

该项目一共分四步走

1.首先进行数据处理

数据处理，我们用到的是pandas常用库

先进行去重drop_duplicates，把重复的内容处理好之后，我们开始删除一些无效的内容，例如表情包，无效词等

这里首先就是先去掉表情包，然后再判断该文本是否为中文，接着再去用停用词文本，去除无效词

```
def emoji_tihuan(x):
    x1 = str(x)
    x2 = re.sub('(\\.|.+?\\))', '', x1)
    x3 = re.sub(r'@[\w\u2E80-\u9FFF]+?:?/[\\w+\\_]', '', x2)
    x4 = re.sub(r'\n', '', x3)
    return x4

def is_all_chinese(strs):
    for _char in strs:
        if not '\u4e00' <= _char <= '\u9fa5':
            return False
    return True

def get_cut_words(content_series):
    # 读入停用词表
    stop_words = ["请问", "真的", "支持", "不用", "昨天", "只能", "感觉", "谢谢", "安排", "你好", "告诉", "好像", "带你去", "这是", "一点", "明天", "刚刚", "一年", "东西", "不到", "我刚", "不上"]

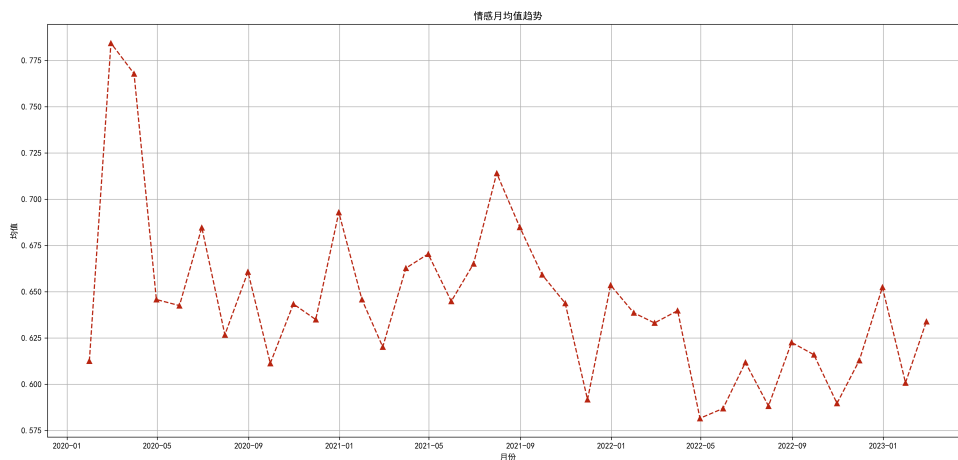
    with open("stopwords.cn.txt", 'r', encoding='utf-8') as f:
        lines = f.readlines()
        for line in lines:
            stop_words.append(line.strip())
```

接着我们stylecloud进行词云图，这样方便我们看看整体的分词效果如何，是否有一些词要不要去掉

最后的结果如下：

	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
	序号	来源	封面昵称	封面号	UID	sec_uid	IP归属	认证	点赞	回复	评论	评论时间	情感分数												
2	0	1	https://www.see1.com/	看见熊类	3.42E+09	6.77E+14	WS641_JaB6	厦门市唐某	42735	2226	在价格上, 2020-04-2	0.75207													
3	1	2	https://www.see1.com/	看见熊类	2.79E+09	3.81E+15	WS641_JaB6		7494	763	爱看美女的 2020-04-0	0.773564													
4	2	3	https://www.see1.com/	看见熊类	1.11E+09	7.39E+10	WS641_JaB6		5722	585	爱看美女的 2020-04-0	0.296296													
5	3	4	https://www.see1.com/	伍月	3.81E+08	8.22E+10	WS641_JaB6		4058	221	河南南阳的 2020-04-0	0.033333													
6	4	5	https://www.see1.com/	繁华落尽	1.24E+09	1E+11	WS641_JaB6		2983	82	云南丽江的 2020-04-0	0.545455													
7	5	6	https://www.see1.com/	人生如梦	2.19E+09	6.2E+10	WS641_JaB6		1108	41	海南有吗 2020-04-0	0.555556													
8	6	7	https://www.see1.com/	用户496077	5.82E+08	2.22E+10	WS641_JaB6		592	24	河南大部市的 2020-04-0	0.490384													
9	7	8	https://www.see1.com/	熊类熊类	1.1E+09	4.01E+11	WS641_JaB6		60	17	福建莆田的 2020-04-0	0.517575													
10	11	12	https://www.see1.com/	骨科骨科科	91318018	8.52E+10	WS641_JaB6		77	0	不不错	2020-04-0	0.863252												
11	12	13	https://www.see1.com/	济南济南市	2.84E+09	4.09E+15	WS641_JaB6	济南唐唐唐	90	4	我就想知道 2020-04-0	0.908795													
12	15	16	https://www.see1.com/	素颜熊	6.17E+08	5.97E+10	WS641_JaB6		1	1	请问熊们 2021-09-0	0.912378													
13	16	17	https://www.see1.com/	怎么又胖了	3.44E+09	3.61E+15	WS641_JaB6		1	0	怎么胖了 2021-11-0	0.333333													
14	18	19	https://www.see1.com/	怎么又胖了	5.85E+08	9.55E+10	WS641_JaB6		0	15	怎么胖了 2022-03-0	0.333333													
15	20	21	https://www.see1.com/	网东	79625603	7.01E+10	WS641_JaB6		0	0	怎么胖了 2021-12-0	0.333333													
16	21	22	https://www.see1.com/	心静安然	1.9E+09	1.09E+11	WS641_JaB6		0	0	这个怎么胖 2021-11-0	0.333333													
17	22	23	https://www.see1.com/	114374821	1.14E+08	7.43E+10	WS641_JaB6		0	0	厦门那个 2021-11-0	0.543982													
18	23	24	https://www.see1.com/	余生安好	3.44E+09	9.23E+14	WS641_JaB6		0	0	厦门那个 2021-11-0	0.543982													
19	26	27	https://www.see1.com/	山北	76811150	6.96E+10	WS641_JaB6		0	0	已打, 己 2021-10-0	0.890175													
20	27	28	https://www.see1.com/	君不见	6.25E+08	8.55E+10	WS641_JaB6		0	0	在哪里胖了 2021-09-0	0.333333													
21	28	29	https://www.see1.com/	小李好物种	5.97E+08	9.55E+10	WS641_JaB6		0	0	打开抖音看 2021-08-0	0.370378													
22	29	30	https://www.see1.com/	小流量高手	1.07E+08	7.35E+10	WS641_JaB6		0	0	谷说医院呀 2021-08-0	0.417629													
23	30	31	https://www.see1.com/	易97啊	1.78E+09	1.08E+11	WS641_JaB6		1	0	猪草呀 2021-07-0	0.588860													
24	31	32	https://www.see1.com/	奥旦熊	1.19E+09	1.02E+11	WS641_JaB6		0	0	西有呀 2021-07-0	0.75													
25	32	33	https://www.see1.com/	熊类熊类	1.1E+09	7.21E+10	WS641_JaB6		0	0	凡尔 2021-07-0	0.77146													
26	33	34	https://www.see1.com/	柚子茶	9.73E+08	1.4E+10	WS641_JaB6		0	0	怎么胖了呀 2021-07-0	0.333333													
27	34	35	https://www.see1.com/	曹诗伟及儿	1.23E+09	1.02E+11	WS641_JaB6		1	0	有九价的呀 2021-07-0	0.840619													
28	35	36	https://www.see1.com/	小辣辣	8.95E+08	6.15E+10	WS641_JaB6		1	0	这个怎么胖 2021-06-10	0.333333													
29	36	37	https://www.see1.com/	Y	2.8E+09	2.21E+15	WS641_JaB6		0	0	河南新蔡县的 2021-06-10	0.707992													
30	37	38	https://www.see1.com/	熊类熊类	2.7E+08	6.18E+10	WS641_JaB6		0	0	在哪里胖 2021-06-10	0.420101													
31	38	39	https://www.see1.com/	梦婷	2.7E+08	8.81E+10	WS641_JaB6		2	0	郑州怎么胖 2021-05-10	0.380952													
32	39	40	https://www.see1.com/	Z	1.58E+09	1.46E+10	WS641_JaB6		0	1	厦门那个 2021-04-0	0.793983													
33	40	41	https://www.see1.com/	旧人`	3.82E+09	1.83E+15	WS641_JaB6		0	0	广西南宁市的 2021-03-0	0.635404													
34	41	42	https://www.see1.com/	叶子`	3.81E+08	9.39E+10	WS641_JaB6		0	0	廊坊这边 2021-03-0	0.656384													
35	42	43	https://www.see1.com/	平静的心	1.08E+09	8.54E+10	WS641_JaB6		0	0	怎么胖了呀 2021-03-0	0.676389													
36	43	44	https://www.see1.com/		3.98E+09	1E+11	WS641_JaB6		0	0	位置在哪里呀 2021-03-0	0.656384													
37	44	45	https://www.see1.com/	啾啾熊	8.73E+08	9.77E+10	WS641_JaB6		0	0	这个位置 2021-03-0	0.656384													
38	45	46	https://www.see1.com/	浩浩	1.99E+08	7.69E+10	WS641_JaB6		0	0	不错不错, 2021-03-0	0.994863													
39	47	48	https://www.see1.com/	木西	3.73E+09	1.93E+15	WS641_JaB6		0	0	那个熊大呀 2021-03-0	0.822518													
40	48	49	https://www.see1.com/	奶里	2.37E+09	8.92E+14	WS641_JaB6		0	0	这谁接龙 2021-03-0	0.917102													

接着我们获取到对应的分值之后，我们可以根据数据来做一个时间趋势图，从而得知，在每个月的一个分值走向，这里才去的是均值处理，把每个月的所有分值相加求平均值，所以这里还是有一定的参考价值，可以作为正确的评判标准，从这里我们可以得知，在2020-01：2020-04,情绪都较为积极比较偏向，这里分值是从0-1直接的，接近0则是负面，接近1则是正面



3、TF-IDF计算

1.TF-IDF计算

TF-IDF (Term Frequency-InversDocument Frequency) 是一种常用于信息处理和数据挖掘的加权技术。该技术采用一种统计方法，根据字词的在文本中出现的次数和在整个语料中出现的文档频率来计算一个字词在整个语料中的重要程度。它的优点是能过滤掉一些常见的却无关紧要本的词语，同时保留影响整个文本的重要字词。计算方法如下面公式所示：

$$TF-IDF = TF * IDF$$

TF (Term Frequency) 表示某个关键词在整篇文章中出现的频率。IDF (InversDocument Frequency) 表示计算倒文本频率。文本频率是指某个关键词在整个语料所有文章中出现的次数。倒文档频率又称为逆文档频率，它是文档频率的倒数，主要用于降低所有文档中一些常见但对文档影响不大的词语的作用。

主要实现代码在这一块

```
# 将文本中的词语转换为词频矩阵 矩阵元素a[i][j] 表示j词在i类文本下的词频
vectorizer = CountVectorizer()

# 该类会统计每个词语的tf-idf权值
transformer = TfidfTransformer()

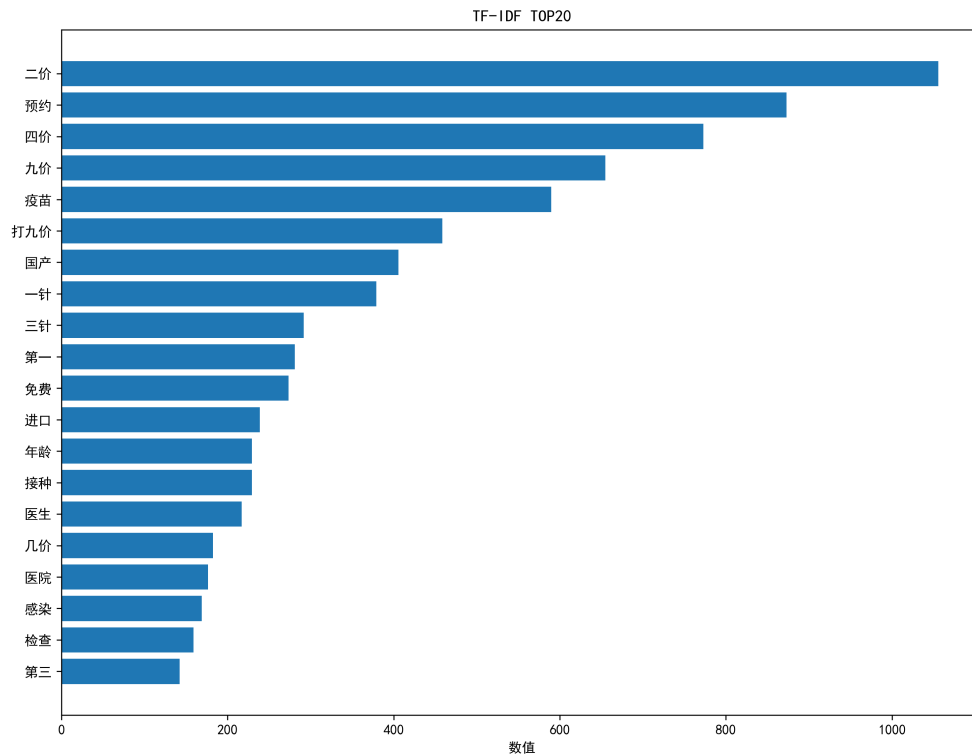
# 第一个fit_transform是计算tf-idf 第二个fit_transform是将文本转为词频矩阵
tfidf = transformer.fit_transform(vectorizer.fit_transform(corpus))
# 获取词袋模型中的所有词语
word = vectorizer.get_feature_names()

# 将tf-idf矩阵抽取出来 元素w[i][j]表示j词在i类文本中的tf-idf权重
weight = tfidf.toarray()

data = {'word': word,
        'tfidf': weight.sum(axis=0).tolist()}
df2 = pd.DataFrame(data)
```

```
df2['tfidf'] = df2['tfidf'].astype('float64')
df2 = df2.sort_values(by=['tfidf'], ascending=False)
```

我们从该文档中，抽取前20的词，会发现这里的词和高频词是一样的，但是它们的权重却是不一样的，这就是高频词和TF-IDF之间的差异了



4、LDA主题建模

LDA参考文献: <https://zhuanlan.zhihu.com/p/75222819>

<https://zhuanlan.zhihu.com/p/76636216>

这里我们的步骤和上面和一样

首先我们还是先分词，把无效词给去掉

接着我们开始构造主题数，寻找最优主题数，这里采用困惑度严格来说，判断标准并不合适，基于此我们这里采用的是另一种方式，也就是通过各个主题间的余弦相似度来衡量主题间的相似程度

(2) 寻找最优主题数

基于相似度的自适应最优LDA模型选择方法，确定主题数并进行主题分析。实验证明该方法可以在不需要人工调试主题数目的情况下，用相对少的迭代，找到最优的主题结构。具体步骤如下。

- ① 取初始主题数k值，得到初始模型，计算各主题之间的相似度（平均余弦距离）。
- ② 增加或减少k值，重新训练模型，再次计算各主题之间的相似度。
- ③ 重复步骤②直到得到最优k值。

利用各主题间的余弦相似度来度量主题间的相似程度。从词频入手，计算它们的相似度，用词越相似，则内容越相近。假定A和B是两个n维向量，A是，B是，则A与B的夹角 θ 的余弦值通过式（4）计算。

$$P(w_j | d_j) = \sum_{s=1}^K P(w_j | z = s) \times P(z = s | d_j) \quad (3)$$

$$\cos\theta = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \sqrt{\sum_{i=1}^n (B_i)^2}} = \frac{AB}{|AB|} \quad (4)$$

使用LDA主题模型，找出不同主题数下的主题词；每个模型各取出若干个主题词（比如前100个），合并成一个集合；生成任何两个主题间的词频向量；计算两个向量的余弦相似度，值越大就表示越相似；计算个主题数的平均余弦相似度，寻找最优主题数，如以下代码清单所示。

具体代码实现方式

```
# 构造主题数寻优函数
def cos(vector1, vector2): # 余弦相似度函数
    dot_product = 0.0
    normA = 0.0
    normB = 0.0
    for a, b in zip(vector1, vector2):
        dot_product += a * b
        normA += a ** 2
        normB += b ** 2
    if normA == 0.0 or normB == 0.0:
        return (None)
    else:
        return (dot_product / ((normA * normB) ** 0.5))

# 主题数寻优

def lda_k(x_corpus, x_dict):
    # 初始化平均余弦相似度
    mean_similarity = []
    mean_similarity.append(1)

    # 循环生成主题并计算主题间相似度
    for i in np.arange(2, 11):
        lda = models.LdaModel(x_corpus, num_topics=i, id2word=x_dict) # LDA
        # 模型训练

        for j in np.arange(i):
            term = lda.show_topics(num_words=30)

        # 提取各主题词
        top_word = []
```

```

for k in np.arange(i):

    top_word.append([''.join(re.findall('"(.*?)"', i)) \
                     for i in term[k][1].split('+')]) # 列出所有词

# 构造词频向量
word = sum(top_word, []) # 列出所有的词
unique_word = set(word) # 去除重复的词

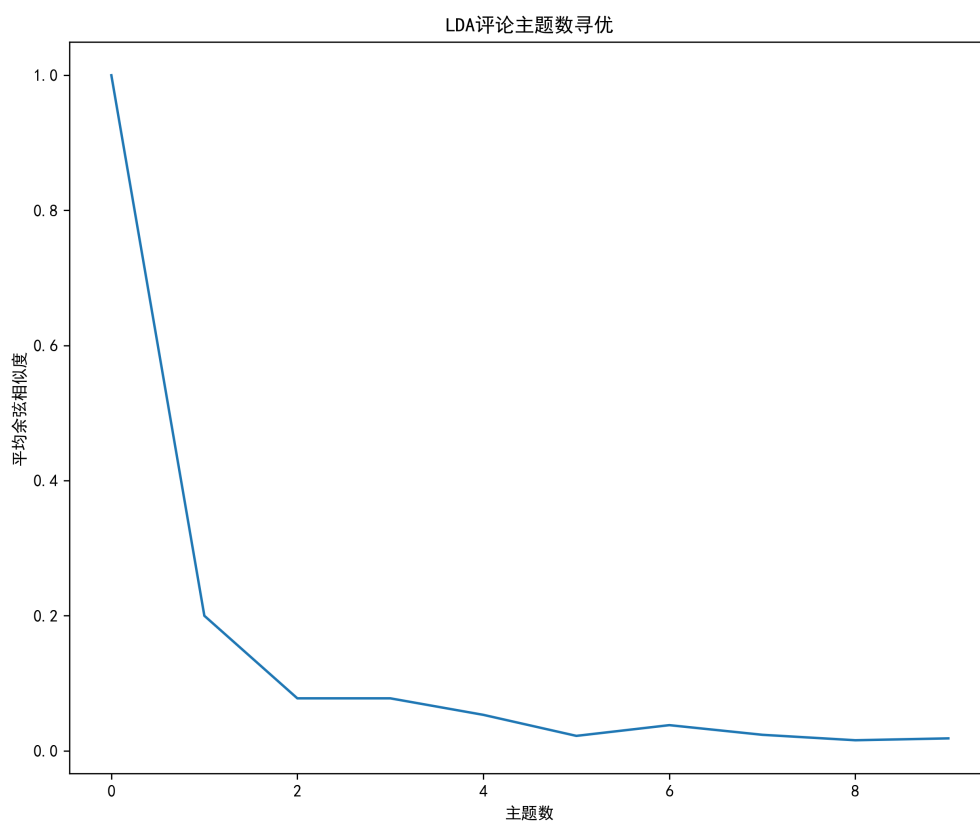
# 构造主题词列表，行表示主题号，列表示各主题词
mat = []
for j in np.arange(i):
    top_w = top_word[j]
    mat.append(tuple([top_w.count(k) for k in unique_word]))

p = list(itertools.permutations(list(np.arange(i)), 2))
l = len(p)
top_similarity = [0]
for w in np.arange(l):
    vector1 = mat[p[w][0]]
    vector2 = mat[p[w][1]]
    top_similarity.append(cos(vector1, vector2))

# 计算平均余弦相似度
mean_similarity.append(sum(top_similarity) / l)
return (mean_similarity)

```

处理好之后，再通过matplotlib来进行作图，在这里明显在5的时候处于谷底，因此根据这一点，我们可以选取主题数为5



从而进行建模，最后呈现的效果图如下：

可以通过点击不同的圆圈，来查看不同主题下，不同主题词的权重

