# Data Analytics Case Study: King County Housing Prices

Holly Ringsak

# Introduction

Overview: Discover variables that influence the housing market in King County, Washington, and create a predictive model to accurately predict housing prices.

Purpose & Context: This project was created as part of my data analytics course through CareerFoundry to demonstrate my abilities in exploratory visual analysis in Python, as well developing a hypothesis and testing them through various analytical approaches, and presenting my findings through Tableau.

Objective: To build an interactive dashboard visually showcasing well-curated results of an advanced exploratory analysis conducted in Python.
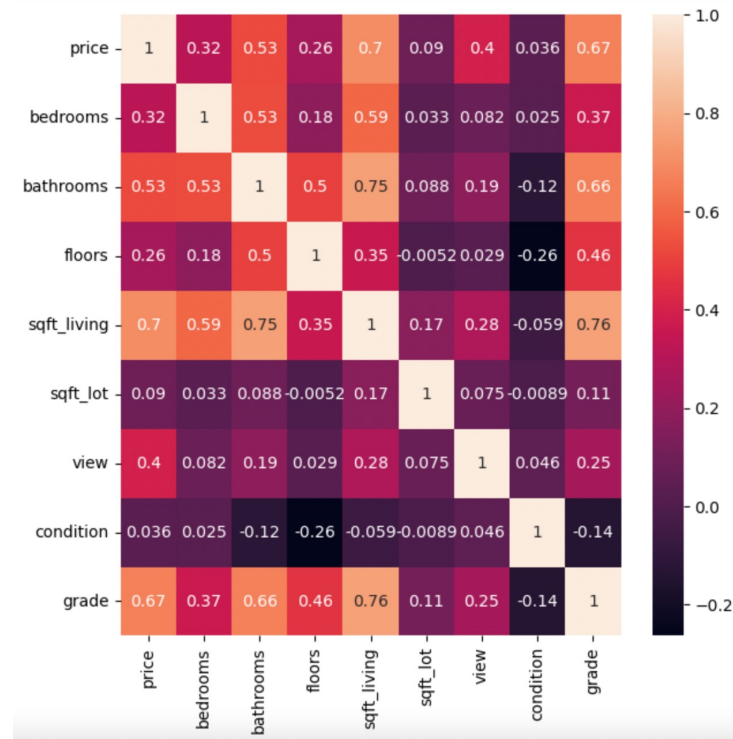
Role: Data Analyst

Project Duration: 2 weeks

Tools Used:
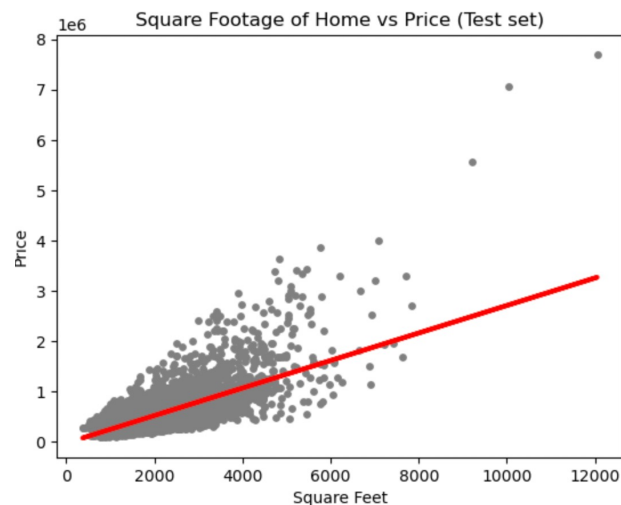- Excel
- Python
- Tableau

# Analysis Process:

- Sourcing Open Data : I selected the [House Sales in King County, USA](#) dataset for my analysis

- Exploring Relationships: Conducted exploratory analysis in Python to test the correlation of different variables with the pricing of homes.

- Geographical Visualizations with Python: Sourced a shapefile with the same location as my data set to perform geospatial analysis.

- Supervised Machine Learning, Regression: Prepare data for regression analysis, split data into training and test sets, run linear regression and analyze the model performance.
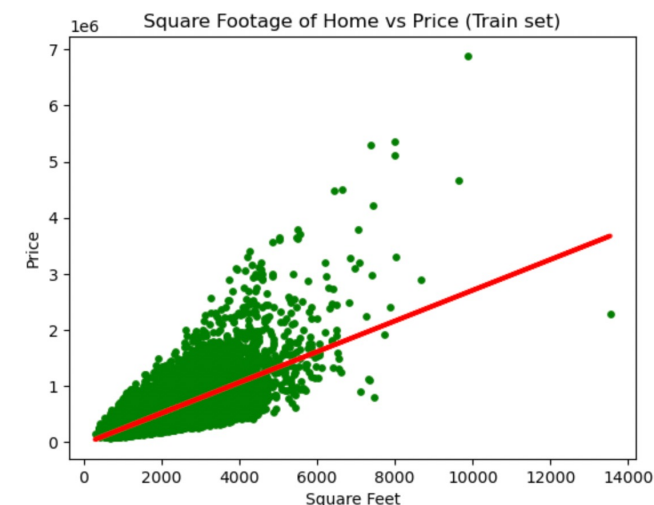


*Correlation heatmap showing the correlation between different variables in the dataset*

*Results of linear regression test on the price and square-footage of a home. Model did not perform well.*
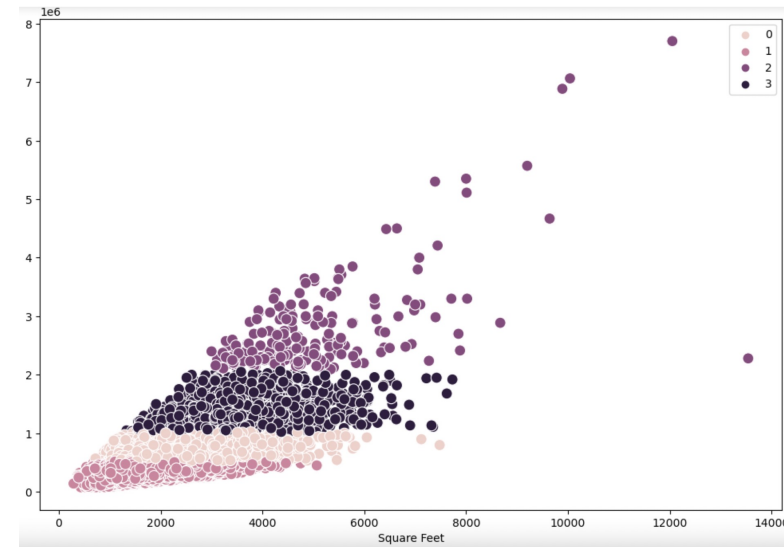


Slope: [[273.70843969]]
Mean squared error:  72796296390.32614
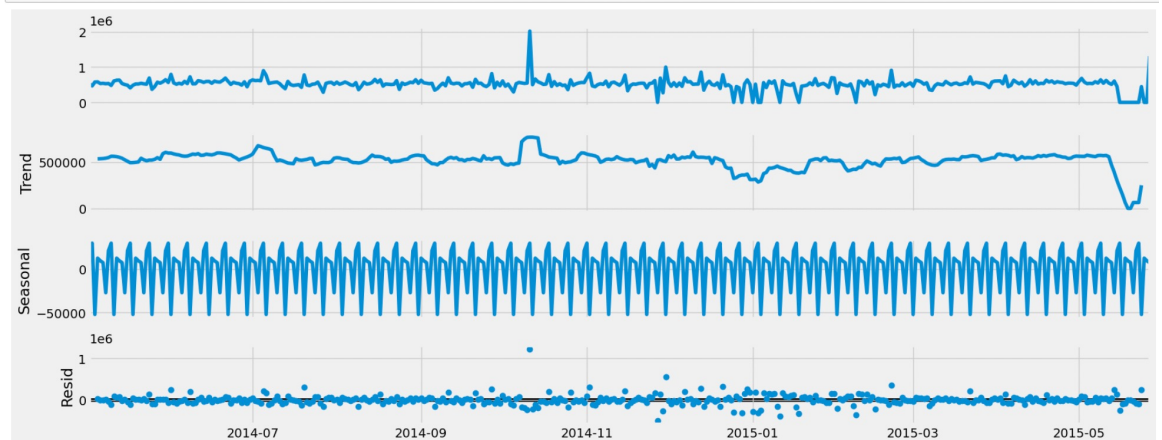R2 score:  0.5057892303219628

Slope: [[273.70843969]]
Mean squared error:  66507995357.54318
R2 score:  0.4860368981218999

# Analysis Process:

- Unsupervised Machine Learning, Clustering: Used the elbow technique to determine optimal number of clusters, ran k-means algorithm, created a variety of visualizations using clustered data, calculated descriptive statistics.

- Sourcing & Analyzing Time Series Data: Prepared data for time series analysis, split the time series into 4 components (as seen on the right), tested for stationarity using the Dickey-Fuller test and an autocorrelation test, stationarized the data through differencing.

- Creating Data Dashboards: Created a storyboard in Tableau to present findings made in Python. *Link to Tableau Presentation*



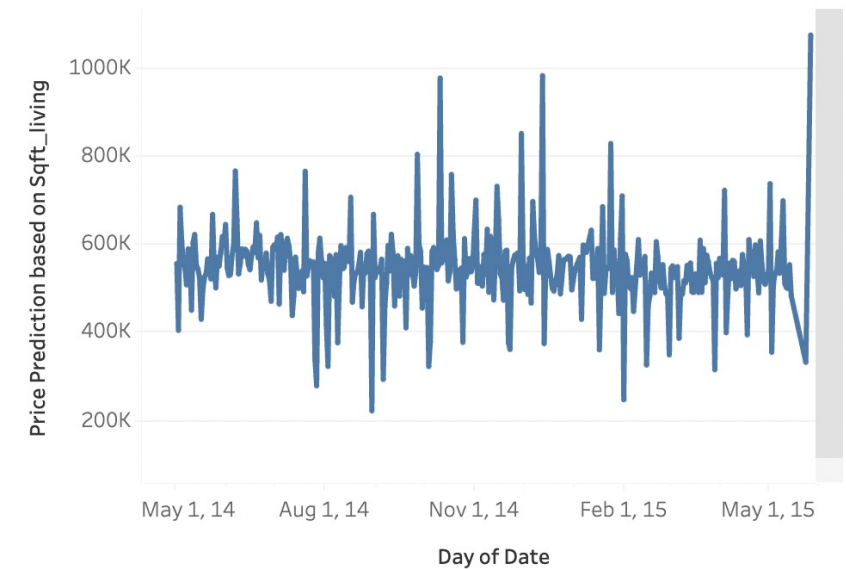*Plotted 4 clusters based on results of the elbow technique*



*Applied an additive decomposition model to split the time series to show the underlying trend, seasonality, and noise.*

# Reflection:

- Although the grade and square footage of the homes were the variables with the greatest correlation to the price, the correlation is not strong enough to build an accurate predictive model.

- Next steps: In order to make an accurate predictive model there need to be more variables with a higher correlation to price or to each other. The dataset I analyzed was sourced from King County's government website which has many other public datasets that could potentially be merged with this one to provide more insights.



Prediction of Price based on Sqft_living variable



Actual Average Price