

# Module2

February 25, 2023

<IPython.core.display.HTML object>

For module 2 we'll be looking at techniques for dealing with big data. In particular binning strategies and the datashader library (which possibly proves we'll never need to bin large data for visualization ever again.)

To demonstrate these concepts we'll be looking at the PLUTO dataset put out by New York City's department of city planning. PLUTO contains data about every tax lot in New York City.

PLUTO data can be downloaded from [here](#). Unzip them to the same directory as this notebook, and you should be able to read them in using this (or very similar) code. Also take note of the data dictionary, it'll come in handy for this assignment.

```
/var/folders/f4/v17stl257nn9w40qytyt496380000gn/T/ipykernel_24847/1232822874.py:1  
1: DtypeWarning:
```

```
Columns (21,22,24,26,28) have mixed types. Specify dtype option on import or set  
low_memory=False.
```

I'll also do some prep for the geographic component of this data, which we'll be relying on for datashader.

You're not required to know how I'm retrieving the latitude and longitude here, but for those interested: this dataset uses a flat x-y projection (assuming for a small enough area that the world is flat for easier calculations), and this needs to be projected back to traditional latitude and longitude.

## 0.1 Part 1: Binning and Aggregation

Binning is a common strategy for visualizing large datasets. Binning is inherent to a few types of visualizations, such as histograms and [2D histograms](#) (also check out their close relatives: [2D density plots](#) and the more general form: [heatmaps](#)).

While these visualization types explicitly include binning, any type of visualization used with aggregated data can be looked at in the same way. For example, lets say we wanted to look at building construction over time. This would be best viewed as a line graph, but we can still think of our results as being binned by year:

Something looks off... You're going to have to deal with this imperfect data to answer this first question.

But first: some notes on pandas. Pandas dataframes are a different beast than R dataframes, here are some tips to help you get up to speed:

---

Hello all, here are some pandas tips to help you guys through this homework:

**Indexing and Selecting:** `.loc` and `.iloc` are the analogs for base R subsetting, or `filter()` in dplyr

**Group By:** This is the pandas analog to `group_by()` and the appended function the analog to `summarize()`. Try out a few examples of this, and display the results in Jupyter. Take note of what's happening to the indexes, you'll notice that they'll become hierarchical. I personally find this more of a burden than a help, and this sort of hierarchical indexing leads to a fundamentally different experience compared to R dataframes. Once you perform an aggregation, try running the resulting hierarchical dataframe through a [reset\\_index\(\)](#).

**Reset\_index:** I personally find the hierarchical indexes more of a burden than a help, and this sort of hierarchical indexing leads to a fundamentally different experience compared to R dataframes. `reset_index()` is a way of restoring a dataframe to a flatter index style. Grouping is where you'll notice it the most, but it's also useful when you filter data, and in a few other split-apply-combine workflows. With pandas indexes are more meaningful, so use this if you start getting unexpected results.

Indexes are more important in Pandas than in R. If you delve deeper into the using python for data science, you'll begin to see the benefits in many places (despite the personal gripes I highlighted above.) One place these indexes come in handy is with time series data. The pandas docs have a [huge section](#) on datetime indexing. In particular, check out [resample](#), which provides time series specific aggregation.

**Merging, joining, and concatenation:** There's some overlap between these different types of merges, so use this as your guide. `Concat` is a single function that replaces `cbind` and `rbind` in R, and the results are driven by the indexes. Read through these examples to get a feel on how these are performed, but you will have to manage your indexes when you're using these functions. Merges are fairly similar to merges in R, similarly mapping to SQL joins.

**Apply:** This is explained in the "group by" section linked above. These are your analogs to the `plyr` library in R. Take note of the lambda syntax used here, these are anonymous functions in python. Rather than predefining a custom function, you can just define it inline using `lambda`.

Browse through the other sections for some other specifics, in particular reshaping and categorical data (pandas' answer to factors.) Pandas can take a while to get used to, but it is a pretty strong framework that makes more advanced functions easier once you get used to it. Rolling functions for example follow logically from the apply workflow (and led to the best google results ever when I first tried to find this out and googled "pandas rolling")

Google Wes McKinney's book "Python for Data Analysis," which is a cookbook style intro to pandas. It's an O'Reilly book that should be pretty available out there.

---

### 0.1.1 Question

After a few building collapses, the City of New York is going to begin investigating older buildings for safety. The city is particularly worried about buildings that were unusually tall when they were built, since best-practices for safety hadn't yet been determined. Create a graph that shows how many buildings of a certain number of floors were built in each year (note: you may want to use a log scale for the number of buildings). Find a strategy to bin buildings (It should be clear 20-29-story buildings, 30-39-story buildings, and 40-49-story buildings were first built in large numbers, but does it make sense to continue in this way as you get taller?)

---

### 0.1.2 Answer

First we can check what the range of values for numfloors is:

```
min      1.0
max     104.0
Name: numfloors, dtype: float64
```

Then we can visualize the distribution of numfloors using a histogram:

We can see the vast majority of buildings have 10 or fewer floors, but there are a few outliers. We can use a logarithmic scale to better visualize the distribution of the data:

From the early plot of buildings built by year it seems that certain years had few or no buildings built. Let's check the number of buildings built in each year using a subset of the data:

|        |       |
|--------|-------|
| 1900.0 | 6440  |
| 1901.0 | 22230 |
| 1902.0 | 476   |
| 1903.0 | 460   |
| 1904.0 | 522   |
| 1905.0 | 6914  |
| 1906.0 | 1026  |
| 1907.0 | 996   |
| 1908.0 | 839   |
| 1909.0 | 1540  |
| 1910.0 | 41983 |
| 1911.0 | 1083  |
| 1912.0 | 883   |
| 1913.0 | 732   |
| 1914.0 | 668   |
| 1915.0 | 15363 |
| 1916.0 | 687   |
| 1917.0 | 534   |
| 1918.0 | 291   |
| 1919.0 | 369   |
| 1920.0 | 87703 |
| 1921.0 | 1118  |
| 1922.0 | 1246  |

```

1923.0    1641
1924.0    2246
1925.0   69677
1926.0    3256
1927.0    3598
1928.0    4346
1929.0    2189
1930.0   74907
1931.0   31000
1932.0    1224
1933.0     946
1934.0     374
1935.0   25085
1936.0     643
1937.0     701
1938.0     764
1939.0     929
1940.0   37904
Name: yearbuilt, dtype: int64

```

We can see while there are no years with 0 buildings built in this range there is a lot of variance in the data. For example 1933 only has 946 buildings built, while 1931 has 31,000. This gives us good reason to aggregate the data by decade.

Now we can create a chart that shows the number of buildings built in each decade by number of floors. We will bin using the min and max floor numbers in increments of 10. We will also make sure to use the log scale for the y-axis to better visualize the data.

It is clear from the plot above that even using a logarithmic scale, the 1-20 floor buildings are dominating the plot. Also, the line scatter chart format gives us a lot of noise in the form of data resolution that we do not need in order to visualize the data. Let's filter out buildings with less than 20 floors and replot using stacked bar charts:

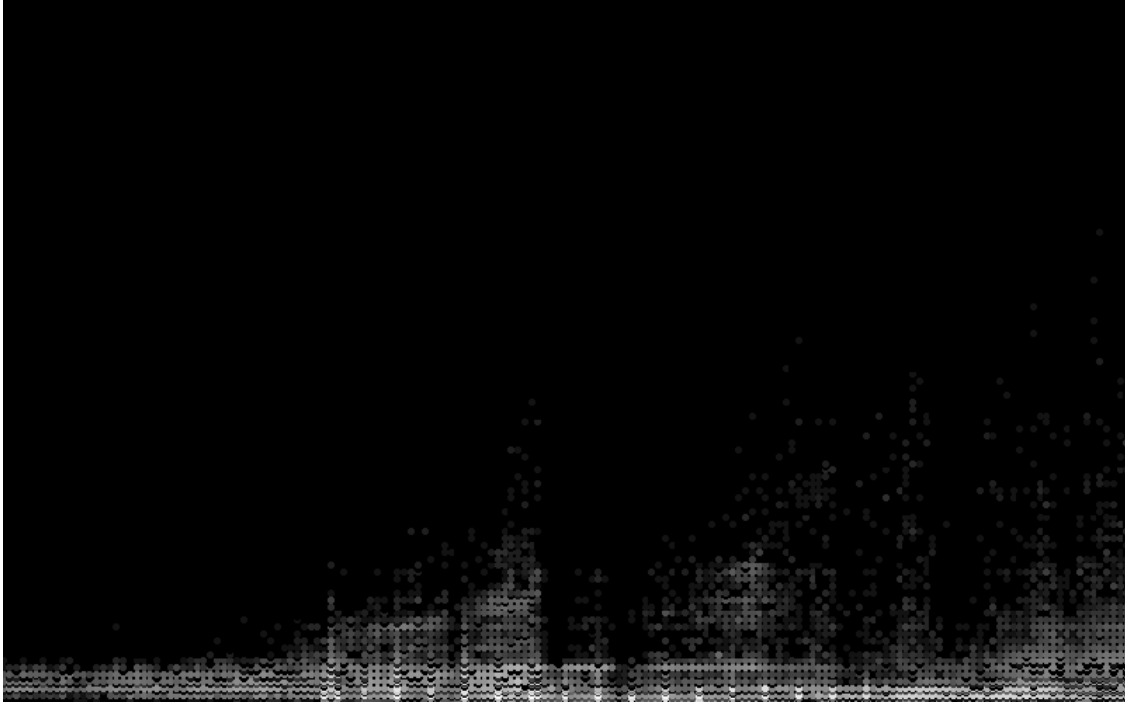
## 0.2 Part 2: Datashader

Datashader is a library from Anaconda that does away with the need for binning data. It takes in all of your datapoints, and based on the canvas and range returns a pixel-by-pixel calculations to come up with the best representation of the data. In short, this completely eliminates the need for binning your data.

As an example, lets continue with our question above and look at a 2D histogram of YearBuilt vs NumFloors:

This shows us the distribution, but it's subject to some biases discussed in the Anaconda notebook [Plotting Perils](#).

Here is what the same plot would look like in datashader:



That's technically just a scatterplot, but the points are smartly placed and colored to mimic what one gets in a heatmap. Based on the pixel size, it will either display individual points, or will color the points of denser regions.

Datashader really shines when looking at geographic information. Here are the latitudes and longitudes of our dataset plotted out, giving us a map of the city colored by density of structures:



Interestingly, since we're looking at structures, the large buildings of Manhattan show up as less dense on the map. The densest areas measured by number of lots would be single or multi family townhomes.

Unfortunately, Datashader doesn't have the best documentation. Browse through the examples from their [github repo](#). I would focus on the [visualization pipeline](#) and the [US Census Example](#) for the question below. Feel free to use my samples as templates as well when you work on this problem.

### 0.2.1 Question

You work for a real estate developer and are researching underbuilt areas of the city. After looking in the [Pluto data dictionary](#), you've discovered that all tax assessments consist of two parts: The assessment of the land and assessment of the structure. You reason that there should be a corre-

lation between these two values: more valuable land will have more valuable structures on them (more valuable in this case refers not just to a mansion vs a bungalow, but an apartment tower vs a single family home). Deviations from the norm could represent underbuilt or overbuilt areas of the city. You also recently read a really cool blog post about [bivariate choropleth maps](#), and think the technique could be used for this problem.

Datashader is really cool, but it's not that great at labeling your visualization. Don't worry about providing a legend, but provide a quick explanation as to which areas of the city are overbuilt, which areas are underbuilt, and which areas are built in a way that's properly correlated with their land value.

## 0.2.2 Answer

There is no separated value of the structural assessment in the dataset but we have the total and land assessment values so we can subtract the land value from the total and get what we assume to be the structural assessment value.

|   | borough | block | lot | cd    | bct2020   | bctcb2020    | ct2010 | cb2010 | \ |
|---|---------|-------|-----|-------|-----------|--------------|--------|--------|---|
| 0 | SI      | 1597  | 125 | 502.0 | 5029104.0 | 5.029104e+10 | 291.04 | 3007.0 |   |
| 2 | BK      | 4794  | 1   | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 3 | BK      | 1488  | 105 | 303.0 | 3037500.0 | 3.037500e+10 | 375.00 | 1001.0 |   |
| 4 | BK      | 4794  | 17  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 5 | BK      | 4794  | 78  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |

|   | schooldist | council | ... | plutomapid | firm07_flag | pfirm15_flag | version | \ |
|---|------------|---------|-----|------------|-------------|--------------|---------|---|
| 0 | 31.0       | 50.0    | ... | 1          | NaN         | NaN          | 22v2    |   |
| 2 | 17.0       | 41.0    | ... | 1          | NaN         | NaN          | 22v2    |   |
| 3 | 16.0       | 41.0    | ... | 1          | NaN         | NaN          | 22v2    |   |
| 4 | 17.0       | 41.0    | ... | 1          | NaN         | NaN          | 22v2    |   |
| 5 | 17.0       | 41.0    | ... | 1          | NaN         | NaN          | 22v2    |   |

|   | dcpedited | latitude  | longitude  | notes | decade | assessstruct |
|---|-----------|-----------|------------|-------|--------|--------------|
| 0 | NaN       | 40.611140 | -74.164376 | NaN   | 1960.0 | 46020.0      |
| 2 | NaN       | 40.661794 | -73.942532 | NaN   | 1890.0 | 308700.0     |
| 3 | NaN       | 40.686484 | -73.920169 | NaN   | 1990.0 | 40740.0      |
| 4 | NaN       | 40.661859 | -73.941991 | NaN   | 1990.0 | 46380.0      |
| 5 | NaN       | 40.661517 | -73.942539 | NaN   | 2000.0 | 78480.0      |

[5 rows x 94 columns]

It wasn't clear to me how to generate the color palettes, so I used this tool <https://colorbrewer2.org/#type=diverging&scheme=PiYG&n=9/>

Going to use **Quantile** segmentation for the data which because we are segmenting into 3 distinct groups means simply dividing the data set into three parts and labeling it.

```
[0.0, 11580.0]
[11580.0, 18120.0]
```

[18120.0, 3205633833.0]

|   | borough | block | lot | cd    | bct2020   | bctcb2020    | ct2010 | cb2010 | \ |
|---|---------|-------|-----|-------|-----------|--------------|--------|--------|---|
| 0 | SI      | 1597  | 125 | 502.0 | 5029104.0 | 5.029104e+10 | 291.04 | 3007.0 |   |
| 2 | BK      | 4794  | 1   | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 3 | BK      | 1488  | 105 | 303.0 | 3037500.0 | 3.037500e+10 | 375.00 | 1001.0 |   |
| 4 | BK      | 4794  | 17  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 5 | BK      | 4794  | 78  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |

|   | schooldist | council | ... | firm07_flag | pfirm15_flag | version | dcpedited | \ |
|---|------------|---------|-----|-------------|--------------|---------|-----------|---|
| 0 | 31.0       | 50.0    | ... | NaN         | NaN          | 22v2    | NaN       |   |
| 2 | 17.0       | 41.0    | ... | NaN         | NaN          | 22v2    | NaN       |   |
| 3 | 16.0       | 41.0    | ... | NaN         | NaN          | 22v2    | NaN       |   |
| 4 | 17.0       | 41.0    | ... | NaN         | NaN          | 22v2    | NaN       |   |
| 5 | 17.0       | 41.0    | ... | NaN         | NaN          | 22v2    | NaN       |   |

|   | latitude  | longitude  | notes | decade | assessstruct | colorOne |
|---|-----------|------------|-------|--------|--------------|----------|
| 0 | 40.611140 | -74.164376 | NaN   | 1960.0 | 46020.0      | 1        |
| 2 | 40.661794 | -73.942532 | NaN   | 1890.0 | 308700.0     | 1        |
| 3 | 40.686484 | -73.920169 | NaN   | 1990.0 | 40740.0      | 2        |
| 4 | 40.661859 | -73.941991 | NaN   | 1990.0 | 46380.0      | 3        |
| 5 | 40.661517 | -73.942539 | NaN   | 2000.0 | 78480.0      | 3        |

[5 rows x 95 columns]

[0.0, 33300.0]

[33300.0, 60000.0]

[60000.0, 4343286717.0]

|   | borough | block | lot | cd    | bct2020   | bctcb2020    | ct2010 | cb2010 | \ |
|---|---------|-------|-----|-------|-----------|--------------|--------|--------|---|
| 0 | SI      | 1597  | 125 | 502.0 | 5029104.0 | 5.029104e+10 | 291.04 | 3007.0 |   |
| 2 | BK      | 4794  | 1   | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 3 | BK      | 1488  | 105 | 303.0 | 3037500.0 | 3.037500e+10 | 375.00 | 1001.0 |   |
| 4 | BK      | 4794  | 17  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 5 | BK      | 4794  | 78  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |

|   | schooldist | council | ... | pfirm15_flag | version | dcpedited | latitude  | \ |
|---|------------|---------|-----|--------------|---------|-----------|-----------|---|
| 0 | 31.0       | 50.0    | ... | NaN          | 22v2    | NaN       | 40.611140 |   |
| 2 | 17.0       | 41.0    | ... | NaN          | 22v2    | NaN       | 40.661794 |   |
| 3 | 16.0       | 41.0    | ... | NaN          | 22v2    | NaN       | 40.686484 |   |
| 4 | 17.0       | 41.0    | ... | NaN          | 22v2    | NaN       | 40.661859 |   |
| 5 | 17.0       | 41.0    | ... | NaN          | 22v2    | NaN       | 40.661517 |   |

|   | longitude  | notes | decade | assessstruct | colorOne | colorTwo |
|---|------------|-------|--------|--------------|----------|----------|
| 0 | -74.164376 | NaN   | 1960.0 | 46020.0      | 1        | B        |
| 2 | -73.942532 | NaN   | 1890.0 | 308700.0     | 1        | C        |
| 3 | -73.920169 | NaN   | 1990.0 | 40740.0      | 2        | B        |
| 4 | -73.941991 | NaN   | 1990.0 | 46380.0      | 3        | B        |



5 -73.942539 NaN 2000.0 78480.0 3 C

[5 rows x 96 columns]

|   | borough | block | lot | cd    | bct2020   | bctcb2020    | ct2010 | cb2010 | \ |
|---|---------|-------|-----|-------|-----------|--------------|--------|--------|---|
| 0 | SI      | 1597  | 125 | 502.0 | 5029104.0 | 5.029104e+10 | 291.04 | 3007.0 |   |
| 2 | BK      | 4794  | 1   | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 3 | BK      | 1488  | 105 | 303.0 | 3037500.0 | 3.037500e+10 | 375.00 | 1001.0 |   |
| 4 | BK      | 4794  | 17  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |
| 5 | BK      | 4794  | 78  | 309.0 | 3080600.0 | 3.080600e+10 | 806.00 | 2000.0 |   |

|   | schooldist | council | ... | version | dcpedited | latitude  | longitude  | notes | \ |
|---|------------|---------|-----|---------|-----------|-----------|------------|-------|---|
| 0 | 31.0       | 50.0    | ... | 22v2    | NaN       | 40.611140 | -74.164376 | NaN   |   |
| 2 | 17.0       | 41.0    | ... | 22v2    | NaN       | 40.661794 | -73.942532 | NaN   |   |
| 3 | 16.0       | 41.0    | ... | 22v2    | NaN       | 40.686484 | -73.920169 | NaN   |   |
| 4 | 17.0       | 41.0    | ... | 22v2    | NaN       | 40.661859 | -73.941991 | NaN   |   |
| 5 | 17.0       | 41.0    | ... | 22v2    | NaN       | 40.661517 | -73.942539 | NaN   |   |

|   | decade | assessstruct | colorOne | colorTwo | combined |
|---|--------|--------------|----------|----------|----------|
| 0 | 1960.0 | 46020.0      | 1        | B        | 1B       |
| 2 | 1890.0 | 308700.0     | 1        | C        | 1C       |
| 3 | 1990.0 | 40740.0      | 2        | B        | 2B       |
| 4 | 1990.0 | 46380.0      | 3        | B        | 3B       |
| 5 | 2000.0 | 78480.0      | 3        | C        | 3C       |

[5 rows x 97 columns]

