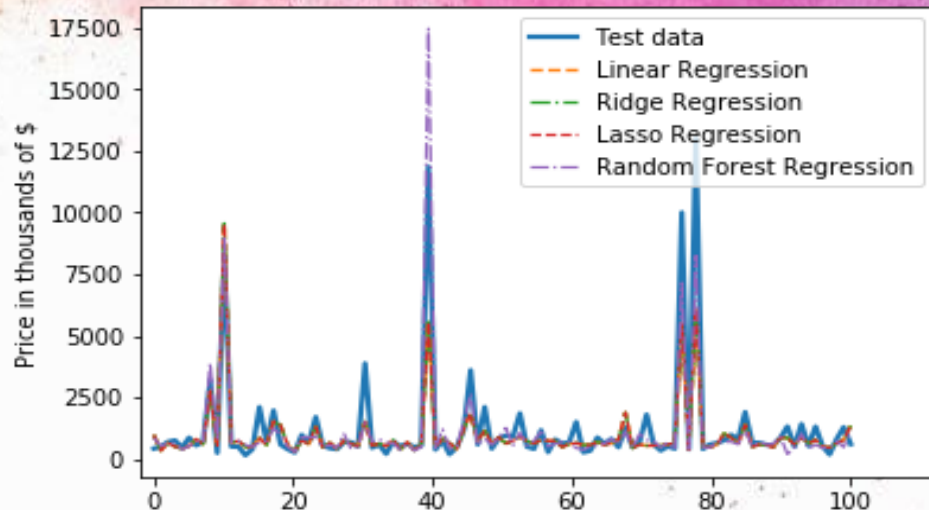


# Predicting the price of properties sell in NYC: Using Machine Learning algorithms.

Full code in:

<https://github.com/Thaleia18>



*The prediction task is to determine the price of a property in sale in NYC.*

*We will create our models using real state data of properties sold in NYC between September 2016 to September 2017.*

The machine learning algorithms that I used are:

- Linear regression.
- Ridge regression.
- Lasso regression.
- Random forest regression.

The metrics to evaluate my predictions are:

- $R^2$  Coefficient of determination.
- RMSE Root Mean Square Error



# THE DATA

This data was extracted from

<https://www.kaggle.com/new-york-city/nyc-property-sales>

Consists on information of properties sold in New York City over a 12-month period from September 2016 to September 2017.

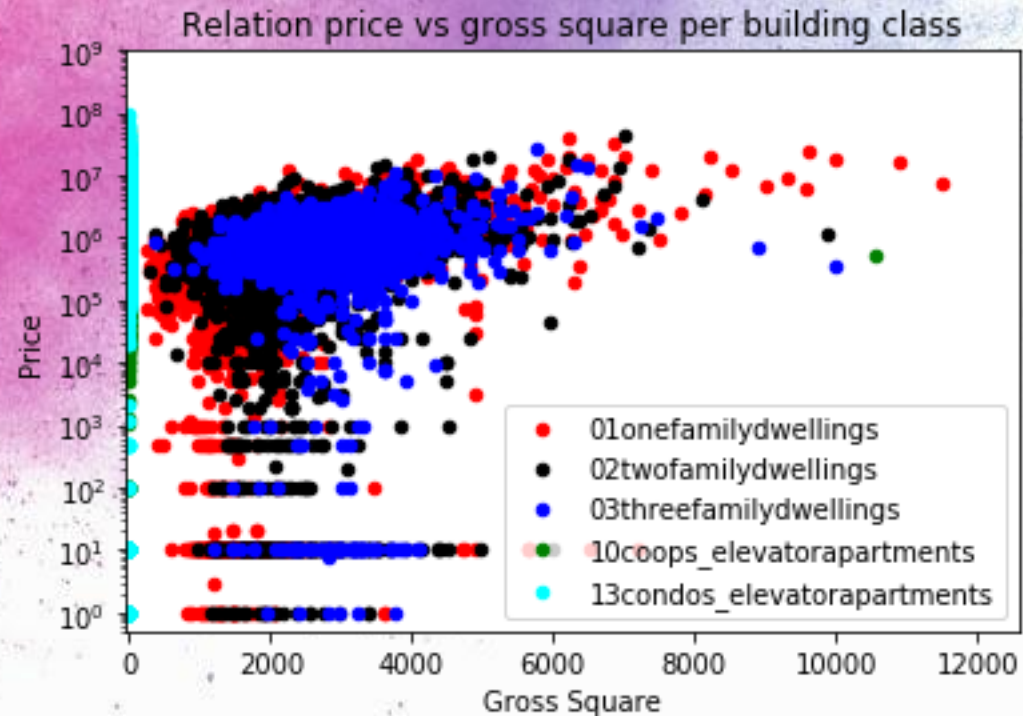
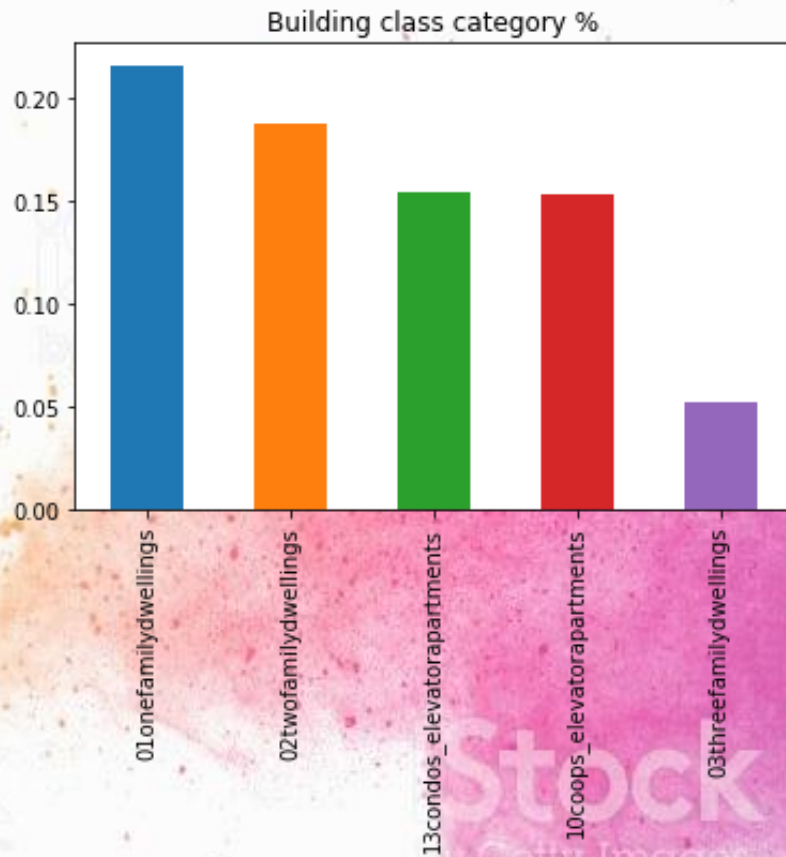
## Attributes

- |                              |                               |
|------------------------------|-------------------------------|
| 1. BOROUGH                   | 13. COMMERCIAL UNITS,         |
| 2. NEIGHBORHOOD              | 14. TOTAL UNITS               |
| 3. BUILDING CLASS CATEGORY,  | 15. LAND SQUARE FEET          |
| 4. TAX CLASS AT PRESENT      | 16. GROSS SQUARE FEET         |
| 5. BLOCK,                    | 17. YEAR BUILT,               |
| 6. LOT,                      | 18. TAX CLASS AT TIME OF SALE |
| 7. EASE-MENT                 | 19. BUILDING CLASS AT TIME OF |
| 8. BUILDING CLASS AT PRESENT | SALE                          |
| 9. ADDRESS                   | 20. SALE PRICE,               |
| 10. APARTMENT NUMBER         | 21. SALE DATE                 |
| 11. ZIP CODE                 |                               |

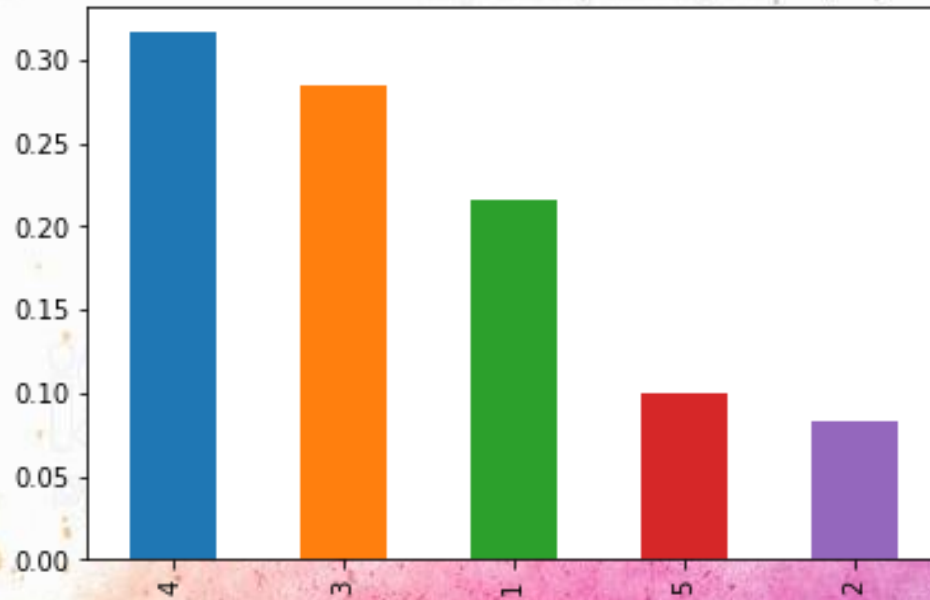
[illegible]



# SOME DATA VISUALIZATIONS.



Borough %



Borough codes:

Manhattan (1),

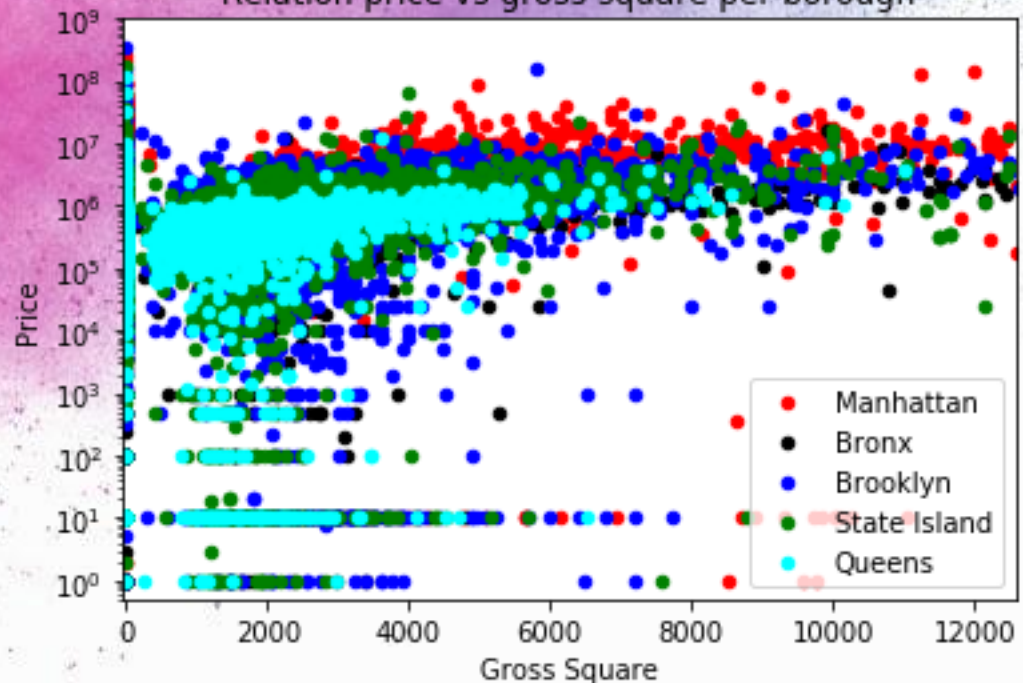
Bronx (2)

Brooklyn (3)

Queens (4)

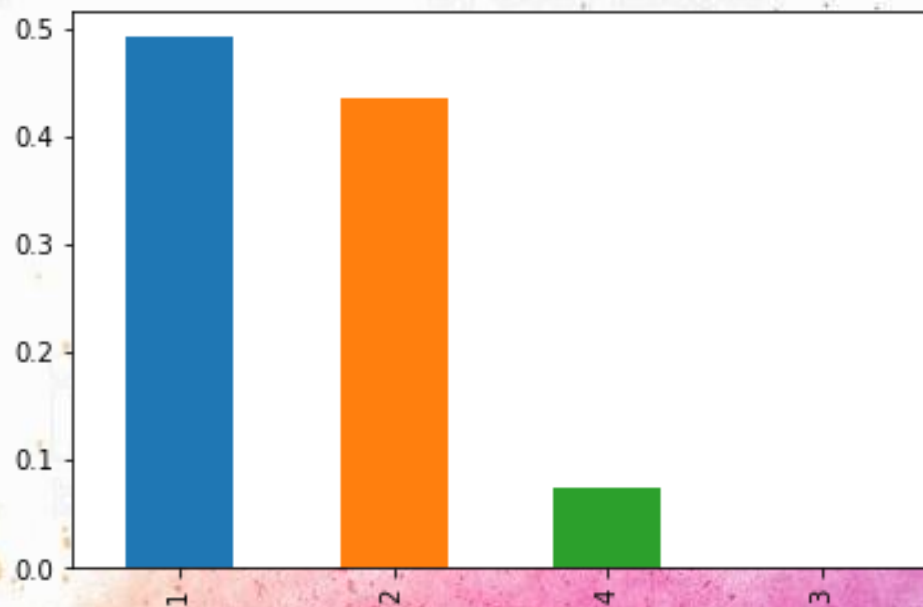
State Island (5)

Relation price vs gross square per borough

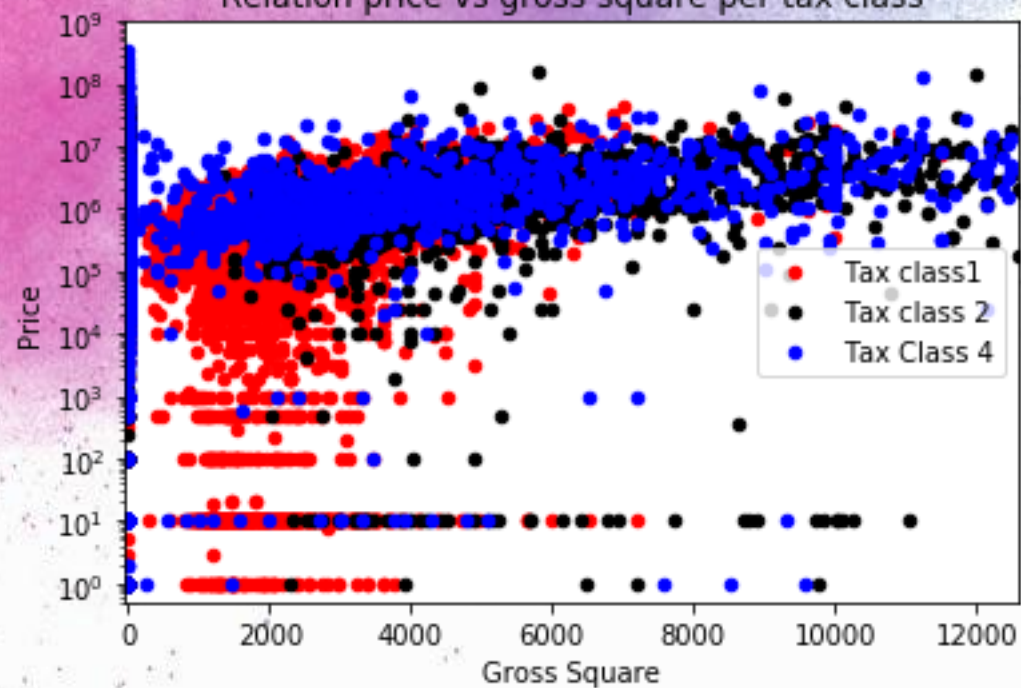




Tax Class at time of sale %

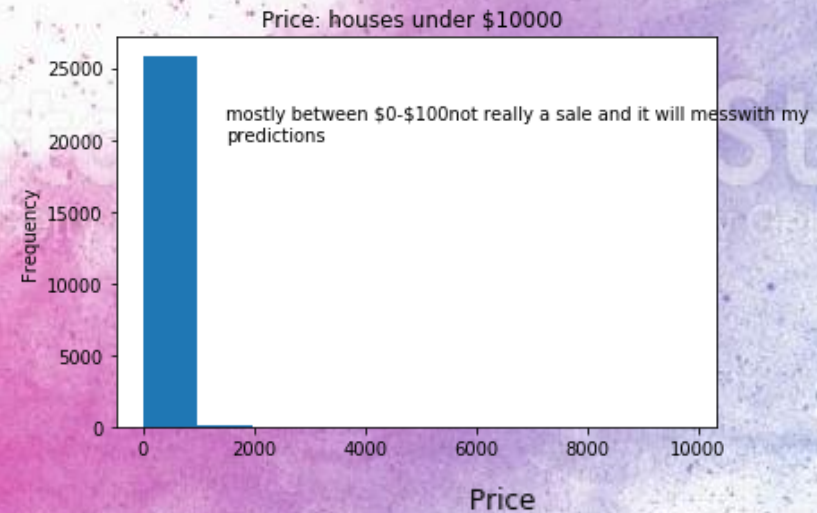
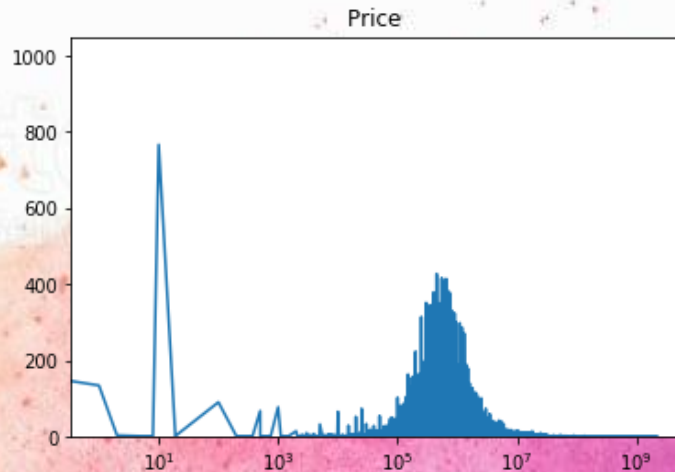


Relation price vs gross square per tax class

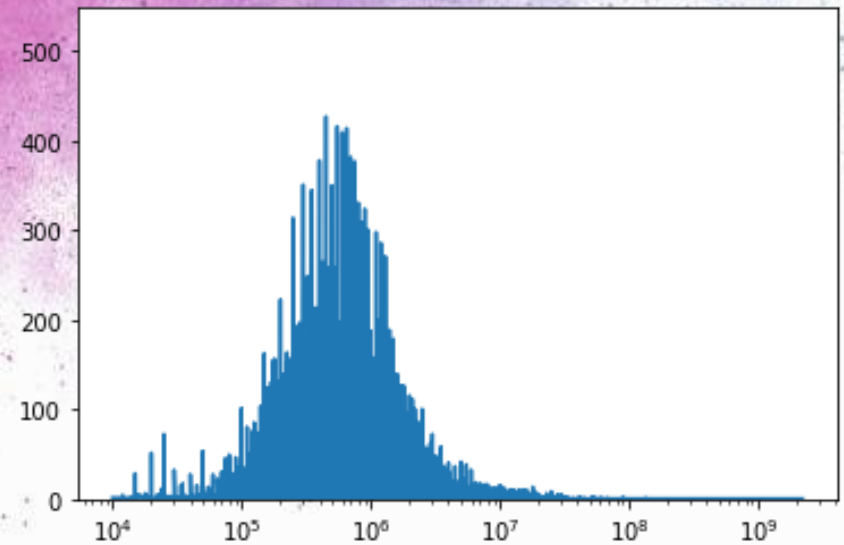


I plotted the frequency of the different sale prices for the properties sold.

A lot of the houses were sold for low prices, under \$10,000. These were not sales, just transfers between family members.

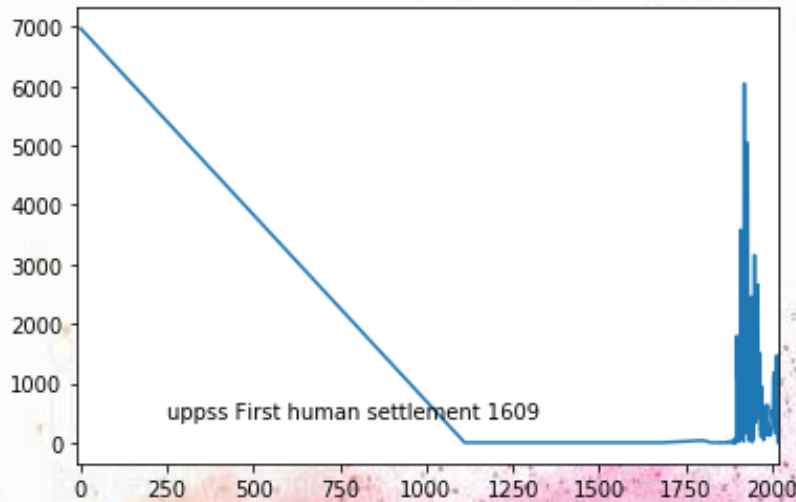


I used data with prices over \$10,000. I consider noise that will mess my models all the data related with transfers.





Year Built



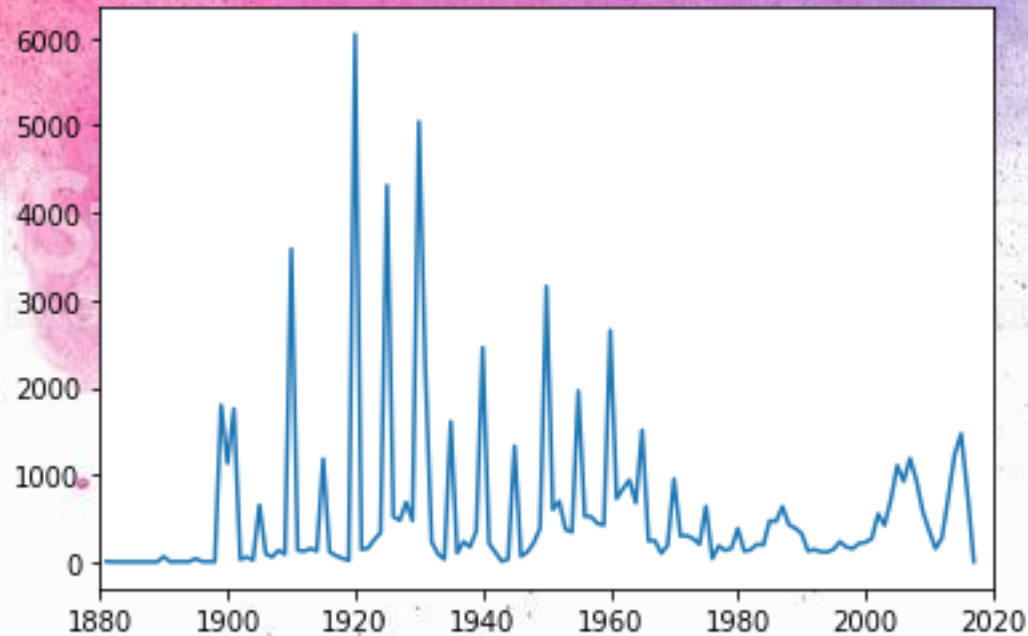
I have data about the year built for each property sold.

I plotted the frequency of houses built every year.

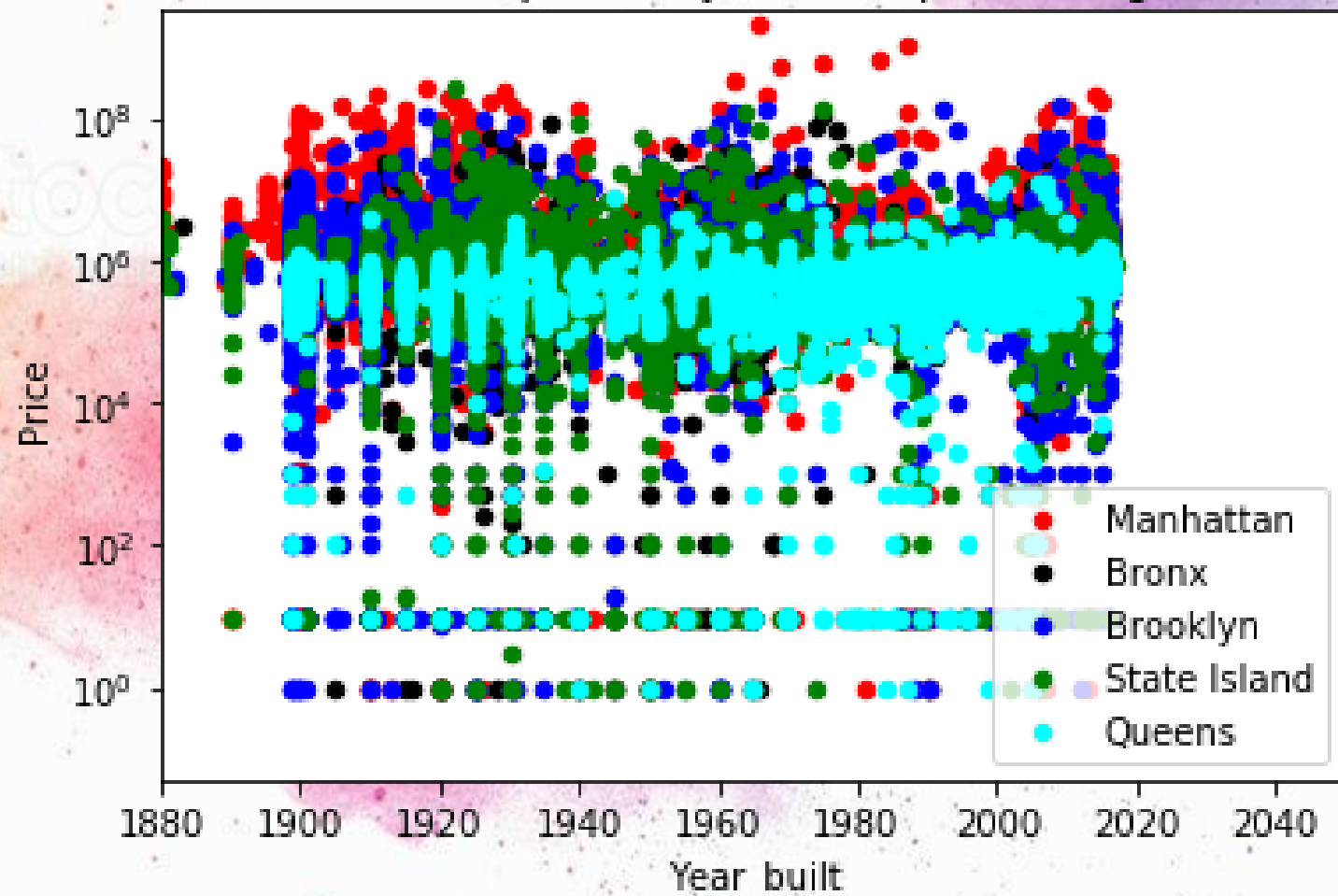
We have some data that says that the properties were built in year 0, this seems wrong.

I removed all the data with houses built before 1880.

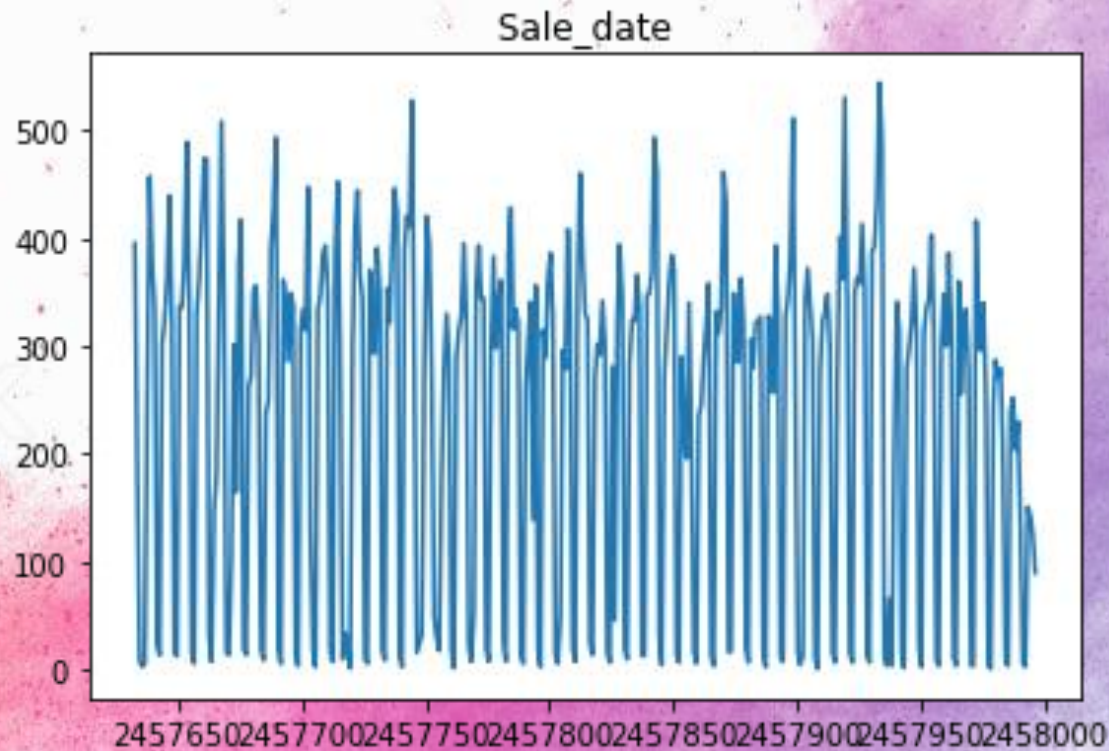
Year Built



Relation price vs year built per borough







I transform the date to Julian date. I plot the frequencies to see if in some point of the year more houses were sold.

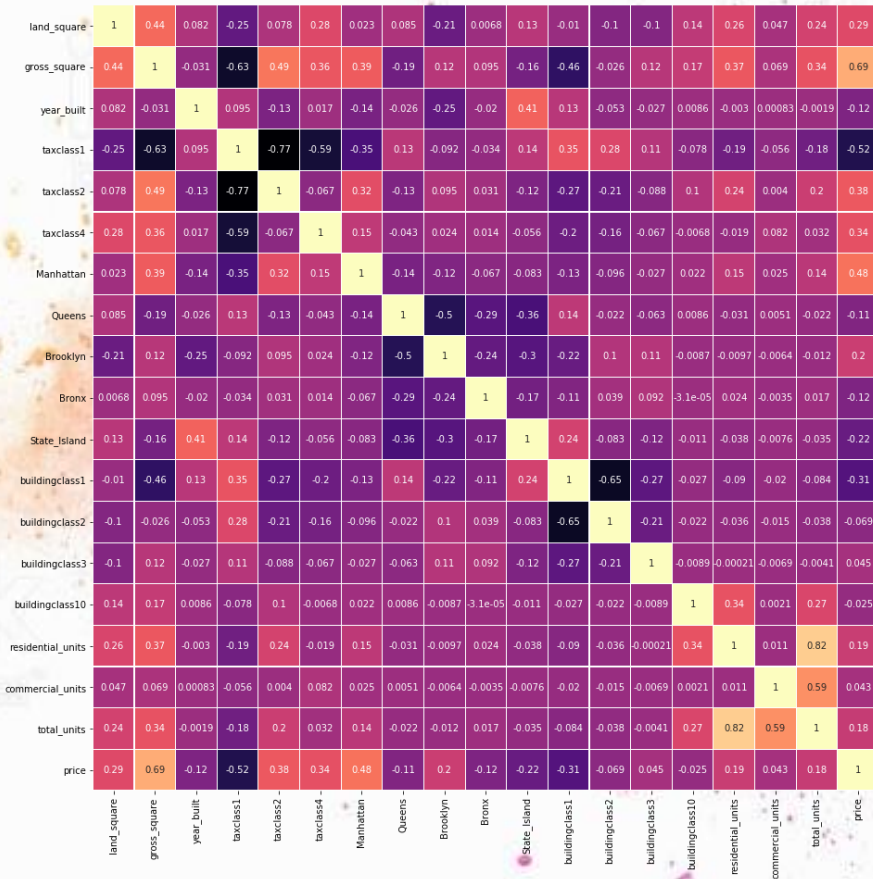
I didn't find any pattern I decided no to use the sale date in my regression.

## models:



# The correlation of these features with the sale price:

Pearson Correlation of Features



```
price 1.000000
gross_square 0.689004
Manhattan 0.475596
taxclass2 0.376039
taxclass4 0.335969
land_square 0.292000
Brooklyn 0.200388
residential_units 0.192468
total_units 0.180459
buildingclass3 0.045443
commercial_units 0.042755
buildingclass10 -0.025421
buildingclass2 -0.068948
Queens -0.109829
Bronx -0.120081
year_built -0.123245
State_Island -0.216972
buildingclass1 -0.309559
taxclass1 -0.520487
Name: price, dtype: float64
```

# MODEL RESULTS

Linear Regression.

```
r2 0.6111805534821673  
rmse 0.5340183237978753
```

Ridge regression.

```
r2 0.6111805532023231  
rmse 0.5340183239900492
```



# Lasso Regression.

```
r2 0.6113140862445937  
rmse 0.5339266166236538
```

# Random forest regression.

```
r2 0.6233097547347948  
rmse 0.5256229959990711
```

The best results are for the random forest regression.

