# Predicting Craigslist NYC Apartment Rental Prices with Linear Regression

Ami Tian

# Project Goal: Do an apartment's features help predict its rental price?



$1650 Per Month!!!

No Oven
No Freezer

The small counter. Cameron Knowlton

https://www.insider.com/worst-apartment-new-york-city-tiktok-reactions-2021-2



7:13
58°
NY
1 HD

Jimmy McMillan
Rent is Too Damn High Party

# Web Scraping and Data Cleaning

- Selenium/BeautifulSoup for web scraping
- Features:
  - \# of bedrooms (BR)
  - \# of bathrooms (BA)
  - Location
  - Square footage (SQFT)
  - Pets allowed
  - Application Fee
  - Broker's Fee
  - Laundry
  - Parking
- Cleaning: Dealing with outliers, incorrect/misplaced information, removing duplicate posts



$1 / 1br - 800ft$^2$ - $1,5oo+ (3-Units) w/Elevator+Laundry (near Queens Cente
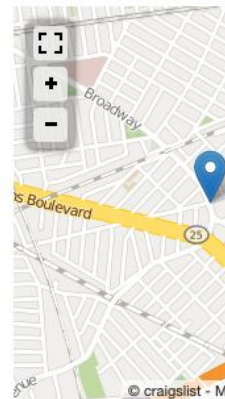
image 1 of 4

Good Income Required

STUDIO $1,500
1BED 1BATH $1,700
1BED 1BATH $1,800
Move in ASAP

51st Avenue off Broadway
Elmhurst, NY 11373
2 Blocks to R+M trains

© craigslist - Ma

1BR / 1Ba    800ft$^2$

cats are OK - purrr

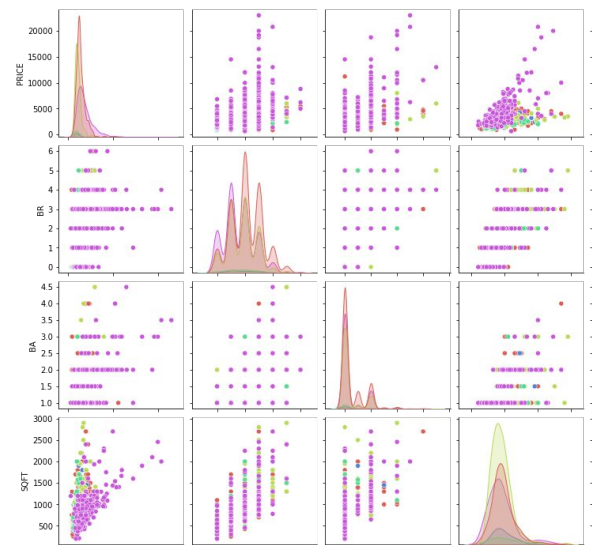flooring: **wood**

apartment

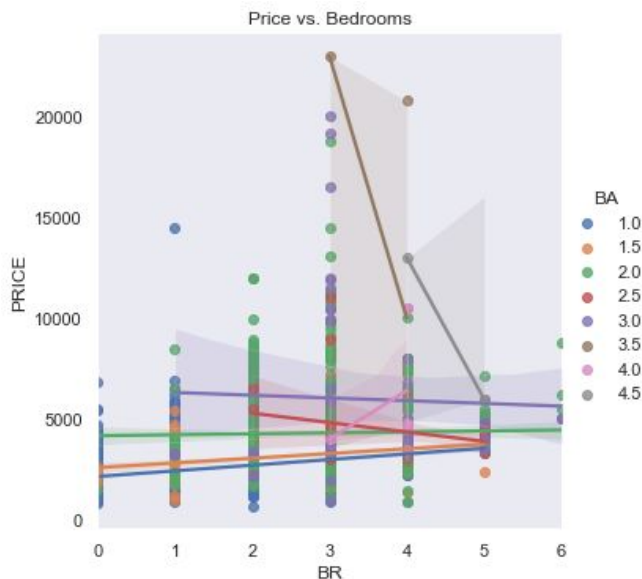laundry on site

street parking

rent period: **yearly**

# Data Prep and Feature Engineering

- Out of 4609 listings, only 969 listed square footage.
  - Imputed missing values based on median, borough, bedrooms.
- Created dummy variables for categoricals
  - Boolean: pets, application fee, broker's fee
  - Washer/Dryer: none, in-unit, in-building
  - Parking: none, off-street (includes carports and garages), valet
  - Borough: Brooklyn, Manhattan, Queens, the Bronx, Staten Island
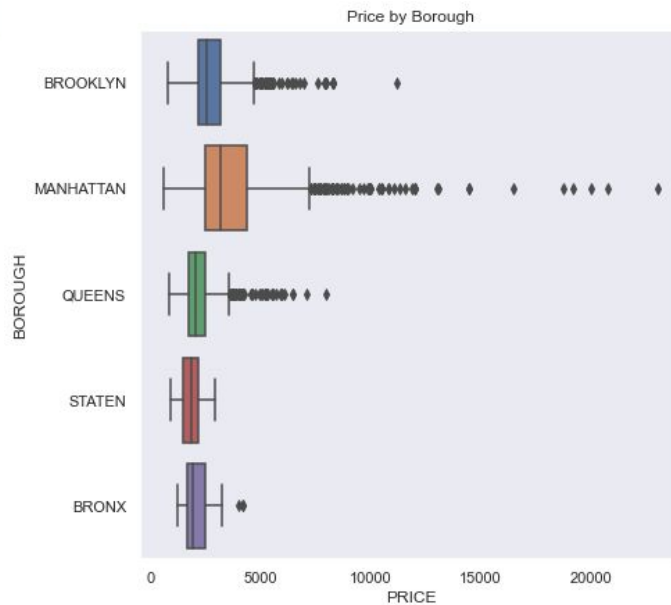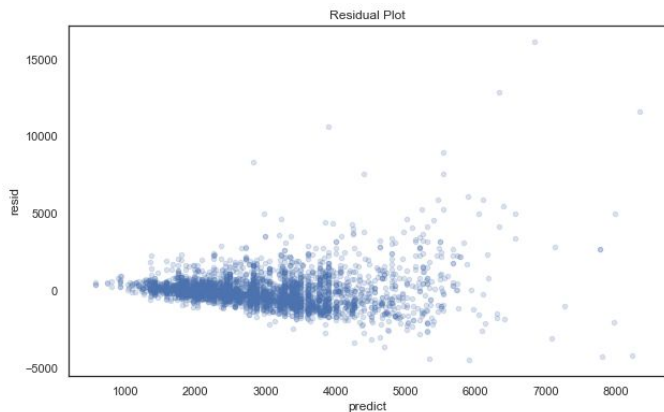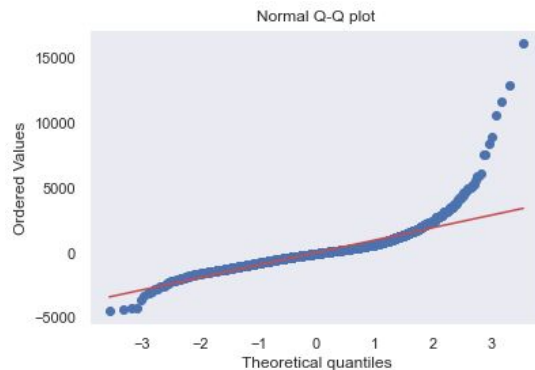




# of Listings per Borough

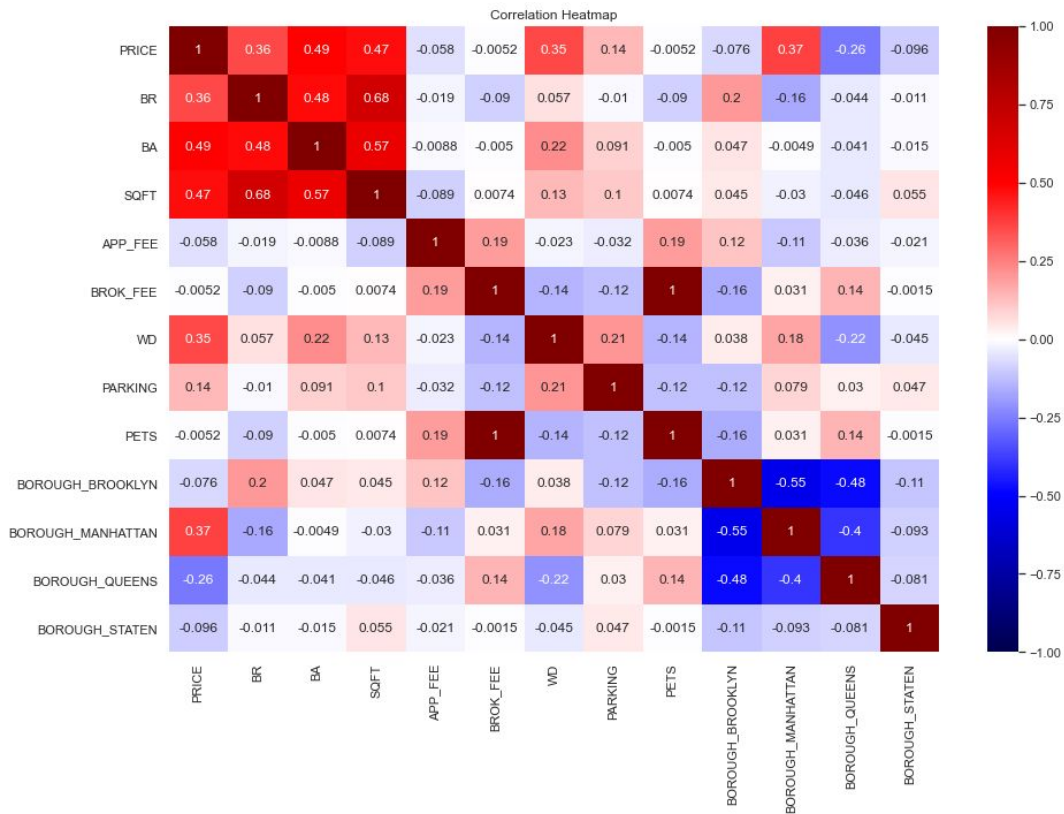# More Data Prep and Feature Engineering

- Removed some features (application fee, pets)
- Added interaction terms
  - BR / BA
  - SQFT * MANHATTAN

# Issues

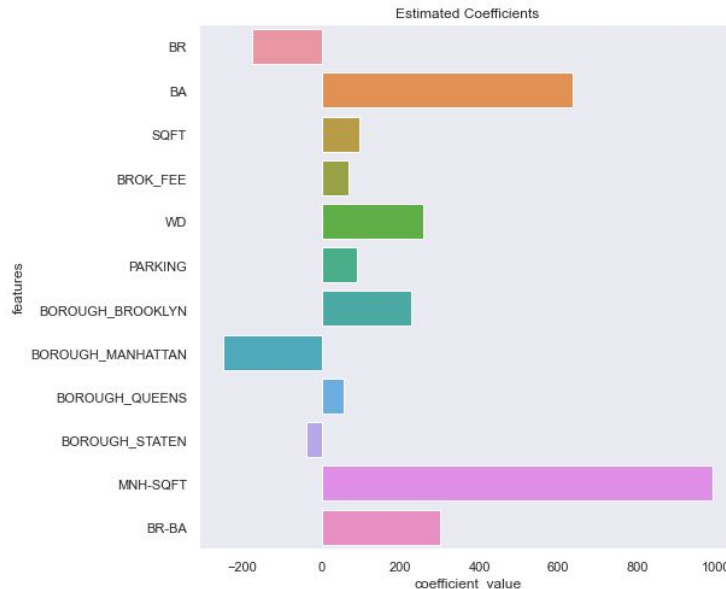- Multicollinearity
- Heteroskedasticity

Correlation Heatmap

| | variables | vif |
|---|---|---|
| 0 | BR | 24.626210 |
| 1 | BA | 8.557337 |
| 2 | SQFT | 4.497295 |
| 3 | BROK_FEE | 1.078687 |
| 4 | WD | 1.184902 |
| 5 | PARKING | 1.087153 |
| 6 | BOROUGH_BROOKLYN | 19.939513 |
| 7 | BOROUGH_MANHATTAN | 29.935476 |
| 8 | BOROUGH_QUEENS | 16.101196 |
| 9 | BOROUGH_STATEN | 2.444164 |
| 10 | MNH-SQFT | 12.682606 |
| 11 | BR-BA | 17.592696 |

# Model Selection and Results

- Used kfold cross-validation on training set to get the R^2 score, mean absolute error (MAE) and root mean squared error (RMSE)
- Applied Ridge and Lasso regularization, compared scores
- Lasso regularization improved scores slightly



## TRAINING DATA

- R^2 value: 0.526
- Mean Absolute Error: $648.65
- Root Mean Squared Error: $1,027.45

## TEST DATA

- R^2 value: 0.519
- Mean Absolute Error: $671.15
- Root Mean Squared Error: $1,127.72

# Conclusions

- The model can only explain about 50% of an apartment's rental price
- Other possible factors: neighborhood, amenities, utilities
- External factors: housing policies, market forces





Real Estate
## New York Renters Face 70% Increases as Pandemic Discounts Expire

The era of widespread Covid concessions for apartment hu[...]

NEW YORK
## New York City, Income Inequality Capital Of America, Now Also Facing Soaring Rent Prices: Report

# Future Steps



- Get more data, use information from other sites (e.g. StreetEasy, apartments.com)
- Look into more features
- Figure out more accurate interaction terms
- Categorize by neighborhood instead of borough, get more information about location
- Try log transformation or use weighted least squares regression to deal with heteroskedasticity