

“智慧政务”中的文本挖掘与应用

摘要

近年来，随着网络的快速发展，微信，微博，市长信箱，阳光热线逐渐成为政府了解民意，汇集民智的重要渠道。因此运用文本分析和数据挖掘技术对提升政府管理水平和施政效率具有极大的推动作用。

对于问题 1，利用 python 读取附件 2，获取留言详情作为新的列表。运用 jieba 分词对留言详情分词，并调用停用词去掉部分分词。用 txt 文件写入各个一级标签的相关词汇，利用 python 读取留言详情并将文本向量化，计算每一条留言中包含的各个一级标签相关词汇的频数，将包含词汇频数最高的词汇所属的一级标签定义为该留言的一级标签，运用循环对每条留言都做如此处理，保存为新的列表。计算各个一级标签的查重率和查全率，利用 F-score 评价该模型。

对于问题 2，利用 python 读取附件 3，运用 jieba 分词对留言主题分词，并调用停用词去掉部分分词。调用 gensim 分析文本相似度，构建语料库，并通过 TF-IDF 算法得到前五个热点，去除热点的留言，得到排名前五的热点的位置。使用位置数组构建成新的数组 Grid 以及 Form 分别是热点问题表以及热点问题留言明细表，再将他们保存为 xlsx 文件格式。

对于问题 3，利用 python 读取附件 4，将时间统一格式为年/月/日，时：分：秒，并存入一个数组中。将元组中的前 4 个元素重组为字符串，元组元素分别为年，月，日，时，分，秒，使用正则表达式提取。计算回复与留言的时间差，自定义时效的划分等级以及对应的分数。从时间上进行评判。提取回复内容以及留言内容，建立语料库，利用 TF-IDF 进行相似度比较，得到相似度的数值，基于两种标准建立模型。

关键词：jieba 分词，TF-IDF 算法，热点问题；

Text mining and application in "intelligent government affairs"

Abstract:In recent years, with the rapid development of the Internet, WeChat, weibo, mayor's mailbox and sunshine hotline have gradually become important channels for the government to understand public opinion and gather people's wisdom. Therefore, the use of text analysis and data mining technology to improve the level of government management and governance efficiency has a great role in promoting.

For question 1, use python to read attachment 2 and get the message details as a new list. Use jieba word segmentation for message details, and call stop words to remove part of the word segmentation. With TXT file is written to the level of vocabulary related to labels, use python to read the message details and the text to quantify, calculate each message contained in each level 1 label related words frequency, will contain the highest frequency vocabulary's primary tag defined as the level of the message, using the loop to do so with every message, save for the new list. The recall and recall rates of each level tag were calculated, and the model was evaluated with f-score.

For question 2, use python to read attachment 3, use jieba word segmentation on the subject of the message, and call the stop word to remove part of the word segmentation. Gensim was called to analyze the text similarity, build the corpus, and obtain the top five hot spots through the tf-idf algorithm, remove the comments of hot spots, and obtain the position of the top five hot spots. Use the position number fabric to build a new array Grid and Form is the hot issues table and hot issues message detail table, and then save them as XLSX file format.

For question 3, use python to read attachment 4, and format the time as year/month/day, hour: minute: second, and store it in an array. The first four elements in the tuple are reorganized into strings. The tuple elements are year, month, day, hour, minute and second respectively. Calculate the time difference between the reply and the message, customize the classification of time and the corresponding score. Judge in terms of time. The reply content and message content were extracted, the corpus was established, the similarity value was obtained by using tf-idf for similarity comparison, and the model was established based on the two standards.

Key words: jieba word segmentation, tf-idf algorithm, hot issues

目录

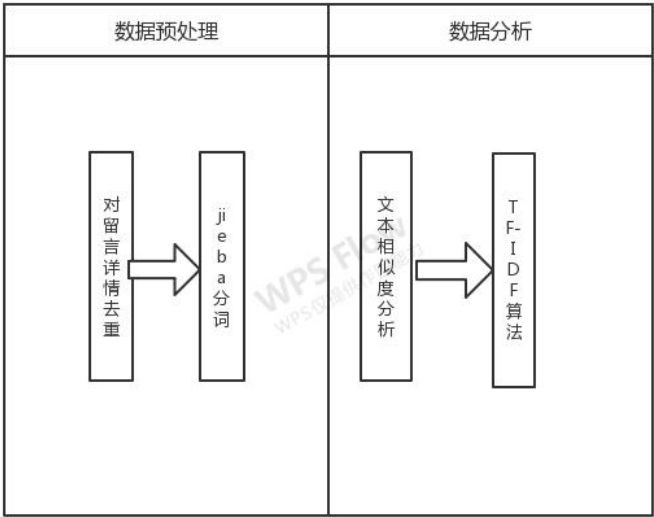
1、 挖掘目标.....	4
2、 分析方法与过程.....	4
2.1 问题一分析方法与过程.....	5
2.1.1 数据处理.....	5
2.1.2 文本空间的向量表示.....	7
2.2 问题 2 分析方法与过程.....	8
2.2.1 问题 2 流程图	8
2.2.2 对留言主题分词.....	8
2.2.3 制作语料库.....	8
2.2.4 文本相似度分析.....	9
2.2.5 TF-IDF 算法分析.....	9
2.2.6 建立热点模型.....	10
2.3 问题 3 分析过程与方法.....	11
2.3.1 问题 3 流程图	11
2.3.2 提取时间.....	11
2.3.3 及时性模型.....	11
2.3.4 相关性模型.....	12
3. 结果分析.....	14
3.1 问题 1 结果分析.....	14
3.1.1 一级标签分析.....	14
3.2 问题 2 结果分析.....	14
3.3 问题 3 结果分析.....	15
4、 结论.....	16
5. 参考文献.....	17

1、挖掘目标

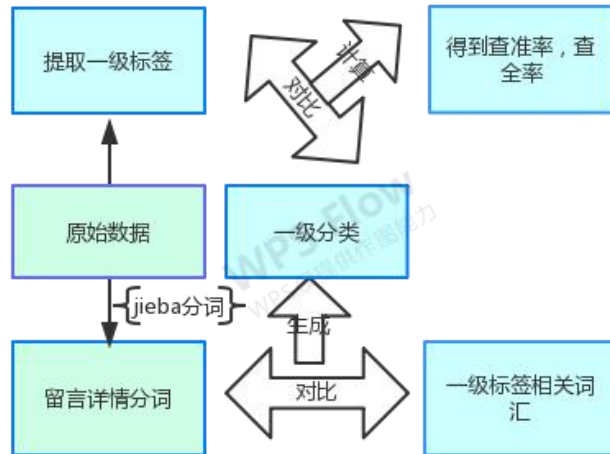
本次建模目标是利用微信微博，市长信箱，阳光热线等网络信息平台的数据，利用 jieba 中文分词工具对留言进行分词，相似度比较，命名体识别，达到以下三个目标：

- （1）利用文本分词和文本向量化进行文本挖掘，自动将留言按所属一级标签分类
- （2）根据附件 3 的数据，分析某一时段的热点问题，了解热点问题的分布。
- （3）根据附件 4 的数据，对相关部门对留言的答复情况，从答复的相关性，完整性，可解释性等角度评价回复情况。

2、分析方法与过程



2.1 问题一分析方法与过程



2.1.2 数据处理

2.1.2.1 对留言详情进行分词

在对留言信息进行数据挖掘的时候，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 2 中，以中文文本的方式给出了数据，为了方便转化，将数据进行分词。采用 python 的中文分词 jieba 进行分词。jieba 支持三种分词模式：精确模式、全模式和搜索引擎模式。精确模式：试图将语句最精确的切分，不存在冗余数据。全模式：将语句中所有可能是词的词语都切分出来，速度很快，但是存在冗余数据。搜索引擎模式：在精确模式的基础上，对长词再次进行切分[1]

下图为部分分词结果:

A3区米兴海路
 一黄中A3A2市梓县
 A7
 阳光路
 合成步行街
 法果行际泉绿路立
 合成品国麓富路坡特
 纳纳古道社改新村公与路
 税示道区变村车
 纳纳古道社改新村公与路
 芝城巷三期产意路四
 木城巷三期产意路四
 摄影
 咨询牌户期质道汇口
 住任四谷性变交路
 A6A7区县生地珠行铁高峰
 A7卫空明市通地晚
 道春粪夜小地A3凌太
 路华便间区铁区晨堵
 镇
 名鼎排施工规利议
 命金外施栋速保点建
 规划泥市噪声用麓施调
 市音空工工整
 路
 步婚A3区民质疑林民号
 初A3扰性质公台扰信
 婚A3扰性质公台扰信
 语灯

2.1.2.2 正则表达式

['A3区', '大道', '西行', '便道', '未管', '路口', '加油站', '路段', '人行道', '包括', '路灯', '杆', '圈', '西湖', '建筑', '集团', '燕子', '山', '安置', '房', '项目', '施工', '围墙', '上下班', '期间', '条', '路上', '人流', '车流', '安全隐患', '请求', '文明城市', 'A市', '整改', '文明', '路段']

['位于', '书院', '路', '主干道', '在水一方', '大厦', '一楼', '四楼', '人为', '拆除', '水', '电等', '设施', '烂尾', '多年', '护栏', '围着', '占用', '人行道', '路', '护栏', '锈迹斑斑', '倒塌', '危机', '过往行人', '车辆', '请求', '部门', '牵头']

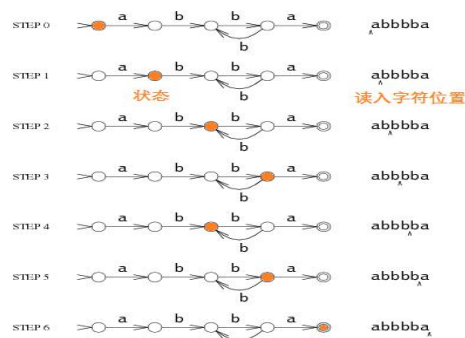
['尊敬', '领导', 'A1区', '苑', '小区', '位于', 'A1区', '火炬', '路', '小区', '物业', 'A市', '程明', '物业管理', '有限公司', '未经', '小区业主', '同意', '利用', '业主', '公摊', '公共', '面积', '业主', '滥收', '停车费', '收费', '车辆', '区内', '损坏', '拒', '承担责任', '业主', '物业', '社区', '物业', '置之不理', '妄自', '违规', '收费', '社区', '不闻不问', '今特请', '领导', '明查', '事实', '老百姓', '做主', '督办', '社区', '物业', '依法', '依规', '小区', '居民', '应享', '正当权益', '甚', '感谢', '物权法', '73', '条', '区内', '道路', '绿地', '公用设施', '物业', '服务', '用房', '业主', '共有', '74', '条', '占用', '业主', '共有', '道路', '场地', '用于', '停放', '汽车', '车位', '业主', '共有', '车位', '业主', '缴纳', '停车费', '用于', '出租', '收益', '业主', '物业管理', '条例', '54', '条', '利用', '物业', '共用', '部位', '共用', '设施', '设备', '经营', '征得', '相关', '业主', '业主大会', '物业管理', '企业', '同意', '办理', '手续', '而本', '小区', '地面', '停车位', '程明', '物业', '征得', '业主', '意见', '强制', '违规', '收费', '管控', '物业', '南门', '商业', '活动', '楼', '违规', '改造', '转租', '用于', '商业', '车位', '变更', '性质', '商户', '专用车', '位', '影响', '相关', '消防通道', '功能', '小区', '业主', '隐患', '类', '业主', '物业', '社区', '提出', '诉求', '物业', '并未', '解决', '不知', '程明', '物业', '解释', '给予', '整改', '地下', '人防', '车']

在计算一级分类相关词汇在留言详情中出现的频率，需要使用正则表达式进行计算。正则表达式的算法具体如下：

(1) 给定正则表达式，例如： $a(bb)^+a$ ，其中 $+$ 表示前面的子表达式一次或多次。所以 $abba$ 或者 $abbbba$ 都能被该表达式匹配。

(2) 读入字符串，此时一个状态转换为另一个状态。FSM 有开始和匹配两种特殊状态，分别位于头部和尾部。[2]

(3) 进行匹配，下图为进行匹配的示意图：



特别的，当状态机结束于最后一个状态则匹配成功，否则匹配失败。

2.1.3 文本的空间向量模型

由于计算机不能够直接处理文本信息，我们需要对文本进行处理，将文本表示成为计算机能够直接处理的形式，即文本数字化。文本表示也称为文本特征表达，它不仅要求能够真实准确的反映文档的内容，而且要对不同的文档具有区分能力。目前常用的文本表示模型有向量空间模型、布尔模型和概率模型等。[3] 这里使用向量空间模型。[4] 向量空间模型最早是由 Salton 和 McGill 于 20 世纪 60 年代末提出的，是目前在文本挖掘技术中最常用的表示模型。其主要思想：将每一个文本表示为向量空间的一个向量，并以每一个不同的特征项对应为向量空间中的一个维度，而每一个维的值就是对应的特征项在文本中的权重，这里的权重可以由 TF-IDF 等算法得到。向量空间模型就是将文本表示成为一个特征向量：

$$V(d) = (t_1, w_1(d), t_2, w_2(d), \dots, t_n, w_n(d))$$

其中， $t_i (i=1, 2, 3, \dots, n)$ 为文档 d 的一个特征项， $w_i (i=1, 2, \dots, n)$ 为特征值的权值。

2.1.3.1 文本的向量化表示

根据留言详情的特征项对应词袋的位置，组成统一维数的向量

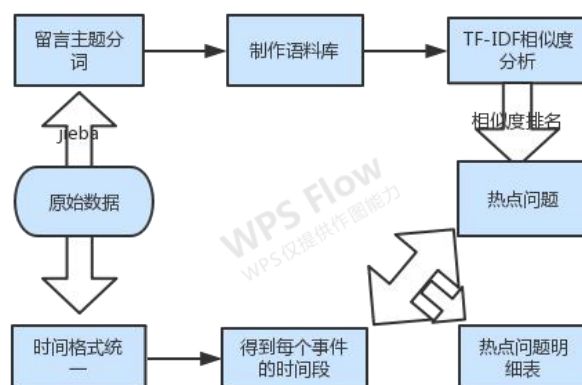
$C = (t_1, t_2, \dots, t_n)$ ，其中 C 为词袋集合， t_n 为每个词在向量对应的位置

下图为部分向量化结果：

卫生计生	[0, 2, 0, 2, 1, 1, 0]
城乡建设	[0, 2, 0, 0, 2, 0, 0]
城乡建设	[0, 0, 0, 0, 2, 0, 1]
卫生计生	[0, 0, 1, 0, 1, 0, 2]
劳动和社会保障	[0, 0, 1, 0, 2, 1, 0]
卫生计生	[0, 0, 0, 1, 1, 0, 0]
商贸旅游	[0, 0, 0, 0, 1, 0, 0]
卫生计生	[0, 0, 0, 0, 3, 0, 0]
卫生计生	[0, 0, 0, 1, 3, 0, 1]
卫生计生	[0, 1, 0, 0, 1, 0, 0]
城乡建设	[0, 1, 0, 0, 1, 0, 0]
城乡建设	[0, 1, 0, 0, 1, 0, 0]
城乡建设	[0, 0, 1, 1, 1, 1, 1]
教育文体	[0, 0, 0, 0, 1, 0, 1]
卫生计生	[0, 0, 0, 0, 0, 0, 0]

2.2 问题 2 分析方法与过程

2.2.1 问题 2 流程图



2.2.2 对留言主题分词

在对留言内容进行挖掘分析之前，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 给出的数据是由中文文本的形式给出，为了便于转换，先要对这些留言详情进行中文分词。这里采用 python 的中文分词包 jieba 进行分词。jieba 采用了基于前缀词典实现的高效词图扫描，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)，同时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，使得能更好的实现中文分词效果。

2.2.3 制作语料库

用 dictionary 方法获取词袋 (bag-of-words)，在词袋中用数字对所有词进行了编号，一个编号对应一个词。使用 doc2bow 制作语料库，语料库是一组向量，向量中的元素是一个二元组 (编号、频次数)，对应分词后的文档中的每一个词。用同样的方法，把测试文档也转换为二元组的向量。

下图为部分语料库：


```

[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1)], [(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, 1), (22, 1), (23, 1), (24, 1)], [(25, 1),
(26, 1), (27, 1), (28, 1), (29, 1), (30, 1), (31, 1), (32, 1), (33, 1)], [(0,
1), (34, 1), (35, 1), (36, 1), (37, 1), (38, 1), (39, 1), (40, 1), (41, 1),
(42, 1), (43, 1), (44, 1)], [(0, 1), (43, 1), (45, 1), (46, 1), (47, 1), (48,
1), (49, 1), (50, 1), (51, 1), (52, 1), (53, 1), (54, 2)], [(25, 1), (47, 1),
(55, 1), (56, 1), (57, 1)], [(34, 1), (58, 1), (59, 1), (60, 1), (61, 1)],
[(34, 1), (62, 1), (63, 1), (64, 1), (65, 1), (66, 1)], [(0, 1), (41, 1),
(42, 1), (54, 1), (58, 1), (67, 1), (68, 1), (69, 1), (70, 1), (71, 1), (72,
1), (73, 1), (74, 1), (75, 1)], [(16, 1), (18, 1), (64, 1), (76, 1), (77, 1),
(78, 1), (79, 1), (80, 1), (81, 1), (82, 1), (83, 1), (84, 1), (85, 1)], [(0,
1), (41, 1), (46, 1), (86, 1), (87, 1), (88, 1), (89, 1), (90, 1), (91, 1),
(92, 1), (93, 1), (94, 1), (95, 1)], [(96, 1), (97, 1), (98, 1), (99, 1),
(100, 1), (101, 1), (102, 1), (103, 1), (104, 1), (105, 1)], [(34, 1), (40,
1), (41, 1), (106, 1), (107, 1), (108, 1), (109, 1), (110, 1)], [(34, 1),
(58, 2), (111, 1), (112, 1), (113, 1), (114, 1), (115, 1), (116, 1), (117,
1)], [(46, 1), (118, 1), (119, 1), (120, 1), (121, 1), (122, 1), (123, 1),
(124, 1)], [(18, 1), (125, 1), (126, 1), (127, 1), (128, 1), (129, 1), (130,
1), (131, 1), (132, 1)], [(16, 1), (18, 1), (19, 1), (24, 1), (133, 1), (134,
1), (135, 1), (136, 1), (137, 1)], [(10, 1), (138, 1), (139, 1), (140, 1),
(141, 1)], [(34, 1), (142, 1), (143, 1), (144, 1), (145, 1), (146, 1), (147,
1), (148, 1)], [(6, 1), (39, 1), (149, 1), (150, 1), (151, 1), (152, 1),
(153, 1), (154, 1)], [(16, 1), (18, 1), (155, 1), (156, 1), (157, 1), (158,
1), (159, 1)], [(25, 1), (46, 1), (160, 1), (161, 1), (162, 1), (163, 1)],
[(34, 1), (164, 1), (165, 1), (166, 1), (167, 1), (168, 1), (169, 1), (170,
1), (171, 1), (172, 1)], [(34, 1), (173, 1), (174, 1), (175, 1), (176, 1),
(177, 1), (178, 1)], [(34, 1), (35, 1), (46, 1), (157, 1), (179, 1), (180,
1), (181, 1), (182, 1)], [(34, 1), (51, 1), (183, 1), (184, 1), (185, 1),
(186, 1), (187, 1), (188, 1)], [(34, 1), (123, 1), (189, 1), (190, 1), (191,
1), (192, 1), (193, 1), (194, 1), (195, 1), (196, 1), (197, 1), (198, 1)],
[(34, 1), (142, 1), (199, 1), (200, 1), (201, 1), (202, 1), (203, 1), (204,
1)], [(25, 1), (41, 1), (193, 1), (205, 1), (206, 1), (207, 1), (208, 1),
(209, 1), (210, 1), (211, 1)], [(16, 1), (18, 1), (123, 1), (132, 1), (212,
1), (213, 1), (214, 1), (215, 1), (216, 1), (217, 1), (218, 1)], [(16, 1),
(18, 1), (46, 1), (219, 1), (220, 1), (221, 1), (222, 1), (223, 1), (224,
1)], [(144, 1), (225, 1), (226, 1), (227, 1), (228, 1), (229, 1), (230, 1)],
[(34, 1), (41, 1), (231, 1), (232, 1), (233, 1), (234, 1), (235, 1), (236,
1), (237, 1)], [(34, 1), (110, 1), (115, 1), (238, 1), (239, 1), (240, 1),
(241, 1)], [(25, 1), (141, 1), (242, 1), (243, 1), (244, 1), (245, 1), (246,

```

2.2.4 文本相似度分析

这里利用 TF-IDF 进行文本相似度分析。获取文档中，每个词的 TF-IDF 值，对每个项目，分析测试文档的相似度，根据相似度排序。

2.2.5 TF-IDF 算法分析

TF-IDF 是一种用于信息检索的常用加权技术。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。字词的重要性随着它在文件中出现的次数成正比增加，但同时会随着它在语料库中出现的频率成反比下降。

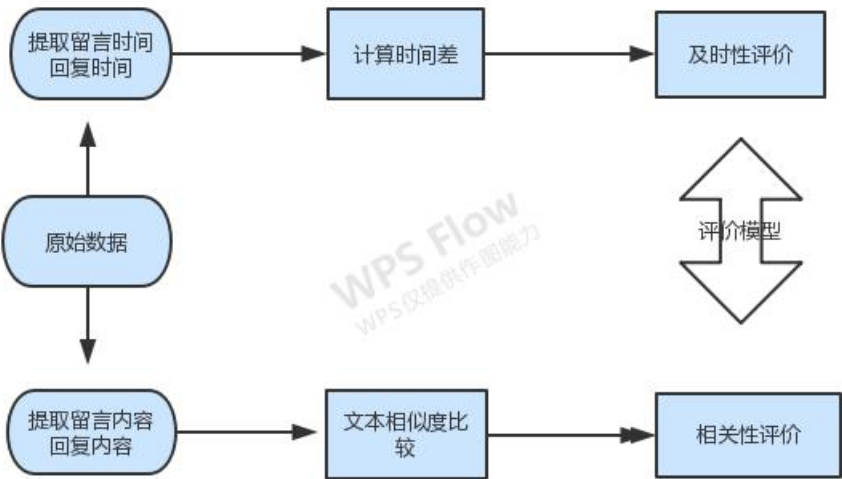
TF 表示词条（关键字）在文本中出现的频率。这个数字通常会被归一化（一般是词频除以文章总词数），以防止它偏向长的文件。公式： $tf_{ij} = \frac{n_{i,j}}{\sum_k n_{k,j}}$ 即：

通过词云图显示，出现最多的为 A 市，A7 县等多为地名，因此利用词频方式效果不好，需进一步建立热点模型

由 TF-IDF 模型得到相似度排序，计算每条相似度大于 0.3 的数目将其与该留言放入一个数组里，遍历循环附件 3，得到每条留言以及与他相似的数目和该留言的反对数和点赞数，计算权重，将相似的数目的百分之八十五加上点赞与反对数的总和的百分之十五，作为最后的热点度，按照热点度的排行从大到小生成 excel 表格。

2.3 问题 3 分析过程与方法

2.3.1 问题 3 流程图



2.3.2 提取时间

回复的是否及时有时也是评判回复的标准，利用 python 读取留言时间以及回复时间，统一时间的格式为年/月/日 时：分：秒，存入数组，将元组前 4 个元素形成字符串。使用正则查询。

2.3.3 及时性模型

计算回复与留言的时间差，自定义以 10 天为一个范围，例如当时间差为 10

'诸多', '问题', '的', '投诉'], ['投诉', 'K9', '县', '二中', '强制', '补考', '并', '收取', '补课费'], ['投诉', 'K9', '县', '二中', '重点班', '强制', '补课'], ['举报', '江', '花县', '教育局', '营养餐', '采购', '项目', '违法', '制作', '招标', '文件'], ['投诉', 'K6', '县绍基', '学校', '周末', '有偿', '补课'], ['投诉', 'K3', '县', '二中', '强制', '高三', '学生', '补课'], ['投诉', 'K9', '县', '二中', '变相', '要求', '学生', '补课'], ['投诉', 'K', '市', '明德', '楚南', '学校', '乱收费'], ['对', 'K6', '县', '小升初', '政策', '的', '质疑'], ['对', 'K7', '县', '公费', '定向生', '录取', '结果', '表示', '质疑'], ['呼吁', '政府部门', '及', '教育局', '为', 'K1', '区', '孩子', '们', '多建', '学校'], ['请求', '放宽', '入学', '界限', '!']]

利用留言详情与回复详情建立语料库，如下图所示：

```
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1), (7, 1), (8, 1), (9, 1), (10, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1), (7, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1), (7, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1), (7, 1), (8, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1), (7, 1) ], []]
[[ (0, 1), (1, 1), (2, 1), (3, 1), (4, 1), (5, 1), (6, 1), (7, 1), (8, 1) ], []]
```

利用 IF-IDF 进行相似度比较，得到对应的数值，如下图所示：

```
[0.9622505, 0.9486833, 0.7537783, 0.6904309, 0.7745967, 0.95618284,
0.82823646, 0.650814, 0.62994075, 0.5604485, 0.70710677, 0.8140999,
0.7071067, 0.0, 0.70710677, 0.7493075, 0.37796447, 0.7302968,
0.43519413, 0.35355338, 0.66943866, 0.53452253, 0.35355338, 0.4447496,
0.65465367, 0.7106691, 0.6634888, 0.57735026, 0.75592893, 0.71330416,
0.40824828, 0.6, 0.0, 0.606977, 0.31622776, 0.57735026, 0.66943866,
0.6681531, 0.42640144, 0.7893522, 0.4082483, 0.6666667, 0.59999996,
0.8838923, 0.5477226, 0.6324555, 0.8164966, 0.71713716, 0.79471934,
0.86154974, 0.8087458, 0.6459422, 0.30151135, 0.45834926, 0.60385966,
0.62994075, 0.65997505, 0.6446584, 0.808122, 0.5111657, 0.57735026,
0.7161149, 0.43495885, 0.5500914, 0.65465367, 0.0, 0.0, 0.80178374,
0.59999996, 0.35355338, 0.5, 0.0, 0.72547626, 0.69985425, 0.8333334,
0.6761234, 0.31622776, 0.7745967, 0.75384617, 0.88741195, 0.5659165,

0.66226614, 0.630126, 0.33333334, 0.57735026, 0.4472136, 0.25,
0.54494923, 0.57735026, 0.63997465, 0.4082483, 0.30151135, 0.35355338,
0.35355338, 0.7337994, 0.4082483, 0.5163978, 0.70710677, 0.37796447,
0.75793684, 0.7947194, 0.73333335, 0.57735026, 0.57735026, 0.5656854,
0.6324555, 0.48038444, 0.6172134, 0.64820373, 0.6324555, 0.65465367,
0.7167132, 0.30151135, 0.8944272, 0.5, 0.86193424, 0.88388336,
0.35355338, 0.5345225, 0.6324555, 0.70710677, 0.8944272, 0.49999997,
0.700649, 0.6459752, 0.51425946, 0.7119914, 0.6324555, 0.6847368,
0.5163978, 0.0, 0.0, 0.7715168, 0.86970985, 0.5378679, 0.6172134,
0.6363636, 0.5163978, 0.6823231, 0.70710677, 0.4767313, 0.848528,
0.8340577, 0.40509576, 0.49999997, 1.0, 0.5556624, 0.5500914,
0.5163978, 0.6981871, 0.35355338, 0.7745967, 0.49999997, 0.49999997,
0.6324555, 0.81649655, 0.5345225, 0.35355338, 0.4685213, 0.70710677,
```


3. 结果分析

3.1 问题 1 结果分析

3.1.1 一级标签分析

利用正则表达式对 jieba 分词后的留言详情与录入的一级分类的相关词汇进行匹配，将所含数目最多的关键词定义为该留言详情的一级标签。得到如“一级标签 txt”的一级标签，再将一级标签 txt 与读取并保存的一级分类 txt 进行对比，计算每一类的查准率以及查全率。如下表所示：

一级标签	查全率	查准率
城乡建设	100%	71.34%
环境保护	100%	84.77%
交通运输	100%	85.99%
教育文体	100%	86.54%
劳动与社会保障	100%	58.27%
商贸旅游	100%	42.76%
卫生计生	100%	76.99%

通过 F-Score 公式 $\frac{1}{n} \sum_{i=1}^n \frac{2P_iR_i}{P_i+R_i}$ 计算得 0.8297，其中 P_i 为第 i 类的查准率， R_i 为第 i 类的查全率。

3.2 问题 2 结果分析

通过 TF-IDF 对数据的相似度分析以及排序，得到排名前 5 的热点问题如下图所示：

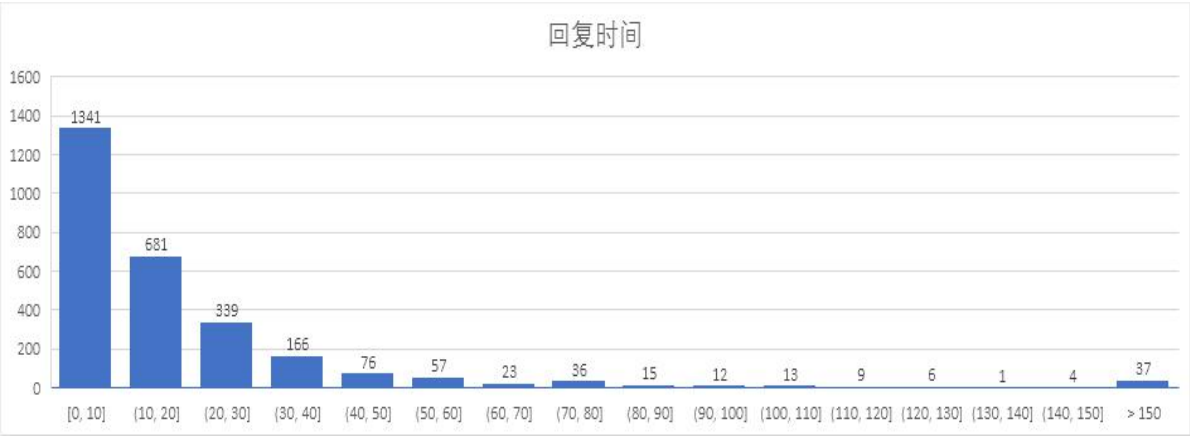
热度排名	问题ID	热度指数	时间范围	地点/人群	问题描述
1	1	54	2019/07/03至2020/01/26	新城小区新城	丽发新城小区附近揽拌站的一些问题
2	2	50	2019/07/07至2019/09/01	伊滨河伊滨河个人	伊景园滨河苑车位捆绑销售！广铁集团做个人吧！
3	3	48	2019/01/02至2020/01/06	A市	请A市加快地铁的建设
4	4	48	2018/11/15至2019/12/24	A市人才	咨询A市人才购房补贴政策
5	5	46	2019/02/27至2019/12/25	A市	询问A市住房公积金贷款的相关问题



由表格可以得到群众反映的热点问题为丽发新城小区附近搅拌站扰民占比 22%，伊景园滨河苑车位捆绑销售占比 20%，A 市地铁问题占比 20%，A 市购房补贴问题占比 20%，A 市住房公积金占比 19%，且每个热点问题所占的比例近似相等，由此可以得到百姓所反映的问题围绕生活，出行，住房。

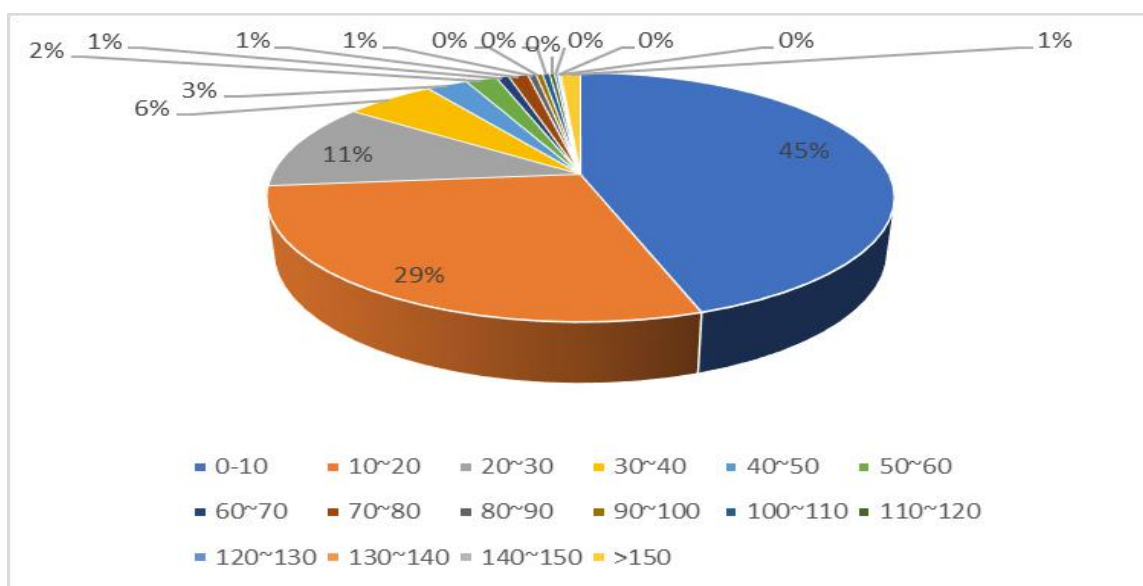
3.3 问题 3 结果分析

通过对每条留言的及时性分析，我们绘制了频数直方图：



整理成表格得到
绘制饼图得到

0-10	10-20	20-30	30-40	40-50	50-60	60-70	70-80	80-90	90-100	100-110	110-120	120-130	130-140	140-150	>150
1341	681	339	166	76	57	23	36	15	12	13	9	6	1	4	37



由以上数据可以得到，回复时间在 10 天之内的占比最多达到百分之四十五，随着时间的增加，数目也在逐渐减少，而当回复时间超过 140 天时，数目开始增多

通过综合考虑相关性以及及时性得到

该模型的得分为(100分制度)： 64.02172372189604
该模型评价等级为： 合格

也就是说市长的回复较为及时，与留言内容具有一定相关性

4、结论

对于留言信息进行分析研究，对提升政府管理水平和施政效率具有极大的推动作用。本文采用 jieba 中文分词以及 TF-IDF 对数据进行研究。统计目前市民反映最多的问题，同时定义及时性以及相关性评分标准，对市长的回复进行评价。

由分析结果可得，目前市民最关注的问题为丽发新城小区附近搅拌站扰民，伊景园滨河苑车位捆绑销售，A 市地铁问题，A 市购房补贴问题，A 市住房公积金，围绕着市民的起居生活，交通出行。

市长回复消息的时间在 10 天之内居多，占比百分之四十五，但是也有超过 150 的天，占比百分之一，因此可以说市长回复消息较为及时。根据自定义的分数，综合考察及时性以及相关性，最后得分 64 分。因此，市长的回复与留言内容具有一定相关性。也就是说，市民反映的问题基本得到回复解决。

参考文献:

- [1] 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(4).
- [2] [2] 详解正则表达式匹配算法原理_网络_ybdesire 的专栏-CSDN 博客
<https://blog.csdn.net/ybdesire/article/details/78255427>
- [3] 周昭涛. 文本聚类分析效果评价及文本表示研究[D]. 中国科学院研究生院(计算技术研究所), 2005.
- [4] 胡学钢, 董学春, 谢飞. 基于词向量空间模型的中文文本分类方法[J]. 合肥工业大学学报:自然科学版, 2007, 30(10):1261-1264.
- [5] TF-IDF 算法介绍及实现_人工智能_Asia-Lee 的博客-CSDN 博客
https://blog.csdn.net/asialeee_bird/article/details/81486700