# Allegheny County Home Value Index Case Study

For a Fidelity Management and Research Interview Matheus C. Fernandes (PhD Candidate - Harvard University) 2/11/2021

Please find more information and full code at <u>https://git.fer.me/fidelity-interview</u>

### About Allegheny County

- Located in southwest Pennsylvania
- State's second-most populated county
- City of Pittsburgh is in center
- With 446 bridges, Pittsburgh has more bridges than any other city in the world\*

\*https://uncoveringpa.com/facts-about-pittsburgh



### About the Allegheny County Dataset

- 86 columns (potential features)
- 580,997 property assessments
- Columns contain information on:
  - Property features: bedrooms, bathrooms, fireplace etc.
  - Location: physical address and location codes
  - Different levels of assessment values county, local for building and land

### About the Allegheny County Dataset

- Has numerical, categorical, and binary variables
- Dataset contains repetitive information i.e. descriptions of codes
- Contains administrative information i.e. deed recording information and legal descriptions

### Goal of Analysis

#### Develop a monthly "Allegheny County Home Value Index" (HVI) to understand key features of the market.

Create model to gain insights for investment opportunities.

### Structure of Data Science Workflow



### Data Cleaning



**Exploratory Data Analysis** 



Model Exploration and Selection



Computing Home Value Index



#### Removing false data





- Remove data missing important information such as:
  - Sale price
  - Sale date
  - Sale price is 0 or unreasonably low (<\$1000)



- Imputed missing data depending on datatype
  - Creating a new category of unknown, zero, or boolean
  - Replacing mean of data for continuous variables
  - Imputing information from different column

![](_page_9_Picture_0.jpeg)

- Converted variable types to increase bit efficiency
  - From 64 bit to 32bit and 8 bit
  - Reduced memory ~380MB to ~120MB without loosing information
- Converted date types to numerical
- Standardize input data

![](_page_10_Picture_0.jpeg)

### Feature Engineering

#### **Geolocation Exploration**

#### Address Granularity

![](_page_10_Figure_4.jpeg)

#### Zip code Granularity

![](_page_10_Figure_6.jpeg)

11

### Structure of Data Science Workflow

![](_page_11_Picture_1.jpeg)

### Data Cleaning

![](_page_11_Picture_3.jpeg)

**Exploratory Data Analysis** 

![](_page_11_Figure_5.jpeg)

Model Exploration and Selection

![](_page_11_Figure_7.jpeg)

Computing Home Value Index

![](_page_12_Picture_0.jpeg)

### Is housing prices in Allegheny county a martingale?

 $\mathbf{E}(|X_n|) < \infty$  $\mathbf{E}(X_{n+1} \mid X_1, \dots, X_n) = X_n$ 

By fitting an exponential line, we see that the expectation follows an exponential growth.

![](_page_12_Figure_4.jpeg)

![](_page_13_Picture_0.jpeg)

### **Exploratory Data Analysis**

### **Correlation Matrix**

- Provides information on correlation of variables
- No correlation does not mean no useful information
- High correlation between variables means potential redundancy between those variables.

![](_page_13_Figure_6.jpeg)

![](_page_14_Picture_0.jpeg)

### **Exploratory Data Analysis**

### A deeper dive into the data

- How does each variable depends on the other
- Look for trends in the data
- How does property sale price depends on each feature

![](_page_14_Figure_6.jpeg)

![](_page_15_Picture_0.jpeg)

#### How do house properties impact pricing?

![](_page_15_Figure_2.jpeg)

![](_page_16_Picture_0.jpeg)

#### How does location of properties impact pricing?

![](_page_16_Figure_2.jpeg)

![](_page_17_Picture_0.jpeg)

### **Exploratory Data Analysis**

#### What is the distribution of the data across different locations?

![](_page_17_Figure_3.jpeg)

![](_page_18_Picture_0.jpeg)

#### How do the assessments impact pricing?

![](_page_18_Figure_2.jpeg)

### Structure of Data Science Workflow

![](_page_19_Picture_1.jpeg)

### Data Cleaning

![](_page_19_Picture_3.jpeg)

**Exploratory Data Analysis** 

![](_page_19_Figure_5.jpeg)

Model Exploration and Selection

![](_page_19_Figure_7.jpeg)

Computing Home Value Index

![](_page_20_Picture_0.jpeg)

### **Model Information**

# **Model Goal:** Predict valuation of existing homes for a variable sale date.

Target variable: Sale Price

**Input variables:** Sale date and important features that provide information on the valuation of a property at a certain date

![](_page_21_Picture_0.jpeg)

### Model Information

#### Model Assumptions:

- No information on the buyers side (demand)
- No listing prices or spread of ask/bid
- No information on interest rates
- No demographic information
- No information on the economy
- No refined information on location
- Based on assessments only from 2021 (dataset)
- Discrete daily sampling

![](_page_22_Picture_0.jpeg)

#### Seek these regression model characteristics:

- Deal with sparse data
- Good for dealing with categorical and numerical data
- Efficient at training (limited computational resources on my end)
- Scalable to potentially adding more data in the future

![](_page_23_Picture_0.jpeg)

### Model Choices

Model	Train Score	Test Score
Linear Model: LassoCV	0.573	0.495
Support Vector Machine (SD)	0.201	0.007
Ensemble: Random Forest	0.900	0.523
Ensemble: Bagging	0.930	0.588
Ensemble: Adaptive Boosting (SD)	0.829	0.380
Ensemble: Extreme Gradient Boosting	0.814	0.767

\*Scores are measured using R2 Score: 1-(sum of square residuals/total sum of squares) SD = sampled dataset

![](_page_24_Picture_0.jpeg)

- Decision-tree-based ensemble model
- Uses gradient boosting framework
  - Converting weak to strong learner through sequential learning
  - Gradient descent algorithm
- Great for small-to-medium structured/tabular data
- Boosting optimized for software and hardware parallelization

![](_page_25_Picture_0.jpeg)

### **Another Level of Feature Selection**

#### **Feature Importance**

Based on Gini importance

Gini Index = 
$$1 - \sum_{i=1}^{n} (P_i)^2$$

• The higher the importance the more crucial it is for prediction

![](_page_25_Figure_6.jpeg)

![](_page_26_Picture_0.jpeg)

### **Another Level of Feature Selection**

### **Permutation Importance**

- Compute score of model
- For each feature shuffle column and compute score for corrupted dataset
- The higher the importance the more crucial a particular feature

![](_page_26_Figure_6.jpeg)

![](_page_27_Picture_0.jpeg)

### XGBoost Hyperparameter Tuning

- For computational efficiency, only tuned 1 param – number of estimators
- Use cross-validation with 4-way split
- A bias-variance balance is obtained at n=110, with underfitting before and overfitting after
- As n increases computational time increases

![](_page_27_Figure_6.jpeg)

### Structure of Data Science Workflow

![](_page_28_Picture_1.jpeg)

### Data Cleaning

![](_page_28_Picture_3.jpeg)

**Exploratory Data Analysis** 

![](_page_28_Figure_5.jpeg)

Model Exploration and Selection

![](_page_28_Figure_7.jpeg)

Computing Home Value Index

![](_page_29_Picture_0.jpeg)

### What does it mean? (based on Zillow's definition)

- Insight on typical expected home values
- Insight on housing market at a given time
- Appreciation over time

![](_page_30_Picture_0.jpeg)

How is it defined:

![](_page_30_Figure_3.jpeg)

![](_page_31_Picture_0.jpeg)

![](_page_31_Figure_2.jpeg)

![](_page_32_Picture_0.jpeg)

![](_page_32_Figure_2.jpeg)

![](_page_33_Picture_0.jpeg)

![](_page_33_Figure_2.jpeg)

![](_page_34_Picture_0.jpeg)

- Monthly time series of index
- Extract key economic features
  - Late 1930s housing boom
  - 1990 housing crisis
  - Dot-com bubble
  - 2007 housing bubble
  - 2007 housing recovery
  - 2020 pandemic

![](_page_34_Figure_9.jpeg)

### **Conclusions and Recommendations**

- Model needs further development and validation
- Must relax certain model assumptions
- Allegheny county real estate has historically been a good longterm investment
- Based on recent trends I would advise not investing in the housing stock of Allegheny county

### Model Extension

- Create monthly HVI for each separate Zip code and create an index based on that
- Implement a better geolocation scheme to refine location
- Account for economic data such as interest rates into the model
- Include additional sale data and historical assessment data from multiple listing service (MLS)

Thank you for the opportunity!

For a Fidelity Management and Research Interview Matheus C. Fernandes (PhD Candidate - Harvard University)

Please find more information and full code at <a href="https://git.fer.me/fidelity-interview">https://git.fer.me/fidelity-interview</a>