**初始化：**

```python
import re
from time import time
import jieba
import numpy as np
import pandas as pd
from datetime import datetime
import itertools
from geopy.distance import geodesic, geodesic
```

```python
angerfile=r'C:\Users\86186\Desktop\现代程序设计\anger.txt'
disgustfile=r'C:\Users\86186\Desktop\现代程序设计\disgust.txt'
fearfile=r'C:\Users\86186\Desktop\现代程序设计\fear.txt'
joyfile=r'C:\Users\86186\Desktop\现代程序设计\joy.txt'
sadnessfile=r'C:\Users\86186\Desktop\现代程序设计\sadness.txt'
stopwordsfile=r"C:\Users\86186\Desktop\现代程序设计\stopwords_list.txt"
filename=r'C:\Users\86186\Desktop\现代程序设计\3\weibo.txt'
```

**创建情绪词典：**

```python
anger=[]
disgust=[]
fear=[]
joy=[]
sadness=[]
```

```python
with open(angerfile,'r',encoding='utf-8') as f:
    for line in f.readlines():
        anger.append(line.strip())

with open(disgustfile,'r',encoding='utf-8') as f:
    for line in f.readlines():
        disgust.append(line.strip())

with open(fearfile,'r',encoding='utf-8') as f:
    for line in f.readlines():
        fear.append(line.strip())

with open(joyfile,'r',encoding='utf-8') as f:
    for line in f.readlines():
        joy.append(line.strip())

with open(sadnessfile,'r',encoding='utf-8') as f:
    for line in f.readlines():
        sadness.append(line.strip())
```

1. 实现一个函数，对微博数据进行清洗，去除噪声（如 url 等），过滤停用词。注意分词的时候应该将情绪词典加入 Jieba 的自定义词典，以提高这些情绪词的识别能力。

**创建停用词列表，并且事先实现去除噪声的功能，这步只考虑分词，过去除噪声只保留中文**

```python
def clean(line):#去除文本中的噪声(只保留中文)
    pattern = r'[^\u4e00-\u9fa5]'
    line = re.sub(pattern,'',line) #将其中所有非中文字符替换
    return line

def stopwordslist():
    stopwords=[]
    with open(stopwordsfile,"r",encoding='utf-8') as sw:
        for line in sw.readlines():
            stopwords.append(line.strip())
    stopwords.append(' ')
    return stopwords
```

实现分词的函数:

```python
def worddepart():
    jieba.load_userdict(angerfile)
    jieba.load_userdict(disgustfile)
    jieba.load_userdict(fearfile)
    jieba.load_userdict(joyfile)
    jieba.load_userdict(sadnessfile)
    res=[]
    stoplist=stopwordslist()

    with open(filename , 'r' , encoding='utf-8') as f:
        for line in f.readlines():
            midres=[]
            s=clean(line)
            s=jieba.lcut(s)
            for i in s:
                if i not in stoplist:
                    midres.append(i)
            res.append(midres)
    return res
```

运行结果:



2. 实现两个函数，实现一条微博的情绪分析，返其情绪向量或情绪值。目前有两种方法，一是认为一条微博的情绪是混合的，即一共有 n 个情绪词，如果 joy 有 n1 个，则 joy 的比例是 n1/n；二是认为一条微博的情绪是唯一的，即 n 个情绪词里，anger 的情绪词最多，则该微博的情绪应该为 angry。注意，这里要求用闭包实现，尤其是要利用闭包实现一次加载情绪词典且局部变量持久化的特点。同时，也要注意考虑一些特别的情况，如无情绪词出现，不同情绪的情绪词出现数目一样等，并予以处理（如定义为无情绪，便于在后面的分析中去除）。

**方法一：情绪向量**

```python
def getvector(anger,disgust,fear,joy,sadness):
    emodict={}
    emodict['anger']=anger
    emodict['disgust']=disgust
    emodict['fear']=fear
    emodict['joy']=joy
    emodict['sadness']=sadness
    def inget(comment):
        vector=[0,0,0,0,0]
        for i in comment:
            if i in emodict['anger']:
                vector[0] +=1
            elif i in emodict['disgust']:
                vector[1] +=1
            elif i in emodict['fear']:
                vector[2] +=1
            elif i in emodict['joy']:
                vector[3] +=1
            elif i in emodict['sadness']:
                vector[4] +=1
        return vector
    return inget
```

**运行结果：**

```
PS C:\Users\86186\Desktop\现代程序设计\3> c:; cd 'c:\Users\86186\Desktop\现代程序
86\.vscode\extensions\ms-python.python-2021.10.1317843341\pythonFiles\lib\python\
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\86186\AppData\Local\Temp\jieba.cache
Loading model cost 0.717 seconds.
Prefix dict has been built successfully.
[0, 0, 1, 1, 0]
```

**方法二：唯一值**

是认为一条微博的情绪是唯一的，即 n 个情绪词里，anger 的情绪词最多，则该微博的情绪应该为 angry，如果没有情绪词则为 nonemotion，如果最大情绪词数目有多个则为 complexemotion

```python
def getvalue(anger,disgust,fear,joy,sadness):
    emodict={}
    emodict['anger']=anger
    emodict['disgust']=disgust
    emodict['fear']=fear
    emodict['joy']=joy
    emodict['sadness']=sadness
    def inget(comment):
        vector=[0,0,0,0,0]
        for i in comment:
            if i in emodict['anger']:
                vector[0] +=1
            elif i in emodict['disgust']:
                vector[1] +=1
            elif i in emodict['fear']:
                vector[2] +=1
            elif i in emodict['joy']:
                vector[3] +=1
            elif i in emodict['sadness']:
                vector[4] +=1
```

```python
        vector = np.array(vector)
        pos = np.where(vector == vector.max())
        length = np.size(pos)
        if length == 1:
            if 0 in pos:
                return 'anger'
            elif 1 in pos:
                return 'disgust'
            elif 2 in pos:
                return 'fear'
            elif 3 in pos:
                return 'joy'
            elif 4 in pos:
                return 'sadness'
        elif length == 5:
            return 'noneemotion'
        else:
            return 'complexemotion'
    return inget
```

**运行结果：**



```
PS C:\Users\86186\Desktop\现代程序设计\3> c:; cd 'c:\Users\86186\Desktop\
86\.vscode\extensions\ms-python.python-2021.10.1317843341\pythonFiles\lib\
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\86186\AppData\Local\Temp\jieba.cache
Loading model cost 0.721 seconds.
Prefix dict has been built successfully.
complexemotion
```

3. 微博中包含时间，可以讨论不同时间情绪比例的变化趋势，实现一个函数，可以通过参数来控制并返回对应情绪的时间模式，如 joy 的小时模式，sadness 的周模式等。

**为了实现这个功能，首先要将每条微分类，一共有七种分别是 anger，disgust，fear，joy，sadness，noneemotion，complexemotion。用列表储存微博的原文本。然后利用正则表达式将某个情绪的时间提取出来，并利用列表储存**

```python
def get_time(emotionlist):
    time=[]
    i=0
    length=len(emotionlist)
    model_time = re.compile(r'\d\d\s\d\d:\d\d:\d\d')

    for i in range(length):
        content=emotionlist[i]
        t ="2013 11 "+ model_time.search(content).group()
        mid=[]
        mid.append(t)
        time.append(mid)

    return time
```

然后利用 pandas 库来统计某个情绪在不同时间段的频数，model 为不同的模式，0 为每半小时，1 为每小时，2 为每天，emotype 为不同情绪类型，0 为 anger，1 为 disgust，2 为 fear，3 为 joy，4 为 sadness

```python
def timesum(model,weibolist,emotype):
    emotionlist=weibolist[emotype]
    time=get_time(emotionlist)



    time = list(itertools.chain.from_iterable(time))
    time = pd.to_datetime(pd.Series(time),format='%Y %m %d %H:%M:%S')
    diction = {'time':time}
    data = pd.DataFrame(diction)
    if model == 0:
        m = data.resample('30T',on='time').count()
    elif model == 1:
        m = data.resample('60T',on='time').count()
    elif model == 2:
        m = data.resample('24h',on='time').count()
    return m
```

运行结果：

| time | | time | |
|---|---|---|---|
| time | | 2013-11-12 00:00:00 | 10 |
| 2013-11-11 00:00:00 | 9 | 2013-11-12 01:00:00 | 2 |
| 2013-11-11 01:00:00 | 6 | 2013-11-12 02:00:00 | 1 |
| 2013-11-11 02:00:00 | 0 | 2013-11-12 03:00:00 | 1 |
| 2013-11-11 03:00:00 | 3 | 2013-11-12 04:00:00 | 0 |
| 2013-11-11 04:00:00 | 0 | 2013-11-12 05:00:00 | 0 |
| 2013-11-11 05:00:00 | 3 | 2013-11-12 06:00:00 | 1 |
| 2013-11-11 06:00:00 | 6 | 2013-11-12 07:00:00 | 2 |
| 2013-11-11 07:00:00 | 3 | 2013-11-12 08:00:00 | 4 |
| 2013-11-11 08:00:00 | 21 | 2013-11-12 09:00:00 | 5 |
| 2013-11-11 09:00:00 | 15 | 2013-11-12 10:00:00 | 5 |
| 2013-11-11 10:00:00 | 9 | 2013-11-12 11:00:00 | 5 |
| 2013-11-11 11:00:00 | 9 | 2013-11-12 12:00:00 | 6 |
| 2013-11-11 12:00:00 | 18 | 2013-11-12 13:00:00 | 7 |
| 2013-11-11 13:00:00 | 15 | 2013-11-12 14:00:00 | 6 |
| 2013-11-11 14:00:00 | 18 | 2013-11-12 15:00:00 | 8 |
| 2013-11-11 15:00:00 | 6 | 2013-11-12 16:00:00 | 10 |
| 2013-11-11 16:00:00 | 30 | 2013-11-12 17:00:00 | 5 |
| 2013-11-11 17:00:00 | 6 | 2013-11-12 18:00:00 | 9 |
| 2013-11-11 18:00:00 | 9 | 2013-11-12 19:00:00 | 8 |
| 2013-11-11 19:00:00 | 18 | 2013-11-12 20:00:00 | 7 |
| 2013-11-11 20:00:00 | 30 | 2013-11-12 21:00:00 | 10 |
| 2013-11-11 21:00:00 | 12 | 2013-11-12 22:00:00 | 7 |
| 2013-11-11 22:00:00 | 6 | 2013-11-12 23:00:00 | 15 |
| 2013-11-11 23:00:00 | 30 | 2013-11-13 00:00:00 | 2 |

4. 微博中包含空间，可以讨论情绪的空间分布，实现一个函数，可以通过参数来控制并返回对应情绪的空间分布，即围绕某个中心点，随着半径增加该情绪所占比例的变化，中心点可默认值可以是城市的中心位置。

**首先获取每条评论的坐标，储存在列表中**

```python
def get_xy(emotionlist):

    xy=[]
    i=0
    length=len(emotionlist)
    model_x = re.compile(r'116\.\d+')
    model_y = re.compile(r'39\.\d+|40\.\d+')


    for i in range(length):
        content=emotionlist[i]
        x = model_x.search(content).group()
        y = model_y.search(content).group()

        xy.append((float(y),float(x)))

    return xy
```

定义北京市中心的坐标：

```
midxy=(39.9299857781,116.395645038,)
```

然后计算每条评论的坐标距离中心的位置，小于半径 r 的进行统计，emotype 为不同情绪类型，0 为 anger，1 为 disgust，2 为 fear，3 为 joy，4 为 sadness

```python
def xysum(r,weibolist,emotype):
    xycount=0
    i=0
    emotionlist=weibolist[emotype]

    xy=get_xy(emotionlist)
    for i in xy:
        if geodesic(i,midxy) <= r:
            xycount += 1

    return xycount
```

运行结果：

```
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\86186\AppData\Local\Temp\jieba.cache
Loading model cost 0.741 seconds.
Prefix dict has been built successfully.
2764
```

主函数：

```python
def main():
    res=worddepart()
    #print(res)

    #f1=getvector(anger,disgust,fear,joy,sadness)
    #print(f1(res[5]))

    #f2=getvalue(anger,disgust,fear,joy,sadness)
    #print(f2(res[5]))

    weibolist=weiboclass()
    #print(timesum(1,weibolist,0))
    print(xysum(10,weibolist,3))


if __name__=='__main__':
    main()
```