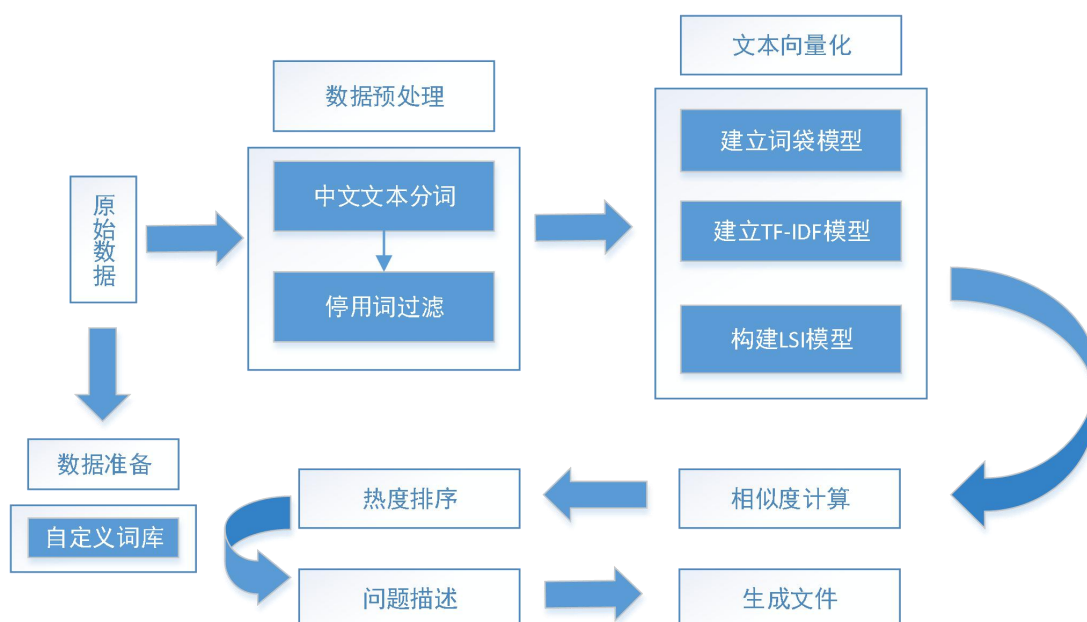


图 5 问题二流程图

4.2.2 算法描述

（一）命名实体识别

（1）命名实体识别（Named Entity Recognition，简称 NER），又称作“专名识别”，是指识别文本中具有特定意义的实体，主要包括人名、地名、机构名、专有名词等。

（2）命名实体一般包括三大类（实体类、时间类、数字类），七小类（人名、机构名、地名、时间、日期、货币和百分比）。

（3）命名实体识别过程：确定实体的边界，即确定哪些词属于实体；确定实体的类别，即确定实体属于人名或者机构名等。

（4）主要方法

● 基于规则的方法

主要通过手工构造规则模板，选用特征，如关键字、标点符号、位置词、指

向词等，基于不同的规则权值进行判断。基于规则的方法性能上可解释性强，对于 badcase 的解决能力高，但构建规则库需要语言学专家且耗费时间长。

● 基于统计的方法

基于统计机器学习的方法主要包括：隐马尔科夫模型、最大熵模型、条件随机场等。实际上是将命名实体识别转化为一个序列标注任务，这部分工具与分词及词性标注有一定的重合之处。

(二) Python+Gensim-文本相似度分析

使用 Gensim 计算文本相似度一般是先用 corpora 模块把文档转为简单的稀疏矩阵；然后用 models 模块得到符合需要的向量模型；最后用 similarities 模块计算相似度。具体流程如图 6 所示。

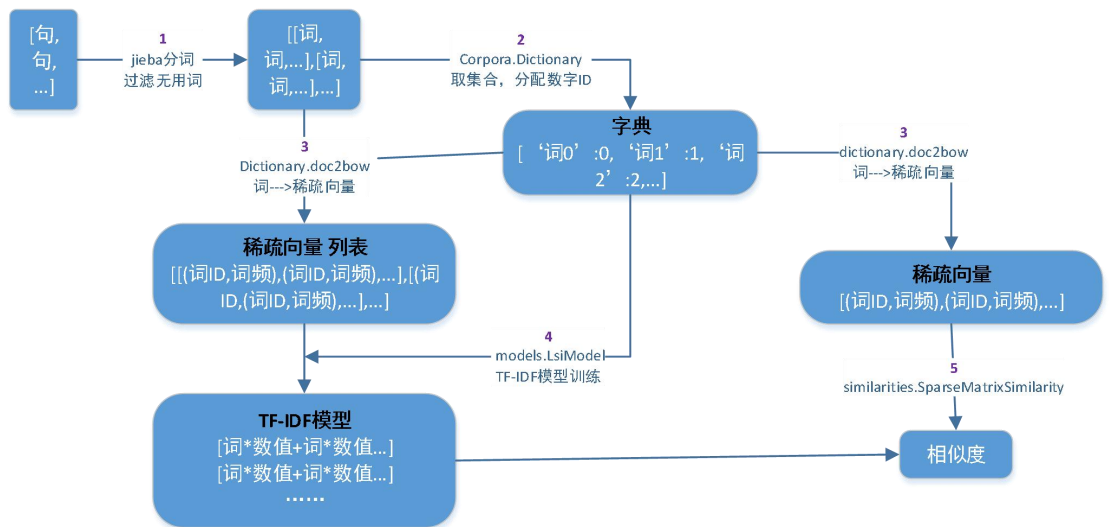


图 6 Gensim 使用流程图

4.2.3 解题方法及过程

(一) 数据描述

通过观察所给数据，可以发现全部数据量共计 4326 条，且附件 3.xlsx 中的字段大多为文本格式，需要将其量化成数值形式才能对其进行分析。该题主要通过通过对数据进行相似度比较后再进行热度排序，由于文本数据存在大量复杂地名，所以可单独通过“留言主题”列将特定地点、特定人群等命名实体提取出来，形成自定义词库，以便提高中文分词的准确性，进而提高文本之间的相似性。同时，因已知文本数据包含留言主题及留言详情，两者可能存在一定的误差，为保证最终得到较为精准的热点问题，可分别通过对留言主题和留言详情进行相似度比较、热度排序后取二者较为精确的热点问题挖掘结果。（由于两者解决算法基本

一致，所以下面过程分析以“留言详情”为例）

（二）数据准备

根据已知数据（附件 3），通过命名实体识别基于规则的方法提取相应实体（主要以地名为主），并保存为文件，以便后续使用。具体过程如图 7 所示。



图 7 命名实体识别（地点识别）步骤

（三）数据预处理（中文分词）

在对热点问题挖掘分析之前，先要把非结构化的文本信息转换为计算机能够识别的结构化信息。在附件 3 中，以中文文本的方式给出了数据。为了便于后续转换，先要对附件 3 中“留言详情”列进行中文分词。这里采用 Python 开发的一个中文分词模块——jieba 分词^[4]的方法。jieba 分词用到的算法：

- 基于 Trie 树结构^[5]实现高效的词图扫描，生成句子中汉字所有可能成词情况所构成的有向无环图（DAG）；
- 采用了动态规划查找最大概率路径，找出基于词频的最大切分组合；
- 对于未登录词，采用了基于汉字成词能力的 HMM 模型，使用了 Viterbi 算法

同时，为节省存储空间和提高搜索效率，在处理文本之前将停用词进行过滤，且为提高分词的准确性，并导入自定义词典，将相应的命名实体合在一起。部分分词结果示例如图 8 所示。

0 [座落在, A市, A3区, 联丰路, 米兰, 春天, G2, 栋, 320, 一家, 名叫...
 1 [A市, A6区, 道路, 命名, 规划, 初步, 成果, 公示, 文件, 转化, 正式, ...
 2 [系, 春华镇, 金鼎村, 七里, 组, 村民, 不知, 相关, 水泥路, 到户, 政策, ...
 3 [靠近, 黄兴路, 步行街, 城南路, 街道, 古道, 巷, 一步, 两, 搭桥, 小区, ...
 4 [A市, A3区, 中海国际社区, 三期, 四期, 蓝天, 璞, 洲, 幼儿园, 旁边, 块...
 5 [麓泉社区, 麓谷明珠小区, 栋, 居民, 近期, 感觉, 震惊, 伤心, 购房, 签合同,...
 6 [“, 二高, 一部, 发出, 非法, 集资, 打击, 通知, 中是, 金融, 犯罪, 通知...
 7 [一名, A市, 地铁站, 上班, 安检员, 中介, 公司, 介绍, 上班, 安检员, 岗位...
 8 [12, 21, 下午, 17, 52, 分许, 6路公交车, 司机, 座位, 旁边, 汽车...
 9 [保利, 麓, 谷林语, 桐梓坡路, 麓松路, 交汇处, 地铁, 凌晨, 点, 施工, 噪音...
 10 [下午, 晚, 高峰, 五点, 半左右, 特立路, 东四路, 口时东, 往西方, 向车, 越...
 11 [宁静, 我要, 复习, 迎考, 大半, 年底, 商, 空调, 冰柜, 外机, 轰响, 扰民...
 12 [桐梓, 坡, 589, 号, 白鹤, 咀, 停车场, 由聚, 美龙楚, 新能源, 公司, ...
 13 [想, 举报, A市, 利保, 壹号, 公馆, 项目, 夜间, 噪声, 扰民, A市, 利保...
 14 [A7县, 星沙街道, 两个, 孩子, 上有老下有, 居民, 居民, 直斥, A市, 地铁...
 15 [北辰, D1, 区, 一名, 业主, 小区, 住, 改商, 数量, 达上, 百家, 安全隐...
 16 [K3县, 乡村, 卫生室, 处于, 无证, 行医, 状态, 情况, 原因, 卫生室, 医疗...
 17 [春华镇, 一名, 村民, 接到, 政府, 通知, 029, 县, 道旁, 开, 麻将馆, ...
 18 [书记, 外地人, 2018, 年, 毕业, 落户, A市, 申请, 短期, 出国, 旅游...
 19 [退费, 之日起, 之内, 退还, 学习, 费用, 今日, 超过, 合同期, 未, 退费, ...
 20 [A市, A6区, 月亮岛街道, 乾源, 国际, 广场, 停车场, 违章, 乱建, 现象, ...
 21 [A7县, 时代星城小区, 幢, 1321, 室, 非法, 一家, 家庭旅馆, 住宅小区, ...
 22 [敬爱, 佳兆业, 水新, 一期, 小区业主, 小区, 43, 栋, 101, 户, 违建...
 23 [沙坪, 街上, 一家, 无证, 经营, 理疗, 馆, 打着, 免费, 做, 健康, 理疗...
 24 [a, 我姐, 杜四娇, 贵区, 管辖, 内德鸡, 餐饮, 打工, 五年, 2019, 年,...
 25 [长房云时代小区, 月底, 开盘, 月初, 摇号, 选房, 10, 号, 业主, 观察, 1...
 26 [谢谢您, 这份, 非常感谢, 松, 雅西, 地省, 几个, 小区, 传销, 记得, 特别...
 27 [2019, 年, 12, 18, 上午, 九点, A市, 01, 年, 11, 所下区, ...
 28 [黄兴路, 步行街, 仕弘, 教育, A市, 晚报, 大厦, 八楼, 前程, 优学, 确保...
 29 [多年, A2区, 迎新路, 区政府, 东门, 万芙路, 段, 改装车, 飙车, 高考, 期...

图 8 “留言详情”部分分词结果

(四) TF-IDF 文本相似度分析

- (1) 使用 gensim 中的 corpora 模块, 将分词形成后的二维数组生成词典;
- (2) 将二维数组通过 doc2bow 稀疏向量, 形成语料库 (如图 9);

语料库:
 [[(0, 1)], [(1, 1), (2, 1), (3, 1), (4, 1), (5, 2), (6, 1), (7, 1), (8, 1), (9, 1), (10, 1), (11, 1), (12, 1), (13, 1), (14, 1), (15, 1), (16, 1), (17, 1), (18, 1), (19, 1), (20, 1), (21, 2), (22, 2), (23, 1), (24, 1), (25, 2), (26, 1), (27, 1), (28, 1), (29, 1), (30, 1), (31, 1), (32, 1), (33, 1), (34, 1), (35, 1), (36, 1), (37, 1), (38, 1), (39, 1), (40, 1), (41, 1), (42, 1), (43, 1), (44, 2), (45, 2), (46, 1), (47, 2), (48, 1), (49, 1), (50, 1), (51, 1), (52, 1), (53, 1), (54, 1), (55, 1)], [(56, 1), (57, 1), (58, 1), (59, 1), (60, 1), (61, 1), (62, 1), (63, 1), (64, 1), (65, 1), (66, 1), (67, 1), (68, 1), (69, 1), (70, 1), (71, 1), (72, 1), (73, 1), (74, 3), (75, 1), (76, 1), (77, 1), (78, 1), (79, 1), (80, 1), (81, 1)], [(58, 1), (67, 1), (68, 1), (69, 1), (71, 1), (74, 1), (82, 1), (83, 1), (84, 1), (85, 1), (86, 1), (87, 1), (88, 1), (89, 1), (90, 1), (91, 1)], [(1, 1), (2, 1), (5, 2), (13, 2), (14, 1), (19, 1), (22, 1), (28, 4), (38, 2), (42, 2), (43, 2), (44, 1), (45, 1), (47, 1), (48, 3), (52, 1), (72, 1), (80, 1), (92, 1), (93, 1), (94, 2), (95, 1), (96, 1), (97, 1), (98, 1), (99, 2), (100, 1), (101, 2), (102, 1), (103, 1), (104, 1), (105, 1), (106, 1), (107, 2), (108, 1), (109, 1), (110, 1), (111, 1), (112, 2), (113, 1), (114, 3)], [(1, 1), (2, 1), (5, 1), (13, 1), (16, 1), (19, 1), (22, 1), (26, 2), (30, 3), (38, 1), (39, 1), (41, 3), (42, 1), (43, 1), (45, 2), (47, 1), (48, 3), (52, 1), (54, 1), (56, 2), (94, 2), (102, 2), (110, 2), (113, 1), (114, 1), (115, 1), (116, 1), (117, 2), (118, 1), (119, 1), (120, 1), (121, 1), (122, 1), (123, 1), (124, 1), (125, 1), (126, 1), (127, 1), (128, 1), (129, 1), (130, 3), (131, 1), (132, 1), (133, 1), (134, 1), (135, 1), (136, 1), (137, 1), (138, 1), (139, 1), (140, 1), (141, 1), (142, 1), (143, 1), (144, 1), (145, 1), (146, 1), (147, 1), (148, 1), (149, 1), (150, 1), (151, 1), (152, 1), (153, 1), (154, 2), (155, 1), (156, 1), (157, 1), (158, 1), (159, 1), (160, 1), (161, 1), (162, 1), (163, 1), (164, 1), (165, 5), (166, 1), (167, 1)], [(28, 1), (38, 1), (42, 1), (43, 1), (48, 1), (52, 1), (94, 2), (98, 2), (110, 3), (117, 1), (119, 1), (136, 2), (154, 2), (168, 1), (169, 1), (170, 1), (171, 1), (172, 2), (173, 1), (174, 1), (175, 1), (176, 1), (177, 1), (178, 1), (179, 1), (180, 1), (181, 1), (182, 1), (183, 1), (184, 1), (185, 1), (186, 1)], [(30, 1), (59, 1), (63, 1), (80, 1), (148, 1), (187, 1), (188, 1), (189, 1), (190, 1), (191, 1), (192, 1), (193, 1), (194, 1), (195, 1), (196, 1), (197, 1), (198, 1), (199, 1), (200, 1), (201, 1), (202, 1), (203, 1), (204, 1), (205, 1), (206, 1), (207, 1), (208, 1), (209, 1), (210, 1), (211, 1), (212, 1), (213, 1), (214, 1), (215, 1), (216, 1), (217, 1), (218, 1), (219, 1), (220, 1), (221, 1), (222, 1), (223, 1), (224, 1), (225, 1), (226, 1), (227, 1), (228, 1), (229, 1), (230, 1), (231, 1), (232, 1), (233, 1), (234, 1), (235, 1), (236, 1), (237, 1), (238, 1), (239, 1), (240, 1), (241, 1), (242, 1), (243, 1), (244, 1), (245, 1), (246, 1), (247, 1), (248, 1), (249, 1), (250, 1), (251, 1), (252, 1), (253, 1), (254, 1), (255, 1), (256, 1), (257, 1), (258, 1), (259, 1), (260, 1), (261, 1), (262, 1), (263, 1), (264, 1), (265, 1), (266, 1), (267, 1), (268, 1), (269, 1), (270, 1), (271, 1), (272, 1), (273, 1), (274, 1), (275, 1), (276, 1), (277, 1), (278, 1), (279, 1), (280, 1), (281, 1), (282, 1), (283, 1), (284, 1), (285, 1), (286, 1), (287, 1), (288, 1), (289, 1), (290, 1), (291, 1), (292, 1), (293, 1), (294, 1), (295, 1), (296, 1), (297, 1), (298, 1), (299, 1), (300, 1), (301, 1), (302, 1), (303, 1), (304, 1), (305, 1), (306, 1), (307, 1), (308, 1), (309, 1), (310, 1), (311, 1), (312, 1), (313, 1), (314, 1), (315, 1), (316, 1), (317, 1), (318, 1), (319, 1), (320, 1), (321, 1), (322, 1), (323, 1), (324, 1), (325, 1), (326, 1), (327, 1), (328, 1), (329, 1), (330, 1), (331, 1), (332, 1), (333, 1), (334, 1), (335, 1), (336, 1), (337, 1), (338, 1), (339, 1), (340, 1), (341, 1), (342, 1), (343, 1), (344, 1), (345, 1), (346, 1), (347, 1), (348, 1), (349, 1), (350, 1), (351, 1), (352, 1), (353, 1), (354, 1), (355, 1), (356, 1), (357, 1), (358, 1), (359, 1), (360, 1), (361, 1), (362, 1), (363, 1), (364, 1), (365, 1), (366, 1), (367, 1), (368, 1), (369, 1), (370, 1), (371, 1), (372, 1), (373, 1), (374, 1), (375, 1), (376, 1), (377, 1), (378, 1), (379, 1), (380, 1), (381, 1), (382, 1), (383, 1), (384, 1), (385, 1), (386, 1), (387, 1), (388, 1), (389, 1), (390, 1), (391, 1), (392, 1), (393, 1), (394, 1), (395, 1), (396, 1), (397, 1), (398, 1), (399, 1), (400, 1), (401, 1), (402, 1), (403, 1), (404, 1), (405, 1), (406, 1), (407, 1), (408, 1), (409, 1), (410, 1), (411, 1), (412, 1), (413, 1), (414, 1), (415, 1), (416, 1), (417, 1), (418, 1), (419, 1), (420, 1), (421, 1), (422, 1), (423, 1), (424, 1), (425, 1), (426, 1), (427, 1), (428, 1), (429, 1), (430, 1), (431, 1), (432, 1), (433, 1), (434, 1), (435, 1), (436, 1), (437, 1), (438, 1), (439, 1), (440, 1), (441, 1), (442, 1), (443, 1), (444, 1), (445, 1), (446, 1), (447, 1), (448, 1), (449, 1), (450, 1), (451, 1), (452, 1), (453, 1), (454, 1), (455, 1), (456, 1), (457, 1), (458, 1), (459, 1), (460, 1), (461, 1), (462, 1), (463, 1), (464, 1), (465, 1), (466, 1), (467, 1), (468, 1), (469, 1), (470, 1), (471, 1), (472, 1), (473, 1), (474, 1), (475, 1), (476, 1), (477, 1), (478, 1), (479, 1), (480, 1), (481, 1), (482, 1), (483, 1), (484, 1), (485, 1), (486, 1), (487, 1), (488, 1), (489, 1), (490, 1), (491, 1), (492, 1), (493, 1), (494, 1), (495, 1), (496, 1), (497, 1), (498, 1), (499, 1), (500, 1), (501, 1), (502, 1), (503, 1), (504, 1), (505, 1), (506, 1), (507, 1), (508, 1), (509, 1), (510, 1), (511, 1), (512, 1), (513, 1), (514, 1), (515, 1), (516, 1), (517, 1), (518, 1), (519, 1), (520, 1), (521, 1), (522, 1), (523, 1), (524, 1), (525, 1), (526, 1), (527, 1), (528, 1), (529, 1), (530, 1), (531, 1), (532, 1), (533, 1), (534, 1), (535, 1), (536, 1), (537, 1), (538, 1), (539, 1), (540, 1), (541, 1), (542, 1), (543, 1), (544, 1), (545, 1), (546, 1), (547, 1), (548, 1), (549, 1), (550, 1), (551, 1), (552, 1), (553, 1), (554, 1), (555, 1), (556, 1), (557, 1), (558, 1), (559, 1), (560, 1), (561, 1), (562, 1), (563, 1), (564, 1), (565, 1), (566, 1), (567, 1), (568, 1), (569, 1), (570, 1), (571, 1), (572, 1), (573, 1), (574, 1), (575, 1), (576, 1), (577, 1), (578, 1), (579, 1), (580, 1), (581, 1), (582, 1), (583, 1), (584, 1), (585, 1), (586, 1), (587, 1), (588, 1), (589, 1), (590, 1), (591, 1), (592, 1), (593, 1), (594, 1), (595, 1), (596, 1), (597, 1), (598, 1), (599, 1), (600, 1), (601, 1), (602, 1), (603, 1), (604, 1), (605, 1), (606, 1), (607, 1), (608, 1), (609, 1), (610, 1), (611, 1), (612, 1), (613, 1), (614, 1), (615, 1), (616, 1), (617, 1), (618, 1), (619, 1), (620, 1), (621, 1), (622, 1), (623, 1), (624, 1), (625, 1), (626, 1), (627, 1), (628, 1), (629, 1), (630, 1), (631, 1), (632, 1), (633, 1), (634, 1), (635, 1), (636, 1), (637, 1), (638, 1), (639, 1), (640, 1), (641, 1), (642, 1), (643, 1), (644, 1), (645, 1), (646, 1), (647, 1), (648, 1), (649, 1), (650, 1), (651, 1), (652, 1), (653, 1), (654, 1), (655, 1), (656, 1), (657, 1), (658, 1), (659, 1), (660, 1), (661, 1), (662, 1), (663, 1), (664, 1), (665, 1), (666, 1), (667, 1), (668, 1), (669, 1), (670, 1), (671, 1), (672, 1), (673, 1), (674, 1), (675, 1), (676, 1), (677, 1), (678, 1), (679, 1), (680, 1), (681, 1), (682, 1), (683, 1), (684, 1), (685, 1), (686, 1), (687, 1), (688, 1), (689, 1), (690, 1), (691, 1), (692, 1), (693, 1), (694, 1), (695, 1), (696, 1), (697, 1), (698, 1), (699, 1), (700, 1), (701, 1), (702, 1), (703, 1), (704, 1), (705, 1), (706, 1), (707, 1), (708, 1), (709, 1), (710, 1), (711, 1), (712, 1), (713, 1), (714, 1), (715, 1), (716, 1), (717, 1), (718, 1), (719, 1), (720, 1), (721, 1), (722, 1), (723, 1), (724, 1), (725, 1), (726, 1), (727, 1), (728, 1), (729, 1), (730, 1), (731, 1), (732, 1), (733, 1), (734, 1), (735, 1), (736, 1), (737, 1), (738, 1), (739, 1), (740, 1), (741, 1), (742, 1), (743, 1), (744, 1), (745, 1), (746, 1), (747, 1), (748, 1), (749, 1), (750, 1), (751, 1), (752, 1), (753, 1), (754, 1), (755, 1), (756, 1), (757, 1), (758, 1), (759, 1), (760, 1), (761, 1), (762, 1), (763, 1), (764, 1), (765, 1), (766, 1), (767, 1), (768, 1), (769, 1), (770, 1), (771, 1), (772, 1), (773, 1), (774, 1), (775, 1), (776, 1), (777, 1), (778, 1), (779, 1), (780, 1), (781, 1), (782, 1), (783, 1), (784, 1), (785, 1), (786, 1), (787, 1), (788, 1), (789, 1), (790, 1), (791, 1), (792, 1), (793, 1), (794, 1), (795, 1), (796, 1), (797, 1), (798, 1), (799, 1), (800, 1), (801, 1), (802, 1), (803, 1), (804, 1), (805, 1), (806, 1), (807, 1), (808, 1), (809, 1), (810, 1), (811, 1), (812, 1), (813, 1), (814, 1), (815, 1), (816, 1), (817, 1), (818, 1), (819, 1), (820, 1), (821, 1), (822, 1), (823, 1), (824, 1), (825, 1), (826, 1), (827, 1), (828, 1), (829, 1), (830, 1), (831, 1), (832, 1), (833, 1), (834, 1), (835, 1), (836, 1), (837, 1), (838, 1), (839, 1), (840, 1), (841, 1), (842, 1), (843, 1), (844, 1), (845, 1), (846, 1), (847, 1), (848, 1), (849, 1), (850, 1), (851, 1), (852, 1), (853, 1), (854, 1), (855, 1), (856, 1), (857, 1), (858, 1), (859, 1), (860, 1), (861, 1), (862, 1), (863, 1), (864, 1), (865, 1), (866, 1), (867, 1), (868, 1), (869, 1), (870, 1), (871, 1), (872, 1), (873, 1), (874, 1), (875, 1), (876, 1), (877, 1), (878, 1), (879, 1), (880, 1), (881, 1), (882, 1), (883, 1), (884, 1), (885, 1), (886, 1), (887, 1), (888, 1), (889, 1), (890, 1), (891, 1), (892, 1), (893, 1), (894, 1), (895, 1), (896, 1), (897, 1), (898, 1), (899, 1), (900, 1), (901, 1), (902, 1), (903, 1), (904, 1), (905, 1), (906, 1), (907, 1), (908, 1), (909, 1), (910, 1), (911, 1), (912, 1), (913, 1), (914, 1), (915, 1), (916, 1), (917, 1), (918, 1), (919, 1), (920, 1), (921, 1), (922, 1), (923, 1), (924, 1), (925, 1), (926, 1), (927, 1), (928, 1), (929, 1), (930, 1), (931, 1), (932, 1), (933, 1), (934, 1), (935, 1), (936, 1), (937, 1), (938, 1), (939, 1), (940, 1), (941, 1), (942, 1), (943, 1), (944, 1), (945, 1), (946, 1), (947, 1), (948, 1), (949, 1), (950, 1), (951, 1), (952, 1), (953, 1), (954, 1), (955, 1), (956, 1), (957, 1), (958, 1), (959, 1), (960, 1), (961, 1), (962, 1), (963, 1), (964, 1), (965, 1), (966, 1), (967, 1), (968, 1), (969, 1), (970, 1), (971, 1), (972, 1), (973, 1), (974, 1), (975, 1), (976, 1), (977, 1), (978, 1), (979, 1), (980, 1), (981, 1), (982, 1), (983, 1), (984, 1), (985, 1), (986, 1), (987, 1), (988, 1), (989, 1), (990, 1), (991, 1), (992, 1), (993, 1), (994, 1), (995, 1), (996, 1), (997, 1), (998, 1), (999, 1), (1000, 1), (1001, 1), (1002, 1), (1003, 1), (1004, 1), (1005, 1), (1006, 1), (1007, 1), (1008, 1), (1009, 1), (1010, 1), (1011, 1), (1012, 1), (1013, 1), (1014, 1), (1015, 1), (1016, 1), (1017, 1), (1018, 1), (1019, 1), (1020, 1), (1021, 1), (1022, 1), (1023, 1), (1024, 1), (1025, 1), (1026, 1), (1027, 1), (1028, 1), (1029, 1), (1030, 1), (1031, 1), (1032, 1), (1033, 1), (1034, 1), (1035, 1), (1036, 1), (1037, 1), (1038, 1), (1039, 1), (1040, 1), (1041, 1), (1042, 1), (1043, 1), (1044, 1), (1045, 1), (1046, 1), (1047, 1), (1048, 1), (1049, 1), (1050, 1), (1051, 1), (1052, 1), (1053, 1), (1054, 1), (1055, 1), (1056, 1), (1057, 1), (1058, 1), (1059, 1), (1060, 1), (1061, 1), (1062, 1), (1063, 1), (1064, 1), (1065, 1), (1066, 1), (1067, 1), (1068, 1), (1069, 1), (1070, 1), (1071, 1), (1072, 1), (1073, 1), (1074, 1), (1075, 1), (1076, 1), (1077, 1), (1078, 1), (1079, 1), (1080, 1), (1081, 1), (1082, 1), (1083, 1), (1084, 1), (1085, 1), (1086, 1), (1087, 1), (1088, 1), (1089, 1), (1090, 1), (1091, 1), (1092, 1), (1093, 1), (1094, 1), (1095, 1), (1096, 1), (1097, 1), (1098, 1), (1099, 1), (1100, 1), (1101, 1), (1102, 1), (1103, 1), (1104, 1), (1105, 1), (1106, 1), (1107, 1), (1108, 1), (1109, 1), (1110, 1), (1111, 1), (1112, 1), (1113, 1), (1114, 1), (1115, 1), (1116, 1), (1117, 1), (1118, 1), (1119, 1), (1120, 1), (1121, 1), (1122, 1), (1123, 1), (1124, 1), (1125, 1), (1126, 1), (1127, 1), (1128, 1), (1129, 1), (1130, 1), (1131, 1), (1132, 1), (1133, 1), (1134, 1), (1135, 1), (1136, 1), (1137, 1), (1138, 1), (1139, 1), (1140, 1), (1141, 1), (1142, 1), (1143, 1), (1144, 1), (1145, 1), (1146, 1), (1147, 1), (1148, 1), (1149, 1), (1150, 1), (1151, 1), (1152, 1), (1153, 1), (1154, 1), (1155, 1), (1156, 1), (1157, 1), (1158, 1), (1159, 1), (1160, 1), (1161, 1), (1162, 1), (1163, 1), (1164, 1), (1165, 1), (1166, 1), (1167, 1), (1168, 1), (1169, 1), (1170, 1), (1171, 1), (1172, 1), (1173, 1), (1174, 1), (1175, 1), (1176, 1), (1177, 1), (1178, 1), (1179, 1), (1180, 1), (1181, 1), (1182, 1), (1183, 1), (1184, 1), (1185, 1), (1186, 1), (1187, 1), (1188, 1), (1189, 1), (1190, 1), (1191, 1), (1192, 1), (1193, 1), (1194, 1), (1195, 1), (1196, 1), (1197, 1), (1198, 1), (1199, 1), (1200, 1), (1201, 1), (1202, 1), (1203, 1), (1204, 1), (1205, 1), (1206, 1), (1207, 1), (1208, 1), (1209, 1), (1210, 1), (1211, 1), (1212, 1), (1213, 1), (1214, 1), (1215, 1), (1216, 1), (1217, 1), (1218, 1), (1219, 1), (1220, 1), (1221, 1), (1222, 1), (1223, 1), (1224, 1), (1225, 1), (1226, 1), (1227, 1), (1228, 1), (1229, 1), (1230, 1), (1231, 1), (1232, 1), (1233, 1), (1234, 1), (1235, 1), (1236, 1), (1237, 1), (1238, 1), (1239, 1), (1240, 1), (1241, 1), (1242, 1), (1243, 1), (1244, 1), (1245, 1), (1246, 1), (1247, 1), (1248, 1), (1249, 1), (1250, 1), (1251, 1), (1252, 1), (1253, 1), (1254, 1), (1255, 1), (1256, 1), (1257, 1), (1258, 1), (1259, 1), (1260, 1), (1261, 1), (1262, 1),

(7) 通过 doc2bow 计算测试数据的稀疏向量;

(8) 求得测试数据与样本数据的相似度。

根据上述得出的热度指数进行降序处理, 即得到题目所求的热点问题。

由于最终结果是根据留言详情进行热度排序后输出, 而题目要求按照特定格式给出排名前 5 的热点问题及其对应的留言信息, 并保存为相应的文件, 则后续还需利用算法概括对应的热点问题进行简单描述, 填入对应表中, 此处就便不予详细说明。

4.2.4 方法所存在的不足

不足点 1: 使用基于规则的方法提取命名实体, 方法虽简单, 但并不能很好的从大量文本中提取出特定地点和特定人群, 同时会影响到后面的相似度比较。

不足点 2: 使用“Gensim-文本相似度分析”的方法对进行相似度计算, 对大文本而言, 相似度的相似性较低, 需不断进行测试, 提高相似性, 才能更好的解决问题, 将热点问题提取出来。

4.3 问题三: 制定评价方案

(1) 相关性

相关性是指两个变量的关联程度, 也指在网络搜索中, 关键字被搜索引擎收录所依据的相关指标。而本题的题干要求从相关性分析答复意见, 根据是否与留言详情相关作为判断, 只要留言详情中的关键字和答复意见的关键字相似, 就可以判断答复意见的有效性。

拟操作: 先在附件 4 分别提取留言详情和答复意见, 分别对他们进行分词, 去除停用词, 统计每一列的词频, 并把词频放在向量中, 用 TF-IDF 模型计算相似度。

(2) 完整性

包括两个方面, 第一是格式完整性, 第二是内容完整性。

格式完整性: 开头要有尊称, 要问好。要用标准词不能使用口头语, 避免群众误会理解错误。结尾要附上时间。

内容完整性: 首先先明确群众提出的问题, 大致的复述一遍, 让群众明白这是哪个问题的答复意见, 还要表示已收到, 对于询问相关的法律法规或对法律法规有所疑惑, 可以直接在答复中给出相应的法律法规条例或仔细解析。对于环境

还是怀疑税务之类的，则说明已经反馈到相应的职能部门，正在等待处理，请群众耐心等待。

(3) 合理性

回复的内容是否客观，没有激烈的言辞，能否让人理解。

5 结果分析

5.1 问题一结果分析

问题一主要通过 TF-IDF 算法和 SVM 方法对留言进行分类。通过结果可以发现方法中仍存在某些不足之处，从而导致结果产生偏差。例如：在使用 TF-IDF 算法过程中可能会因为前期的数据清洗的不到位导致特征值不准确，而特征值主要是用以评估一字词对于一个文件集或一个语料库中的其中一份文件的重要程度，使得分类的准确度降低；通常使用 SVM 方法实现构造并训练分类器，但该方法的执行需要大量的存储空间，若硬件配置较差，也会对结果产生一定的影响，甚至使得程序无法正常运行。综上所述，若想结果准确性较高，硬件、软件设备都要跟得上。具体结果如图 10 所示。

accuracy 0.9903279697529236				
	precision	recall	f1-score	support
交通运输	0.99	1.00	1.00	1562
劳动和社会保障	0.98	0.98	0.98	1621
卫生计生	0.99	0.99	0.99	1372
商贸旅游	0.99	0.99	0.99	1599
城乡建设	0.99	0.98	0.98	1648
教育文体	0.99	0.99	0.99	2100
环境保护	0.99	1.00	1.00	1471
avg / total	0.99	0.99	0.99	11373

图 10 题一模型评价结果

5.2 问题二结果分析

通过对民众留言信息进行热度排序计数，选取热度排名前 5 的热点问题（结

果见表 2 排名前 5 的热点问题)并结合通过“留言主题”形成的词云图(见图 11)进行分析,可以发现现阶段的热点问题还是主要以环境污染、噪音污染等为主,总的来说就是与人们日常生活息息相关,这也说明目前政府对环境管理这方面所采取的措施仍有限,后期需对这方面加强管理力度,更好的服务民众。

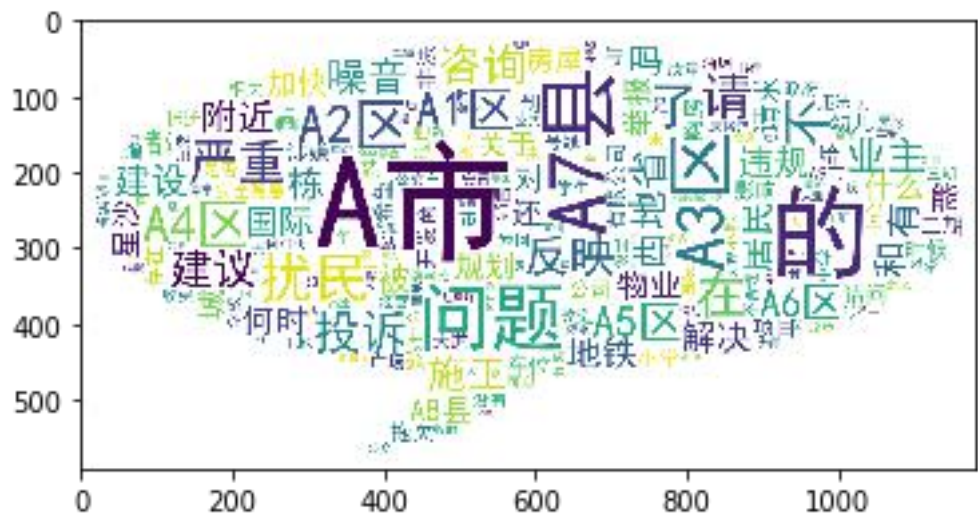


图 11 “留言主题”词云图

表 2 排名前 5 的热点问题

热度排名	问题ID	热度指数	时间范围	地点/人群	问题描述
1	1	55	2019/07/11 至 2019/12/11	A 市伊景园滨河苑	开发商捆绑销售车位
2	2	52	2019/07/21 至 2020/01/15	A 市 A2 区丽发新城小区	小区附近搅拌厂灰尘噪音污染严重
3	3	21	2019/07/21 至 2019/12/04	A 市 A5 区魅力之城小区	小区临街餐饮店油烟噪音扰民
4	4	12	2017/06/08 至 2019/11/27	A 市经济学院学生	学校强制学生去定点企业学习
5	5	12	2019/01/13 至 2019/05/22	A1 区辉煌国际城二期	居民楼下商铺违法开饭店

参考文献

- [1] 胡国胜,钱玲,张国红.支持向量机的多分类算法[J].系统工程与电子技术,2006,28(1):127-132.
- [2] 大漠帝国.文本数据增强二(EDA、同义词替换-新增-交换-删除-生成同义句).CSDN 博客,2019-04-27 19:32:13.
- [3] -派神-.使用python和sklearn的中文文本多分类实战开发.CSDN 博客,2019-03-02 00:32:15.
- [4] 陶伟.警务应用中基于双向最大匹配法的中文分词算法实现[J].电子技术与软件工程,2016(04):153-155.
- [5] 刘志.基于用户兴趣的协同过滤算法的广告推荐研究[D].昆明理工大学,2014.