

# 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+\\u]', '', 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 = []

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

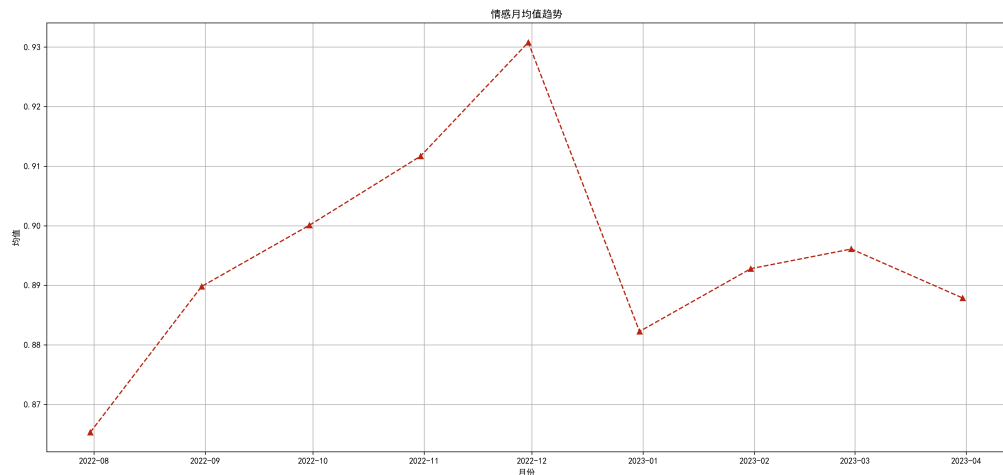
接着我们stylecloud进行词云图，这样方便我们看看整体的分词效果如何，是否有一些词要不要去掉

最后的结果如下：

就是snownlp毕竟是广义的，无法做到精准判断该文本的正确评分，只能给出大概的数值，这也是所有机器学习的通病，毕竟是机器，NLP还是主观占比居多，所以提供的数值仅供参考

	A	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
1	序号	链接	用户	ID	昵称	等级	评分	环境	论题	图片数	评论时间	人均	浏览数	点赞量	回复数	评论内容	情感分数	情感类型								
2	1	<a href="https://www.guoguo.com/2022-03-16/17561441.html">https://www.guoguo.com/2022-03-16/17561441.html</a>	故宫博物院	8.67E+08	奥兹国小凉	v6	5		18	2023-03-16 17:56:14	4	0	0	0	0	看故宫高	1	pos								
3	2	<a href="https://www.guoguo.com/2022-03-16/17561442.html">https://www.guoguo.com/2022-03-16/17561442.html</a>	故宫博物院	1.41E+09	dpuar	v5	5		6	2023-03-16 17:54:01	1	0	0	0	0	故宫博物院	0.999788	pos								
4	3	<a href="https://www.guoguo.com/2022-03-16/17561443.html">https://www.guoguo.com/2022-03-16/17561443.html</a>	故宫博物院	1.53E+08	爱吃玩	v5	5		19	2023-03-16 15:49:01	53	0	0	0	0	虽然不是	0.999999	pos								
5	4	<a href="https://www.guoguo.com/2022-03-16/17561444.html">https://www.guoguo.com/2022-03-16/17561444.html</a>	故宫博物院	4.88E+08	周四是去月	v5	5		19	2023-03-16 15:49:01	53	0	0	0	0	谁来故宫	0.999974	pos								
6	5	<a href="https://www.guoguo.com/2022-03-16/17561445.html">https://www.guoguo.com/2022-03-16/17561445.html</a>	故宫博物院	1.58E+09	大故	v1	5		0	2023-03-16 14:53:12	15	0	0	0	0	带我弟弟	0.999338	pos								
7	6	<a href="https://www.guoguo.com/2022-03-16/17561446.html">https://www.guoguo.com/2022-03-16/17561446.html</a>	故宫博物院	-1	匿名用户	v6	5		0	2023-03-16 14:53:12	8	0	0	0	0	占地面积	0.999998	pos								
8	7	<a href="https://www.guoguo.com/2022-03-16/17561447.html">https://www.guoguo.com/2022-03-16/17561447.html</a>	故宫博物院	2.84E+09	全国最采	v1	5		2	2023-03-16 14:30:11	2	0	0	0	0	寄波阿哥	0.999154	pos								
9	8	<a href="https://www.guoguo.com/2022-03-16/17561448.html">https://www.guoguo.com/2022-03-16/17561448.html</a>	故宫博物院	35914421	SarimaZe	v5	5		1	2023-03-16 14:30:11	16	0	0	0	0	京波阿哥	0.9969	pos								
10	9	<a href="https://www.guoguo.com/2022-03-16/17561449.html">https://www.guoguo.com/2022-03-16/17561449.html</a>	故宫博物院	1.81E+08	已候周士	v1	5		9	2023-03-16 14:06:11	11	0	0	0	0	一天根本	0.920265	neg								
11	10	<a href="https://www.guoguo.com/2022-03-16/17561450.html">https://www.guoguo.com/2022-03-16/17561450.html</a>	故宫博物院	-1	匿名用户	v4	5		4	2023-03-16 12:55:11	1	0	0	0	0	40块钱的	0.951396	pos								
12	11	<a href="https://www.guoguo.com/2022-03-16/17561451.html">https://www.guoguo.com/2022-03-16/17561451.html</a>	故宫博物院	7.9E+08	瓜子牙	v5	4.5		10	2023-03-16 12:49:11	10	0	0	0	0	淡李游故	1	pos								
13	12	<a href="https://www.guoguo.com/2022-03-16/17561452.html">https://www.guoguo.com/2022-03-16/17561452.html</a>	故宫博物院	1.58E+09	小幸运	v4	5		2	2023-03-16 12:39:11	4	0	0	0	0	和朕去一	0.998123	pos								
14	13	<a href="https://www.guoguo.com/2022-03-16/17561453.html">https://www.guoguo.com/2022-03-16/17561453.html</a>	故宫博物院	-1	匿名用户	v2	5		0	2023-03-16 12:11:11	3	0	0	0	0	双双导游	0.991485	pos								
15	14	<a href="https://www.guoguo.com/2022-03-16/17561454.html">https://www.guoguo.com/2022-03-16/17561454.html</a>	故宫博物院	74096280	后曹博	v6	5		8	2023-03-16 11:58:11	5	0	0	0	0	天家门就	1	pos								
16	15	<a href="https://www.guoguo.com/2022-03-16/17561455.html">https://www.guoguo.com/2022-03-16/17561455.html</a>	故宫博物院	8.4E+08	D.Mr.	v1	5		3	2023-03-16 11:45:11	22	0	0	0	0	导游小郭	1	pos								
17	16	<a href="https://www.guoguo.com/2022-03-16/17561456.html">https://www.guoguo.com/2022-03-16/17561456.html</a>	故宫博物院	2128966	是个吊吊	v8	4.5		10	2023-03-16 11:02:11	236	0	0	0	0	故宫真的	1	pos								
18	17	<a href="https://www.guoguo.com/2022-03-16/17561457.html">https://www.guoguo.com/2022-03-16/17561457.html</a>	故宫博物院	5.51E+08	Piggy小	v7	5		13	2023-03-16 09:16:11	88	5	5	5	5	魏紫紫紫	1	pos								
19	18	<a href="https://www.guoguo.com/2022-03-16/17561458.html">https://www.guoguo.com/2022-03-16/17561458.html</a>	故宫博物院	7.64E+08	草原_7588	v4	5		9	2023-03-16 08:52:11	173	0	0	0	0	上次去故	0.999917	pos								
20	19	<a href="https://www.guoguo.com/2022-03-16/17561459.html">https://www.guoguo.com/2022-03-16/17561459.html</a>	故宫博物院	1.12E+09	Jas和正	v5	5		9	2023-03-16 01:04:11	21	1	0	0	0	喜欢中国	1	pos								
21	20	<a href="https://www.guoguo.com/2022-03-16/17561460.html">https://www.guoguo.com/2022-03-16/17561460.html</a>	故宫博物院	401949	Angelchar	v2	5		4	2023-03-15 21:19:11	72	1	1	1	1	一个有趣	0.99084	pos								
22	21	<a href="https://www.guoguo.com/2022-03-16/17561461.html">https://www.guoguo.com/2022-03-16/17561461.html</a>	故宫博物院	1.55E+08	超沐沐	v7	5		13	2023-03-15 20:20:11	89	0	0	0	0	这个还用	0.999999	pos								
23	22	<a href="https://www.guoguo.com/2022-03-16/17561462.html">https://www.guoguo.com/2022-03-16/17561462.html</a>	故宫博物院	13926099	飞若帝	v8	5		9	2023-03-15 20:12:11	41	0	0	0	0	已经几次	0.999998	pos								
24	23	<a href="https://www.guoguo.com/2022-03-16/17561463.html">https://www.guoguo.com/2022-03-16/17561463.html</a>	故宫博物院	25645024	阔月	v5	4.5		15	2023-03-15 19:31:11	41	1	1	1	1	哈哈哈本	0.995993	pos								
25	24	<a href="https://www.guoguo.com/2022-03-16/17561464.html">https://www.guoguo.com/2022-03-16/17561464.html</a>	故宫博物院	1.1E+08	松子儿	v4	5		6	2023-03-15 19:31:11	175	2	2	2	2	这个地方	0.997969	pos								
26	25	<a href="https://www.guoguo.com/2022-03-16/17561465.html">https://www.guoguo.com/2022-03-16/17561465.html</a>	故宫博物院	23399809	圆圆和	v5	4		4	2023-03-15 19:24:11	76	0	0	0	0	我们是	0.999939	pos								
27	26	<a href="https://www.guoguo.com/2022-03-16/17561466.html">https://www.guoguo.com/2022-03-16/17561466.html</a>	故宫博物院	2.8E+09	富家子	v2	5		13	2023-03-15 17:23:11	37	0	0	0	0	真的很好	1	pos								
28	27	<a href="https://www.guoguo.com/2022-03-16/17561467.html">https://www.guoguo.com/2022-03-16/17561467.html</a>	故宫博物院	1.45E+09	吃嘎玩	v1	5		5	2023-03-15 16:39:11	42	1	0	0	0	在没进去	1	pos								
29	28	<a href="https://www.guoguo.com/2022-03-16/17561468.html">https://www.guoguo.com/2022-03-16/17561468.html</a>	故宫博物院	7.63E+08	晚霜	v6	0.5		4	2023-03-15 16:09:11	20	0	0	0	0	给出来的	2.13E-06	neg								
30	29	<a href="https://www.guoguo.com/2022-03-16/17561469.html">https://www.guoguo.com/2022-03-16/17561469.html</a>	故宫博物院	1.53E+09	Irislab	v5	5		4	2023-03-15 16:08:11	54	4	0	0	0	故宫真的	0.9981	pos								
31	30	<a href="https://www.guoguo.com/2022-03-16/17561470.html">https://www.guoguo.com/2022-03-16/17561470.html</a>	故宫博物院	9.23E+08	KLXNAA	v4	5		6	2023-03-15 15:16:11	78	0	0	0	0	墙头上	1	pos								
32	31	<a href="https://www.guoguo.com/2022-03-16/17561471.html">https://www.guoguo.com/2022-03-16/17561471.html</a>	故宫博物院	84216809	象朝星	v5	4.5		1	2023-03-15 15:14:11	51	0	0	0	0	去年带孩	0.986747	pos								
33	32	<a href="https://www.guoguo.com/2022-03-16/17561472.html">https://www.guoguo.com/2022-03-16/17561472.html</a>	故宫博物院	1.51E+09	从小就吃	v4	5		4	2023-03-15 13:00:11	74	0	0	0	0	从午门进	1	pos								
34	33	<a href="https://www.guoguo.com/2022-03-16/17561473.html">https://www.guoguo.com/2022-03-16/17561473.html</a>	故宫博物院	2.43E+09	Tenaru	v5	4.5		6	2023-03-15 12:43:11	54	0	0	0	0	文通坐地	0.997328	pos								
35	34	<a href="https://www.guoguo.com/2022-03-16/17561474.html">https://www.guoguo.com/2022-03-16/17561474.html</a>	故宫博物院	3.21E+09	名小	v5	5		3	2023-03-15 11:58:11	48	1	0	0	0	故宫真的	1	pos								
36	35	<a href="https://www.guoguo.com/2022-03-16/17561475.html">https://www.guoguo.com/2022-03-16/17561475.html</a>	故宫博物院	1.98E+09	群群群	v7	5		1	2023-03-15 11:30:11	32	0	0	0	0	讲解非常	0.998151	pos								
37	36	<a href="https://www.guoguo.com/2022-03-16/17561476.html">https://www.guoguo.com/2022-03-16/17561476.html</a>	故宫博物院	1.83E+09	小明	v7	5		19	2023-03-15 09:41:11	63	0	0	0	0	还是去年	0.960972	pos								
38	37	<a href="https://www.guoguo.com/2022-03-16/17561477.html">https://www.guoguo.com/2022-03-16/17561477.html</a>	故宫博物院	1.33E+08	Cynthia	v5	4.5		9	2023-03-15 08:48:11	45	0	0	0	0	历史的博	1	pos								
39	38	<a href="https://www.guoguo.com/2022-03-16/17561478.html">https://www.guoguo.com/2022-03-16/17561478.html</a>	故宫博物院	-1	匿名用户	v5	5		5	2023-03-15 07:21:11	18	0	0	0	0	香花并	1	pos								
40	39	<a href="https://www.guoguo.com/2022-03-16/17561479.html">https://www.guoguo.com/2022-03-16/17561479.html</a>	故宫博物院	7.65E+08	一级吃	v7	4.5		4	2023-03-15 02:19:11	88	18	18	18	18	2 香风十里	1	pos								

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，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、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一样

首先我们还是先分词，把无效词给除去掉

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 $\theta$ 的余弦值通过式（4）计算。

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

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

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

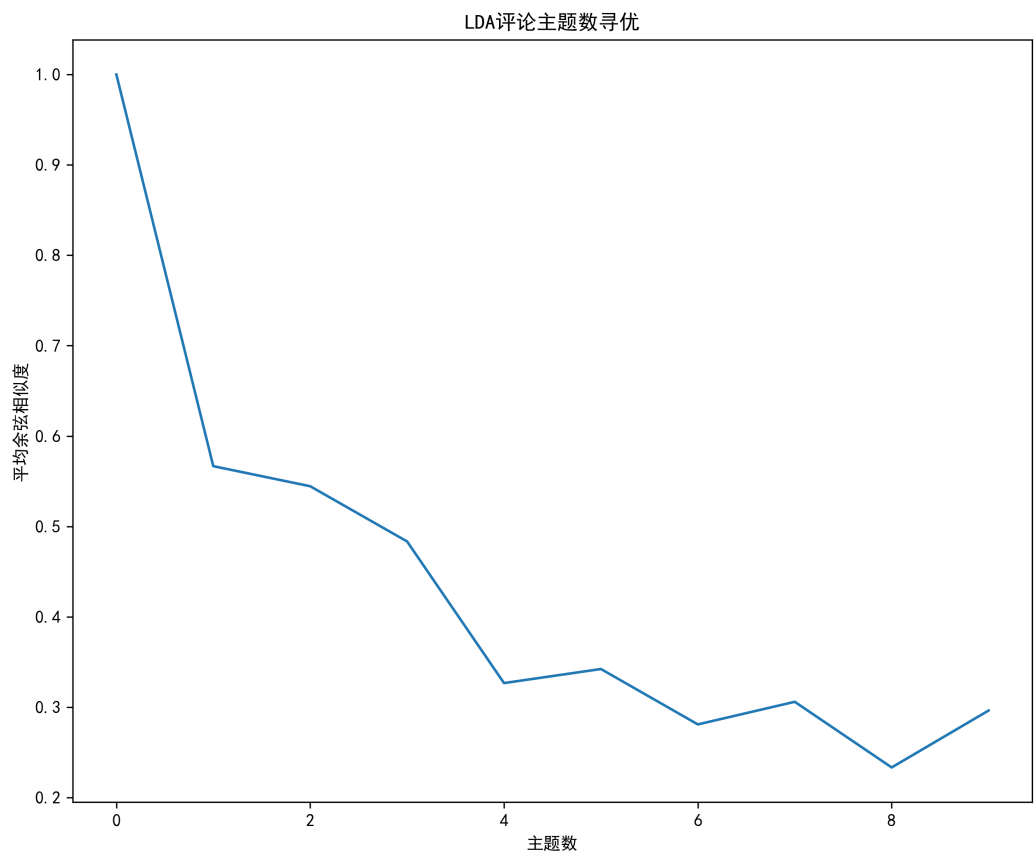
# 构造主题词列表，行表示主题号，列表示各主题词
mat = []
for j in np.arange(i):
    top_w = top_word[j]
    mat.append(tuple([top_w.count(k) for k in unique_word]))

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，在这里有3个低谷，可以选择4或者6或者8，根据最后lda呈现的模型去判断，效果好与坏，这样才用4



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，这里可以看出来选择4的效果挺不错的，圈圈分的比较开，不会重叠在一起

