

# “智慧政务”中的文本挖掘应用

## 摘要

本文旨在设计运用自然语言处理和文本挖掘等方法使计算机能够解决群众留言自动分类，热点问题挖掘，答复意见的质量评价等任务。近年来，随着互联网的广泛应用和网络问政的迅速发展，网络问政平台已经成为政府获取群众信息的主要渠道。因此，运用网络文本分析和数据挖据对群众问政留言记录及相关部门对部分群众留言的答复意见的研究具有重大意义。

**针对问题一**，首先，将附件 2 中列名为“留言详情”的内容作为特征，列名为“一级分类”的内容作为标签，建立关于留言内容分类的 **BERT-wwm 模型**，得到每个测试数据属于所有类别的概率值矩阵，矩阵中概率值最大的列即为所属类别，使用**查准率、查全率、F<sub>1</sub>-score** 对分类结果进行评价，计算出**查准率**的值为 0.92，**查全率**的值为 0.91，**F<sub>1</sub>-score** 的值为 0.91。

**针对问题二**，本文通过建立热度评价模型： $\text{热度} = w_1 * \text{该热点留言总数} + w_2 * (\text{总点赞数} - \text{总反对数}) / \text{该热点留言总数}$ ，根据附件 3 提供的内容，使用 **DBSCAN 聚类模型**根据留言主题和留言详情摘要对留言内容进行聚类，得到 27 个聚类簇，对每个聚类簇进行排名，得到前 5 个**热点问题**，提取时间范围和地点人群的信息，使用 **TextRank 算法**提取每一簇里面得分最高的留言主题作为该热点的问题描述，然后分类汇总每条留言。

**针对问题三**，本文根据附件 4 所给数据，对答复的**相关性、完整性、可解释性、及时性**进行定义，建立对答复意见质量的综合评价模型。

**关键词：**智慧政务；BERT-wwm 模型；热点问题；DBSCAN 密度聚类；TextRank 算法；Doc2vec 模型；余弦距离。

## Abstract

This article aims to design the use of natural language processing and text mining methods to enable the computer to solve tasks such as automatic classification of mass messages, mining of hot issues, and quality evaluation of answers. In recent years, with the widespread use of the Internet and the rapid development of online questioning, online questioning platforms have become the main channel for the government to obtain mass information. Therefore, it is of great significance to use network text analysis and data mining to study the records of the masses' political messages and the relevant departments' responses to some masses' messages.

For question one, first, take the content listed in Annex 2 as "message details" as features, and the content listed as "first-class classification" as tags, establish a BERT-wwm model for the classification of message content, and get each test data belonging to all categories. The probability value matrix of the matrix, the column with the largest probability value in the matrix is the category. Use the precision rate, recall rate, and F1-score to evaluate the classification results, calculate the precision rate value is 0.92, the recall rate value is 0.91, and the value of F1-score is 0.91.

For question two, In this paper, through the establishment of a heat evaluation model,  $\text{heat} = w_1 * \text{total number of hotspot messages} + w_2 * (\text{total likes} - \text{total opposition}) / \text{total number of hotspots}$ . According to the content provided in Annex 3, use the DBSCAN clustering model according to the subject of the message and The message details summary clusters the content of the message to get 27 clusters, ranks each cluster, gets the top 5 hotspot questions, extracts the time range and location crowd information, and uses the TextRank algorithm to extract each cluster. The subject of the message with the highest score is used as the problem description of the hotspot, and then each message is classified and summarized.

In response to question three, This article defines the relevance, completeness, interpretability, and timeliness of the response based on the data given in Annex 4, and establishes a comprehensive evaluation model for the quality of the response opinion.

**Key word:** Smart government affairs; BERT-wwm model; DBSCAN clustering;  
TextRank algorithm; Doc2vec model; cosine distance

# 目录

|                                      |    |
|--------------------------------------|----|
| 1 问题重述.....                          | 6  |
| 1.1 问题背景.....                        | 6  |
| 1.2 要解决的问题.....                      | 6  |
| 2 问题一：群众留言分类.....                    | 7  |
| 2.1 问题分析.....                        | 7  |
| 2.2 模型建立.....                        | 7  |
| 2.2.1 BERT-www 模型.....               | 7  |
| 2.2.2 Bert 网络框架.....                 | 8  |
| 2.2.3 Bert 所用 Transformer 内部结构图..... | 9  |
| 2.3 数据预处理.....                       | 10 |
| 2.3.1 数据清洗.....                      | 10 |
| 2.3.2 数据集划分.....                     | 10 |
| 2.4 模型求解.....                        | 11 |
| 2.4.1 实验条件.....                      | 11 |
| 2.4.2 数据选择.....                      | 11 |
| 2.4.3 算法流程.....                      | 12 |
| 2.4.4 程序框架.....                      | 12 |
| 2.4.5 模型的训练.....                     | 13 |
| 2.4.6 模型的训练结果.....                   | 15 |
| 2.5 模型评价.....                        | 15 |
| 3 问题二： .....                         | 18 |
| 3.1 问题分析.....                        | 18 |
| 3.2 模型建立.....                        | 18 |
| 3.2.1 热度评价模型.....                    | 18 |
| 3.3 数据预处理.....                       | 18 |
| 3.3.1 数据清洗.....                      | 18 |
| 3.3.2 特征提取.....                      | 19 |

|                            |    |
|----------------------------|----|
| 3.4 算法介绍.....              | 19 |
| 3.4.1 DSSCAN 聚类算法.....     | 19 |
| 3.4.2 TextRank 自动文本摘要..... | 21 |
| 3.5 模型求解.....              | 22 |
| 3.5.1 求解流程.....            | 22 |
| 3.5.2 结果展示.....            | 25 |
| 4 问题三： .....               | 28 |
| 4.1 问题分析.....              | 28 |
| 4.2 解题思路.....              | 28 |
| 4.3 模型建立.....              | 28 |
| 5 模型的局限性.....              | 28 |
| 参考文献.....                  | 29 |

# 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、凝聚民气的重要渠道，各类社情民意相关的文本数据量不断攀升，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了极大挑战。同时，随着大数据、云计算、人工智能等技术的发展，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，对提升政府的管理水平和施政效率具有极大的推动作用。

## 1.2 要解决的问题

1) 根据附件2给出的数据，建立有关于留言内容的一级标签分类模型。通过建立好的模型，对留言内容进行分类，使用F<sub>1</sub>-score对模型的分类效果进行评价。

2) 根据附件3对留言内容进行聚类，再对聚好类的结果进行热度分析，通常是定义一个热度指数，得到排名前五的热点问题，并给出结果，汇总留言信息，整理成相关表格。

3) 根据附件4，对留言的答复意见进行相关性、完整性、可解释性和及时性方面的评价，从而得出对答复意见的质量的评价方案。

## 2 问题一：群众留言分类

### 2.1 问题分析

首先，需要提取附件 2 中的留言内容及其对应的标签内容，对数据进行预处理，整理成可用于 bert 模型的数据集，按 8:1:1 的比例进行分层随机抽样将数据集划分为 train(训练集)，test(测试集)，dev(验证集)，利用该数据集建立有关于留言内容的一级标签分类 bert 模型，并对模型效果进行评价。

### 2.2 模型建立

#### 2.2.1 BERT-wwm 模型

BERT-wwm 模型，即中文全词覆盖的（Whole Word Masking）BERT 预训练模型，它在 BERT 模型的基础上进行了小升级，主要更改了原预训练阶段的训练样本生成策略。BERT(Bidirectional Encoder Representations from Transformers)模型<sup>[1]</sup>是由 Google 于 2018 年提出的，模型采用了表义能力更强的 Transformer 网络结构<sup>[2]</sup>来对语言模型进行训练，是一种通过大量语言资料训练得来的一个模型，是第一个用于在预训练 NLP 上的无监督的、深度双向系统。

Devlin J 等<sup>[1]</sup>设计的 BERT 模型主要分为两大部分，主要分为输入层和双向 Transformer 编码层，如下图 2-1 所示：

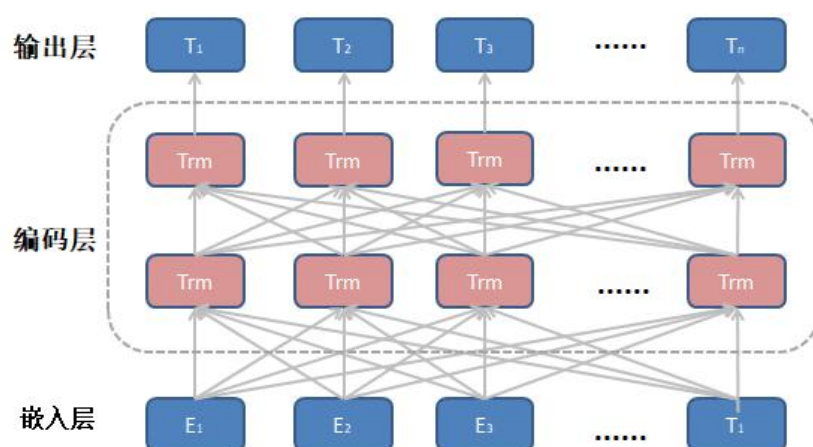


图 2-1 BERT 模型基本结构示意图

### 2.2.2 Bert 网络框架

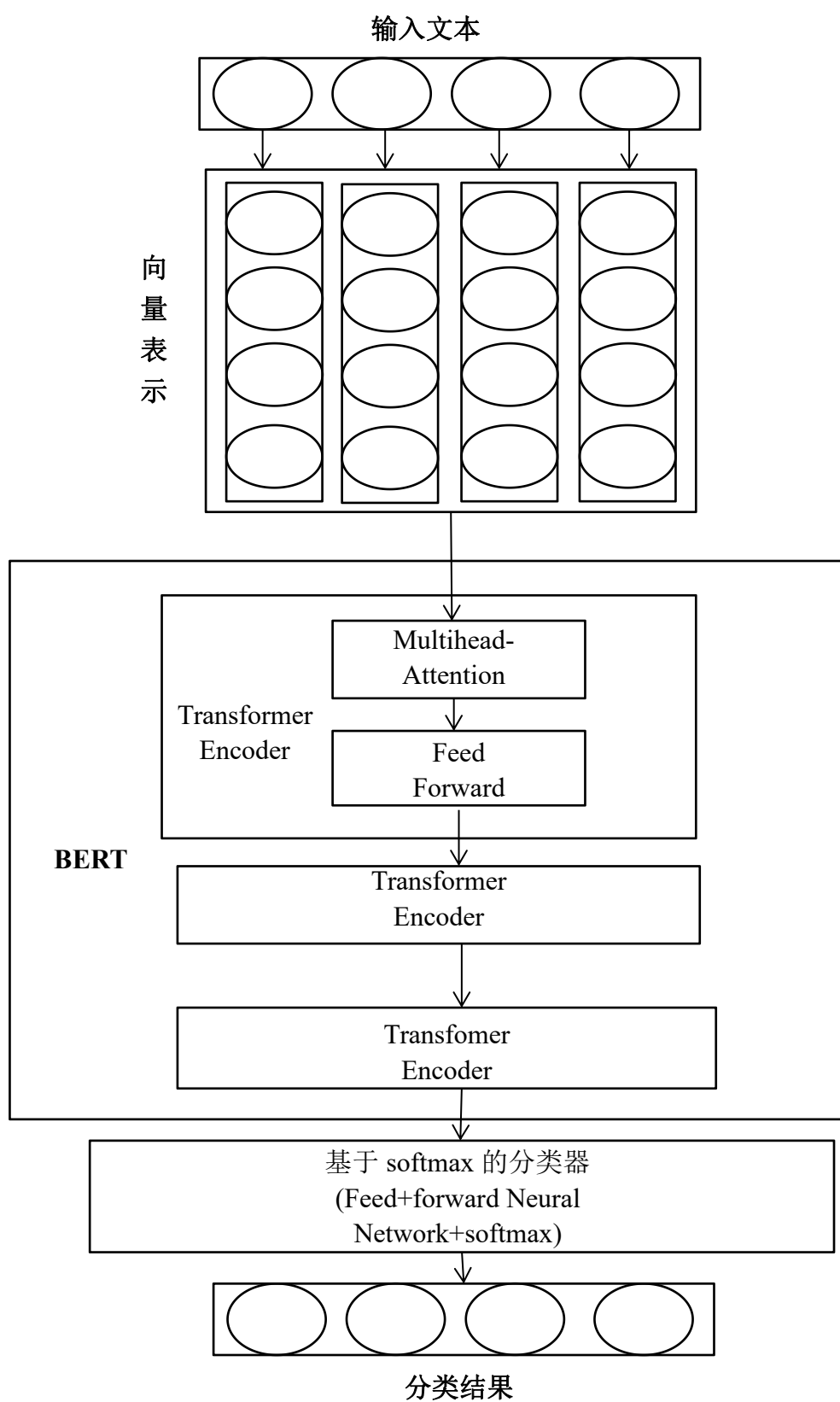


图 2-2 BERT 模型网络框架



### 2.2.3 Bert 所用 Transformer 内部结构图

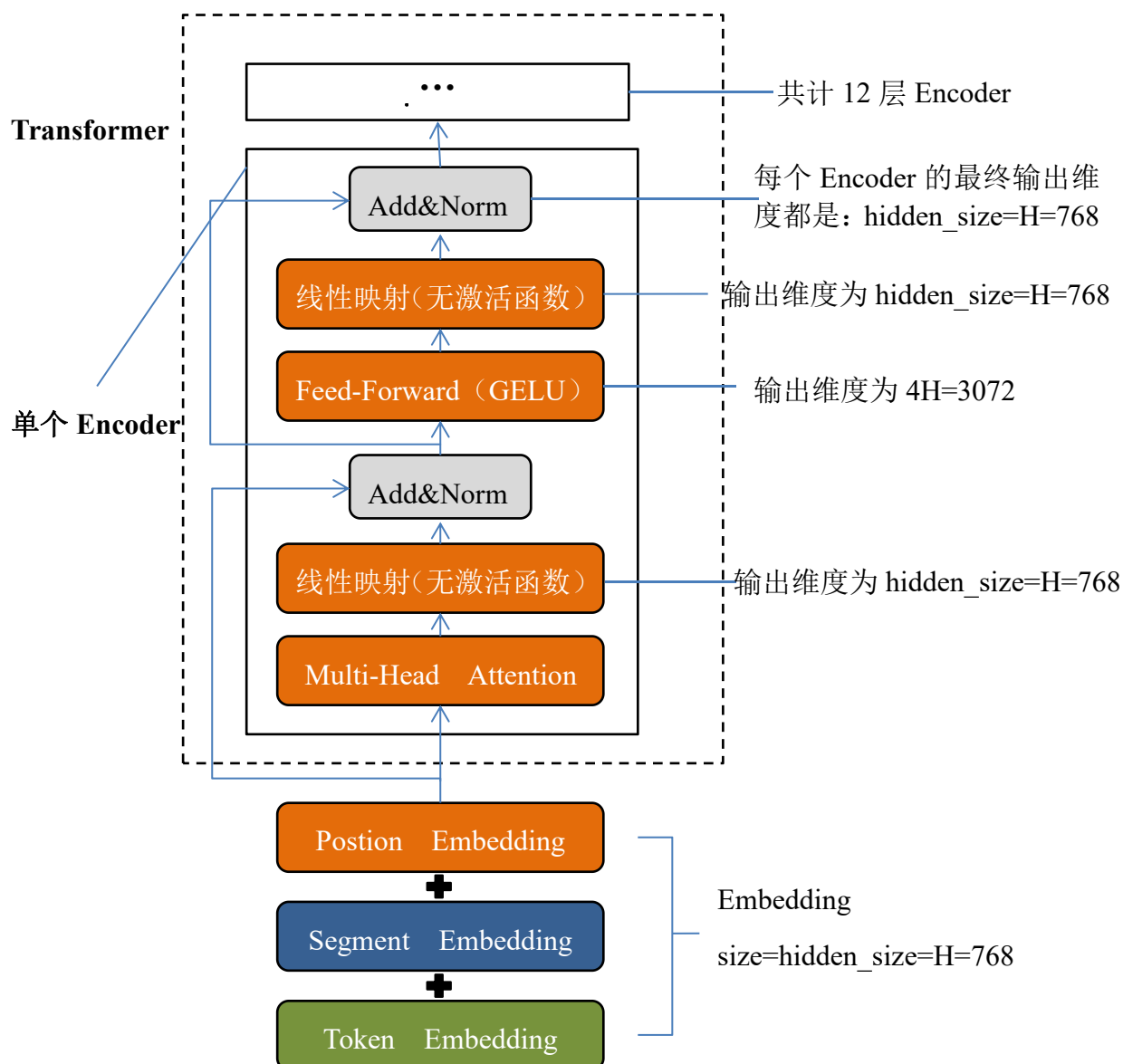


图 2-3 Transformer 内部结构图

由于谷歌官方发布的 BERT-base (Chinese) 中, 中文是以字为粒度进行切分, 没有考虑中文需要分词的特点。应用全词 mask, 而非字粒度的中文 BERT 模型可能有更好的表现, 因此我们在这里采用, 哈工大讯飞联合实验室发布的全词覆盖的中文 BERT 预训练模型(中文 BERT-wwm)。

下面展示了全词 Mask 的生成样例。 注意: 为了方便理解, 下述例子中只考虑替换成[MASK]标签的情况:

表 1-1 全词 Mask 的生成样例

| 说明         | 样例   |
|------------|--|
| 原始文本       | 使用语言模型来预测下一个词的 probability。  |
| 分词文本       | 使用 语言 模型 来 预测 下 一个 词 的 probability。                                  |
| 原始 Mask 输入 | 使 用 语 言 [MASK] 型 来 [MASK] 测 下 一 个 词 的 pro [MASK] ##lity。             |
| 全词 Mask 输入 | 使 用 语 言 [MASK] [MASK] 来 [MASK] [MASK] 下 一 个 词 的[MASK] [MASK] [MASK]。 |

本文模型在预训练的 BERT-wwm 模型的基础上进行了结构的微调，使其能更好的应用于本次文本分类任务。

## 2.3 数据预处理

### 2.3.1 数据清洗

(1) 读取数据并去除文本两端空格：运用 pandas 包中 read\_excel 导入附件 2 中的 4, 5 列，即留言详情与一级分类，由于留言详情中内容两端存在大量无意义的空格占据大量字符长度，因此对留言详情中的内容运用 str.strip() 去除两端空格。

(2) 文本去重：经过排查，留言详情中存在重复行，且对应的标签内容不一致，因此对于重复的数据进行删除

(3) 对数据进行空值和异常值的排查

### 2.3.2 数据集划分

对清洗完的数据按照 8: 1: 1 的比例按标签类别进行分层随机抽样，将数据集划分为训练集(train.tsv)，测试集(test.tsv)，开发集(dev.tsv)。(设置随机种子，保证每次随机的结果都是一样的)

将数据集导出为以下形式：

表 2-2 数据集示例

| Label   | content   |
|---------|---|
| 教育文体    | L3 县七甲坪镇的赶尸、傩文化、哭嫁等风俗风情即将流失。七甲坪正在筹建新开发区。七甲坪离 H 市... |
| 卫生计生    | 网络订餐这一新兴的餐饮营销迎合了消费者的需求，但是，作为一种新型的网络购物形式，网络订餐还存在...  |
| 劳动和社会保障 | 领导们：你们好！我们都是好润佳员工都属防损员共 11 人，我们在好润佳工作六七个月了他们没有...   |
| 卫生计生    | 领导，您好，我是病人家属，病人在 D7 县人民医院住院期间由于医护人员的违规操作在灌肠的时候把...  |
| 教育文体    | 局长你好，想问一下 B 市民办学校贷款，政府是否有贴息补助？                      |

## 2.4 模型求解

### 2.4.1 实验条件

表 2-3 实验环境配置情况

| 实验工具   | 型号或版本                |
|--------|----------------------|
| 显卡     | Tesla P100-PCIE-16GB |
| 操作系统   | Linux                |
| 编程语言   | Python 3.6.9         |
| 深度学习框架 | Tensorflow 1.15.0    |

### 2.4.2 数据选择

对于数据预处理后生成的 3 组数据集，train.tsv 作为训练集进行模型训练。dev.tsv 作为训练过程中的验证集，验证模型拟合效果。test.tsv 作为模型训练完成后的测试集，对模型的分类效果进行评价。

### 2.4.3 算法流程

整体算法流程图

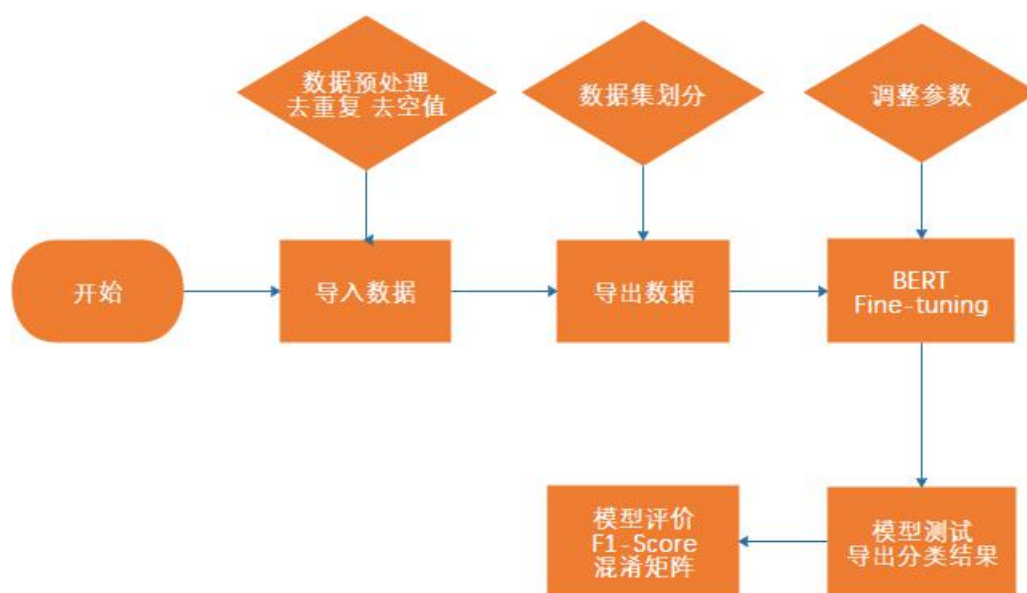


图 2-4 算法流程图

### 2.4.4 程序框架

BERT 本质上是一个两段式的 NLP 模型。第一个阶段叫做：Pre-training，跟 WordEmbedding 类似，利用现有无标记的语料训练一个语言模型。第二个阶段叫做：Fine-tuning，利用预训练好的语言模型，完成具体的 NLP 下游任务。Google 已经投入了大规模的语料和昂贵的机器帮我们完成了 Pre-training 过程

在 fine-tuning 过程这里使用 `run_classifier.py` 进行句子分类任务。从主函数开始，可以发现它指定一些必须的参数：“data\_dir”，“task\_name”，“vocab\_file”，“bert\_config\_file”，“output\_data”。

```
if __name__ == "__main__":
    flags.mark_flag_as_required("data_dir")
    flags.mark_flag_as_required("task_name")
    flags.mark_flag_as_required("vocab_file")
    flags.mark_flag_as_required("bert_config_file")
    flags.mark_flag_as_required("output_dir")
    tf.app.run()
```

BERT 代码中 `processor` 就是负责对模型的输入进行处理。这样，我们需要在原本 `main` 函数的 `processor` 字典里，加入自定义的 `processor` 类，取名为

MyTaskProcessor，即可在运行参数里指定调用该 processor。

```
def main(_):
    tf.logging.set_verbosity(tf.logging.INFO)
    processors = {
        "cola": ColaProcessor,
        "mnli": MnliProcessor,
        "mrpc": MrpcProcessor,
        "xnli": XnliProcessor,
        "mytask": MyTaskProcessor,#自定义的 Processor
    }
```

自定义的 processor 里需要继承 DataProcessor，并重载获取 label 的 get\_labels 和获取单个输入的 get\_train\_examples，get\_dev\_examples 和 get\_test\_examples 函数。其分别会在 main 函数的 FLAGS.do\_train、FLAGS.do\_eval 和 FLAGS.do\_predict 阶段被调用。

之后就可以直接运行 run\_classifier.py 进行模型的训练。在运行时需要制定一些参数。

### 2.4.5 模型的训练

经过分析，由于留言详情字符串长度百分之 75 分位数为: 465.0，为了模型有更好的分类效果,因此 max\_seq\_length 设置为 bert 模型所能设置的最大值 512，learning\_rate 则设置为由许多专业人士测试过的篇章级文本分类最佳学习率 2e-5

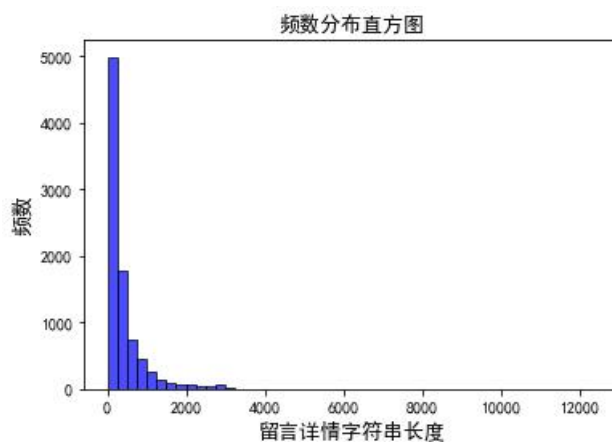


图 2-5 频数分布直方图



图 2-6 中 Num examples 为加载训练数据样本 6123 个，batch\_size: 批大小，即 1 次迭代所使用的样本量。以及对数据进行标注和向量化，并且随机 mask 后训练模型的过程。

## 2.4.6 模型的训练结果

```
***** Eval results *****  
  
eval_accuracy = 0.92005074  
  
eval_loss = 0.4481641
```

图 2-7 验证集评估结果

由于参数 do\_eval （是否在训练结束后验证），这里为 True，所以在模型训练完成后会对验证集进行准确率的评估，eval\_accuracy 即为准确率，评估结果显示:模型对于验证集的预测准确率为 0.92005074。

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (2-1)$$

True Positive (TP)：预测为正例，实际为正例，False Positive (FP)：预测为正例，实际为负例，True Negative (TN)：预测为负例，实际为负例，False Negative (FN)：预测为负例，实际为正例。

## 2.5 模型评价

使用建立的留言内容分类的 BERT-wwm 模型，对测试数据的类别进行预测，得到每个测试数据属于所有类别的概率值矩阵，矩阵中概率值最大的列即为所属类别，使用 F<sub>1</sub>-score 对分类结果进行评价。

### ■ 查准率 precision

$$precision = \frac{TP}{TP + FP} \quad (2-2)$$

### ■ 查全率 recall

$$recall = \frac{TP}{TP + FN} \quad (2-3)$$

## ■ F<sub>1</sub>-score

通常我们使用 F1-score 对分类方法进行评价，其中 F1-score 的式子如式 2-4 所示：

$$F_1 - score = \frac{1}{n} \sum_{i=1}^n \frac{2P_i R_i}{P_i + R_i} \quad (2-4)$$

其中为  $P_i$  第  $i$  类的查准率，为  $R_i$  第  $i$  类的查全率。F1-score 为 0 到 1 之间的值，越接近于 1，模型的效果越好。

True Positive (TP)：预测为正例，实际为正例，False Positive (FP)：预测为正例，实际为负例，True Negative (TN)：预测为负例，实际为负例，False Negative (FN)：预测为负例，实际为正例。

## ■ 查准率、查全率、F<sub>1</sub>-score

表 2-5 查准率、查全率、F<sub>1</sub>-score 结果

|           | precision | recall | F <sub>1</sub> -score | support |
|-----------|-----------|--------|-----------------------|---------|
| 城乡建设      | 0.86      | 0.92   | 0.89                  | 196     |
| 环境保护      | 0.96      | 0.94   | 0.95                  | 90      |
| 交通运输      | 0.82      | 0.89   | 0.86                  | 57      |
| 教育文体      | 0.96      | 0.91   | 0.94                  | 150     |
| 劳动和社会保障   | 0.97      | 0.94   | 0.96                  | 188     |
| 商贸旅游      | 0.86      | 0.81   | 0.83                  | 109     |
| 卫生计生      | 0.92      | 0.95   | 0.94                  | 85      |
| avg/total | 0.92      | 0.91   | 0.91                  | 875     |



## ■ 混淆矩阵

表 2-6 混淆矩阵

|         | 城乡建设 | 环境保护 | 交通运输 | 教育文体 | 劳动和社会保障 | 商贸旅游 | 卫生计生 |
|---------|------|------|------|------|---------|------|------|
| 城乡建设    | 181  | 2    | 5    | 2    | 1       | 5    | 0    |
| 环境保护    | 5    | 85   | 0    | 0    | 0       | 0    | 0    |
| 交通运输    | 5    | 0    | 51   | 0    | 0       | 1    | 0    |
| 教育文体    | 6    | 0    | 0    | 137  | 3       | 4    | 0    |
| 劳动和社会保障 | 4    | 1    | 2    | 0    | 177     | 1    | 3    |
| 商贸旅游    | 8    | 1    | 4    | 3    | 1       | 88   | 4    |
| 卫生计生    | 1    | 0    | 0    | 0    | 0       | 3    | 81   |

从结果上看，BERT-wwm 模型分类的  $F_1$ -score 达到了 0.91，准确率较高。

## 3 问题二：

### 3.1 问题分析

为了及时发现热点问题，让相关部门进行有针对性地处理，提升服务效率，需要做热点问题挖，根据附件 3 提供的内容，对数据进行清洗和提取特征后，采用 **DBSCAN 聚类模型**对留言主题和留言详情的摘要进行聚类，得到多个聚类簇，通过自定义热度指数对每个聚类簇进行排名，得到前 5 个热点问题，提取时间范围和地点人群的信息，使用 **TextRank 算法**提取每一簇里面得分最高的留言主题作为该热点的问题描述，并保存为文件“热点问题表.xlsx”。最后给出相应热点问题对应的留言信息，并保存为“热点问题留言明细表.xlsx”。

### 3.2 模型建立

#### 3.2.1 热度评价模型

经过聚类之后，表述同一个问题的留言被聚成了一簇，形成了若干问题簇，需要判断哪些簇算得上是热点，因此我们定义一个热度指数来将其进行排名，本文考虑到的第一个指标便是反应同一个问题的留言数量，第二个便是留言的点赞数和反对数，因为点赞数和反对数同样能反应民众的心声，点赞表示对问题有着同样的看法，反对表示此问题持否定态度。因此，本文将热度指标设定为：

$$\text{热度指数} = w_1 \times \text{该热点留言总数} + w_2 \times \frac{\text{该热点总点赞数} - \text{该热点总反对数}}{\text{该热点留言总数}} \quad (3-1)$$

其中  $w_1$  和  $w_2$  为自定义权重，具体取值可见模型求解。

### 3.3 数据预处理

#### 3.3.1 数据清洗

对留言进行去重、然后对新闻内容文本进行 jieba 分词并词性标注，过滤出名词、动词、简称等词性，分词前使用自定义的用户词词典增加分词的准确性，分词后使用停用词词典、消歧词典、保留单字词典过滤掉对话题无关并且影响聚

类准确性的词，建立每篇留言的词库。

### 3.3.2 特征提取

利用 TF-IDF<sup>[3]</sup>特征提取，得到每篇留言的特征向量。

TF 算法：统计一个词在一篇文档中出现的频次。基本思想是一个词在文档中出现的次数越多，则其对文档的表达能力也就越强。公式如下：

$$tf_{ij} = \frac{n_{ij}}{\sum_k n_{kj}} \quad (3-2)$$

IDF 算法：统计一个词在文档集的多少个文档中出现。基本的思想是一个词在越少的文档中出现，其对文档的区分能力也就越强。公式如下：

$$idf_i = \log\left(\frac{|D|}{1+|D_i|}\right) \quad (3-3)$$

TF-IDF 算法就是 TF 算法与 IDF 算法的综合，计算公式如下：

$$tf \times idf(i, j) = tf_{ij} \times idf_i = \frac{n_{ij}}{\sum_k n_{kj}} \times \log\left(\frac{|D|}{1+|D_i|}\right) \quad (3-4)$$

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，可以产生出高权重的 TF-IDF。因此，TF-IDF 倾向于过滤掉常见的词语，保留重要的词语。

## 3.4 算法介绍

### 3.4.1 DSSCAN 聚类算法

利用 TF-IDF 特征提取之后对新闻进行 DBSCAN 聚类<sup>[4]</sup>，并对每个类的大小进行排序。基于密度的 DBSCAN 算法进行新闻聚类是最为准确的做法，这种算法的特点可以发现任意形状的簇，同时可以过滤离群点，而不必把离群点分在某一个簇中，增加聚类的偏差。

DBSCAN 聚类的原理很简单：由密度可达关系导出最大密度相连的样本集合（聚类）。这样的集合中有一个或多个核心对象，如果只有一个核心对象，则簇中其他非核心对象都在这个核心对象的  $\varepsilon$  邻域内；如果是多个核心对象，那么任意一个核心对象的  $\varepsilon$  邻域内一定包含另一个核心对象（否则无法密度可达）。这些核心对象以及包含在它  $\varepsilon$  邻域内的所有样本构成一个类。算法流程图如下图

所示：

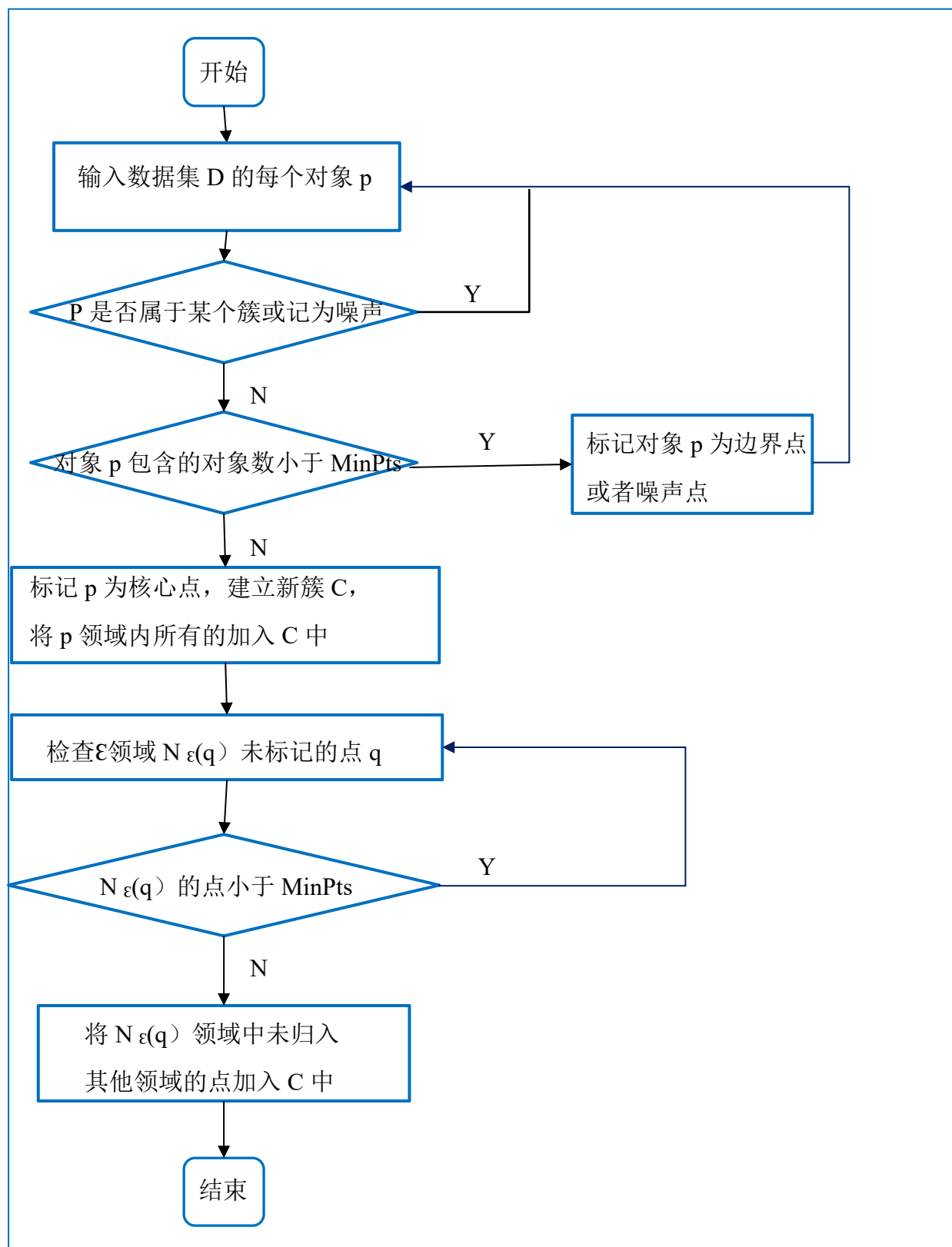


图 3-1 DBSCAN 流程图

在聚类的过程中需要输入三个参数：

(1)  $D$ : 一个包含  $n$  个对象的数据集

(2)  $\epsilon$ : 半径参数

(3)  $MinPts$ : 领域密度阈值

在计算过程中, 使用“余弦相似度”来计算距离, 所以本文在设置  $\epsilon$  参数时, 一般都设置为 0.3-0.5 之间, 因为超过 0.5 的半径值会使不属于同类新闻聚在一起, 小于 0.3 则无法识别相同事件的新闻。设置  $MinPts$  参数时, 可以人为的假定一个阈值, 假设有 10 条新闻报道同一事件就认为它就是一个热点, 那么可以设置  $MinPts$  值为 10。

### 3.4.2 TextRank 自动文本摘要

针对聚类后的每一类热点, 为了得到该处热点的话题信息, 还需要提取它们的标题, 利用 TextRank 算法<sup>[5]</sup>, 对标题的重要程度进行排序, 用重要性最高的标题来描述该处热点的话题。

TextRank 是一种基于图的用于文本的排序算法, 基本思想来自于 Google 的 PageRank 算法<sup>[3]</sup>。类似于网页的排名, 对于词语可得到词语的排行, 对于句子也可得到句子的排名, 所以 TextRank 可以进行关键词提取, 也可以进行自动文摘。其用于自动文摘时的思想是: 将每个句子看成 PageRank 图中的一个节点, 若两个句子之间的相似度大于设定的阈值, 则认为这两个句子之间有相似联系, 对应的这两个节点之间便有一条无向有权边, 边的权值是相似度, 接着利用 PageRank 算法即可得到句子的得分, 把得分较高的句子作为文章的摘要。

TextRank 算法的主要步骤如下:

(1) 预处理: 分割原文本中的句子得到一个句子集合, 然后对句子进行分词以及去停用词处理, 筛选出候选关键词集。

(2) 计算句子间的相似度: 采用如下公式进行计算句子  $i$  和句子  $j$  的相似度:

$$Similarity(S_i, S_j) = \frac{n}{\log(S_i) + \log(S_j)} \quad (3-5)$$

其中,  $n$  为两个句子都出现的词的数目,  $S_i$  为句子  $i$  中的词的数目,  $S_j$  为句子  $j$  中的词的数目。对于两个句子之间的相似度大于设定的阈值的两个句子节点用边连接起来, 设置其边的权重为两个句子的相似度。

(3) 计算句子权重:

$$WS(V_i) = (1-d) + d \times \sum_{V_j \in In(V_i)} \left( \frac{w_{ji}}{\sum_{V_k \in Out(V_j)} w_{jk}} \times WS(V_j) \right) \quad (3-6)$$

其中， $WS(V_i)$  表示句子  $i$  的权重， $d$  为阻尼系数， $w_{ji}$  表示句子  $i$  和句子  $j$  的相似度， $WS(V_j)$  表示句子  $j$  的权重， $V_j \in In(V_i)$  表示与句子  $i$  相连的句子， $\sum_{V_k \in Out(V_j)} w_{jk}$  表示所有与句子  $j$  相连的句子的边的权重和。

(4) 形成文摘：将句子按照句子得分进行倒序排序，抽取得分排序最前的几个句子作为候选文摘句，再依据字数或句子数量要求筛选出符合条件的句子组成文摘。

## 3.5 模型求解

### 3.5.1 求解流程

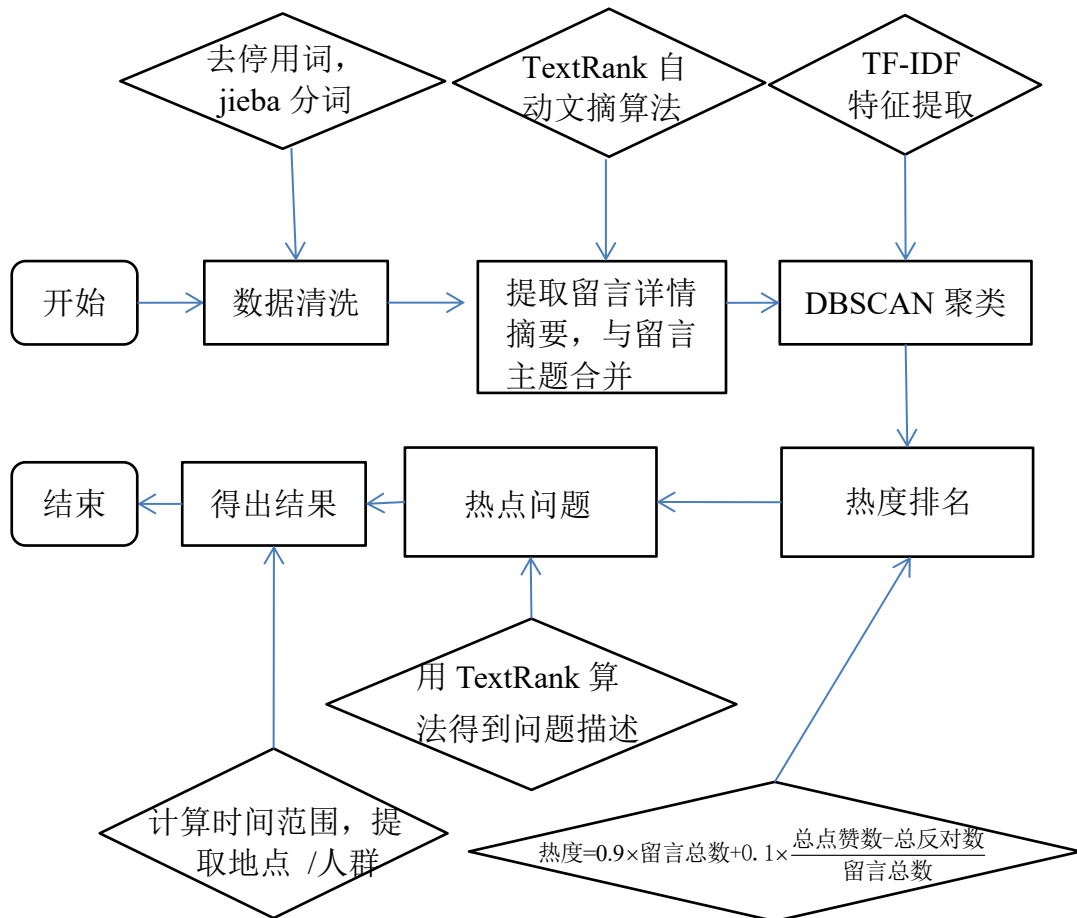


图 3-2 求解流程图

第一步：对留言内容进行清洗并进行结巴分词及词性筛选。

第二步：由于单纯使用留言主题进行聚类 and 单纯使用留言详情进行聚类，效果都不是很好。于是本文决定采用留言主题和留言详情相结合的方法进行聚类，使用 TextRank 自动文本摘要算法提取每条留言详情里得分最高的 3 个句子作为每条留言的摘要，并将其与留言主题合并。

这样做是因为留言主题是对留言详情的高度概括，所以要利用起来，而留言详情内容太多又会影响聚类效果，使用 TextRank 算法提取留言摘要，这样就间接地利用了留言详情，能够拥有较好的聚类效果。

第三步：使用 TF-IDF 算法提取特征，得到每篇留言的特征向量。

第四步：使用 DBSCAN 密度聚类模型根据留言主题和留言详情的摘要结合起来的文本对留言内容进行聚类，得到多个聚类簇。离群点会被筛选掉。关于 DBSCAN 密度聚类需要设置的参数：本文设置  $\epsilon$  参数为 0.49，MinPts 参数为 5；一共聚成 27 个聚类簇（不包含离群点）。

第五步：通过式 3-1 所建立的热度评价模型对每个聚类簇进行排名，令其  $w_1 = 0.9$ ， $w_2 = 0.1$ ；使用 TextRank 算法提取每一簇里面得分比较高的留言主题作为该问题描述，从而得热点问题的排名。

$$\text{热度指数} = 0.9 \times \text{该热点留言总数} + 0.1 \times \frac{\text{该热点总点赞数} - \text{该热点总反对数}}{\text{该热点留言总数}} \quad (3-7)$$

结果如表 3-1：

表 3-1 聚类和排名结果

| 热度排名 | 热点内容                          |
|------|-------------------------------|
| 热点 1 | 投诉 A 市伊景园滨河苑捆绑车位销售            |
| 热点 2 | 丽发新城小区附近的搅拌站噪声严重扰民            |
| 热点 3 | A4 区绿地海外滩小区距渝长夏高铁太近了          |
| 热点 4 | A 市人才新政落户后屡次申请购房补贴不成功         |
| 热点 5 | 咨询 A 市公积金贷款买房的问题              |
| 热点 6 | 请 A 市加快国家中心城市建设力度             |
| 热点 7 | 咨询 A3 区西湖街道茶场村五组的拆迁规划         |
| 热点 8 | 何时发布 2019 年 A7 县普惠性幼儿园清单及收费标准 |

|       |                                 |
|-------|---------------------------------|
| 热点 9  | A 市高新区金领公寓每天停水                  |
| 热点 10 | A 市万科魅力之城小区近百户楼板开裂墙面开裂          |
| 热点 11 | A7 县星沙四区凉塘路旧城改造要拖到何年何月才能动工      |
| 热点 12 | A 市很多驾校有违规收取求职者押金的行为            |
| 热点 13 | A 市 805 路公交车改到能不能走木莲中路          |
| 热点 14 | 请问 A 市地铁 3 号线什么时候开通             |
| 热点 15 | A4 区洪山公园何时开工建设                  |
| 热点 16 | 对西地省高速公路建设开发总公司薪酬改革合理性的质疑       |
| 热点 17 | A7 县第一中学违规补课收费                  |
| 热点 18 | 举报湖楚财富及两大国资委侵吞百姓血汗钱             |
| 热点 19 | 请解决 A4 区凯乐国际城周边太平路的路灯问题         |
| 热点 20 | A 市普通上班族的生育津贴和医院生产费一份都不能报销      |
| 热点 21 | 咨询 A 市转业士官异地安置问题                |
| 热点 22 | 建议 A 市经开区泉星公园项目规划进一步优化          |
| 热点 23 | A3 区保利西海岸配套幼儿园迟迟不交付             |
| 热点 24 | A 市 A3 区兰亭湾畔小区违法开餐厅             |
| 热点 25 | 在 A8 县连续缴纳社保 36 个月，是否拥有 A 市购房资格 |
| 热点 26 | A7 县榔梨领东汇小区非法住改商，开设麻将馆扰民严重      |
| 热点 27 | A3 区青青家园小区违规住改商存隐患              |

第六步：提取排名前 5 的热点问题，统计热点问题的留言时间的始末，采用词性标注并提取特定词性的方法提取留言的地点人群，并保存为文件“热点问题表.xlsx”。最后给出相应热点问题对应的留言信息，并保存为“热点问题留言明细表.xlsx”。



### 3.5.2 结果展示

#### ■ 热点问题表

表 3-2 热点问题表

| 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    | 37.82381 | 2019-07-07至2020-01-06 | A市伊景园  | 投诉A市伊景园滨河苑捆绑车位销售     |
| 2    | 2    | 36.99756 | 2019-11-02至2020-01-25 | 丽发新城小区 | 丽发新城小区附近的搅拌站噪音严重扰民   |
| 3    | 3    | 15.8375  | 2019-08-23至2019-09-06 | A4区绿地  | A4区绿地海外滩小区距渝长厦高铁太近了  |
| 4    | 4    | 14.5     | 2018-11-15至2019-12-02 | A市人才新政 | A市人才新政落户后屡次申请购房补贴不成功 |
| 5    | 5    | 10.96667 | 2019-01-02至2019-12-24 | A市公积金  | 咨询A市公积金贷款买房的问题       |

#### ■ 热点问题留言明细表

表 3-3 热点问题明细表部分结果示例

| 问题ID | 留言编号   | 留言用户      | 留言主题                | 留言时间                | 留言详情   | 点赞数 | 反对数 |
|------|--------|-----------|---------------------|---------------------|--|-----|-----|
| 1    | 188801 | A909180   | 投诉滨河苑针对广铁职工购房的霸王规定  | 2019-08-01 00:00:00 | 尊敬的张市长，您好！我叫李建义，来自湖北仙桃，虽然已经在广铁集团A市公司工作了十几年，但是依旧没有一套自己的房子。所以当我听说A市政府和集团牵头让我们可以定向购买商品房的伊景园滨河苑……                        | 0   | 0   |
| 1    | 190337 | A00090519 | 关于伊景园滨河苑捆绑销售车位的维权投诉 | 2019-08-23 12:22:00 | 投诉伊景园.滨河苑开发商捆绑销售车位！A市武广新城片区下的伊景园.滨河苑是广铁集团铁路职工的定向商品房，之前已经统一交了18.5万的认购款，但没有正规的合同，现在集团下发文件，强制要求职工再交12万的车位费，不交就取消购房资格，…… | 0   | 0   |
| 1    | 195511 | A90923    | 车位捆绑违规              | 2019-08-01          | 对于伊景园滨河苑商品房，A市广铁集团违规捆  | 0   | 0   |

|       |            |                   |  |                                    |   |       |       |
|-------|------------|-------------------|--|------------------------------------|---|-------|-------|
|       |            | 7                 | 销售   | 8-16<br>14:20:<br>26               | 绑车位销售至今,买房必须<br>买车位我们反映多次一直<br>没有人彻查处理,到底是什<br>么样的力量在支持他们。现<br>在各部门都是推来推去。<br>他们难道不违法吗?   |       |       |
| ..... | .....      | .....             | .....  | .....                              | .....   | ..... | ..... |
| 2     | 188<br>809 | A<br>90913<br>9   | A 市<br>万家丽南<br>路丽发新<br>城居民区<br>附近搅拌<br>站扰民                   | 2<br>019/1<br>1/19<br>18:07:<br>54 | A 市万家丽南路丽发<br>新城居民区,开发商在小区<br>旁 50 米处建搅拌站, 运渣<br>车吵得人精神崩溃,灰尘满<br>天飞, 我们该怎么办?  | 1     | 0     |
| 2     | 189<br>950 | A<br>90920<br>4   | 投诉<br>A2 区丽发<br>新城附近<br>建搅拌站<br>噪音扰民                         | 2<br>019-1<br>1-13<br>11:20:<br>21 | 我是 A2 区丽发新城小<br>区的一名业主,我要投诉同<br>发投资有限公司在未经小<br>区业主同意的情况下,在离<br>小区不到百米的地方建搅<br>拌站。可想而知,一个大型<br>搅拌站每天的噪音输出有<br>多严重! 还有扬尘污染,<br>.....  | 0     | 0     |
| 2     | 190<br>108 | A<br>90924<br>0   | 丽发<br>新城小区<br>旁边建搅<br>拌站                                     | 2<br>019-1<br>2-21<br>15:11:<br>29 | 丽发新城小区旁边建<br>的搅拌站几百米外就是小<br>学,扬尘严重影响几千名学<br>生的健康,很多业主反应家<br>里小孩各种皮肤病症状<br>.....   | 1     | 0     |
| 3     | 263<br>672 | A<br>00041<br>448 | A4 区<br>绿地海外<br>滩小区距<br>长赣高铁<br>最近只有<br>30 米不<br>到, 合理<br>吗? | 2<br>019/9/<br>5<br>13:06:<br>55   | 您好,近日看到了渝长厦高<br>铁最新的红线征地范围以<br>及走向经过,其经过北三环<br>的地方紧挨着绿地海外滩<br>小区二期,我测算了一下距<br>离,最近的位置只有 30 米<br>不到,这严重不符合我国于<br>1988 年颁布的国家标准<br>gb8702-88《电磁辐射防护<br>规定》。按照设计要求, ”<br>铁路两侧 30 米内严禁新建<br>住宅, ..... | 669   | 0     |
| ..... | .....      | .....             | .....  | .....                              | .....   | ..... | ..... |

|       |            |                   |   |                                   |  |       |       |
|-------|------------|-------------------|---|-----------------------------------|--|-------|-------|
| 4     | 205<br>771 | A<br>00020<br>115 | 夫妻<br>共同买的<br>房为何申<br>请 A 市人<br>才购房补<br>贴不通<br>过? | 2<br>019/5/<br>29<br>17:12:<br>29 | 年初申请人才购房补贴不<br>通过,原因是购房合同上必<br>须有申请人姓名。但是这套<br>房是我在读研期间和我老<br>公共同购买,当时由于在外<br>地上学,购房合同都是老公<br>一人处理, 只签了他的名<br>字.....       | 0     | 0     |
| 4     | 206<br>983 | A<br>00049<br>301 | A 市<br>人才新政<br>补贴最近<br>两个月的<br>怎么还没<br>发?         | 2<br>019/5/<br>28<br>16:18:<br>02 | A 市人才新政的补贴已经<br>两个月没发了,想知道具体<br>是什么原因,不发也没给我<br>们一个解释说明,就突然之<br>间就停了,大家都没收到。<br>诉求: 1, 具体是什么原因<br>不发? 2, 什么时间补<br>发? ..... | 8     | 0     |
| 5     | 261<br>570 | A<br>00042<br>277 | 关于<br>A 市住房<br>贷款商转<br>公问题的<br>建议                 | 2<br>019/3/<br>22<br>16:50:<br>16 | 目前 A 市的商转公, 需要<br>公民自己先借钱还清商业<br>银行贷款, 拿到产权证后,<br>再到公积金中心办理,然后<br>公积金中心再将贷款给个<br>人,让个人去还借款。为何<br>不能向其他省份一样.....            | 26    | 0     |
| ..... | .....      | .....             | .....   | .....                             | .....  | ..... | ..... |
| 5     | 265<br>551 | A<br>00076<br>292 | 咨询<br>A 市公积<br>金贷款买<br>房的问题                       | 2<br>019/5/<br>23<br>15:06:<br>05 | 去年购买了首套房,但是当<br>时刚参加工作,还未缴纳公<br>积金, 故使用的是商业贷<br>款,现在单位已经缴纳了公<br>积金,请问是否可以将商业<br>贷款转化.....                                  | 0     | 0     |

具体的见“热点问题明细表.xlsx”。

## 4 问题三：

### 4.1 问题分析

针对附件 4 中的留言详情和答复意见，对答复的相关性、完整性、可解释性，及时性进行定义，建立对答复意见的质量评价模型，从而完成对答复质量的评价。

### 4.2 解题思路

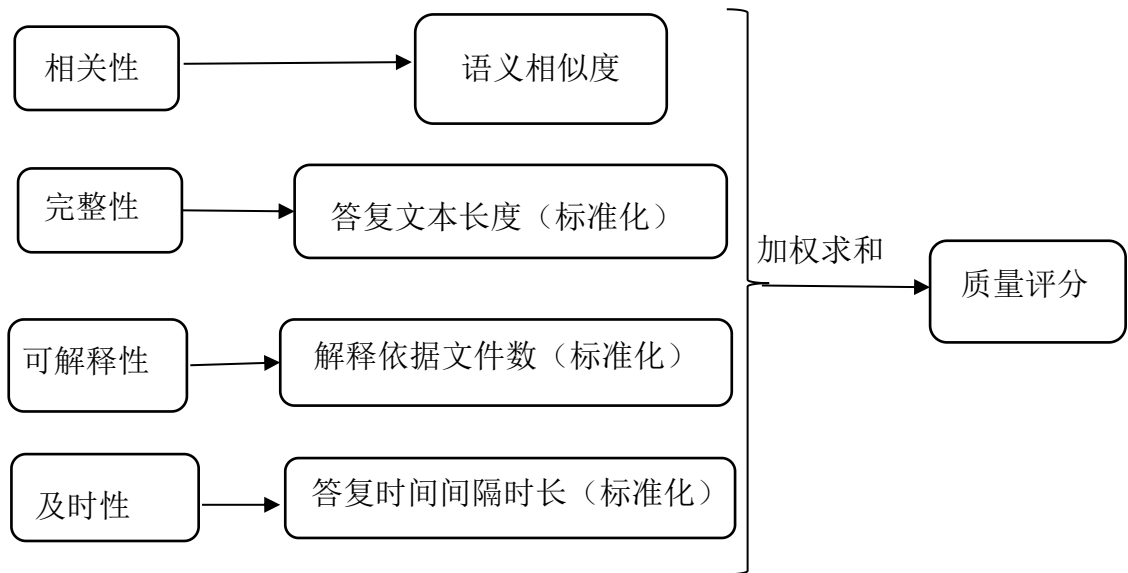


图 4-1 解题思路示意图

### 4.3 模型建立

## 5 模型的局限性

本文问题二的热度指数和问题三的质量评分中各指标的权重参数都是人为假设，难以把握一个最佳值，对模型的推广有一定的难度，需重新人工调参。

## 参考文献

- [1] DEVLIN J, CHANG M W, LEE K, et al. Bert:Pre-training of deep bidirectional transformers for language understanding[J/OL]. 2018, arXiv:1810. 04805, (2018-10-11) [2019-06-01]. <https://arxiv.org/abs/1810.04805>
- [2] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need[C]//Advances in Neural Information Processing Systems. 2017:5998-6008.
- [3] 涂铭、刘祥、刘树春. Python 自然语言处理实战核心技术与算法[M]. 机械工业出版社, 86-91
- [4] Clustering 聚类 密度聚类——DBSCAN[z]. <https://www.cnblogs.com/PJQ000/p/11838288.html>, 2020. 4. 4
- [5] 叶建成. 利用文本挖掘技术进行新闻热点关注问题分析[D]. 广州: 广州大学, 2018: 14-20
- [6] Benabderrahmane S, Mellouli N, Lamolle M. On the predictive analysis of behavioral massive job data using embedded clustering and deep recurrent neural networks[J]. Knowledge-Based Systems, 2018, 151: 95-113.
- [7] Johnson0722.doc2vec 原理及实践[z]. [https://blog.csdn.net/John\\_xyz/article/details/79208564](https://blog.csdn.net/John_xyz/article/details/79208564), 2020. 4. 23