# FINDING THE 5 BEST ZIP CODES FOR REAL ESTATE INVESTMENT IN THE NYC SUBURBS

# Business Understanding

Thanks for hiring me . It's May 2018 and you, a Private Equity company, want to invest in NYC suburban real estate.

You have narrowed your focus down to 2 counties in NY: Westchester and Nassau and 2 counties in New Jersey: Bergen and Hudson.

The question you have is: **which zip codes should you invest in**?

Investment Horizon: 5 Years

$$ROI\% = \frac{\text{final predicted data price} - \text{final observed data price}}{\text{final observed data price}} \times 100$$

# Business Understanding

- You are seeking the 5 zip codes with the highest ROI%.

## Data Understanding

I used an outstanding dataset from Zillow that contains monthly average sales data from almost every zip code in the United States from 1996-2018.

I narrowed the data down to the 201 zip codes in the 4 chosen counties and used that data to predict which would be the 5 most profitable in the next 5 years.

# Modeling

I used an auto.arima model to predict the prices.

# Results— 5 Highest ROI% Zip Codes

| Zipcode | Projected 5 Year ROI% |
|---------|----------------------|
| 10590 | 371.187770 |
| 10553 | 255.495333 |
| 11804 | 190.458716 |
| 10536 | 188.628919 |
| 10504 | 180.919491 |



Top 5 Zipcodes by Projected 5 Year ROI%

# Results

- 4 of the zip codes are in Westchester County.

- All of the zip codes, except 10553, are considered high income areas.
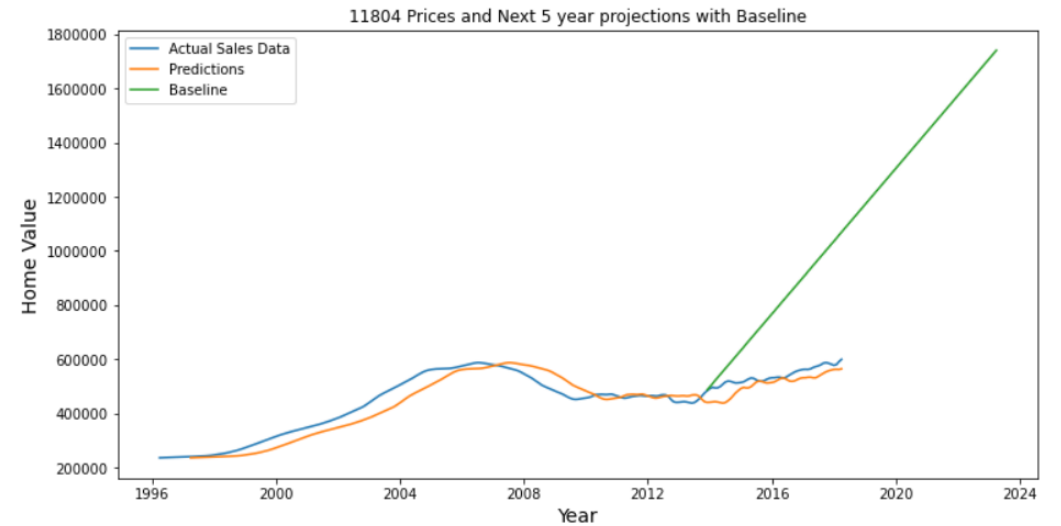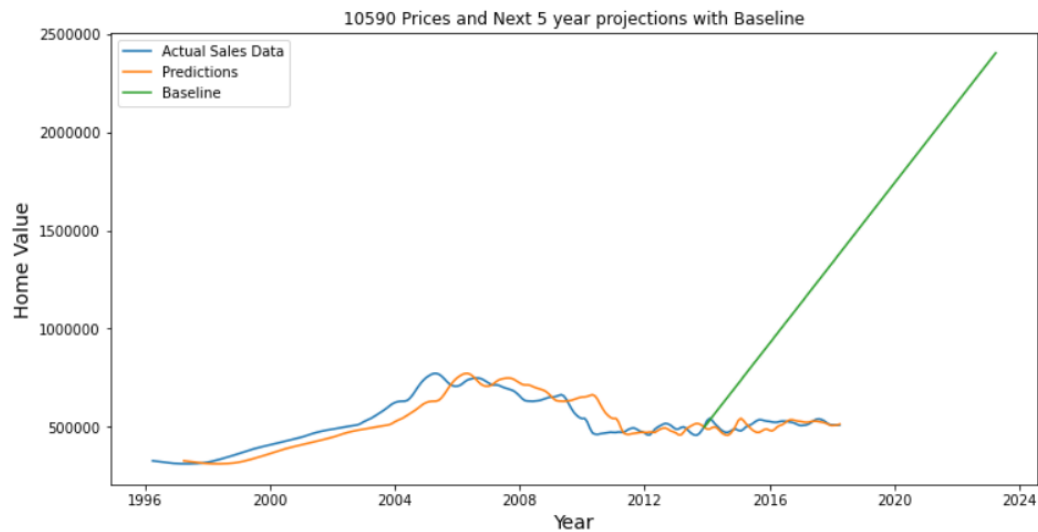
- 10590, 10536, and 10504 are all in Upper Westchester County and are close to each other.
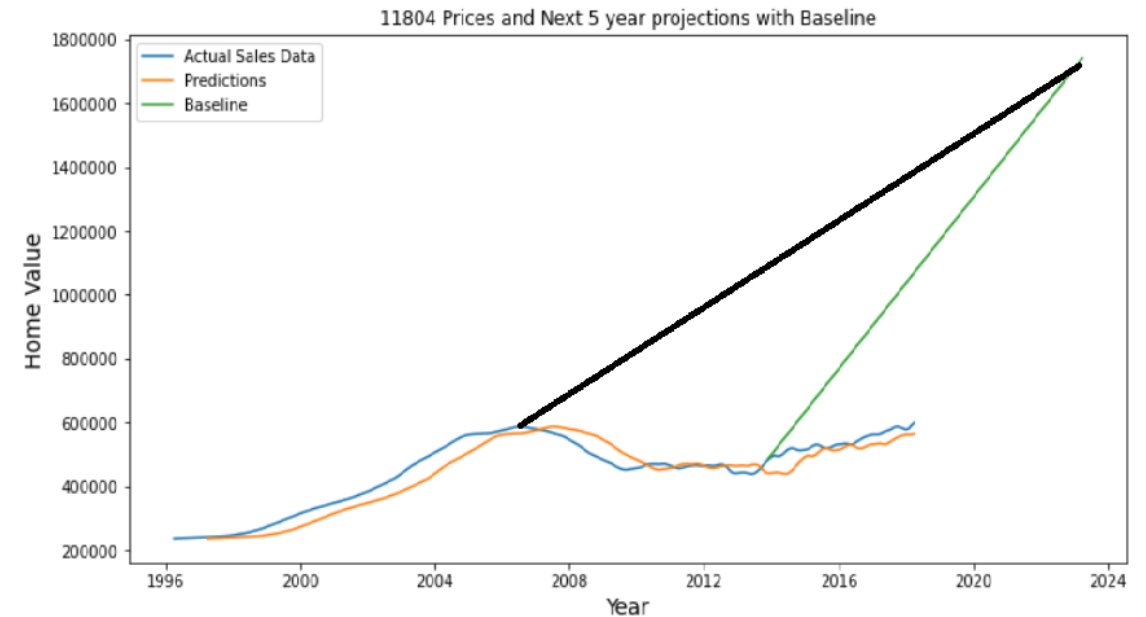
# Results-Prediction Problems

The predictions for 2013-2018 varied greatly from the actual sales data from 2018-2023.

# Results-Predictions Problems

- The discrepancy between the predicted and actual prices may be a product of the Great Financial Crash. The model may be expecting the prices to increase after the financial crash. Prices did rebound but the model may be expecting prices to reach the levels they would have had there been no downturn.



11804 Prices and Next 5 year projections with Baseline

# Recommendations

- The top 5 zip codes with the highest projected ROI% were 10590, 10553, 11804, 10536, and 10504. Invest in those zip codes.

- The least expensive of the top 5 is 10553 and it has an expected ROI% of 255%. So, if you are targeting homes under $1 million then focus on that zip code.

# Next Steps

- More data on Fairfield County, CT would also be very useful to investors and to the model. Unfortunately, the Zillow dataset omitted many zip codes in Fairfield County, CT(parts of which are only a 40-minute train ride to Midtown, Manhattan and contains some of the wealthiest areas in the country) and omitted data many zip codes in the city itself.

- Another ML model may yield more fruitful results. An XGBoost Regressor or a neural network model may yield better and more meaningful results.

# Thanks!

Github:https://github.com/icapeli/Phase_4_Time_Series

Image from:https://www.youtube.com/watch?v=1TFzXBLBsAM