
✓ Hands-on Activity 8.1: Aggregating Data with Pandas

8.1.1 Intended Learning Outcomes

After this activity, the student should be able to:

- Demonstrate 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

✓ Weather Data Collection

```

1 import requests
2 import datetime
3 import pandas as pd
4 import sqlite3
5 from IPython import display
6
7 class weather_data():
8     def __init__(self):
9         self.make_request()
10        self.date()
11        self.results = []
12        self.response()
13        self.dataframe()
14        self.connection()
15
16    def make_request(self, endpoint = 'data', payload = None):
17        return requests.get(f'https://www.ncdc.noaa.gov/cdo-web/api/v2/{endpoint}',
18                            headers={'token': 'uPbSRvXwGYFwftSwWzZNLZsxpPKvvaYN'},
19                            params=payload)
20
21    def date(self):
22        self.current = datetime.date(2018, 1, 1)
23        self.end = datetime.date(2019, 1, 1)
24
25    def response(self):
26        while self.current < self.end:
27            display.clear_output(wait=True)
28            display.display(f'Gathering data for {str(self.current)}')
29
30            response = self.make_request(
31                'data',
32                {
33                    'datasetid' : 'GHCND',
34                    'locationid' : 'CITY:US360019',
35                    'startdate' : self.current,
36                    'enddate' : self.current,
37                    'units' : 'metric',
38                    'limit' : 1000
39                }
40            )
41
42            if response.ok:
43                self.results.extend(response.json()['results'])
44
45            self.current += datetime.timedelta(days=1)
46
47    def dataframe(self):
48        self.df = pd.DataFrame(self.results)
49        self.df.to_csv('/content/weather_data/nyc_weather_2018.csv', index = False)
50
51    def connection(self):
52        with sqlite3.connect('/content/weather_data/weather.db') as connection:
53            self.df.to_sql('weather', connection, index = False, if_exists = 'replace')
54
55        self.response = self.make_request('stations',{'datasetid':'GHCND',
56                                                    'locationid':'CITY:US360019',
57                                                    'limit':1000}
58        )
59
60        self.stations = pd.DataFrame(self.response.json()['results'])[['id', 'name', 'latitude', 'longitude', 'elevation']]
61        self.stations.to_csv('/content/weather_data/weather_stations.csv', index=False)
62
63        with sqlite3.connect('/content/weather_data/weather.db') as connection:
64            self.stations.to_sql('stations', connection, index=False, if_exists='replace')
65
66        weather = weather_data()
67
68        'Gathering data for 2018-12-31'

```

Querying and Merging

```

1 import pandas as pd
2 import sqlite3
3
4 class query_and_merge():
5     def __init__(self):
6         self.weather()
7         self.query_weather()
8         self.stationinfo()
9         self.dirty_data()
10
11     while True:
12         try:
13             x = input("\n0. Exit \n"\
14                 "1. See the head of the weather dataframe \n"\
15                 "2. See the head of the queried weather dataframe \n"\
16                 "3. See the head of the station_info dataframe \n"\
17                 "4. Describe unique values of station_info df \n"\
18                 "5. Describe unique values of weather df \n"\
19                 "6. Get info on station_info and weather dataframes \n"\
20                 "7. Get sample on merged weather and station dfs \n"\
21                 "8. Get sample on inner join \n"\
22                 "9. Get tail on right join \n"\
23                 "10. Get tail on left join \n"\
24                 "11. Get sample on outer join \n"\
25                 "12. Rows and columns after ij, lj, and rj \n"\
26                 "13. See the head of the dirty data dataframe \n"\
27                 "14. Get sample on valid stations \n"\
28                 "15. Get sample on invalid stations \n"\
29                 "16. See the head on the merged valid and invalid stations' dataframes \n"\
30                 "17. See the head on the joined valid and invalid stations' dataframes \n"\
31                 "18. See the intersection of the weather index \n"\
32                 "19. See the difference of the weather index \n"\
33                 "20. See the difference of the station info index \n"\
34                 "21. ny_in_name == weather? \n"\
35                 "22. All unique indexes of unioned weather and station_info dataframes \n"\
36                 "> ")
37         if x == "0":
38             break
39         elif x == "1":
40             print(self.weather.head())
41             continue
42         elif x == "2":
43             print(self.snow_data.head())
44             continue
45         elif x == "3":
46             print(self.station_info.head())
47             continue
48         elif x == "4":
49             print(self.station_info.id.describe())
50             continue
51         elif x == "5":
52             print(self.weather.station.describe())
53             continue
54         elif x == "6":
55             print(self.get_info('shape', self.station_info, self.weather))
56             continue
57         elif x == "7":
58             print(self.weather.merge(self.station_info.rename(dict(id='station')), axis = 1, on = 'station').sample(5, random_state=0))
59             continue
60         elif x == "8":
61             self.inner_join()
62             print(self.inner_join.sample(5, random_state=0))
63             continue
64         elif x == "9":
65             self.right_join()
66             print(self.right_join.tail())
67             continue
68         elif x == "10":
69             self.left_join()
70             print(self.left_join.tail())
71             continue
72         elif x == "11":
73             self.outer_join()
74             print(self.outer_join.sample(4, random_state=0).append(self.outer_join[self.outer_join.station.isna()].head(2)))
75             continue
76         elif x == "12":
77             print(self.get_info('shape', self.inner_join, self.left_join, self.right_join))

```



```

78         continue
79     elif x == "13":
80         print(self.dirty_data.head())
81         continue
82     elif x == "14":
83         self.valid()
84         print(self.valid_station.sample(5, random_state = 0))
85         continue
86     elif x == "15":
87         self.invalid()
88         print(self.station_with_wesf.sample(5, random_state = 0))
89         continue
90     elif x == "16":
91         print(self.valid_station.merge(self.station_with_wesf, left_index=True, right_index=True).query('WESF > 0').head())
92         continue
93     elif x == "17":
94         print(self.valid_station.join(self.station_with_wesf, rsuffix='_?').query('WESF > 0').head())
95         continue
96     elif x == "18":
97         self.set_index()
98         print(self.weather.index.intersection(self.station_info.index))
99         continue
100    elif x == "19":
101        print(self.weather.index.difference(self.station_info.index))
102        continue
103    elif x == "20":
104        print(self.station_info.index.difference(self.weather.index))
105        continue
106    elif x == "21":
107        self.symmetric_difference()
108        print(self.ny_in_name.index.difference(self.weather.index).shape[0]\
109        + self.weather.index.difference(self.ny_in_name.index).shape[0]\
110        == self.weather.index.symmetric_difference(self.ny_in_name.index).shape[0])
111        continue
112    elif x == "22":
113        print(self.weather.index.unique().union(self.station_info.index))
114        continue
115
116    except ValueError:
117        print("Error Occured. Please try again")
118        continue
119
120    def weather(self):
121        self.weather = pd.read_csv('/content/query_merge/nyc_weather_2018.csv')
122
123    def query_weather(self):
124        self.snow_data = self.weather.query('datatype == "SNOW" and value > 0')
125
126    def stationinfo(self):
127        self.station_info = pd.read_csv('/content/query_merge/weather_stations.csv')
128
129    def inner_join(self):
130        self.inner_join = self.weather.merge(self.station_info, left_on = 'station', right_on = 'id')
131
132    def left_join(self):
133        self.left_join = self.station_info.merge(self.weather, left_on = 'id', right_on = 'station', how = 'left')
134
135    def right_join(self):
136        self.right_join = self.weather.merge(self.station_info, left_on = 'station', right_on = 'id', how = 'right')
137
138    def outer_join(self):
139        self.outer_join = self.weather.merge(self.station_info[self.station_info.name.str.contains('NY')],
140        left_on = 'station', right_on = 'id', how = 'outer', indicator = True)
141
142    def dirty_data(self):
143        self.dirty_data = pd.read_csv('/content/query_merge/dirty_data.csv', index_col='date').drop_duplicates().drop(columns='SNWD')
144
145    def valid(self):
146        self.valid_station = self.dirty_data.query('station != "?"').copy().drop(columns=['WESF', 'station'])
147
148    def invalid(self):
149        self.station_with_wesf = self.dirty_data.query('station == "?"').copy().drop(columns=['station', 'TOBS', 'TMIN', 'TMAX'])
150
151    def set_index(self):
152        self.weather.set_index('station', inplace=True)
153        self.station_info.set_index('id', inplace=True)
154

```

```

155 def symmetric_difference(self):
156     self.ny_in_name = self.station_info[self.station_info.name.str.contains('NY')]
157
158 def station_describe(self):
159     self.station_info.id.describe()
160
161 def get_info(self, attr, *dfs):
162     return list(map(lambda x: getattr(x, attr), dfs))
163
164
165
166 QueryAndMerge = query_and_merge()

```

```

0. Exit
1. See the head of the weather dataframe
2. See the head of the queried weather dataframe
3. See the head of the station_info dataframe
4. Describe unique values of station_info df
5. Describe unique values of weather df
6. Get info on station_info and weather dataframes
7. Get sample on merged weather and station dfs
8. Get sample on inner join
9. Get tail on right join
10. Get tail on left join
11. Get sample on outer join
12. Rows and columns after ij, lj, and rj
13. See the head of the dirty data dataframe
14. Get sample on valid stations
15. Get sample on invalid stations
16. See the head on the merged valid and invalid stations' dataframes
17. See the head on the joined valid and invalid stations' dataframes
18. See the intersection of the weather index
19. See the difference of the weather index
20. See the difference of the station info index
21. ny_in_name == weather?
22. All unique indexes of unioned weather and station_info dataframes
> 1

```

	attributes	datatype	date	station	value
0	,,N,	PRCP	2018-01-01T00:00:00	GHCND:US1CTFR0039	0.0
1	,,N,	PRCP	2018-01-01T00:00:00	GHCND:US1NJBG0015	0.0
2	,,N,	SNOW	2018-01-01T00:00:00	GHCND:US1NJBG0015	0.0
3	,,N,	PRCP	2018-01-01T00:00:00	GHCND:US1NJBG0017	0.0
4	,,N,	SNOW	2018-01-01T00:00:00	GHCND:US1NJBG0017	0.0

```

0. Exit
1. See the head of the weather dataframe
2. See the head of the queried weather dataframe
3. See the head of the station_info dataframe
4. Describe unique values of station_info df
5. Describe unique values of weather df
6. Get info on station_info and weather dataframes
7. Get sample on merged weather and station dfs
8. Get sample on inner join
9. Get tail on right join
10. Get tail on left join
11. Get sample on outer join
12. Rows and columns after ij, lj, and rj
13. See the head of the dirty data dataframe
14. Get sample on valid stations
15. Get sample on invalid stations
16. See the head on the merged valid and invalid stations' dataframes
17. See the head on the joined valid and invalid stations' dataframes
18. See the intersection of the weather index
19. See the difference of the weather index
20. See the difference of the station info index
21. ny_in_name == weather?
22. All unique indexes of unioned weather and station_info dataframes
> 2

```

	attributes	datatype	date	station	value
--	------------	----------	------	---------	-------

Dataframe Operations

Arithmetic and Statistics

```

1 import numpy as np
2 import pandas as pd
3 weather = pd.read_csv('/content/Arithmetic_and_Stat/nyc_weather_2018.csv', parse_dates=['date'])
4 fb = pd.read_csv('/content/Arithmetic_and_Stat/fb_2018.csv', index_col='date', parse_dates=True)
5
6 x = fb.assign(abs_z_score_volume=lambda x: x.volume.sub(x.volume.mean()).div(x.volume.std()).abs())
7 y = fb.assign(volume_pct_change=fb.volume.pct_change(),
8               pct_change_rank=lambda x: x.volume_pct_change.abs().rank(ascending=False))
9 print("abs_z_score_volume > 3")
10 print(x.query('abs_z_score_volume > 3'))
11 print("\n\nnpct_change_rank")
12 print(y.nsmallest(5, 'pct_change_rank'))
13 print("\n\n2018-01-11':'2018-01-12")
14 print(fb['2018-01-11':'2018-01-12'])
15 print("\n\nThroughout 2018, Facebook's stock price never had a low above $215:")
16 print((fb > 215).any())
17 print("\n\nFacebook's OHLC (open, high, low, and close) prices all had at least one day they were at $215 or less:")
18 print((fb > 215).all())

```

```

abs_z_score_volume > 3
   open  high  low  close  volume  abs_z_score_volume
date
2018-03-19  177.01  177.17  170.06  172.56  88140060      3.145078
2018-03-20  167.47  170.20  161.95  168.15  129851768      5.315169
2018-03-21  164.80  173.40  163.30  169.39  106598834      4.105413
2018-03-26  160.82  161.10  149.02  160.06  126116634      5.120845
2018-07-26  174.89  180.13  173.75  176.26  169803668      7.393705

```

```

pct_change_rank
   open  high  low  close  volume  volume_pct_change \
date
2018-01-12  178.06  181.48  177.40  179.37  77551299      7.087876
2018-03-19  177.01  177.17  170.06  172.56  88140060      2.611789
2018-07-26  174.89  180.13  173.75  176.26  169803668      1.628841
2018-09-21  166.64  167.25  162.81  162.93  45994800      1.428956
2018-03-26  160.82  161.10  149.02  160.06  126116634      1.352496

```

```

pct_change_rank
date
2018-01-12      1.0
2018-03-19      2.0
2018-07-26      3.0
2018-09-21      4.0
2018-03-26      5.0

```

```

2018-01-11':'2018-01-12
   open  high  low  close  volume
date
2018-01-11  188.40  188.40  187.38  187.77  9588587
2018-01-12  178.06  181.48  177.40  179.37  77551299

```

Throughout 2018, Facebook's stock price never had a low above \$215:

```

open      True
high      True
low       False
close     True
volume    True
dtype: bool

```

Facebook's OHLC (open, high, low, and close) prices all had at least one day they were at \$215 or less:

```

open      False
high      False
low       False
close     False
volume    True
dtype: bool

```

▼ Binning and Thresholds

```

1 import matplotlib.pyplot as plt
2
3 volume_binned = pd.cut(fb.volume, bins=3, labels=['low', 'med', 'high'])
4 print(volume_binned.value_counts())
5 print("\n\nJuly 25th Facebook announced disappointing user growth and the stock tanked in the after hours:")
6 print(fb['2018-07-25':'2018-07-26'])
7 print("\n\nCambridge Analytica scandal broke on Saturday March 17th, so we look to the Monday for the numbers:")
8 print(fb['2018-03-16':'2018-03-20'])

```

```

low      240
med       8
high      3
Name: volume, dtype: int64

```

July 25th Facebook announced disappointing user growth and the stock tanked in the after hours:

	open	high	low	close	volume
date					
2018-07-25	215.715	218.62	214.27	217.50	64592585
2018-07-26	174.890	180.13	173.75	176.26	169803668

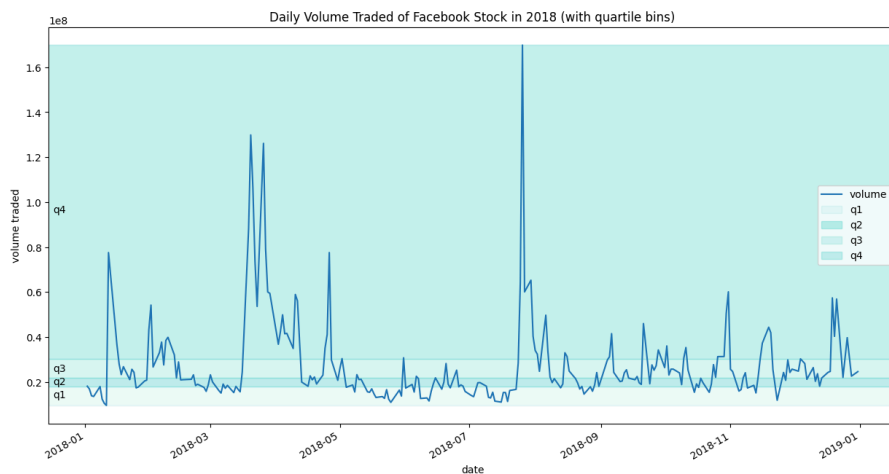
Cambridge Analytica scandal broke on Saturday March 17th, so we look to the Monday for the numbers:

	open	high	low	close	volume
date					
2018-03-16	184.49	185.33	183.41	185.09	24403438
2018-03-19	177.01	177.17	170.06	172.56	88140060
2018-03-20	167.47	170.20	161.95	168.15	129851768

```

1 fb.plot(y='volume', figsize=(15, 8), title='Daily Volume Traded of Facebook Stock in 2018 (with quartile bins)')
2
3 for bin_name, alpha, bounds in zip(
4     ['q1', 'q2', 'q3', 'q4'], [0.1, 0.35, 0.2, 0.3], pd.qcut(fb.volume, q=4).unique().categories.values
5 ):
6     plt.axhspan(bounds.left, bounds.right, alpha=alpha, label=bin_name, color='mediumturquoise')
7     plt.annotate(bin_name, xy=('2017-12-17', (bounds.left + bounds.right)/2.1))
8
9 plt.ylabel('volume traded')
10 plt.legend()
11 plt.show()

```




```

1 central_park_weather = weather.query(
2     'station == "GHCND:USW00094728"'
3     ).pivot(index='date', columns='datatype', values='value')
4 central_park_weather.SNOW.clip(0, 1).value_counts()

0.0    354
1.0     11
Name: SNOW, dtype: int64

```

▼ Applying Functions

```

1 import numpy as np
2 oct_weather_z_scores = central_park_weather.loc['2018-10', ['TMIN', 'TMAX', 'PRCP']].apply(lambda x: x.sub(x.mean()).div(x.std()))
3 print("October 27th rained much more than the rest of the days:")
4 print(oct_weather_z_scores.query('PRCP > 3'))
5 print("\n\nIndeed, this day was much higher than the rest:")
6 print(central_park_weather.loc['2018-10', 'PRCP'].describe())
7 print("\n\n")
8 fb.apply(lambda x: np.vectorize(lambda y: len(str(np.ceil(y))))(x)).astype('int64').equals(fb.applymap(lambda x: len(str(np.ceil(x))))

```

```

October 27th rained much more than the rest of the days:
datatype      TMIN      TMAX      PRCP
date
2018-10-27 -0.751019 -1.201045  3.936167

```

Indeed, this day was much higher than the rest:

```

count    31.000000
mean      2.941935
std       7.458542
min       0.000000
25%       0.000000
50%       0.000000
75%       1.150000
max       32.300000
Name: PRCP, dtype: float64

```

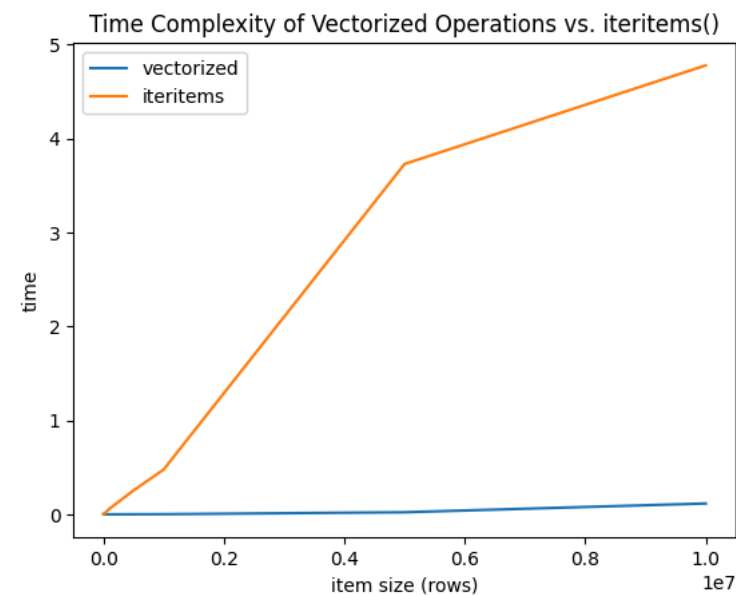
True

```

1 import time
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
6
7 np.random.seed(0)
8
9 vectorized_results = {}
10 iteritems_results = {}
11
12 for size in [10, 100, 1000, 10000, 100000, 500000, 1000000, 5000000, 10000000]:
13     test = pd.Series(np.random.uniform(size=size))
14
15     start = time.time()
16     x = test + 10
17     end = time.time()
18     vectorized_results[size] = end - start
19
20     start = time.time()
21     x = []
22     for i, v in test.iteritems():
23         x.append(v + 10)
24     x = pd.Series(x)
25     end = time.time()
26     iteritems_results[size] = end - start
27
28 pd.DataFrame(
29     [pd.Series(vectorized_results, name='vectorized'), pd.Series(iteritems_results, name='iteritems')]
30 ).T.plot(title='Time Complexity of Vectorized Operations vs. iteritems()')
31
32 plt.xlabel('item size (rows)')
33 plt.ylabel('time')
34 plt.show()

```

<ipython-input-17-b1f8f4319006>:22: FutureWarning: iteritems is deprecated and will be removed in a future version. Use .items() instead.
for i, v in test.iteritems():



Window Calculations

```

1 rolling = central_park_weather['2018-10'].assign(rolling_PRCP=lambda x: x.PRCP.rolling('3D').sum())[['PRCP', 'rolling_PRCP']]
2 print("Rolling Calculation of Central Park Weather for October 2018")
3 print(rolling.head(7).T)
4 print("\n\nPerforming Several Operations on the data in the '2018-10' column")
5 print(central_park_weather['2018-10'].rolling('3D').mean().head(7).iloc[:, :6])
6
7 rolling_func_per_column = central_park_weather['2018-10-01': '2018-10-07'].rolling('3D').agg(
8     {'TMAX': 'max', 'TMIN': 'min', 'AWND': 'mean', 'PRCP': 'sum'}).join(central_park_weather[
9     ['TMAX', 'TMIN', 'AWND', 'PRCP']], lsuffix='_rolling')
10 print("\n\nDifferent Aggregation Functions per Column")
11 print(rolling_func_per_column.sort_index(axis=1))
12
13 isequal = central_park_weather.PRCP.expanding().sum().equals(central_park_weather.PRCP.cumsum())
14 print("\n\nIs the expanding calculation of the precipitation in central park equal to the cumulative sum of it?")
15 print(isequal)
16
17 expanding_agg_per_column = central_park_weather['2018-10-01': '2018-10-07'].expanding().agg({'TMAX': np.max, 'TMIN': np.min, 'AWND': np.me
18     central_park_weather[['TMAX', 'TMIN', 'AWND', 'PRCP']], lsuffix='_expanding')
19 print("\n\nExpanding Aggregations per column")
20 print(expanding_agg_per_column.sort_index(axis=1))
21
22 moving_average = fb.assign(close_ewma=lambda x: x.close.ewm(span=5).mean())
23 print("\n\nExponentially Weighted Moving Average")
24 print(moving_average.tail(10)[['close', 'close_ewma']])

```

2018-10-06	0.5	0.833333	0.0	1.0	20.0	24.4	17.2
2018-10-07	1.1	1.066667	0.0	0.0	26.1	26.1	19.4

datatype	TMIN_rolling
date	
2018-10-01	17.2
2018-10-02	17.2
2018-10-03	17.2
2018-10-04	16.1
2018-10-05	15.6
2018-10-06	15.6
2018-10-07	15.6

```

Is the expanding calculation of the precipitation in central park equal to the cumulative sum of it?
False

```

```

Expanding Aggregations per column
datatype    AWND  AWND_expanding  PRCP  PRCP_expanding  TMAX  TMAX_expanding \
date

```

✓ Pipes

```
1 def get_info(df):
2     return '%d rows and %d columns and max closing z-score was %d' % (df.shape, df.close.max())
3
4 getinfo = fb['2018-Q1'].apply(lambda x: (x - x.mean())/x.std()).pipe(get_info)\
5 == get_info(fb['2018-Q1'].apply(lambda x: (x - x.mean())/x.std()))
6
7 print("Is the result of using the custom 'get_info' function the same as the get_info function using the pipe method?")
8 print(getinfo)
9
10 pipe_rolling = fb.pipe(pd.DataFrame.rolling, '20D').mean().equals(fb.rolling('20D').mean())
11 print("\n\nIs the result of directly using the rolling method for the calculation of the mean the same as using pipe?")
12 print(pipe_rolling)
13
14 def window_calc(df, func, agg_dict, *args, **kwargs):
15     return df.pipe(func, *args, **kwargs).agg(agg_dict)
16
17 print("\n\nExpanding Median of FB using Pipe and Rolling Method (without args and kwargs)")
18 print(window_calc(fb, pd.DataFrame.expanding, np.median).head())
19
20 print("\n\nExponentially Weighted Moving Average (with kwargs)")
21 print(window_calc(fb, pd.DataFrame.ewm, 'mean', span=3).head())
22
23 print("\n\nWindow Calculation using args for the indow size")
24 print(window_calc(central_park_weather['2018-10'], pd.DataFrame.rolling,
25 {'TMAX': 'max', 'TMIN': 'min', 'AWND': 'mean', 'PRCP': 'sum'}, '3D').head())
```

Is the result of using the custom 'get_info' function the same as the get_info function using the pipe method?
True

Is the result of directly using the rolling method for the calculation of the mean the same as using pipe?
True

Expanding Median of FB using Pipe and Rolling Method (without args and kwargs)

	open	high	low	close	volume
date					
2018-01-02	177.68	181.580	177.5500	181.420	18151903.0
2018-01-03	179.78	183.180	179.4400	183.045	17519233.0
2018-01-04	181.88	184.780	181.3300	184.330	16886563.0
2018-01-05	183.39	185.495	182.7148	184.500	15383729.5
2018-01-08	184.90	186.210	184.0996	184.670	16886563.0

Exponentially Weighted Moving Average (with kwargs)

	open	high	low	close	volume
date					
2018-01-02	177.680000	181.580000	177.550000	181.420000	1.815190e+07
2018-01-03	180.480000	183.713333	180.070000	183.586667	1.730834e+07
2018-01-04	183.005714	185.140000	182.372629	184.011429	1.534980e+07
2018-01-05	184.384000	186.078667	183.736560	185.525333	1.440299e+07
2018-01-08	185.837419	187.534839	185.075110	186.947097	1.625679e+07

Window Calculation using args for the indow size

datatype	TMAX	TMIN	AWND	PRCP
date				
2018-10-01	24.4	17.2	0.900000	0.0
2018-10-02	25.0	17.2	0.900000	17.5
2018-10-03	25.0	17.2	0.966667	17.5
2018-10-04	25.0	16.1	0.800000	18.5
2018-10-05	24.4	15.6	1.033333	1.0

```
<ipython-input-30-3fcf3e16cf1b>:4: FutureWarning: Indexing a DataFrame with a datetimelike index using a single string to slice the rows
getinfo = fb['2018-Q1'].apply(lambda x: (x - x.mean())/x.std()).pipe(get_info)\
<ipython-input-30-3fcf3e16cf1b>:5: FutureWarning: Indexing a DataFrame with a datetimelike index using a single string to slice the rows
== get_info(fb['2018-Q1'].apply(lambda x: (x - x.mean())/x.std()))
<ipython-input-30-3fcf3e16cf1b>:24: FutureWarning: Indexing a DataFrame with a datetimelike index using a single string to slice the row
print(window_calc(central_park_weather['2018-10'], pd.DataFrame.rolling,
```

▼ Aggregations

▼ Summarizing DataFrames

```
1 import numpy as np
2 import pandas as pd
3 weather = pd.read_csv('/content/Aggregations/weather_by_station.csv', index_col='date', parse_dates=True)
4 fb = pd.read_csv('/content/Aggregations/fb_2018.csv', index_col='date', parse_dates=True).assign(
5     trading_volume=lambda x: pd.cut(x.volume, bins=3, labels=['low', 'med', 'high'])
6 )
7 pd.set_option('display.float_format', lambda x: '%.2f' % x) #All numerical outputs will be shown in 2 rounded decimal places whilst not a
8 print("Aggregate into a Single Series")
9 agg = fb.agg({
10     'open': np.mean,
11     'high': np.max,
12     'low': np.min,
13     'close': np.mean,
14     'volume': np.sum
15 })
16 print(agg)
17 print("\n\nSnowfall and Precipitation Recorded in Central Park in 2018")
18 print(weather.query('station == "GHCND:USW00094728").pivot(
19     columns='datatype', values='value')[['SNOW', 'PRCP']].agg('sum'))
20 print("\n\nAggregate into a Data Frame")
21 agg_df = fb.agg({
22     'open': 'mean',
23     'high': ['min', 'max'],
24     'low': ['min', 'max'],
25     'close': 'mean'
26 })
27 print(agg_df)
28
```

Aggregate into a Single Series

open	171.45
high	218.62
low	123.02
close	171.51
volume	6949682394.00
dtype:	float64

Snowfall and Precipitation Recorded in Central Park in 2018

datatype	
SNOW	1007.00
PRCP	1665.30
dtype:	float64

Aggregate into a Data Frame

	open	high	low	close
mean	171.45	NaN	NaN	171.51
min	NaN	129.74	123.02	NaN
max	NaN	218.62	214.27	NaN

▼ Using groupby()

```

1 groupby = fb.groupby('trading_volume').mean()
2 agg_groupby = fb.groupby('trading_volume')['close'].agg(['min', 'max', 'mean'])
3
4 print("Mean based on the trading volume")
5 print(groupby)
6 print("\n\nAggregated Values of each row of the trading volume")
7 print(agg_groupby)
8
9 fb_agg = fb.groupby('trading_volume').agg({
10     'open': 'mean',
11     'high': ['min', 'max'],
12     'low': ['min', 'max'],
13     'close': 'mean'
14 })
15
16 print("\n\nGrouped Aggregation of each row of the trading volume")
17 print(fb_agg)
18
19 print("\n\nHierarchical names of the columns in the dataframe")
20 print(fb_agg.columns)
21
22 print("\n\nJoining the hierarchical names into one")
23 fb_agg.columns = ['_'.join(col_agg) for col_agg in fb_agg.columns]
24 print(fb_agg.head())
25
26 print("\n\nAverage Precipitation Per Day during October 2018")
27 print(weather['2018-10'].query('datatype == "PRCP"]').groupby(
28     pd.Grouper(freq='D')
29 ).mean().head())
30
31 print("\n\nTotal Quarterly Precipitation Per Station")
32 print(weather.query('datatype == "PRCP"]').groupby(
33     ['station_name', pd.Grouper(freq='Q')]
34 ).sum().unstack().sample(5, random_state=1))
35
36 print("\n\nTotal snowfall per New York Station (in millimeters)")
37 print(weather.groupby('station').filter(lambda x: 'NY' in x.name).query(
38     'datatype == "SNOW"]').groupby('station_name').sum().squeeze()
39
40 print("\n\n5 months with the most precipitation")
41 print(weather.query('datatype == "PRCP"]').groupby(
42     pd.Grouper(freq='D')
43 ).mean().groupby(pd.Grouper(freq='M')).sum().value.nlargest())
44
45 print("\n\nThe five largest percentage of daily precipitation relative to the total precipitation in that month")
46 print(weather\
47 .query('datatype == "PRCP")\
48 .rename(dict(value='prcp'), axis=1)\
49 .groupby(pd.Grouper(freq='D')).mean()\
50 .assign(
51     total_prctp_in_month=lambda x: x.groupby(
52         pd.Grouper(freq='M')
53     ).transform(np.sum),
54     pct_monthly_prctp=lambda x: x.prctp.div(
55         x.total_prctp_in_month
56     )
57 ).nlargest(5, 'pct_monthly_prctp'))
58
59 print("\n\nStandardized Dataframe using transform")
60 print(fb[['open', 'high', 'low', 'close']].transform(
61     lambda x: (x - x.mean()).div(x.std())
62 ).head())

```

```

Mean based on the trading volume
      open  high  low  close  volume
trading_volume
low      171.36 173.46 169.31 171.43 24547207.71
med      175.82 179.42 172.11 175.14 79072559.12
high     167.73 170.48 161.57 168.16 141924023.33

```

```

Aggregated Values of each row of the trading volume
      min  max  mean
trading_volume
low      124.06 214.67 171.43
med      152.22 217.50 175.14
high     160.06 176.26 168.16

```

Grouped Aggregation of each row of the trading volume

	open	high		low		close
trading_volume	mean	min	max	min	max	mean
low	171.36	129.74	216.20	123.02	212.60	171.43
med	175.82	162.85	218.62	150.75	214.27	175.14
high	167.73	161.10	180.13	149.02	173.75	168.16

Hierarchical names of the columns in the dataframe

```
MultiIndex([( 'open', 'mean'),
              ( 'high', 'min'),
              ( 'high', 'max'),
              ( 'low', 'min'),
              ( 'low', 'max'),
              ('close', 'mean')],
            )
```

Joining the hierarchical names into one

	open_mean	high_min	high_max	low_min	low_max	close_mean
trading_volume						
low	171.36	129.74	216.20	123.02	212.60	171.43
med	175.82	162.85	218.62	150.75	214.27	175.14
high	167.73	161.10	180.13	149.02	173.75	168.16

Average Precipitation Per Day during October 2018

date	value
2018-10-01	0.01
2018-10-02	2.23
2018-10-03	19.69
2018-10-04	0.32
2018-10-05	0.97

▼ Pivot Tables and Crosstabs

```

1 print("Pivot around trading volume as the column")
2 print(fb.pivot_table(columns='trading_volume'))
3 print("\n\nPivot around trading volume as the index")
4 print(fb.pivot_table(index='trading_volume'))
5
6 print("\n\nWeather data in the wide format")
7 print(weather.reset_index().pivot_table(
8     index=['date', 'station', 'station_name'],
9     columns='datatype',
10    values='value',
11    aggfunc='median'
12    ).reset_index().tail())
13
14 from google.colab import drive
15 drive.mount('/content/drive')
16
17 print("\n\nFrequency Table for Trading Volume per Month")
18 print(pd.crosstab(
19     index=fb.trading_volume,
20     columns=fb.index.month,
21     colnames=['month']
22 ))
23
24 print("\n\nPercentage of the value in each row relative to the total value per month")
25 print(pd.crosstab(
26     index=fb.trading_volume,
27     columns=fb.index.month,
28     colnames=['month'],
29     normalize='columns'
30 ))
31
32 print("\n\nAverage trading volume per month")
33 print(pd.crosstab(
34     index=fb.trading_volume,
35     columns=fb.index.month,
36     colnames=['month'],
37     values=fb.close,
38     aggfunc=np.mean
39 ))
40
41 print("\n\nNumber of times each station recorded snowfall per month")
42 snow_data = weather.query('datatype == "SNOW"')
43 print(pd.crosstab(
44     index=snow_data.station_name,
45     columns=snow_data.index.month,
46     colnames=['month'],
47     values=snow_data.value,
48     aggfunc=lambda x: (x > 0).sum(),
49     margins=True,
50     margins_name='total observations of snow'
51 ))

```

```

Pivot around trading volume as the column
trading_volume      low      med      high
close              171.43    175.14    168.16
high              173.46    179.42    170.48
low               169.31    172.11    161.57
open              171.36    175.82    167.73
volume            24547207.71  79072559.12  141924023.33

```

```

Pivot around trading volume as the index
           close  high  low  open  volume
trading_volume
low          171.43  173.46  169.31  171.36  24547207.71
med          175.14  179.42  172.11  175.82  79072559.12
high         168.16  170.48  161.57  167.73  141924023.33

```

```

Weather data in the wide format
datatype  date      station      station_name \
28740    2018-12-31  GHCND:USW00054787  FARMINGDALE REPUBLIC AIRPORT, NY US
28741    2018-12-31  GHCND:USW00094728  NY CITY CENTRAL PARK, NY US
28742    2018-12-31  GHCND:USW00094741  TETERBORO AIRPORT, NJ US
28743    2018-12-31  GHCND:USW00094745  WESTCHESTER CO AIRPORT, NY US
28744    2018-12-31  GHCND:USW00094789  JFK INTERNATIONAL AIRPORT, NY US

```


datatype	AWND	DAPR	MDPR	PGTM	PRCP	SNOW	SNWD	...	WSF5	WT01	WT02	\
28740	5.00	NaN	NaN	2052.00	28.70	NaN	NaN	...	15.70	NaN	NaN	
28741	NaN	NaN	NaN	NaN	25.90	0.00	0.00	...	NaN	1.00	NaN	
28742	1.70	NaN	NaN	1954.00	29.20	NaN	NaN	...	8.90	NaN	NaN	
28743	2.70	NaN	NaN	2212.00	24.40	NaN	NaN	...	11.20	NaN	NaN	
28744	4.10	NaN	NaN	NaN	31.20	0.00	0.00	...	12.50	1.00	1.00	

datatype	WT03	WT04	WT05	WT06	WT08	WT09	WT11
28740	NaN	NaN	NaN	NaN	NaN	NaN	NaN
28741	NaN	NaN	NaN	NaN	NaN	NaN	NaN
28742	NaN	NaN	NaN	NaN	NaN	NaN	NaN
28743	NaN	NaN	NaN	NaN	NaN	NaN	NaN
28744	NaN	NaN	NaN	NaN	NaN	NaN	NaN

[5 rows x 30 columns]

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

Frequency Table for Trading Volume per Month

month	1	2	3	4	5	6	7	8	9	10	11	12
trading_volume												
low	20	19	15	20	22	21	18	23	19	23	21	19
med	1	0	4	1	0	0	2	0	0	0	0	0
high	0	0	2	0	0	0	1	0	0	0	0	0

Percentage of the value in each row relative to the total value per month

month	1	2	3	4	5	6	7	8	9	10	11	12
-------	---	---	---	---	---	---	---	---	---	----	----	----

Time Series

Time-Based Selection and Filtering

```

1 import numpy as np
2 import pandas as pd
3
4 fb = pd.read_csv('/content/Time_Series/fb_2018.csv', index_col='date', parse_dates=True).assign(
5     trading_volume=lambda x: pd.cut(x.volume, bins=3, labels=['low', 'med', 'high'])
6 )
7 print("As the index of the dataframe is the date, it does not include rows with dates under holidays or weekends")
8 print("Weekends (Called October 11-15):\n",fb['2018-10-11':'2018-10-15'])
9 print("\n\nHolidays (Called January 1-7):\n",fb.first('1W')) #first takes from the beginning, last takes from the end
10 stock_data_per_minute = pd.read_csv('/content/Time_Series/fb_week_of_may_20_per_minute.csv',
11     index_col='date', parse_dates=True,
12     date_parser=lambda x: pd.to_datetime(x, format='%Y-%m-%d %H-%M'))
13 print("\n\nStock Data per Minute:\n", stock_data_per_minute.head())
14 print("\n\nStock Data per Day:\n", stock_data_per_minute.groupby(pd.Grouper(freq='1D')).agg({
15     'open': 'first',
16     'high': 'max',
17     'low': 'min',
18     'close': 'last',
19     'volume': 'sum'
20 })))
21 print("\n\nData at a specific time per day (i.e. 9:30, 15:59-16:00, etc.): \n",
22     "At a time\n", stock_data_per_minute.at_time('9:30'),
23     "\n\n",
24     "Between times\n",stock_data_per_minute.between_time('15:59', '16:00'))
25
26 shares_traded_in_first_30_min = stock_data_per_minute\
27     .between_time('9:30', '10:00')\
28     .groupby(pd.Grouper(freq='1D'))\
29     .filter(lambda x: (x.volume > 0).all())\
30     .volume.mean()
31
32 shares_traded_in_last_30_min = stock_data_per_minute\
33     .between_time('15:30', '16:00')\
34     .groupby(pd.Grouper(freq='1D'))\
35     .filter(lambda x: (x.volume > 0).all())\
36     .volume.mean()
37
38 print("\n\nDifference between the average shares traded in the first 30 minutes to the last 30 minutes everyday: \n",
39     shares_traded_in_first_30_min - shares_traded_in_last_30_min,
40     "-- A positive number means that more shares are traded in the first 30 minutes, while a negative number is the opposite")
41
42 print("\n\nUsing Normalize to remove all specific time information:\n",
43     pd.DataFrame(dict(before=stock_data_per_minute.index,
44         after=stock_data_per_minute.index.normalize())).head())

```

As the index of the dataframe is the date, it does not include rows with dates under holidays or weekends
Weekends (Called October 11-15):

	open	high	low	close	volume	trading_volume
date						
2018-10-11	150.13	154.81	149.16	153.35	35338901	low
2018-10-12	156.73	156.89	151.30	153.74	25293492	low
2018-10-15	153.32	155.57	152.55	153.52	15433521	low

Holidays (Called January 1-7):

	open	high	low	close	volume	trading_volume
date						
2018-01-02	177.68	181.58	177.55	181.42	18151903	low
2018-01-03	181.88	184.78	181.33	184.67	16886563	low
2018-01-04	184.90	186.21	184.10	184.33	13880896	low
2018-01-05	185.59	186.90	184.93	186.85	13574535	low

Stock Data per Minute:

	open	high	low	close	volume
date					
2019-05-20 09:30:00	181.62	181.62	181.62	181.62	159049.00
2019-05-20 09:31:00	182.61	182.61	182.61	182.61	468017.00
2019-05-20 09:32:00	182.75	182.75	182.75	182.75	97258.00
2019-05-20 09:33:00	182.95	182.95	182.95	182.95	43961.00
2019-05-20 09:34:00	183.06	183.06	183.06	183.06	79562.00

Stock Data per Day:

	open	high	low	close	volume
date					
2019-05-20	181.62	184.18	181.62	182.72	10044838.00

```

2019-05-21 184.53 185.58 183.97 184.82 7198405.00
2019-05-22 184.81 186.56 184.01 185.32 8412433.00
2019-05-23 182.50 183.73 179.76 180.87 12479171.00
2019-05-24 182.33 183.52 181.04 181.06 7686030.00

```

Data at a specific time per day (i.e. 9:30, 15:59-16:00, etc.):

```

At a time
      open  high  low  close  volume
date
2019-05-20 09:30:00 181.62 181.62 181.62 181.62 159049.00
2019-05-21 09:30:00 184.53 184.53 184.53 184.53 58171.00
2019-05-22 09:30:00 184.81 184.81 184.81 184.81 41585.00
2019-05-23 09:30:00 182.50 182.50 182.50 182.50 121930.00
2019-05-24 09:30:00 182.33 182.33 182.33 182.33 52681.00

```

```

Between times
      open  high  low  close  volume
date
2019-05-20 15:59:00 182.01 182.01 182.01 182.01 124560.00

```

✓ Shifting for Lagged Data

```

1 print("Using shift() to determine the previous day's closing price and change in price:\n",
2       fb.assign(prior_close=lambda x: x.close.shift(),
3                 after_hours_change_in_price=lambda x: x.open - x.prior_close,
4                 abs_change=lambda x: x.after_hours_change_in_price.abs()).nlargest(5, 'abs_change'))
5
6 print("\n\nAligning the starting time to the Market Hours:\n",
7       pd.date_range('2018-01-01', freq='D', periods=5) + pd.Timedelta('9 hours 30 minutes'))
8
9 print("\n\nFinding the valid/non-null/non-missing data in the month of September\n",
10       "First:", fb['2018-09'].first_valid_index(),
11       "\n Last:", fb['2018-09'].last_valid_index())
12
13 print("\n\nIs the market open on Feb. 30?\n",
14       (fb.index == '2018-09-30').any(),
15       "\n\nData on the last time the market was open:",
16       fb.asof('2018-09-30'))

```

Using shift() to determine the previous day's closing price and change in price:

```

      open  high  low  close  volume trading_volume prior_close \
date
2018-07-26 174.89 180.13 173.75 176.26 169803668      high      217.50
2018-04-26 173.22 176.27 170.80 174.16 77556934      med      159.69
2018-01-12 178.06 181.48 177.40 179.37 77551299      med      187.77
2018-10-31 155.00 156.40 148.96 151.79 60101251      low      146.22
2018-03-19 177.01 177.17 170.06 172.56 88140060      med      185.09

```

```

      after_hours_change_in_price  abs_change
date
2018-07-26                    -42.61      42.61
2018-04-26                    13.53     13.53
2018-01-12                    -9.71      9.71
2018-10-31                     8.78      8.78
2018-03-19                     -8.08      8.08

```

Aligning the starting time to the Market Hours:

```

DatetimeIndex(['2018-01-01 09:30:00', '2018-01-02 09:30:00',
              '2018-01-03 09:30:00', '2018-01-04 09:30:00',
              '2018-01-05 09:30:00'],
              dtype='datetime64[ns]', freq='D')

```

Finding the valid/non-null/non-missing data in the month of September

```

First: 2018-09-04 00:00:00
Last: 2018-09-28 00:00:00

```

```

Is the market open on Feb. 30?
False

```

Data on the last time the market was open: open

168.33

```

high            168.79
low             162.56
close           164.46
volume          34265638
trading_volume  low
Name: 2018-09-30 00:00:00, dtype: object
<ipython-input-67-581a5bc14eac>:10: FutureWarning: Indexing a DataFrame with a datetimelike index using a single string to slice the row
"First:", fb['2018-09'].first_valid_index(),
<ipython-input-67-581a5bc14eac>:11: FutureWarning: Indexing a DataFrame with a datetimelike index using a single string to slice the row
"\n Last:", fb['2018-09'].last_valid_index())

```

▼ Differenced Data

```

1 print("Is the result of subtracting each element from its",
2       "\npreceding element using the shift method equal to the\n",
3       "\nresult obtained by using the diff method?\n",
4       (fb.drop(columns='trading_volume') - fb.drop(columns='trading_volume').shift()
5        ).equals(fb.drop(columns='trading_volume').diff()))
6
7 print("\n\nDaily Difference of Facebook stock\n",
8       fb.drop(columns='trading_volume').diff().head(),
9       "\n\nDifference between current value and the values three positions behind it\n",
10      fb.drop(columns='trading_volume').diff(-3).head())
11

```

Is the result of subtracting each element from its
preceding element using the shift method equal to the
result obtained by using the diff method?
True

Daily Difference of Facebook stock

	open	high	low	close	volume
date					
2018-01-02	NaN	NaN	NaN	NaN	NaN
2018-01-03	4.20	3.20	3.78	3.25	-1265340.00
2018-01-04	3.02	1.43	2.77	-0.34	-3005667.00
2018-01-05	0.69	0.69	0.83	2.52	-306361.00
2018-01-08	1.61	2.00	1.40	1.43	4420191.00

Difference between current value and the values three positions behind it

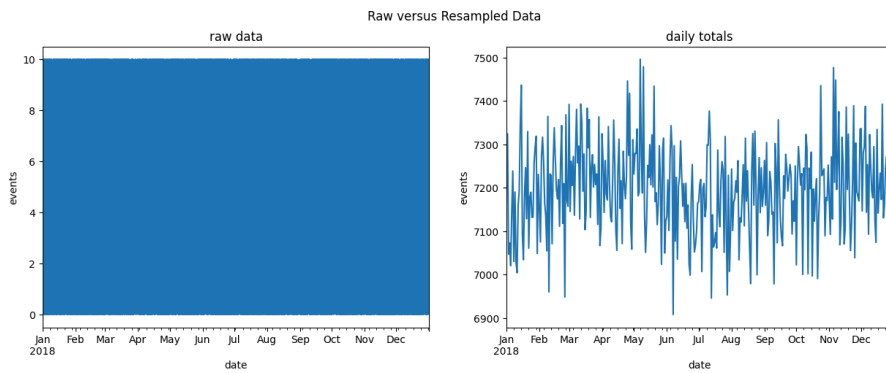
	open	high	low	close	volume
date					
2018-01-02	-7.91	-5.32	-7.38	-5.43	4577368.00
2018-01-03	-5.32	-4.12	-5.00	-3.61	-1108163.00
2018-01-04	-3.80	-2.59	-3.00	-3.54	1487839.00
2018-01-05	-1.35	-0.99	-0.70	-0.99	3044641.00
2018-01-08	-1.20	0.50	-1.05	0.51	8406139.00

▼ Resampling

```

1 import matplotlib.pyplot as plt
2 np.random.seed(0)
3 index = pd.date_range('2018-01-01', freq='T', periods=365*24*60)
4 raw = pd.DataFrame(
5     np.random.uniform(0, 10, size=index.shape[0]), index=index
6 )
7
8 fig, axes = plt.subplots(1, 2, figsize=(15, 5))
9 raw.plot(legend=False, ax=axes[0], title='raw data')
10 raw.resample('1D').sum().plot(legend=False, ax=axes[1], title='daily totals')
11 for ax in axes:
12     ax.set_xlabel('date')
13     ax.set_ylabel('events')
14
15 plt.suptitle('Raw versus Resampled Data')
16 plt.show()

```



```

1 print("Resampling the stock data per minute to get the daily frequency\n",
2       "Head: \n", stock_data_per_minute.head(),
3       "\n\nResampled: \n", stock_data_per_minute.resample('1D').agg({
4         'open': 'first',
5         'high': 'max',
6         'low': 'min',
7         'close': 'last',
8         'volume': 'sum'
9       }))
10 print("\n\n\nDownsampling to quarterly data:\n",
11        fb.resample('Q').mean())
12 print("\n\n\nQuarterly change from start to end:\n",
13        fb.drop(columns='trading_volume').resample('Q').apply(
14          lambda x: x.last('1D').values - x.first('1D').values
15        ))
16
17 melted_stock_data = pd.read_csv('/content/Time_Series/melted_stock_data.csv', index_col='date', parse_dates=True)
18 print("\n\n\nResampling Melted Stock prices:\n",
19        melted_stock_data.resample('1D').ohlc()['price'],
20        "\n\nPer 6 hours\n",
21        fb.resample('6H').asfreq().head(),
22        "\n\nHandle NaN Values using:\n",
23        "\nPad(): \n", fb.resample('6H').pad().head(),
24        "\n\nfillna(): \n", fb.resample('6H').fillna('nearest').head(),
25        "\n\nasfreq() and assign(): \n",
26        fb.resample('6H').asfreq().assign(
27          volume=lambda x: x.volume.fillna(0),
28          close=lambda x: x.close.fillna(method='ffill'),
29          open=lambda x: np.where(x.open.isnull(), x.close, x.open),
30          high=lambda x: np.where(x.high.isnull(), x.close, x.high),
31          low=lambda x: np.where(x.low.isnull(), x.close, x.low)).head())

```

Resampling the stock data per minute to get the daily frequency
Head:

	open	high	low	close	volume
date					
2019-05-20 09:30:00	181.62	181.62	181.62	181.62	159049.00
2019-05-20 09:31:00	182.61	182.61	182.61	182.61	468017.00
2019-05-20 09:32:00	182.75	182.75	182.75	182.75	97258.00
2019-05-20 09:33:00	182.95	182.95	182.95	182.95	43961.00
2019-05-20 09:34:00	183.06	183.06	183.06	183.06	79562.00

Resampled:

	open	high	low	close	volume
date					
2019-05-20	181.62	184.18	181.62	182.72	10044838.00
2019-05-21	184.53	185.58	183.97	184.82	7198405.00
2019-05-22	184.81	186.56	184.01	185.32	8412433.00
2019-05-23	182.50	183.73	179.76	180.87	12479171.00
2019-05-24	182.33	183.52	181.04	181.06	7686030.00

Downsampling to quarterly data:

	open	high	low	close	volume
--	------	------	-----	-------	--------

```

date
2018-03-31 179.47 181.79 177.04 179.55 32926396.70
2018-06-30 180.37 182.28 178.60 180.70 24055317.75
2018-09-30 180.81 182.89 178.96 181.03 27019824.76
2018-12-31 145.27 147.62 142.72 144.87 26974331.73

```

Quarterly change from start to end:

```

date
2018-03-31 [[-22.53, -20.160000000000025, -23.4100000000...
2018-06-30 [[39.509999999999999, 38.399700000000024, 39.84...
2018-09-30 [[-25.039999999999992, -28.659999999999997, -2...
2018-12-31 [[-28.580000000000013, -31.24000000000001, -31...
Freq: Q-DEC, dtype: object

```

Resampling Melted Stock prices:

```

      open  high  low  close
date
2019-05-20 181.62 184.18 181.62 182.72
2019-05-21 184.53 185.58 183.97 184.82
2019-05-22 184.81 186.56 184.01 185.32
2019-05-23 182.50 183.73 179.76 180.87
2019-05-24 182.33 183.52 181.04 181.06

```

Per 6 hours

```

      open  high  low  close  volume trading_volume
date
2018-01-02 00:00:00 177.68 181.58 177.55 181.42 18151903.00      low
2018-01-02 06:00:00   NaN   NaN   NaN   NaN       NaN       NaN
2018-01-02 12:00:00   NaN   NaN   NaN   NaN       NaN       NaN
2018-01-02 18:00:00   NaN   NaN   NaN   NaN       NaN       NaN

```

Merging

```

1 import sqlite3
2
3 with sqlite3.connect('/content/Time_Series/stocks.db') as connection:
4     fb_prices = pd.read_sql(
5         'SELECT * FROM fb_prices', connection,
6         index_col='date', parse_dates=['date']
7     )
8     aapl_prices = pd.read_sql(
9         'SELECT * FROM aapl_prices', connection,
10        index_col='date', parse_dates=['date']
11    )
12
13 print("Merging two data frames' prices column according to the datetime",
14       "index with a 30 seconds tolerance\n",
15       "\nUsing asof():\n",
16       pd.merge_asof(
17         fb_prices, aapl_prices,
18         left_index=True, right_index=True,
19         direction='nearest', tolerance=pd.Timedelta(30, unit='s')).head(),
20       "\n\nUsing merge_ordered():\n",
21       pd.merge_ordered(
22         fb_prices.reset_index(), aapl_prices.reset_index()).set_index('date').head())
23
24 print("\n\nFilling NaN Values\n",
25       pd.merge_ordered(fb_prices.reset_index(), aapl_prices.reset_index(),
26         fill_method='ffill').set_index('date').head())

```

Merging two data frames' prices column according to the datetime index with a 30 seconds tolerance

Using asof():

```

      FB  AAPL
date
2019-05-20 09:30:00 181.62 183.52
2019-05-20 09:31:00 182.61   NaN
2019-05-20 09:32:00 182.75 182.87
2019-05-20 09:33:00 182.95 182.50
2019-05-20 09:34:00 183.06 182.11

```

Using merge_ordered():

```

      FB  AAPL
date
2019-05-20 09:30:00 181.62 183.52

```

```

2019-05-20 09:31:00 182.61    NaN
2019-05-20 09:31:52    NaN 182.87
2019-05-20 09:32:00 182.75    NaN
2019-05-20 09:32:36    NaN 182.50

```

Filling NaN Values

```

           FB    AAPL
date
2019-05-20 09:30:00 181.62 183.52
2019-05-20 09:31:00 182.61 183.52
2019-05-20 09:31:52 182.61 182.87
2019-05-20 09:32:00 182.75 182.87
2019-05-20 09:32:36 182.75 182.50

```

8.1.4 Data Analysis

Provide some comments here about the results of the procedures.

- Doing the procedure took me some effort and time, as I had to painstakingly search for the meaning and logic behind every unfamiliar method and function that I encountered in order for me to understand each and everyone of them effectively and make it worth my time and energy. In the end, after completing everything in the procedures, I left with an ample amount of knowledge and understanding behind the logic and use of every function and method in this Hands-On Activity. But, even though I understand the theory behind them, I am still experiencing a hard time making the codes themselves from scratch without aid, as it is impossible for me to remember every single one of the numerous methods and functions that I understood in the procedures. Nevertheless, the results of the procedures are perfect as according to the one that was presented in the modules 8.1 to 8.5, and I have gained a deeper understanding on what and how one should clean data.

✓ 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.

```

1 import pandas as pd
2 import numpy as np
3 earthquakes = pd.read_csv('/content/Supplementary_Act/earthquakes.csv')
4 faang = pd.read_csv('/content/Supplementary_Act/faang.csv')
5
6 rows_with_japan = len(earthquakes.loc[earthquakes['parsed_place'].str.contains('Japan'), 'place'].value_counts())
7 print("Number of Earthquakes in Japan: ", rows_with_japan, "\n")
8
9 japanEarthquakes = earthquakes[(earthquakes['magType'] == 'mb') & (earthquakes['mag'] >= 4.9) & (earthquakes['place'].str.contains('Japan'))
10 num_earthquakes_in_japan = len(japanEarthquakes)
11 print("Number of earthquakes in Japan with a magnitude of 4.9 or greater: ",
12       num_earthquakes_in_japan, "\n",
13       japanEarthquakes)

```

Number of Earthquakes in Japan: 55

```

Number of earthquakes in Japan with a magnitude of 4.9 or greater: 4
   mag magType      time      place tsunami \
1563  4.9      mb 1538977532250 293km ESE of Iwo Jima, Japan      0
2576  5.4      mb 1538697528010  37km E of Tomakomai, Japan      0
3072  4.9      mb 1538579732490  15km ENE of Hasaki, Japan      0
3632  4.9      mb 1538450871260  53km ESE of Hitachi, Japan      0

   parsed_place
1563      Japan
2576      Japan
3072      Japan
3632      Japan

```

✓ 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.

```
1 magnitude_bins = pd.cut(earthquakes.loc[earthquakes['magType'] == 'ml', 'mag'],
2                           bins=np.arange(0, 11),
3                           labels=[f'{i}-{i+1}' for i in range(10)])
4 earthquakes_per_magnitude = magnitude_bins.value_counts().sort_index()
5 print(earthquakes_per_magnitude)
6
```

0-1	2207
1-2	3105
2-3	862
3-4	122
4-5	2
5-6	1
6-7	0
7-8	0
8-9	0
9-10	0

Name: mag, dtype: int64

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

```
1 faang_agg_groupby = faang.groupby('ticker').resample('M').agg({
2     'open': 'mean',
3     'high': 'max',
4     'low': 'min',
5     'close': 'mean',
6     'volume': 'sum'
7 })
8 print(faang_agg_groupby)
9 #The code above doesn't work because we first have to set the date as the index after converting it into DateTime format
```

```
-----
TypeError                                Traceback (most recent call last)
<ipython-input-5-58c4251796b2> in <cell line: 1>()
----> 1 faang_agg_groupby = faang.groupby('ticker').resample('M').agg({
      2     'open': 'mean',
      3     'high': 'max',
      4     'low': 'min',
      5     'close': 'mean',

-----
↕ 2 frames -----
/usr/local/lib/python3.10/dist-packages/pandas/core/resample.py in _get_resampler(self,
obj, kind)
    1723     )
    1724
-> 1725     raise TypeError(
    1726         "Only valid with DatetimeIndex, "
    1727         "TimedeltaIndex or PeriodIndex, "

TypeError: Only valid with DatetimeIndex, TimedeltaIndex or PeriodIndex, but got an
instance of 'RangeIndex'
```

```
1 faang['date'] = pd.to_datetime(faang['date'])
2 faang.set_index('date', inplace=True)
3 faang_agg_groupby = faang.groupby('ticker').resample('M').agg({
4     'open': 'mean',
5     'high': 'max',
6     'low': 'min',
7     'close': 'mean',
8     'volume': 'sum'
9 })
10 print(faang_agg_groupby)
```


	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	213.6351	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
	2018-12-31	1572.922105	1778.3400	1307.0000	1559.443158	154812304
FB	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	181.910509	192.7200	170.2300	182.930000	401144183
	2018-06-30	194.974067	203.5500	186.4300	195.267619	387265765
	2018-07-31	199.332143	218.6200	166.5600	199.967143	652763259
	2018-08-31	177.598443	188.3000	170.2700	177.491957	549016789
	2018-09-30	164.232895	173.8900	158.8656	164.377368	500468912
	2018-10-31	154.873261	165.8800	139.0300	154.187826	622446235
	2018-11-30	141.762857	154.1300	126.8500	141.635714	518150415
	2018-12-31	137.529474	147.1900	123.0200	137.161053	558786249
GOOG	2018-01-31	1127.200952	1186.8900	1045.2300	1130.770476	28738485
	2018-02-28	1088.629474	1174.0000	992.5600	1088.206842	42384105
	2018-03-31	1096.108095	1177.0500	980.6400	1091.490476	45430049
	2018-04-30	1038.415238	1094.1600	990.3700	1035.696190	41773275
	2018-05-31	1064.021364	1110.7500	1006.2900	1069.275909	31849196
	2018-06-30	1136.396190	1186.2900	1096.0100	1137.626667	32103642
	2018-07-31	1183.464286	1273.8900	1093.8000	1187.590476	31953386
	2018-08-31	1226.156957	1256.5000	1188.2400	1225.671739	28820379
	2018-09-30	1176.878421	1212.9900	1146.9100	1175.808947	28863199
	2018-10-31	1116.082174	1209.9600	995.8300	1110.940435	48496167
	2018-11-30	1054.971429	1095.5700	996.0200	1056.162381	36735570
	2018-12-31	1042.620000	1124.6500	970.1100	1037.420526	40256461
NFLX	2018-01-31	231.269286	286.8100	195.4200	232.908095	238377533
	2018-02-28	270.873158	297.3600	236.1100	271.443684	184585819
	2018-03-31	312.712857	333.9800	275.9000	312.228095	263449491
	2018-04-30	309.129529	338.8200	271.2239	307.466190	262064417
	2018-05-31	329.779759	356.1000	305.7300	331.536818	142051114
	2018-06-30	384.557595	423.2056	352.8200	384.133333	244032001
	2018-07-31	380.969090	419.7700	328.0000	381.515238	305487432
	2018-08-31	345.409591	376.8085	310.9280	346.257826	213144082
	2018-09-30	363.326842	383.2000	335.8300	362.641579	170832156
	2018-10-31	340.025348	386.7999	271.2093	335.445652	363589920
	2018-11-30	290.643333	332.0499	250.0000	290.344762	257126498
	2018-12-31	266.309474	298.7200	231.2300	265.302368	234304628

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.

```
1 earthquake_crosstab = pd.crosstab(
2     earthquakes['tsunami'],
3     earthquakes['magType'],
4     values=earthquakes['mag'],
5     aggfunc='max'
6 )
7
8 print(earthquake_crosstab)
```

magType	mb	mb_lg	md	mh	ml	ms_20	mw	mbw	mwr	mwW
tsunami										
0	5.6	3.5	4.11	1.1	4.2	NaN	3.83	5.8	4.8	6.0
1	6.1	NaN	NaN	NaN	5.1	5.7	4.41	NaN	NaN	7.5

5

Calculate the rolling 60-day aggregations of OHLC data by ticker for the FAANG data. Use the same aggregations as exercise no. 3.

```
1 faang_rolling = faang.groupby('ticker').rolling('60D').agg({
2     'open': 'mean',
3     'high': 'max',
4     'low': 'min',
5     'close': 'mean',
6     'volume': 'sum'
7 })
8 print(faang_rolling)
```

ticker	date	open	high	low	close	volume
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 x 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 volume traded data.

```
1 faang_pivot_table = faang.pivot_table(index='ticker',
2                                     aggfunc='mean',
3                                     values=['open', 'high', 'low', 'close', 'volume'])
4 print(faang_pivot_table)
5
```

ticker	close	high	low	open	volume
AAPL	186.986218	188.906858	185.135729	187.038674	3.402145e+07
AMZN	1641.726175	1662.839801	1619.840398	1644.072669	5.649563e+06
FB	171.510936	173.615298	169.303110	171.454424	2.768798e+07
GOOG	1113.225139	1125.777649	1101.001594	1113.554104	1.742645e+06
NFLX	319.290299	325.224583	313.187273	319.620533	1.147030e+07

7

Calculate the Z-scores for each numeric column of Netflix's data (ticker is NFLX) using apply().

```
1 netflixData = faang[faang['ticker'] == 'NFLX']
2 netflix_numeric_columns = netflixData.select_dtypes(include=[np.number])
3 netflix_ZScores = netflix_numeric_columns.apply(lambda x: (x - x.mean()) / x.std())
4 print("Z-scores for each numeric column of netflix's data:\n", netflix_ZScores)
5 print("\n\nZ-score combined with ticker:\n",
6       pd.concat([netflixData.drop(columns=netflix_numeric_columns.columns), netflix_ZScores], axis=1))
```

Z-scores for each numeric column of netflix's data:

date	open	high	low	close	volume
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
2018-01-04	-2.296272	-2.406077	-2.234616	-2.323429	-0.959287
2018-01-05	-2.275014	-2.345607	-2.202087	-2.234303	-0.782331
2018-01-08	-2.218934	-2.295113	-2.143759	-2.192192	-1.038531
...
2018-12-24	-1.571478	-1.518366	-1.627197	-1.745946	-0.339003
2018-12-26	-1.735063	-1.439978	-1.677339	-1.341402	0.517040
2018-12-27	-1.407286	-1.417785	-1.495805	-1.302664	0.134868
2018-12-28	-1.248762	-1.289018	-1.297285	-1.292137	-0.085164
2018-12-31	-1.203817	-1.122354	-1.088531	-1.055420	0.359444

[251 rows x 5 columns]

Z-score combined with ticker:

	date	ticker	open	high	low	close	volume
2018-01-02	NFLX	-2.500753	-2.516023	-2.410226	-2.416644	-0.088760	
2018-01-03	NFLX	-2.380291	-2.423180	-2.285793	-2.335286	-0.507606	
2018-01-04	NFLX	-2.296272	-2.406077	-2.234616	-2.323429	-0.959287	
2018-01-05	NFLX	-2.275014	-2.345607	-2.202087	-2.234303	-0.782331	
2018-01-08	NFLX	-2.218934	-2.295113	-2.143759	-2.192192	-1.038531	
...	
2018-12-24	NFLX	-1.571478	-1.518366	-1.627197	-1.745946	-0.339003	
2018-12-26	NFLX	-1.735063	-1.439978	-1.677339	-1.341402	0.517040	
2018-12-27	NFLX	-1.407286	-1.417785	-1.495805	-1.302664	0.134868	
2018-12-28	NFLX	-1.248762	-1.289018	-1.297285	-1.292137	-0.085164	
2018-12-31	NFLX	-1.203817	-1.122354	-1.088531	-1.055420	0.359444	

[251 rows x 6 columns]

▼ 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

```
1 events_df = pd.DataFrame({
2     'ticker': ['FB',
3               'FB',
4               'FB'],
5     'date': ['2018-07-25',
6              '2018-03-19',
7              '2018-03-20'],
8     'event': ['Disappointing user growth announced after close.',
9              'Cambridge Analytica story',
10             'FTC investigation']
11 })
12
13 events_df['date'] = pd.to_datetime(events_df['date'])
14 merged_data = faang.merge(events_df, how='outer', on=['date', 'ticker'])
15 print(merged_data)
```


	date	ticker	open	high	low	close	volume	event
0	2018-01-02	FB	177.68	181.58	177.5500	181.42	18151903	NaN
1	2018-01-03	FB	181.88	184.78	181.3300	184.67	16886563	NaN
2	2018-01-04	FB	184.90	186.21	184.0996	184.33	13880896	NaN
3	2018-01-05	FB	185.59	186.90	184.9300	186.85	13574535	NaN
4	2018-01-08	FB	187.20	188.90	186.3300	188.28	17994726	NaN
...
1250	2018-12-24	GOOG	973.90	1003.54	970.1100	976.22	1590328	NaN
1251	2018-12-26	GOOG	989.01	1040.00	983.0000	1039.46	2373270	NaN
1252	2018-12-27	GOOG	1017.15	1043.89	997.0000	1043.88	2109777	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 x 8 columns]

▼ 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 values for the first date in the data for that ticker. This is referred to as an index, and the data for the first date is the base (<https://ec.europa.eu/eurostat/statistics-explained/index.php/Beginners:Statisticalconcept-Indexandbaseyear>). When data is in this format, we can easily see growth over time. Hint: transform() can take a function name.

```
1 faang_transformed = faang.groupby('ticker').transform(lambda x: x / x.iloc[0])
2 print(faang_transformed)
```



	open	high	low	close	volume
date					
2018-01-02	1.000000	1.000000	1.000000	1.000000	1.000000
2018-01-03	1.023638	1.017623	1.021290	1.017914	0.930292
2018-01-04	1.040635	1.025498	1.036889	1.016040	0.764707
2018-01-05	1.044518	1.029298	1.041566	1.029931	0.747830
2018-01-08	1.053579	1.040313	1.049451	1.037813	0.991341
...
2018-12-24	0.928993	0.940578	0.928131	0.916638	1.285047
2018-12-26	0.943406	0.974750	0.940463	0.976019	1.917695
2018-12-27	0.970248	0.978396	0.953857	0.980169	1.704782
2018-12-28	1.001221	0.989334	0.988395	0.973784	1.142383
2018-12-31	1.002499	0.986653	0.979296	0.972404	1.206986

[1255 rows x 5 columns]

Double-click (or enter) to edit