

✓ Aggregations with pandas and numpy

About the Data

In this notebook, we will be working with 2 data sets:

- Facebook's stock price throughout 2018 (obtained using the stock_analysis package).
- daily weather data for NYC from the National Centers for Environmental Information (NCEI) API.

Note: The NCEI is part of the National Oceanic and Atmospheric Administration (NOAA) and, as you can see from the URL for the API, this resource was created when the NCEI was called the NCDC. Should the URL for this resource change in the future, you can search for the NCEI weather API to find the updated one.

Background on the weather data

Data meanings:

- AWND : average wind speed
- PRCP : precipitation in millimeters
- SNOW : snowfall in millimeters
- SNWD : snow depth in millimeters
- TMAX : maximum daily temperature in Celsius
- TMIN : minimum daily temperature in Celsius

> Setup

[] ↳ 4 cells hidden

> Summarizing DataFrames

We learned about `agg()` in the dataframe operations notebook when we learned about window calculations; however, we can call this on the dataframe directly to aggregate its contents into a single series:

[] ↳ 7 cells hidden

> Using `groupby()`

Often we won't want to aggregate on the entire dataframe, but on groups within it. For this purpose, we can run `groupby()` before the aggregation. If we group by the `trading_volume` column, we will get a row for each of the values it takes on:

[] ↳ 23 cells hidden

✓ Pivot tables and crosstabs

We saw pivots in before; however, we weren't able to provide any aggregations. With `pivot_table()`, we get the mean by default as the `aggfunc`. In its simplest form, we provide a column to place along the columns:

```
1 fb.pivot_table(columns='trading_volume')
```

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

By placing the trading volume in the index, we get the aggregation from the first example in the group by section above:

```
1 fb.pivot_table(index='trading_volume')
```

	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

With pivot(), we also weren't able to handle multi-level indices or indices with repeated values. For this reason we haven't been able to put the weather data in the wide format. The pivot_table() method solves this issue:

```
1 weather.reset_index().pivot_table(
2     index=['date', 'station', 'station_name'],
3     columns='datatype',
4     values='value',
5     aggfunc='median'
6 ).reset_index().tail()
```

	datatype	date	station	station_name	AWND	DAPR	MDPR	PGTM	PRCP	SN
28740		2018-12-31	GHCND:USW00054787	FARMINGDALE REPUBLIC AIRPORT, NY US	5.00	NaN	NaN	2052.00	28.70	NaN
28741		2018-12-31	GHCND:USW00094728	NY CITY CENTRAL PARK, NY US	NaN	NaN	NaN	NaN	25.90	0.0
28742		2018-12-31	GHCND:USW00094741	TETERBORO AIRPORT, NJ US	1.70	NaN	NaN	1954.00	29.20	NaN

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Mounted at /content/drive

We can use the pd.crosstab() function to create a frequency table. For example, if we want to see how many low-, medium-, and high-volume trading days Facebook stock had each month, we can use crosstab:

```
1 pd.crosstab(
2     index=fb.trading_volume,
3     columns=fb.index.month,
4     colnames=['month'] # name the columns index
5 )
```

	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

We can normalize with the row or column totals with the normalize parameter. This shows percentage of the total:

```
1 pd.crosstab(
2     index=fb.trading_volume,
3     columns=fb.index.month,
4     colnames=['month'],
5     normalize='columns'
6 )
```

month	1	2	3	4	5	6	7	8	9	10	11	12
trading_volume												
low	0.95	1.00	0.71	0.95	1.00	1.00	0.86	1.00	1.00	1.00	1.00	1.00
med	0.05	0.00	0.19	0.05	0.00	0.00	0.10	0.00	0.00	0.00	0.00	0.00
high	0.00	0.00	0.10	0.00	0.00	0.00	0.05	0.00	0.00	0.00	0.00	0.00

If we want to perform a calculation other than counting the frequency, we can pass the column to run the calculation on to values and the function to use to aggfunc :

```
1 pd.crosstab(
2     index=fb.trading_volume,
3     columns=fb.index.month,
4     colnames=['month'],
5     values=fb.close,
6     aggfunc=np.mean
7 )
```

month	1	2	3	4	5	6	7	8	9	1
trading_volume										
low	185.24	180.27	177.07	163.29	182.93	195.27	201.92	177.49	164.38	154.1
med	179.37	NaN	164.76	174.16	NaN	NaN	194.28	NaN	NaN	Na
high	NaN	NaN	164.11	NaN	NaN	NaN	176.26	NaN	NaN	Na

We can also get row and column subtotals with the margins parameter. Let's count the number of times each station recorded snow per month and include the subtotals:

```
1 snow_data = weather.query('datatype == "SNOW"')
2 pd.crosstab(
3     index=snow_data.station_name,
4     columns=snow_data.index.month,
5     colnames=['month'],
6     values=snow_data.value,
7     aggfunc=lambda x: (x > 0).sum(),
8     margins=True, # show row and column subtotals
9     margins_name='total observations of snow' # name the subtotals
10 )
```

month	1	2	3	4	5	6	7	8	9	10	11	12	total observations of snow
station_name													
ALBERTSON 0.2 SSE, NY US	3.00	1.00	3.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	9
AMITYVILLE 0.1 WSW, NY US	1.00	0.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	3
AMITYVILLE 0.6 NNE, NY US	3.00	1.00	3.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	8
ARMONK 0.3 SE, NY US	6.00	4.00	6.00	3.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	3.00	23
BLOOMINGDALE 0.7 SSE, NJ US	2.00	1.00	3.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	8
...
WESTFIELD 0.6 NE, NJ US	3.00	0.00	4.00	1.00	0.00	NaN	0.00	0.00	0.00	NaN	1.00	NaN	9
WOODBIDGE TWP 1.1 ESE, NJ US	4.00	1.00	3.00	2.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	11
WOODBIDGE TWP 1.1 NNE, NJ US	2.00	1.00	3.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	7
WOODBIDGE TWP 3.0 NNW, NJ US	NaN	0.00	0.00	NaN	NaN	0.00	NaN	NaN	NaN	0.00	0.00	NaN	0
total observations of snow	190.00	97.00	237.00	81.00	0.00	0.00	0.00	0.00	0.00	0.00	49.00	13.00	667

99 rows × 13 columns

Comments and Conclusions

- In this module, I learned a lot about how to summarize all the data or specific groups of data in a dataframe by implementing a specific aggregation operation for each of the columns, rows, or groups using the `.agg()` and `.groupby()` methods, along with the utilization of the `pivot` and `crosstab` method to further arrange the affected data into the specified format.