**concatenation.py**

```python
1  #!/usr/bin/env python2
2  # -*- coding: utf-8 -*-
3  """
4  Created on Fri Nov  9 10:00:13 2018
5
6  @author: fay
7  """
8
9  import pandas as pd
10 import os
11
12 base_path = '/Users/fay/Desktop/3-Tech/AFPD Projects/Data concatenation/data'
13 com1 = pd.read_csv(base_path + os.sep + '2014.csv')
14 com2 = pd.read_csv(base_path + os.sep + '2015.csv')
15
16 com = com1.append(com2)
17
18 output = com.drop(labels = 'Unnamed: 0', axis=1)
19 output.to_csv(base_path + os.sep + 'output.csv', index = False)
```

**datetime.py**

```python
1  #!/usr/bin/env python2
2  # -*- coding: utf-8 -*-
3  """
4  Created on Fri Nov  9 10:21:05 2018
5
6  @author: fay
7  """
8
9  from datetime import datetime
10 import pandas as pd
11 import os
12
13 base_path = '/Users/fay/Desktop/3-Tech/AFPD Projects/Datetime/data'
14 time =pd.read_csv(base_path + os.sep + 'ukpound_exchange.csv')
15
16 # Change the str to datetime in pandas
17 time['Date']= time['Date'].apply(lambda x: datetime.strptime(x, "%m/%d/%Y"))
18
19 # generate year and month from datetime to groupby
20 time['year']= time['Date'].apply(lambda x: x.year)
21 time['month']= time['Date'].apply(lambda x: x.month)
22
23 #| pick the max
24 time_max=time.groupby(['year','month']).max()
25
26 output = time_max.drop(columns=['Unnamed: 0'])
27 output.to_csv(base_path + os.sep + 'output.csv', index = False)
```

```python
transformation.py
1  #!/usr/bin/env python2
2  # -*- coding: utf-8 -*-
3  """
4  Transformation
5  Created on Wed Nov 14 16:26:05 2018
6
7  @author: fay
8  """
9  import pandas as pd
10 import os
11
12 base_path = '/Users/fay/Desktop/3-Tech/AFPD Projects/Data transformation/data'
13
14 # Read files
15 fi = open (base_path + os.sep + 'deal_level_data.csv', 'r')
16 data = fi.readlines()
17 fi = fi.close
18
19 # Use dic to store the data (faster than dataframe)
20 col_name = data[0].split(',')
21 data_dic = []
22
23 for i in range(1, len(data) ):
24     data1 = data[i].split(',')
25     a = {}
26     for j in range (len(col_name)):
27         a[col_name[j]] = data1[j]
28     data_dic.append(a)
29
30 # New columns (col)
31 col_res = col_name[0:14]
32 ## Extract columns from quarter level
33 fi2 = open (base_path + os.sep + 'quarter_level_data.csv', 'r')
34 col_add_ori = fi2.readline().split(',')
35 fi2 = fi2.close
36 col_add = col_add_ori [16:]
37 col_add[25] = 'Tar_TtlSalary_log' #The last word of a line has a '\n'
38 col = col_res + ['quarter_to_the_event_date', 'quarter'] + col_add
39
40 # Lists for "Quarter to event date"
41 num = ['0']
42 suffix = ['']
43 for i in range(12):
44     z = i+1
45     num.append('%d'%z)
46     num.append('-%d'%z)
47     suffix.append('_%d'%z)
48     suffix.append('__%d'%z)
49
50 # Transformation
51 new_data = []
52 for l in range(len(data_dic)):
53     for m in range(25):
54         b={}
55         for n in range(len(col)):
56             if n <= 13:
57                 b[col[n]]=data_dic[l][col[n]]
58             if n == 14:
59                 b[col[n]]=num[m]
60             if n >= 15:
61                 b[col[n]]=data_dic[l][col[n]+suffix[m]]
62         new_data.append(b)
63
64 # Change dicts to dataframe
65 output = pd.DataFrame(new_data)
66 output.columns = col
67 output.to_csv(base_path + os.sep + 'output.csv', index = False)
```

fuzzymatch  multiprocess.p...

```python
#!/usr/bin/env python2
# -*- coding: utf-8 -*-
"""
Created on Sat Nov 10 16:54:32 2018

@author: fay
"""

import time, os
from multiprocessing import Pool
from fuzzywuzzy import process
import pandas as pd

# Read key
base_path = '/Users/fay/Desktop/3-Tech/AFPD Projects/Fuzzy and multiprocessing/data'
key_data = pd.read_excel(base_path + os.sep + 'acquirers.xlsx')
key = list(key_data['Acquirer Name'])

# Read values
values_data = pd.read_csv(base_path + os.sep + 'bank_names.csv')
values = values_data['bank_names']

# Use process.extract to match
def match_name(keyword):
    outcome = process.extract(keyword, values)
    a= [outcome[i][0] for i in range(5)]
    a.insert(0,keyword)
    return a

# Multiprocessing
if __name__ == '__main__':
    start_time = time.time()
    print start_time
    p = Pool(processes=3)
    temp = p.map(match_name, key)
    print("--- %s seconds ---" % (time.time() - start_time))

# print the dataframe
output = pd.DataFrame(temp)
output.columns = ['acquirers','0','1','2','3', '4']
output.to_csv(base_path + os.sep + 'output.csv', index = False)
print('Done!')

# Without multiprocessing: 44.79seconds
#match_name(range(key))
#start_time = time.time()
#print start_time
#map(match_name, key)
#print("--- %s seconds ---" % (time.time() - start_time))
```

**geolocation.py\***

```python
1 #!/usr/bin/env python2
2 # -*- coding: utf-8 -*-
3 """
4 Created on Fri Nov  9 18:56:45 2018
5
6 @author: fay
7 """
8 import time, os
9 import pandas as pd
10 from geopy.distance import vincenty
11 import requests
12
13 # Read key
14 base_path = '/Users/fay/Desktop/3-Tech/AFPD Projects/Google API geolocation/data'
15 key_data = pd.read_excel(base_path + os.sep + 'coname_addresses.xlsx')
16
17 # Get API
18 baselink= 'https://maps.googleapis.com/maps/api/geocode/json?address='
19 api = '&key=AIzaSyDz8wfzYtrecnyn3a1YzwGNjmy0M1Sp25Y'
20
21 # Initialization
22 WhiteHouse=(38.8976763, 77.0387185)
23 lat = []
24 lng = []
25 dis = []
26 Notfound = []
27
28 # Find the coordinate
29 def location(key):
30     try:
31         addr = requests.get(baselink + key + api) # Average request time == 0.59s, max time == 1.57s
32         addr_json = addr.json()
33         loc =addr_json['results'][0]['geometry']['location']
34         key_loc = (loc['lat'], loc['lng'])
35         lat.append(loc['lat'])
36         lng.append(loc['lng'])
37         dis.append(vincenty(key_loc, WhiteHouse).km)
38     except:
39         lat.append('NA')
40         lng.append('NA')
41         dis.append('NA')
42         print key
43         Notfound.append(key) # 1 address cannot find, 2 addresses are empty
45 # Run the function
46 start_time = time.time() # Need 12 mins
47 print start_time
48 map(location, key_data['address'])
49 print("--- %s seconds ---" % (time.time() - start_time))
50 print str(Notfound) + ' Are Not Found. Please check the data!'
51
52 # Print out
53 key_data['lat'] = lat
54 key_data['lng'] = lng
55 key_data['distance'] = dis
56 key_data.to_csv(base_path + os.sep + 'output.csv')
57 print 'Done!'
```