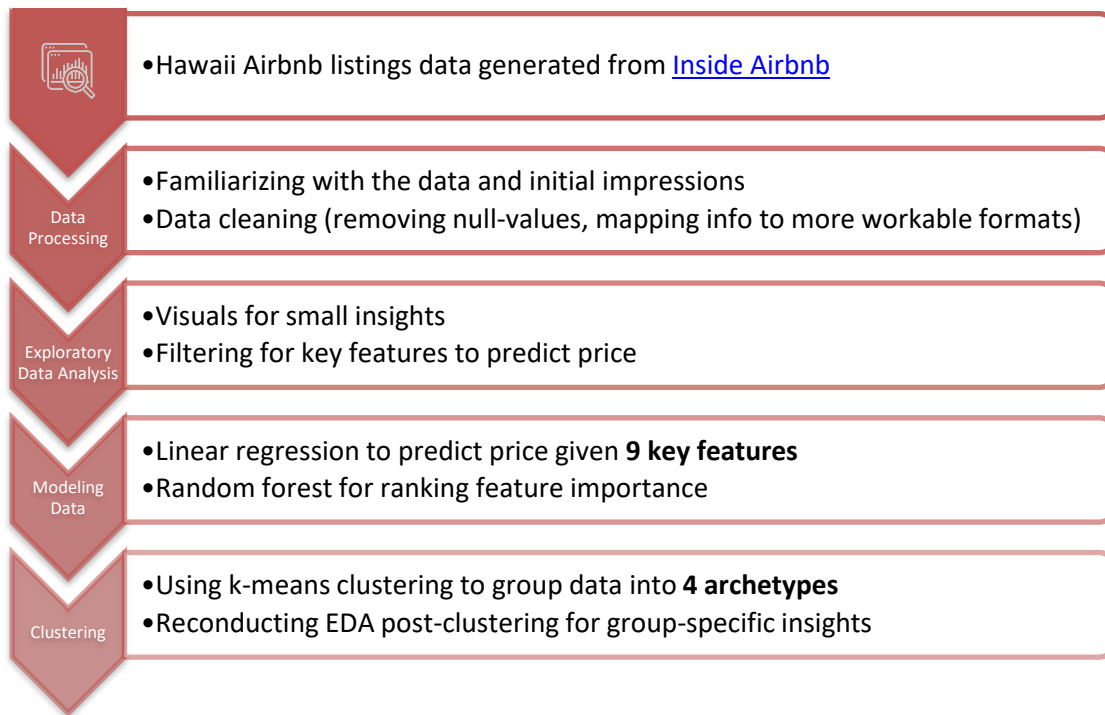


# WHARTON ANALYTICS FELLOWS FALL 2022 DATA CHALLENGE

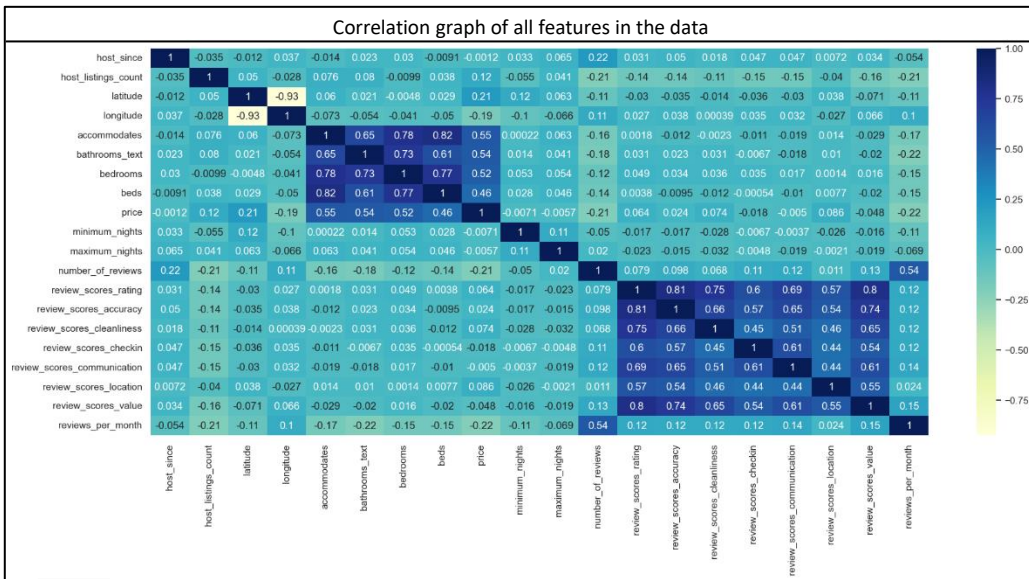
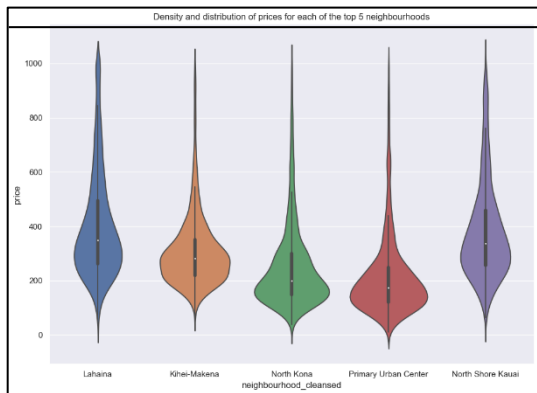
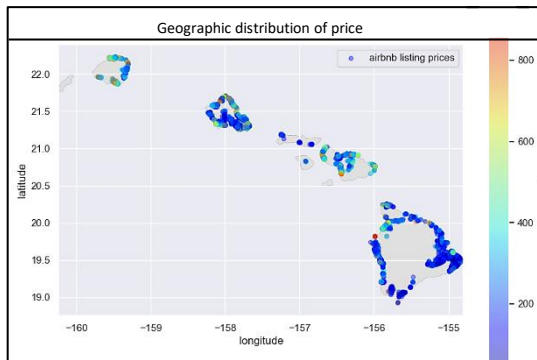
Kyle Liao



# Introduction: Data Analytics Pipeline for Airbnb Data

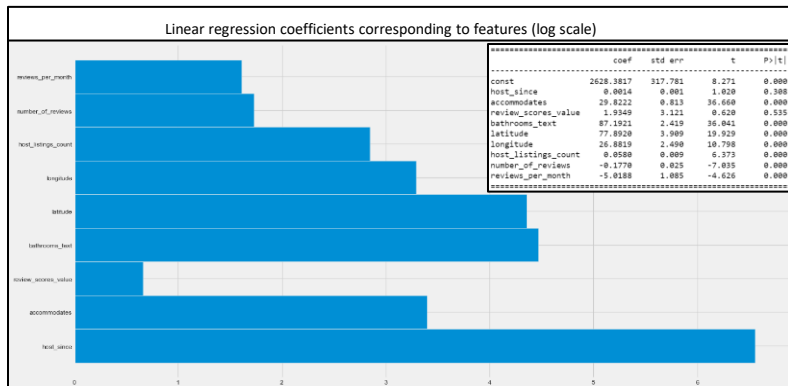


# Part 1. Exploratory Data Analysis



9 key features selected: host experience, # of people accommodated, # of bathrooms, review score, latitude, longitude, host listing count, reviews per month, and number of reviews.

## Part 2. Modeling Data



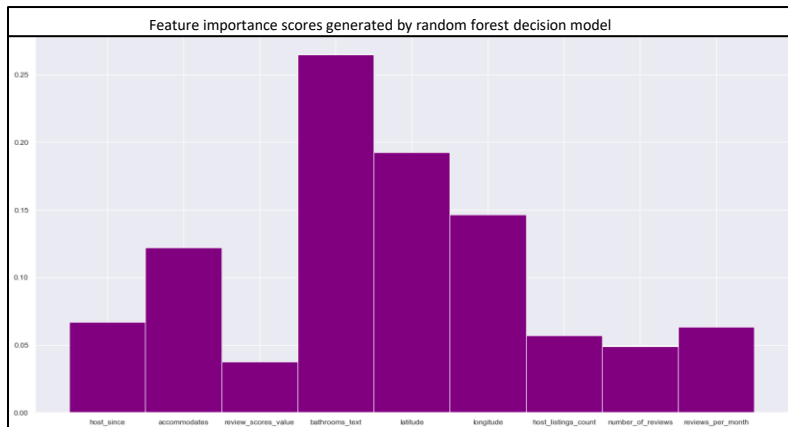
Linear regression coefficients corresponding to features graphed on a logarithmic scale



Larger magnitude implies a larger influence on resulting prediction



Statistics were calculated for each feature; 'host\_since' and 'review\_scores\_value' insignificant



Feature importance scores as determined by random forest decision model

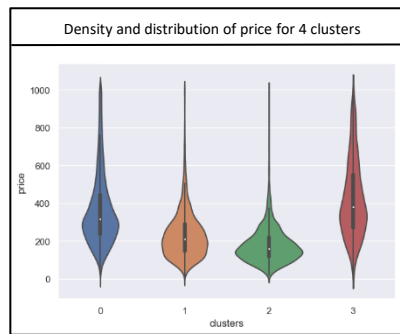
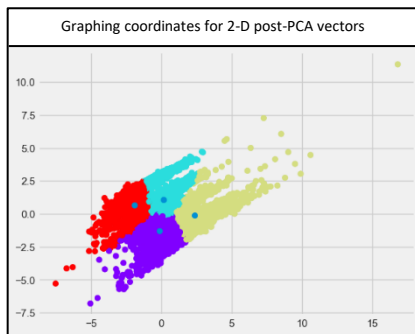
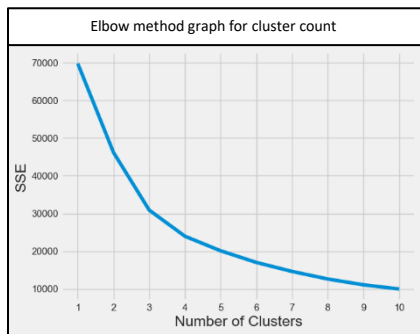


Higher scores (larger magnitude) means that the specific feature will have a larger effect on the price model



# of bathrooms is the most important feature, followed by latitude, longitude, and # of people accommodated

## Part 3. Clustering and Post-EDA



Centroids for four main archetypes of Hawaii Airbnb listings												
Cluster	Days hosting	Host listing count	Latitude	Longitude	Guests	Bathrooms	Bedrooms	Beds	Price	Review Score	Reviews	Reviews/month
1	2259.383	111.5177	20.33584	-156.211	5.463832	1.92021	1.975775	2.988989	363.1719	4.647645	18.05997	0.713817
2	2264.835	35.77911	21.30692	-157.692	3.626732	1.059996	1.195245	1.923624	232.4216	4.73904	42.06496	1.411348
3	2547.212	7.241441	19.95914	-155.826	3.102354	1.085949	1.217546	1.691869	177.5453	4.818006	101.4458	2.615278
4	2182.065	103.6904	21.7785	-158.813	5.822527	1.99726	2.169863	3.218874	423.1081	4.529105	13.3726	0.611513

# Conclusions and Extensions



## Step 1. EDA

- Cleaned and processed data
- Created a geographic distribution for price that displays association between latitude and price
- Found neighborhoods with most listings and graphed price density and distribution
- Created correlation graph to isolate 9 key features for regression model



## Step 2. Modeling

- Trained linear regression model
- Calculated coefficients and corresponding significance values (p-values) for 9 key features identified in EDA
- Trained a random decision forest to rank feature importance



## Step 3. Clustering

- Scaled data down for use in clustering algorithm
- Used PCA to reduce dimension of data from 9 to 2 (for graphing)
- Used K-means clustering algorithm to group data
- Found 4 main archetypes of Hawaii Airbnb listings, as well as centroid for each



## Conclusions

- According to the linear regression model, the top 3 features for predicting price are # bathrooms, # of people accommodated, and the latitude
- According to the random decision forest model, the top 3 most important features are # bathrooms, latitude, and longitude, followed by the # of people accommodated



## Extension and Improvement

- If I had more time, I would do the following:
- See how a gradient-boosted tree would predict prices based on features
- Rigorously check for collinearity before selecting features for models
- Use Tableau for EDA visualizations (ran into pandas export error)
- Learn how to implement NLP on text descriptions description



Thank you for your time and consideration!

