keyb Hadh claron Activity 8.1: Aggregating Data with Pandas

8.1.1 Intended Learning Outcomes

After this activity, the student should be able to:

- · emonstrate querying and merging of dataframes
- · Perform advanced calculations on dataframes
- Aggregate dataframes with pandas and numpy
- Work with time series data

8.1.2 Resources

- · Computing Environment using Python 3.x
- · Attached Datasets (under Instructional Materials)

8.1.3 Procedures

The procedures can be found in the canvas module. Check the following under topics:

- 8.1 Weather Data Collection
- 8.2 Querying and Merging
- 8.3 Dataframe Operations
- 8.4 Aggregations
- 8.5 Time Series

8.1.4 Data Analysis

Provide some comments here about the results of the procedures.

8.1.5 Supplementary Activity

Using the CSV files provided and what we have learned so far in this module complete the following exercises:

- 1. With the earthquakes.csv file, select all the earthquakes in Japan with a magType of mb and a magnitude of 4.9 or greater.
- 2. Create bins for each full number of magnitude (for example, the first bin is 0-1, the second is 1-2, and so on) with a magType of ml and count how many are in each bin.
- 3. Using the faang.csv file, group by the ticker and resample to monthly frequency. Make the following aggregations:
- Mean of the opening price
- · Maximum of the high price
- · Minimum of the low price
- · Mean of the closing price
- · Sum of the volume traded
- 4. Build a crosstab with the earthquake data between the tsunami column and the magType column. Rather than showing the frequency count, show the maximum magnitude that was observed for each combination. Put the magType along the columns.
- 5. Calculate the rolling 60-day aggregations of OHLC data by ticker for the FAANG data. Use the same aggregations as exercise no. 3.
- 6. Create a pivot table of the FAANG data that compares the stocks. Put the ticker in the rows and show the averages of the OHLC and volume traded data.
- 7. Calculate the Z-scores for each numeric column of Netflix's data (ticker is NFLX) using apply().
- 8. Add event descriptions: Create a dataframe with the following three columns: ticker, date, and event. The columns should have the following values: ticker: 'FB' date: ['2018-07-25', '2018-03-19', '2018-03-20'] event: ['Disappointing user growth announced after close.', 'Cambridge Analytica story', 'FTC investigation'] Set the index to ['date', 'ticker'] Merge this data with the FAANG data using an outer join
- 9. Use the transform() method on the FAANG data to represent all the values in terms of the first date in the data. To do so, divide all the values for each ticker by the value

```
import pandas as pd
eq = pd.read_csv('/content/earthquakes.csv')
faang = pd.read_csv('/content/faang (1).csv')
```

eq

```
time
                                                             place tsunami parsed_place
      mag magType
                 ml 1539475168010
                                              9km NE of Aguanga, CA
                                                                                   California
                     1539475129610
                                              9km NE of Aguanga, CA
                                                                                   California
      3.42
                     1539475062610
                                              8km NE of Aguanga, CA
                     1539474978070
                                              9km NE of Aguanga, CA
      0.44
                                                                                   California
                 ml
                     1539474716050
                                              10km NW of Avenal, CA
                                                                                   California
 4
9327 0.62
                     1537230228060 9km ENE of Mammoth Lakes, CA
                                                                                   California
9328
                     1537230135130
                                                 3km W of Julian, CA
                                                                                   California
9329 2.40
                     1537229908180 35km NNE of Hatillo, Puerto Rico
                                                                                 Puerto Rico
9330 1.10
                     1537229545350
                                              9km NE of Aguanga, CA
                                                                                   California
9331 0.66
                     1537228864470
                                              9km NE of Aguanga, CA
                                                                                   California
```

rk

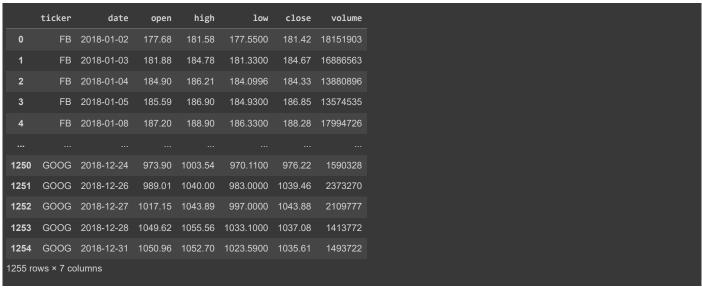
	mag	magType	time	place	tsunami	parsed_place
1563	4.9	mb	1538977532250	293km ESE of Iwo Jima, Japan	0	Japan
2576	5.4	mb	1538697528010	37km E of Tomakomai, Japan	0	Japan
3072	4.9	mb	1538579732490	15km ENE of Hasaki, Japan	0	Japan
3632	4.9	mb	1538450871260	53km ESE of Hitachi, Japan	0	Japan

```
# Create bins for each full number of magnitude (for example, the first bin is 0-1, the second is 1-2, and so on) with a magType of ml and c ml = eq[eq['magType'] == 'ml'] bins = [i for i in range(0, int(ml['mag'].max()) + 4)] counts = pd.cut(ml['mag'], bins=bins, right=False).value_counts().sort_index() counts
```

```
2072
[0, 1)
          3126
[1, 2)
[2, 3)
           985
[3, 4)
           153
[4, 5)
             6
[5, 6)
             2
[6, 7)
             0
             0
[7, 8)
```

Name: mag, dtype: int64

faang



```
# 3. Using the faang.csv file, group by the ticker and resample to monthly frequency. Make the following aggregations:
# Mean of the opening price
# Maximum of the high price
# Minimum of the low price
# Mean of the closing price
\mbox{\em \#} Sum of the volume traded
faang['date'] = pd.to_datetime(faang['date'])
faang.set_index('date', inplace=True)
agg = faang.groupby('ticker').resample('M').agg({
    'open': 'mean',
    'high': 'max',
    'low': 'min',
    'close': 'mean',
    'volume': 'sum'
})
agg
```

```
open
                                    high
                                                low
                                                           close
                                                                     volume
ticker
             date
AAPL 2018-01-31
                    170.714690
                                 176.6782
                                           161.5708
                                                      170.699271 659679440
       2018-02-28
                    164.562753
                                 177.9059
                                           147.9865
                                                      164.921884 927894473
       2018-03-31
                    172.421381
                                 180.7477
                                           162.4660
                                                      171.878919 713727447
       2018-04-30
                    167.332895
                                 176.2526
                                           158.2207
                                                      167.286924 666360147
       2018-05-31
                    182.635582
                                 187.9311
                                           162 7911
                                                      183.207418 620976206
       2018-06-30
                    186.605843
                                 192.0247
                                           178.7056
                                                      186.508652 527624365
       2018-07-31
                    188 065786
                                 193 7650
                                           181 3655
                                                      188 179724 393843881
       2018-08-31
                    210.460287
                                227.1001
                                           195.0999
                                                      211.477743 700318837
       2018-09-30
                    220.611742
                                227 8939
                                                      220.356353 678972040
       2018-10-31
                    219.489426
                                231.6645
                                           204.4963
                                                      219.137822 789748068
       2018-11-30
                    190.828681
                                220.6405
                                           169.5328
                                                      190.246652 961321947
       2018-12-31
                    164.537405
                                 184.1501
                                           145.9639
                                                      163.564732 898917007
AMZN 2018-01-31 1301.377143 1472.5800
                                          1170.5100 1309.010952
                                                                  96371290
       2018-02-28 1447.112632 1528.7000
                                          1265.9300 1442.363158 137784020
       2018-03-31 1542.160476 1617.5400 1365.2000 1540.367619 130400151
       2018-04-30 1475.841905 1638.1000
                                         1352.8800 1468.220476 129945743
       2018-05-31 1590.474545 1635.0000 1546.0200 1594.903636
                                                                  71615299
       2018-06-30 1699.088571 1763.1000
                                          1635.0900 1698.823810
                                                                  85941510
       2018-07-31 1786.305714 1880.0500 1678.0600 1784.649048
                                                                  97629820
       2018-08-31 1891.957826 2025.5700
                                          1776.0200 1897.851304
                                                                  96575676
       2018-09-30 1969.239474 2050.5000 1865.0000 1966.077895
                                                                  94445693
       2018-10-31 1799.630870 2033.1900
                                          1476.3600 1782.058261 183228552
       2018-11-30 1622.323810 1784.0000
                                         1420.0000 1625.483810 139290208
                                          1307.0000 1559.443158 154812304
       2018-12-31 1572.922105 1778.3400
 FΒ
       2018-01-31
                    184.364762
                                 190.6600
                                           175.8000
                                                      184.962857 495655736
       2018-02-28
                    180.721579
                                195.3200
                                           167.1800
                                                      180.269474 516621991
       2018-03-31
                    173.449524
                                 186.1000
                                           149.0200
                                                      173.489524 996232472
       2018-04-30
                    164.163557
                                 177.1000
                                           150.5100
                                                      163.810476 751130388
       2018-05-31
                                 192.7200
                                                      182.930000 401144183
       2018-06-30
                   194.974067
                                203.5500
                                           186.4300
                                                     195.267619 387265765
```

4. Build a crosstab with the earthquake data between the tsunami column and the magType column. Rather than showing the frequency count, s # magnitude that was observed for each combination. Put the magType along the columns. eq['tsunami'].value_counts()

```
magType
        mb mb_lg
                              ml ms_20
tsunami
        5.6
               3.5 4.11
                                   NaN 3.83
                                               5.8
                                                         6.0
              NaN NaN NaN 5.1
        6.1
                                     5.7 4.41 NaN NaN 7.5
       2018-07-31
                  1183.464286
                              1273.8900
                                        1093.8000 1187.590476
                                                               31953386
```

```
# 5. Calculate the rolling 60-day aggregations of OHLC data by ticker for the FAANG data. Use the same aggregations as exercise no. 3.
agg = faang.groupby('ticker').rolling('60D').agg({
    'open': 'mean',
    'high': 'max',
    'low': 'min',
    'close': 'mean',
    'volume': 'sum'
})
agg
```

		open	high	low	close	volume	
ticker	date						
AAPL	2018-01-02	166.927100	169.0264	166.0442	168.987200	25555934.0	
	2018-01-03	168.089600	171.2337	166.0442	168.972500	55073833.0	
	2018-01-04	168.480367	171.2337	166.0442	169.229200	77508430.0	
	2018-01-05	168.896475	172.0381	166.0442	169.840675	101168448.0	
	2018-01-08	169.324680	172.2736	166.0442	170.080040	121736214.0	
NFLX	2018-12-24	283.509250	332.0499	233.6800	281.931750	525657894.0	
	2018-12-26	281.844500	332.0499	231.2300	280.777750	520444588.0	
	2018-12-27	281.070488	332.0499	231.2300	280.162805	532679805.0	
	2018-12-28	279.916341	332.0499	231.2300	279.461341	521968250.0	
	2018-12-31	278.430769	332.0499	231.2300	277.451410	476309676.0	
1255 rows × 5 columns							

6. Create a pivot table of the FAANG data that compares the stocks. Put the ticker in the rows and show the averages of the OHLC and volum table = pd.pivot_table(faang, index='ticker', aggfunc='mean') table

```
close
                         high
                                      low
                                                            volume
                                                 open
ticker
AAPL
        186.986218 188.906858
                                           187.038674 3.402145e+07
AMZN 1641.726175 1662.839801 1619.840398 1644.072669 5.649563e+06
 FΒ
        171.510936 173.615298
                                169.303110 171.454424 2.768798e+07
GOOG 1113.225139 1125.777649 1101.001594 1113.554104 1.742645e+06
NFLX
                   325.224583
                                           319.620533 1.147030e+07
```

```
# 7. Calculate the Z-scores for each numeric column of Netflix's data (ticker is NFLX) using apply()
netflix = faang[faang['ticker'] == 'NFLX']
def z_score(column):
    return (column - column.mean()) / column.std()
z_scores = netflix.select_dtypes(include='number').apply(z_score)
z_scores
```

```
high
                                                   close
                                                            volume
           date
      2018-01-02 -2.500753 -2.516023 -2.410226 -2.416644 -0.088760
      2018-01-03 -2.380291 -2.423180 -2.285793 -2.335286 -0.507606
# 8. Add event descriptions
events_data = pd.DataFrame({
    'ticker': ['FB', 'FB', 'FB'],
    'date': ['2018-07-25', '2018-03-19', '2018-03-20'],
    'event': ['Disappointing user growth announced after close.',
              'Cambridge Analytica story',
              'FTC investigation']
})
events_data['date'] = pd.to_datetime(events_data['date'])
# Merge with FAANG data using outer join
data = pd.merge(faang, events_data, on=['date', 'ticker'], how='outer')
data
                 date ticker
                                          high
                                                     1ow
                                                            close
                                                                     volume event
                                 open
                                        181.58
                                                 177.5500
                                                                              NaN
       0
            2018-01-03
                          FB
                                                 181 3300
                                                           184.67 16886563
                                181 88
                                        184 78
                                                                              NaN
                                184.90
                                        186.21
                                                 184.0996
                                                           184.33 13880896
                                                                              NaN
            2018-01-05
                          FΒ
                                185.59
                                        186.90
                                                 184.9300
                                                           186.85 13574535
                                                                              NaN
            2018-01-08
                                187.20
                                        188.90
                                                 186.3300
                                                           188.28 17994726
                                                                              NaN
      1250 2018-12-24 GOOG
                                973.90 1003.54
                                                           976.22
                                                                    1590328
                                                                              NaN
      1251 2018-12-26 GOOG
                                       1040.00
                                989.01
                                                 983.0000 1039.46
                                                                    2373270
                                                                              NaN
      1252 2018-12-27 GOOG 1017.15 1043.89
                                                 997.0000 1043.88
                                                                              NaN
      1253 2018-12-28 GOOG 1049.62 1055.56 1033.1000 1037.08
                                                                    1413772
                                                                              NaN
      1254 2018-12-31 GOOG 1050.96 1052.70 1023.5900 1035.61
                                                                    1493722
                                                                              NaN
     1255 rows × 8 columns
dataas = data[data['event'] == 'FTC investigation']
dataas
               date ticker
                                    high
                                             low close
                                                             volume
                                                                              event
      53 2018-03-20
                         FB 167.47 170.2 161.95 168.15 129851768 FTC investigation
dataas = data[data['event'] == 'Cambridge Analytica story']
dataas
               date ticker
                                                    close
                                                             volume
                              open
                                      high
                                              low
                                                                                     event
      52 2018-03-19
                         FB 177.01 177.17 170.06 172.56 88140060 Cambridge Analytica story
dataas = data[data['event'] == 'Disappointing user growth announced after close.']
dataas
```