

数据挖掘与最优化: Assignment 1

| | | | | | | |
|-----|---------------|------|---------------|-----|-----|-------------------------------------------------------------------------------------|
| 组长: | 231502004 李子 | 完成内容 | T2, T3 | 贡献度 | 20% |  |
| 组员: | 221300066 季千焜 | 完成内容 | T1 | 贡献度 | 20% |  |
| 组员: | 221098024 李瑛琦 | 完成内容 | T4 | 贡献度 | 20% |  |
| 组员: | 221098145 李傲雪 | 完成内容 | 优化 T2, T3, T4 | 贡献度 | 20% |  |
| 组员: | 221098071 单佳仪 | 完成内容 | T4 | 贡献度 | 20% |  |

目录

| | | |
|---|------------|----|
| 1 | 实验进度 | 3 |
| 2 | T1 | 3 |
| 3 | T2 | 4 |
| 4 | T3 | 6 |
| 5 | T4 | 8 |
| 6 | T2 与 T3 优化 | 12 |
| 7 | T4 优化 | 17 |

实验进度

我们完成了所有内容。

T1

修改 `gradient_descent` 函数，在每次迭代中遍历候选步长并选择使目标函数最小的步长。以下是代码：

```

1 def gradient_descent(max_iterations, threshold, w_init, obj_func, grad_func):
2     w = w_init.copy()
3     w_history = w.reshape(1, -1) # 初始化为二维数组
4     f_history = np.array([obj_func(w)])
5     i = 0
6     diff = 1.0e10
7
8     while i < max_iterations and diff > threshold:
9         grad = grad_func(w)
10        lrs = np.linspace(0, 1, 101)
11        best_lr = 0
12        best_f = np.inf
13        best_w = w
14        # 遍历所有候选步长，寻找最优步长
15        for lr in lrs:
16            w_candidate = w - lr * grad
17            f_candidate = obj_func(w_candidate)
18            if f_candidate < best_f:
19                best_f = f_candidate
20                best_lr = lr
21                best_w = w_candidate
22
23        # 更新参数和记录历史
24        w = best_w
25        w_history = np.vstack((w_history, w.reshape(1, -1)))
26        f_history = np.vstack((f_history, [best_f]))
27        # 计算函数值差异
28        if i == 0:
29            diff = np.abs(f_history[1] - f_history[0])
30        else:
31            diff = np.abs(f_history[-1] - f_history[-2])
32        i += 1
33    return w_history, f_history

```

修改后代码整体运行时间从 0.2 秒左右缩短到 0.083 秒左右，效率显著提高。

T2

我们通过如下代码，首先构建了词典列表：

```

1 def remove_numbers_from_file(input_file, output_file):
2     with open(input_file, 'r', encoding='utf-8') as file:
3         lines = file.readlines()
4
5     with open(output_file, 'w', encoding='utf-8') as file:
6         for line in lines:
7             new_line = ''.join([char for char in line if not char.isdigit()])
8             file.write(new_line)
9
10 if __name__ == "__main__":
11     input_file = 'THUOCL_diming.txt'
12     output_file = 'THUOCL_diming_no_numbers.txt'
13     remove_numbers_from_file(input_file, output_file)

```

然后通过如下代码，将两个文件合并：

```

1 def read_file(file_path):
2     with open(file_path, 'r', encoding='utf-8') as file:
3         return file.readlines()
4
5 def main():
6     file1 = 'THUOCL_diming_no_numbers.txt'
7     file2 = 'THUOCL_caijing_no_numbers.txt'
8
9     content1 = read_file(file1)
10    content2 = read_file(file2)
11
12    combined_content = content1 + content2
13
14    with open('combined_output.txt', 'w', encoding='utf-8') as output_file:
15        output_file.writelines(combined_content)
16
17 if __name__ == "__main__":
18     main()

```

接着修改 cut 代码，使其在仅剩单字的情况下不再切分：

```

1 import numpy as np
2 def cut(dic, text):
3     result = []
4     index = 0
5     text_length = len(text)

```

```

6     length=[]
7     for word in dic:
8         length.append(len(word))
9     m=max(length)
10    while text_length>index:
11        for size in np.arange(m+index,index,-1):
12            piece=text[index:size]
13            if piece in dic or len(piece) == 1:
14                result.append(piece)
15            index=size
16        break
17    print(result)

```

切分效果如下：

['据', '中', '指', '研究院', '统计', ',', ',', '自', '9', '月', '3', '0', '日', '人', '民', '银行', '发', '布', '消', '息', '决', '定', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '以', '来', ',', ',', '截', '至', '1', '0', '月', '9', '日', ',', ',', '已', '有', '杭', '州', '、', '济', '南', '、', '吉', '林', '等', '至', '少', '3', '0', '个', '地', '区', '的', '公', '积', '金', '管', '理', '中', '心', '发', '布', '相', '关', '通', '知', ',', ',', '落', '实', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '0', '、', '1', '5', '个', '百分点', '。', '近期', ',', ',', '包', '括', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '、', '阶', '段', '性', '放', '宽', '部', '分', '城', '市', '首', '套', '住', '房贷', '款', '利率', '下', '限', '、', '支持', '居', '民', '换', '购', '住', '房', '个', '人', '所', '得', '税', '退', '税', '优惠', '政', '策', '等', '接', '连', '出', '台', '。']

显然可见效果并不是很好，因为我们的词典中并没有包含很多常用词汇。并且词典也有一定不全的问题。比如：词典中有“杭州市”，“杭州 xx”，但是没有“杭州”这个词。这就导致无法切分出“杭州”这个词。

T3

修改代码如下：

```

1 import numpy as np
2 def cut(dic,text):
3     result=[]
4     index=0
5     text_length=len(text)
6     length=[]
7     for word in dic:
8         length.append(len(word))
9     m=max(length)
10    while text_length>index:
11        for size in np.arange(m+index,index,-1):
12            piece=text[index:size]
13            if piece in dic or len(piece) == 1:
14                result.append(piece)
15                index=size
16            break
17    return result
18 def cut2(dic, text):
19     result = []
20     index = len(text)
21     length = []
22     for word in dic:
23         length.append(len(word))
24     m = max(length)
25     while index > 0:
26         for size in np.arange(max(0, index - m), index):
27             piece = text[size:index]
28             if piece in dic or len(piece) == 1:
29                 result.append(piece)
30                 index = size
31             break
32     result.reverse()
33     return result
34 def cut3(dic, text):
35     result1 = cut(dic, text)
36     result2 = cut2(dic, text)
37     if len(result1) == len(result2):
38         print("SAME RESULT!")
39     if len(result1) <= len(result2):
40         return result1

```

```

41 else:
42     return result2
43
44 print(cut(dic,text))
45 print(cut2(dic,text))
46 print(cut3(dic,text))

```

得到输出:

['据', '中', '指', '研究院', '统计', ' ', ' ', '自', '9', '月', '3', '0', '日', '人', '民', '银行', '发', '布', '消', '息', '决', '定', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '以', '来', ' ', ' ', '截', '至', '1', '0', '月', '9', '日', ' ', ' ', '已', '有', '杭', '州', ' ', ' ', '济', '南', ' ', ' ', '吉', '林', '等', ' ', '至', ' ', '少', '3', '0', '个', '地', '区', '的', '公', '积', '金', '管', '理', '中', '心', '发', '布', '相', '关', '通', '知', ' ', ' ', '落', '实', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '0', ' ', ' ', '1', '5', '个', '百分点', ' ', ' ', '近期', ' ', ' ', '包', '括', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', ' ', ' ', '阶', '段', '性', '放', '宽', '部', '分', '城', '市', '首', '套', '住', '房贷', '款', '利率', '下', '限', ' ', ' ', '支持', '居', '民', '换', '购', '住', '房', '个', '人', '所', '得', '税', '退', '税', '优惠', '政', '策', '等', '接', '连', '出台', ' ', ' ']

['据', '中', '指', '研究院', '统计', ' ', ' ', '自', '9', '月', '3', '0', '日', '人', '民', '银行', '发', '布', '消', '息', '决', '定', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '以', '来', ' ', ' ', '截', '至', '1', '0', '月', '9', '日', ' ', ' ', '已', '有', '杭', '州', ' ', ' ', '济', '南', ' ', ' ', '吉', '林', '等', ' ', '至', ' ', '少', '3', '0', '个', '地', '区', '的', '公', '积', '金', '管', '理', '中', '心', '发', '布', '相', '关', '通', '知', ' ', ' ', '落', '实', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '0', ' ', ' ', '1', '5', '个', '百分点', ' ', ' ', '近期', ' ', ' ', '包', '括', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', ' ', ' ', '阶', '段', '性', '放', '宽', '部', '分', '城', '市', '首', '套', '住', '房', '贷款', '利率', '下', '限', ' ', ' ', '支持', '居', '民', '换', '购', '住', '房', '个', '人', '所', '得', '税', '退', '税', '优惠', '政', '策', '等', '接', '连', '出台', ' ', ' ']

SAME RESULT!

['据', '中', '指', '研究院', '统计', ' ', ' ', '自', '9', '月', '3', '0', '日', '人', '民', '银行', '发', '布', '消', '息', '决', '定', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '以', '来', ' ', ' ', '截', '至', '1', '0', '月', '9', '日', ' ', ' ', '已', '有', '杭', '州', ' ', ' ', '济', '南', ' ', ' ', '吉', '林', '等', ' ', '至', ' ', '少', '3', '0', '个', '地', '区', '的', '公', '积', '金', '管', '理', '中', '心', '发', '布', '相', '关', '通', '知', ' ', ' ', '落', '实', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', '0', ' ', ' ', '1', '5', '个', '百分点', ' ', ' ', '近期', ' ', ' ', '包', '括', '下调', '首', '套', '个', '人', '住', '房', '公', '积', '金', '贷款', '利率', ' ', ' ', '阶', '段', '性', '放', '宽', '部', '分', '城', '市', '首', '套', '住', '房贷', '款', '利率', '下', '限', ' ', ' ', '支持', '居', '民', '换', '购', '住', '房', '个', '人', '所', '得', '税', '退', '税', '优惠', '政', '策', '等', '接', '连', '出台', ' ', ' ']

说明正向和逆向切分结果一致。自然三种不同分词方法的结果也一致。

T4

(i) 将数据导入 Python，分别统计正面评论和负面评论的条数。以下是代码：

```
1 import re
2 import jieba as jb
3 import matplotlib.pyplot as plt
4 import numpy as np
5 with open('中文酒店评论.txt', 'r', encoding='utf-8') as file:
6     content = file.read()
7 p_content = content.split("\n")
8 gd_num = 0
9 bd_num = 0
10 regex1 = "^1"
11 for line in p_content:
12     if re.search(regex1,line) is not None: # 如果当前行匹配regex的结果不是None
13         gd_num += 1# 就打印这行信息
14     else:
15         bd_num += 1
16 print("正面评论数量: ",gd_num)
17 print("负面评论数量: ",bd_num)
```

得到输出：

正面评论数量：7000 负面评论数量：3001

(ii) 分别对每条评论进行分词并统计词数（或句子长度）。

```
1 jb_cmt = list()#分词结果
2 wd_num = list()#词数
3 len_cmt = list()#句子长度
4 for line in p_content[1:]:
5     cmt = line.split(" ")[1]
6     jb_cut = jb.lcut(cmt, cut_all = False)
7     jb_cmt.append(jb_cut)
8     wd_num.append(len(jb_cut))
9     len_cmt.append(len(cmt))
```

(iii) 分别绘制正面评论和负面评论句子长度的直方图。

```
1 # 绘制直方图
2 gdata = len_cmt[0:gd_num]
3 bdata = len_cmt[gd_num:]
4 plt.hist(gdata, color='blue', alpha=0.7) # bins 表示柱子的数量
5 plt.title('Positive comment histogram') # 设置标题
6 plt.xlabel('Wordage') # 设置 x 轴标签
7 plt.ylabel('Frequency') # 设置 y 轴标签
8 plt.show() # 显示图表
```



```

9 plt.hist(bdata, color='red', alpha=0.7)
10 plt.title('Negative comment histogram') # 设置标题
11 plt.xlabel('Wordage') # 设置 x 轴标签
12 plt.ylabel('Frequency') # 设置 y 轴标签
13 plt.show() # 显示图表

```

得到直方图：

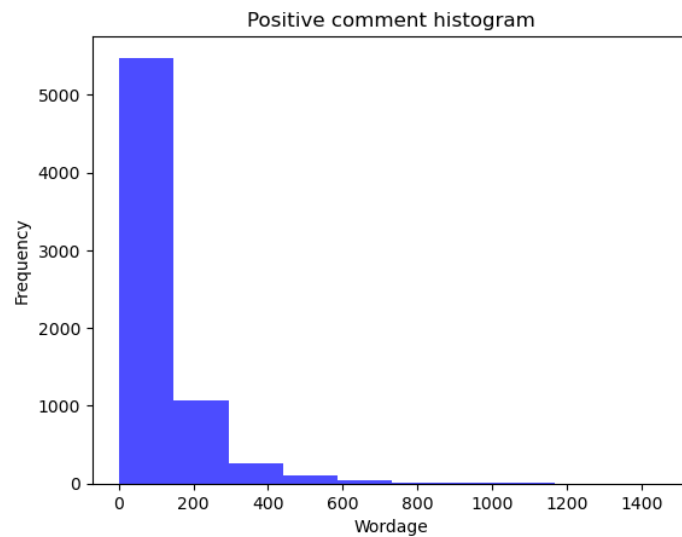


图 1: Positive Comment

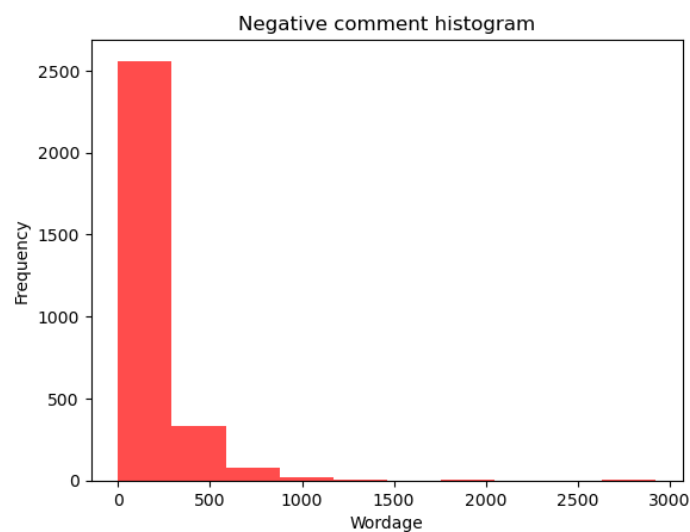


图 2: Negative Comment

(iv) 删除 (ii) 中每条评论分词后得到的非形容词，只保留每条评论分词后得到的形容词。

```

1 import jieba.posseg as pseg
2 def extract_adjectives(text):
3     words = pseg.lcut(text)
4     adjectives = [word.word for word in words if word.flag == 'a']

```

```

5     return adjectives
6     jb_cmtg = list()#正面分形容词结果
7     jbt_cmtb = list()#负面分形容词结果
8     for line in p_content[1:]:
9         cmt = line.split(" ")[1]
10        jb_cut = extract_adjectives(cmt)
11        if line.split(" ")[0]=='1':
12            jb_cmtg.append(jb_cut)
13        else:
14            jbt_cmtb.append(jbt_cut)

```

(v) 分别绘制正面评论和负面评论高频形容词的词云图。

```

1     from wordcloud import WordCloud
2     #创建证明评价词云对象
3     gcloud=list()
4     for gw in jb_cmtg:
5         gwords_str = " ".join(gw)
6         gcloud.append(gwords_str)
7
8     gwordcloud = WordCloud(background_color="white", width=800, height=600, font_path='font.ttf').
        generate(" ".join(gcloud))
9
10    plt.imshow(gwordcloud, interpolation='bilinear')
11    plt.axis("off")
12    plt.show()

```

得到正面评价词云图：



图 3: 正面评价词云图

```

1     #创建负面评价词云对象

```


T2 与 T3 优化

正向最大匹配法优化点：

- 加入结巴词库 (dict.txt.big) 扩大现有词典内容
- 利用正则表达式，完善日期、数字的分词

相关库

```
1 import jieba
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import jieba.posseg as pseg
6 import re
7 import itertools
8 from wordcloud import WordCloud
```

字典与文档

```
1 # 字典
2 dic_financial=pd.read_csv(".\THUOCL_caijing.txt",sep="\t")
3 dic_financial.columns = ['word', 'count']
4 dic_financial=dic_financial['word']
5
6 dic_place=pd.read_csv(".\THUOCL_diming.txt",sep="\t")
7 dic_place.columns = ['word', 'count']
8 dic_place=dic_place['word']
9
10 def load_jieba_dict(path):
11     word_list = [] # 改用列表存储
12     with open(path, 'r', encoding='utf-8') as f:
13         for line in f:
14             word = line.split()[0] # 取每行第一个词
15             word_list.append(word)
16     return word_list # 返回列表
17
18 my_dict = load_jieba_dict("dict.txt.big") #
19 dict=my_dict+list(dic_financial)+list(dic_place)+['吉林']
20
21 # 文档
22 text = "据中指研究院统计，自9月30日人民银行发布消息决定下调首套个人住房公积金贷款利率以来，截至10月9日，已有杭州、济南、吉林等至少30个地区的公积金管理中心发布相关通知，落实下调首套个人住房公积金贷款"
```

利率0.15个百分点。近期，包括下调首套个人住房公积金贷款利率、阶段性放宽部分城市首套住房贷款利率下限、支持居民换购住房个人所得税退税优惠政策等接连出台。"

新的正向最大匹配函数

```

1  def cut_forward(dic,text):
2      result=[]
3      index=0
4      text_length=len(text)
5      max_length = max(len(word) for word in dic) if dic else 1
6
7      # 正则匹配数字、日期、小数等模式
8      patterns = [
9          (r"\d+月\d+日", "DATE"),      # 匹配 "9月30日"
10         (r"\d+\.\d+", "FLOAT"),        # 匹配 "0.15"
11         (r"\d+", "INT")                 # 匹配纯数字
12     ]
13
14     while text_length>index:
15         # 优先用正则匹配数字/日期
16         matched = False
17         for pattern, _ in patterns:
18             regex = re.compile(pattern)
19             match = regex.match(text, index)
20             if match:
21                 result.append(match.group())
22                 index = match.end()
23                 matched = True
24                 break
25         if matched:
26             continue
27
28         current_max = min(max_length, text_length - index)
29         for size in np.arange(current_max,0,-1):
30             piece=text[index:index+size]
31             if piece in dic or len(piece)==1:
32                 result.append(piece)
33                 index+=size
34                 break
35     return(result)
36
37 print(cut_forward(dict,text))

```

得到输出：


```

35     piece = text[index-size:index]
36     if piece in dic or size == 1:
37         result.insert(0, piece) # 插入结果列表头部
38         index -= size
39         found = True
40         break
41     if not found:
42         # 处理未匹配字符
43         result.insert(0, text[index-1])
44         index -= 1
45
46     return result
47
48 print(cut_reverse(dict,text))

```

得到输出：

['据', '中指', '研究院', '统计', ', ', '自', '9月30日', '人民银行', '发布', '消息', '决定', '下调', '首套', '个人住房', '公积金', '贷款', '利率', '以来', ', ', '截至', '10月9日', ', ', '已有', '杭州', ', ', '济南', ', ', '吉林', '等', '至少', '30', '个', '地区', '的', '公积金', '管理中心', '发布', '相关', '通知', ', ', '落实', '下调', '首套', '个人住房', '公积金', '贷款', '利率', '0.15', '个', '百分点', '。', '近期', ', ', '包括', '下调', '首套', '个人住房', '公积金', '贷款', '利率', ', ', '阶段性', '放宽', '部分', '城市', '首套', '住房贷款', '利率', '下限', ', ', '支持', '居民', '换购', '住房', '个人所得税', '退税', '优惠政策', '等', '接连', '出台', '。']

新的双向匹配函数

```

1 # 双向匹配
2 def cut_2direction(dic, text):
3     result1 = cut_forward(dic, text)
4     result2 = cut_reverse(dic, text)
5     if result1 == result2:
6         print("SAME RESULT!")
7     if len(result1) <= len(result2):
8         return result1
9     else:
10         return result2
11
12 print(cut_2direction(dict,text))

```

得到输出：SAME RESULT!

['据', '中指', '研究院', '统计', ', ', '自', '9月30日', '人民银行', '发布', '消息', '决定', '下调', '首套', '个人住房', '公积金', '贷款', '利率', '以来', ', ', '截至', '10月9日', ', ', '已有', '杭州', ', ', '济南', ', ', '吉林', '等', '至少', '30', '个', '地区', '的', '公积金', '管理中心', '发布', '相关', '通知', ', ', '落实', '下调', '首套', '个人住房', '公积金', '贷款', '利率', '0.15', '个', '百分点', '。', '近期', ', ', '包括', '下调', '首套', '个人住房', '公积金', '贷款', '利率', ', ', '阶段性', '放宽', '部分', '城市', '首套', '住房贷款', '利率', '下限', ', ', '支持', '居民', '换购', '住房', '个人所得税', '退税', '优惠政策', '等', '接连', '出台', '。']

包括，' 下调'，' 首套'，' 个人住房'，' 公积金'，' 贷款'，' 利率'，'、'，' 阶段性'，' 放宽'，' 部分'，' 城市'，' 首套'，
' 住房贷款'，' 利率'，' 下限'，'、'，' 支持'，' 居民'，' 换购'，' 住房'，' 个人所得税'，' 退税'，' 优惠政策'，' 等'，'
接连'，' 出台'，'。']

T4 优化

优化点

- 完整所需内容位数据集，方便后续调用
- 部分代码调整

获取文件内容，并转换为 label, sentence 两列

```

1 # 读取整个文件内容
2 with open('中文酒店评论.txt', 'r', encoding='utf-8') as file:
3     lines = file.readlines()
4
5 df_txt = []
6 lines=lines[1:] # 去除掉读入后的第一行数据 "label sentence"，和后续需要处理的格式不相符
7 for line in lines:
8     # 去除行尾的换行符并按空格分割
9     parts = line.strip().split(' ', 1)
10    if len(parts) == 2:
11        label, sentence = parts
12        df_txt.append([int(label), sentence])
13    else:
14        label = parts[0]
15        df_txt.append([int(label), ''])
16
17 df_txt = pd.DataFrame(df_txt, columns=['label', 'sentence'])
18 df_txt['sentence'] = df_txt['sentence'].apply(lambda x: x.strip())
19
20 print(df_txt.head(6))
    
```

只列出代码调整部分，输出图形相同，此处不放置

```

1 # (1)
2 gd_num = len(df_txt[df_txt['label']==1])
3 bd_num = len(df_txt[df_txt['label']==-1])
4
5 # (2)
6 df_txt['length'] = df_txt['sentence'].str.len() # 添加句子长度列
7 df_txt['jieba'] = df_txt['sentence'].apply(lambda x: jieba.lcut(x)) # 添加分词结果列
8 df_txt['count'] = df_txt['jieba'].str.len() # 添加分词数量列
9
10 # (3)
11 gdata = df_txt[df_txt['label']==1]['length']
12 bdata = df_txt[df_txt['label']==-1]['length']
    
```

```

13
14 # (4)
15 # 添加形容词列
16 df_txt['jieba_with_pos'] = df_txt['sentence'].apply(lambda x: pseg.lcut(x))
17 df_txt['adj'] = df_txt['jieba_with_pos'].apply(lambda word_list: [pair.word for pair in word_list
    if pair.flag == 'a'])
18 df_txt.drop(columns=['jieba_with_pos'], inplace=True)
19
20 # (5)
21 g_word = list(itertools.chain(*df_txt[df_txt['label']==1]['adj']))
22 b_word = list(itertools.chain(*df_txt[df_txt['label']==-1]['adj']))

```