

Module2

September 23, 2018

```
In [112]: import datashader as ds
import datashader.transfer_functions as tf
import datashader.glyphs
from datashader import reductions
from datashader.core import bypixel
from datashader.utils import lnglat_to_meters as webm, export_image
from datashader.colors import colormap_select, Greys9, viridis, inferno
import copy

from pyproj import Proj, transform
import numpy as np
import pandas as pd
import urllib
import json
import datetime
import colorlover as cl

import plotly.plotly as py
import plotly.graph_objs as go
from plotly import tools

from shapely.geometry import Point, Polygon, shape
# In order to get shapely, you'll need to run [pip install shapely.geometry] from yo

from functools import partial

from IPython.display import GeoJSON
```

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.

```
In [113]: bk = pd.read_csv('PLUTO18v1/BK_18v1.csv', low_memory=False)
          bx = pd.read_csv('PLUTO18v1/BX_18v1.csv', low_memory=False)
          mn = pd.read_csv('PLUTO18v1/MN_18v1.csv', low_memory=False)
          qn = pd.read_csv('PLUTO18v1/QN_18v1.csv', low_memory=False)
          si = pd.read_csv('PLUTO18v1/SI_18v1.csv', low_memory=False)

          ny = pd.concat([bk, bx, mn, qn, si], ignore_index=True, sort=True)

          # Getting rid of some outliers
          ny = ny[(ny['YearBuilt'] > 1850) & (ny['YearBuilt'] < 2020) & (ny['NumFloors'] != 0)]
```

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.

```
In [114]: wgs84 = Proj("+proj=longlat +ellps=GRS80 +datum=NAD83 +no_defs")
          nyli = Proj("+proj=lcc +lat_1=40.66666666666666 +lat_2=41.03333333333333 +lat_0=40.1")
          ny['XCoord'] = 0.3048*ny['XCoord']
          ny['YCoord'] = 0.3048*ny['YCoord']
          ny['lon'], ny['lat'] = transform(nyli, wgs84, ny['XCoord'].values, ny['YCoord'].values)

          ny = ny[(ny['lon'] < -60) & (ny['lon'] > -100) & (ny['lat'] < 60) & (ny['lat'] > 20)]

          #Defining some helper functions for DataShader
          background = "black"
          export = partial(export_image, background = background, export_path="export")
          cm = partial(colormap_select, reverse=(background!="black"))
```

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:

```
In [115]: tools.set_credentials_file(username='agilharrysr', api_key='zuCRUxfK09iTCU1aIGDE')

          trace = go.Scatter(
              # I'm choosing BBL here because I know it's a unique key.
              x = ny.groupby('YearBuilt').count()['BBL'].index,
              y = ny.groupby('YearBuilt').count()['BBL']
          )
```

```

layout = go.Layout(
    xaxis = dict(title = 'Year Built'),
    yaxis = dict(title = 'Number of Lots Built')
)

fig = go.Figure(data = [trace], layout = layout)

py.iplot(fig, filename = 'ny-year-built')

```

High five! You successfully sent some data to your account on plotly. View your plot in your browser.

Out[115]: <plotly.tools.PlotlyDisplay object>

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 Response (Small Multiples)

This graph needs to visualize 3 dimensions; year, buildings, and number of floors. A heatmap is a possible candidate for this question, however it may be difficult for most users to quickly grasp the information. For this reason I will use the concept of small multiples as a solution to this question. Data will be binned by multiples of ten (up to 50) based on the number of floors. There are not many buildings having more than 50 floors, therefore the final bin will contain buildings with more than 50 floors. The bins will be placed on separate bar charts and clustered on a single image below.

In [116]: *# bin 1: Up to 10 floors*

```
lte10 = pd.DataFrame(ny.loc[ ny['NumFloors'] <= 10 ].groupby('YearBuilt').count()['B
lte10.columns = ['NumBuildings']
lte10['LogNumBuildings'] = lte10['NumBuildings']
lte10.loc[lte10['LogNumBuildings'] > 0, 'LogNumBuildings'] = np.log(lte10.loc[lte10[
trace1 = go.Bar(
    x = lte10.index,
    y = lte10['LogNumBuildings'].values,
    showlegend=False
)
```

bin 2: 11 to 20 floors

```
lte1120 = pd.DataFrame(ny.loc[ (ny['NumFloors'] > 10) & (ny['NumFloors'] <= 20) ].gr
lte1120.columns = ['NumBuildings']
lte1120['LogNumBuildings'] = lte1120['NumBuildings']
lte1120.loc[lte1120['LogNumBuildings'] > 0, 'LogNumBuildings'] = np.log(lte1120.loc[
trace2 = go.Bar(
    # I'm choosing BBL here because I know it's a unique key.
    x = lte1120.index,
```

```

        y = lte1120['LogNumBuildings'].values,
        showlegend=False
    )

    # bin 3: 21 to 30 floors
    lte2130 = pd.DataFrame(ny.loc[ (ny['NumFloors'] > 20) & (ny['NumFloors'] <= 30) ].groupby('YearBuilt').count()['BBL'])
    lte2130.columns = ['NumBuildings']
    lte2130['LogNumBuildings'] = lte2130['NumBuildings']
    lte2130.loc[lte2130['LogNumBuildings'] > 0, 'LogNumBuildings'] = np.log(lte2130.loc[lte2130['LogNumBuildings'] > 0, 'LogNumBuildings'])
    trace3 = go.Bar(
        # I'm choosing BBL here because I know it's a unique key.
        x = lte2130.index,
        y = lte2130['LogNumBuildings'].values,
        showlegend=False
    )

    # bin 4: 31 to 40 floors
    lte3140 = pd.DataFrame(ny.loc[ (ny['NumFloors'] > 30) & (ny['NumFloors'] <= 40) ].groupby('YearBuilt').count()['BBL'])
    lte3140.columns = ['NumBuildings']
    lte3140['LogNumBuildings'] = lte3140['NumBuildings']
    lte3140.loc[lte3140['LogNumBuildings'] > 0, 'LogNumBuildings'] = np.log(lte3140.loc[lte3140['LogNumBuildings'] > 0, 'LogNumBuildings'])
    trace4 = go.Bar(
        # I'm choosing BBL here because I know it's a unique key.
        x = lte3140.index,
        y = lte3140['LogNumBuildings'].values,
        showlegend=False
    )

    # bin 5: 41 to 50 floors
    lte4150 = pd.DataFrame(ny.loc[ (ny['NumFloors'] > 40) & (ny['NumFloors'] <= 50) ].groupby('YearBuilt').count()['BBL'])
    lte4150.columns = ['NumBuildings']
    lte4150['LogNumBuildings'] = lte4150['NumBuildings']
    lte4150.loc[lte4150['LogNumBuildings'] > 0, 'LogNumBuildings'] = np.log(lte4150.loc[lte4150['LogNumBuildings'] > 0, 'LogNumBuildings'])
    trace5 = go.Bar(
        # I'm choosing BBL here because I know it's a unique key.
        x = lte4150.index,
        y = lte4150['LogNumBuildings'].values,
        showlegend=False
    )

    # bin 6: more than 50 floors
    gt50 = pd.DataFrame(ny.loc[ ny['NumFloors'] > 50 ].groupby('YearBuilt').count()['BBL'])
    gt50.columns = ['NumBuildings']
    gt50['LogNumBuildings'] = gt50['NumBuildings']
    gt50.loc[gt50['LogNumBuildings'] > 0, 'LogNumBuildings'] = np.log(gt50.loc[gt50['LogNumBuildings'] > 0, 'LogNumBuildings'])
    trace6 = go.Bar(
        # I'm choosing BBL here because I know it's a unique key.
        x = gt50.index,

```

```

        y = gt50['LogNumBuildings'].values,
        showlegend=False
    )

    # set up clusters
    fig = tools.make_subplots(rows=2, cols=3,
                              subplot_titles=('1 - 10 Floors', '11 - 20 Floors',
                                                '21 - 30 Floors', '31 - 40 Floors',
                                                '41 - 50 Floors', 'Greater Than 50 Floors'))

    fig.append_trace(trace1, 1, 1)
    fig.append_trace(trace2, 1, 2)
    fig.append_trace(trace3, 1, 3)
    fig.append_trace(trace4, 2, 1)
    fig.append_trace(trace5, 2, 2)
    fig.append_trace(trace6, 2, 3)

    fig['layout'].update(title='Building Sizes by Year: Number of Floors')
    fig['layout']['yaxis1'].update(title = 'Log(# of Buildings)', range=[0, 12])
    fig['layout']['yaxis2'].update(range=[0, 12])
    fig['layout']['yaxis3'].update(range=[0, 12])
    fig['layout']['yaxis4'].update(title = 'Log(# of Buildings)', range=[0, 12])
    fig['layout']['yaxis5'].update(range=[0, 12])
    fig['layout']['yaxis6'].update(range=[0, 12])
    fig['layout']['xaxis5'].update(title = 'Year Built')
    py.iplot(fig, filename = 'ny-year-built-floors')

```

This is the format of your plot grid:

```

[ (1,1) x1,y1 ] [ (1,2) x2,y2 ] [ (1,3) x3,y3 ]
[ (2,1) x4,y4 ] [ (2,2) x5,y5 ] [ (2,3) x6,y6 ]

```

Out[116]: <plotly.tools.PlotlyDisplay object>

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:

```

In [7]: yearbins = 200
        floorbins = 200

```

```

yearBuiltCut = pd.cut(ny['YearBuilt'], np.linspace(ny['YearBuilt'].min(), ny['YearBuilt'].max(), yearbins))
numFloorsCut = pd.cut(ny['NumFloors'], np.logspace(1, np.log(ny['NumFloors'].max()), floorbins))

```

```

xlabels = np.floor(np.linspace(ny['YearBuilt'].min(), ny['YearBuilt'].max(), yearbins))
ylabels = np.floor(np.logspace(1, np.log(ny['NumFloors'].max()), floorbins))

data = [
    go.Heatmap(z = ny.groupby([numFloorsCut, yearBuiltCut])['BBL'].count().unstack().f
               colorscale = 'Greens', x = xlabels, y = ylabels)
]

py.iplot(data, filename = 'datashader-2d-hist')

```

Out[7]: <plotly.tools.PlotlyDisplay object>

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

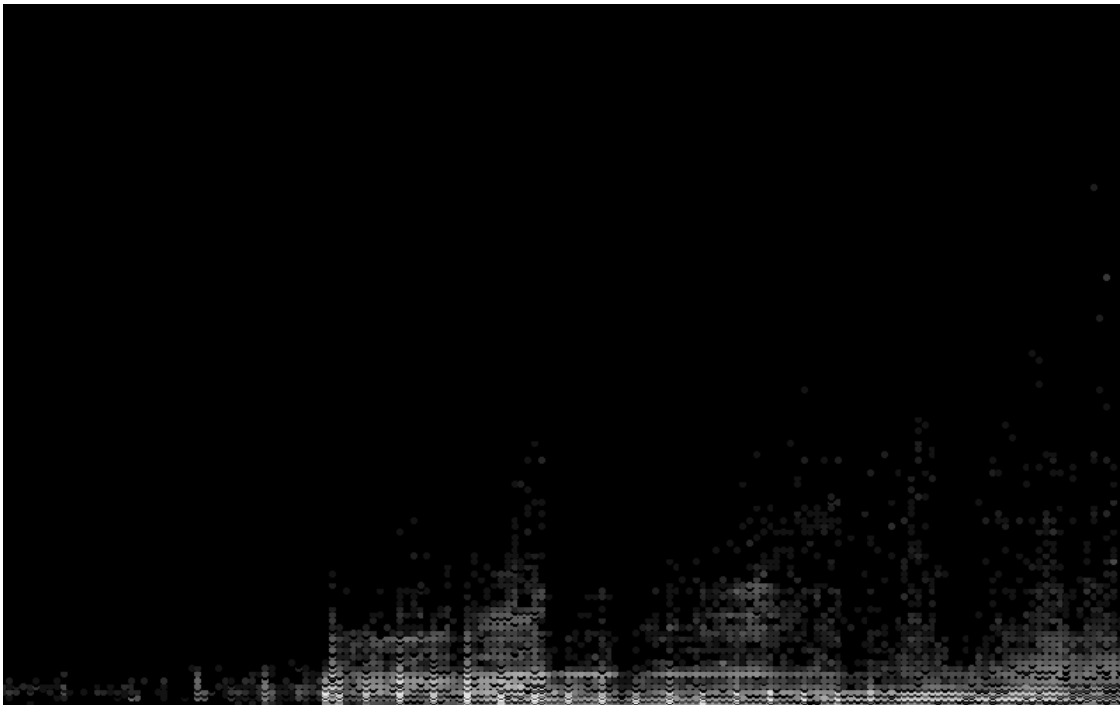
Here is what the same plot would look like in datashader:

```

In [21]: cvs = ds.Canvas(800, 500, x_range = (ny['YearBuilt'].min(), ny['YearBuilt'].max()),
                        y_range = (ny['NumFloors'].min(), ny['NumFloors'].max())
agg = cvs.points(ny, 'YearBuilt', 'NumFloors')
view = tf.shade(agg, cmap = cm(Greys9), how='log')
export(tf.spread(view, px=2), 'yearvsnumfloors')

```

Out[21]:

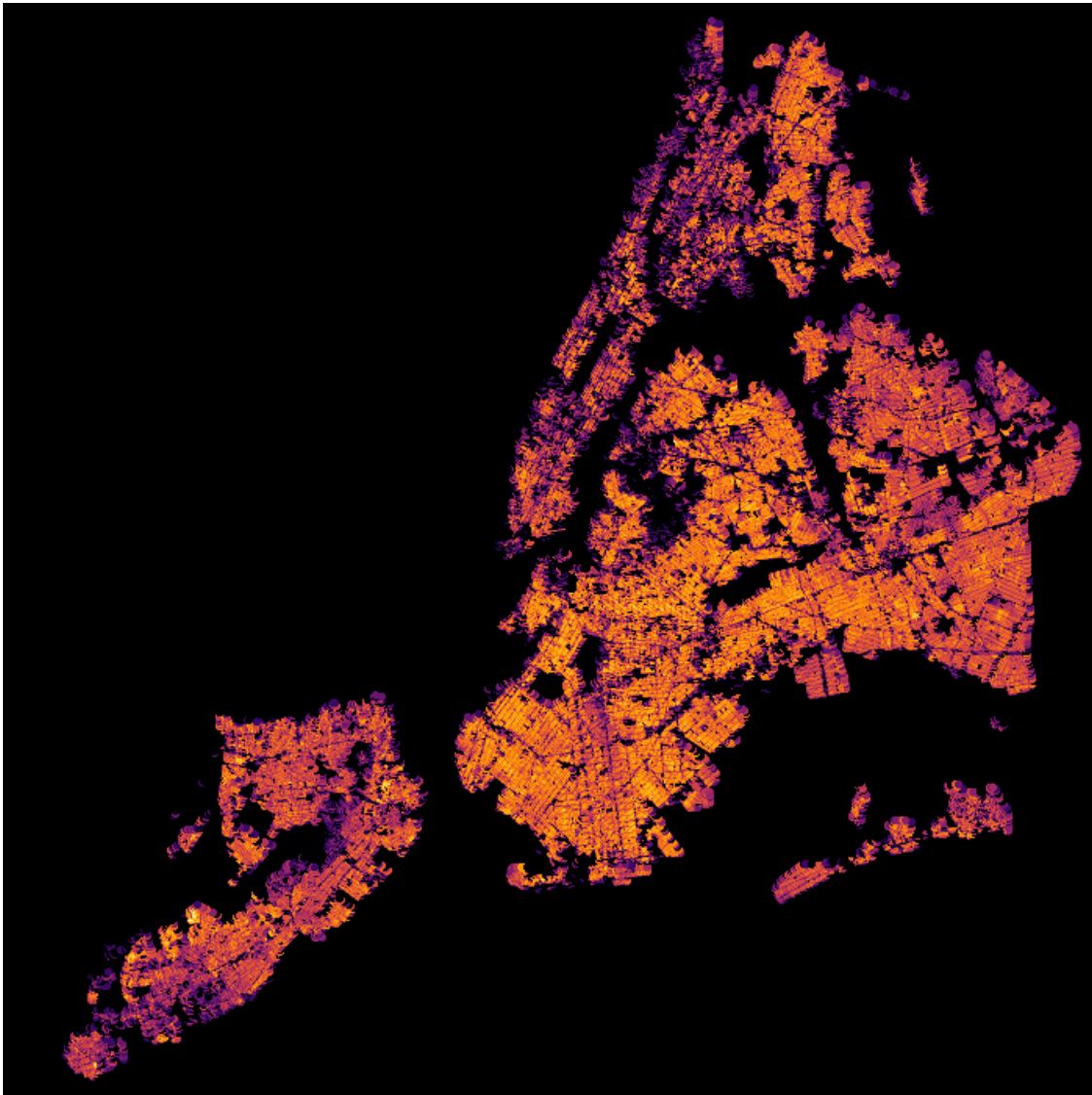


That's technically just a scatterplot, but the points are smartly placed and colored to mimic what one gets in a heatmap. Based on the pixel size, it will either display individual points, or will color the points of denser regions.

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:

```
In [22]: NewYorkCity = (( -74.29, -73.69), (40.49, 40.92))
        cvs = ds.Canvas(700, 700, *NewYorkCity)
        agg = cvs.points(ny, 'lon', 'lat')
        view = tf.shade(agg, cmap = cm.inferno, how='log')
        export(tf.spread(view, px=2), 'firery')
```

Out [22] :



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. [This talk](#) also provides a nice background for datashader. 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 correlation 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 Response

I will use the steps outlined in the above article on creating bivariate choropleth maps. Please note that the results will vary depending the classification strategy. I used a simple classification strategy where the standard deviation of the land value is compared with the standard deviation of the structure value. Having a high land value and low structure value is considered underbuilt. The opposite is considered overbuilt and if they are more or less the same, it is considered normal.

On the map below, white is considered normal, blue is considered overbuilt, and red is considered underbuilt.

```
In [119]: # new dataframe for necessary data.
          shader_df = ny[['AssessLand', 'AssessTot']]

          # calculate the value of the structures
          shader_df = shader_df.assign(AssessStructure = shader_df['AssessTot'] - shader_df['AssessLand'])

          # normalize the data using min-max normalization
          shader_df=(shader_df-shader_df.min()/(2*shader_df.std()-shader_df.min()))

          # load coordinates
          shader_df['lat'] = ny['lat']
          shader_df['lon'] = ny['lon']

          # classify data
          sdStructure = shader_df['AssessStructure'].std()
          sdLand = shader_df['AssessLand'].std()
          shader_df['Classification'] = "normal"
```

```

shader_df.loc[((shader_df['AssessStructure'] <= sdStructure ) & (shader_df['AssessLa
((shader_df['AssessStructure'] <= 1.5 * sdStructure ) & (shader_df['Ass
((shader_df['AssessStructure'] <= 2 * sdStructure ) & (shader_df['Asses
'Classification'] = "underBuilt"

shader_df.loc[((shader_df['AssessStructure'] >= sdStructure ) & (shader_df['AssessLa
((shader_df['AssessStructure'] >= 1.5 * sdStructure ) & (shader_df['Ass
((shader_df['AssessStructure'] >= 2 * sdStructure ) & (shader_df['Asses
'Classification'] = "overBuilt"
# datashader requires a categorical type for categorical data classification
shader_df['Classification'] = shader_df['Classification'].astype('category')

NewYorkCity = (( -74.29, -73.69), (40.49, 40.92))
cvs = ds.Canvas(700, 700, *NewYorkCity)
agg = cvs.points(shader_df, 'lon', 'lat', ds.count_cat('Classification'))
view = tf.shade(agg, color_key={'normal':'white', 'overBuilt':'blue', 'underBuilt':'
export(tf.spread(view, px=3), 'real_estate')

```

Out[119]:

