# How to price your Harlem Airbnb?

Jan-Willem Reijnen July - 1- 2020
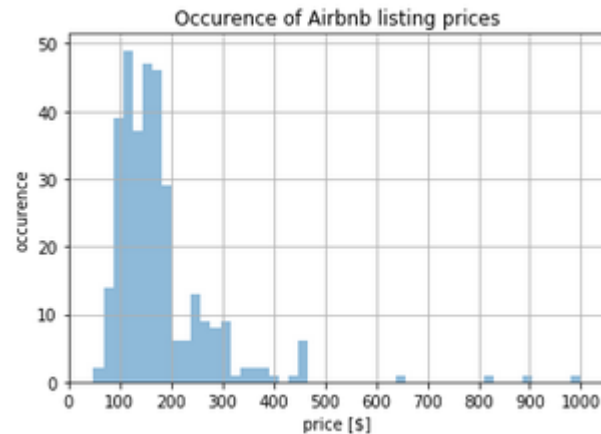
Source: Google Maps 2020

# Overview

- Introduction
- Data description
- Methodology
- Results
- Discussion & Conclusion

# Introduction

- **"Airbnb's** growth is alarming and threatening hoteliers. Having recorded more than **4 million spaces** for rent **across the world** in 65,000 cities and 191 countries, the company is waxing strong in the **United States** with approximately **600,000** listings."

- What can be learnt from current listings for future listings

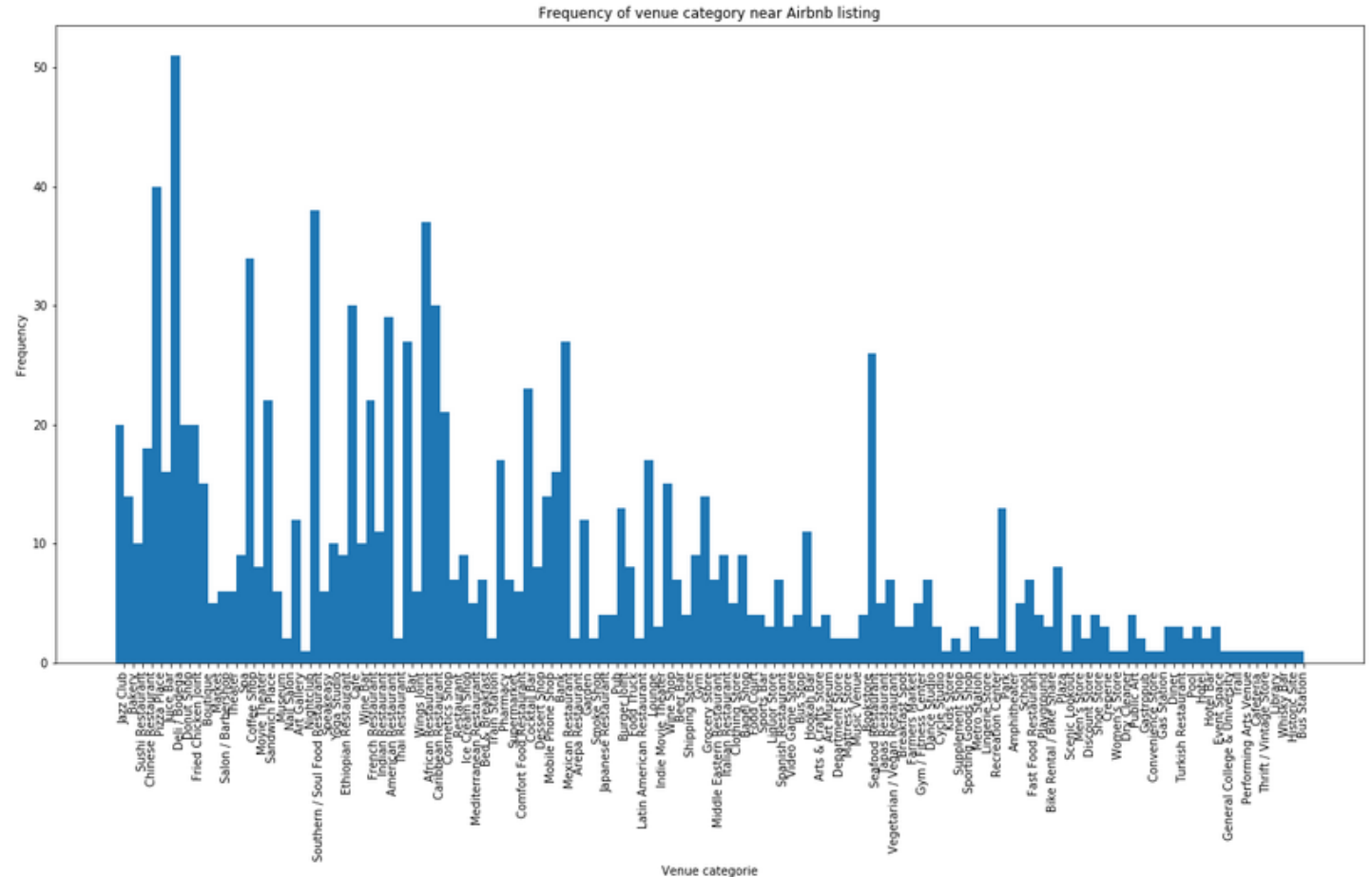- Goal: Can the venues nearby an Airbnb listing be used for future price consulting?

# Introduction

- NYC has 47900 listings in 2019
- Limited to Harlem due to Foursquare request restrictions



Neighbourhood vs. Number of listings

# Data description

- Filtered set contains 263 listings in Harlem
- All listings are full apartments or houses



Occurence of Airbnb listing prices

|       | longitude  | latitude  | minimum_nights | availability_365 | price       |
|-------|------------|-----------|----------------|------------------|-------------|
| count | 333.000000 | 333.000000| 333.000000     | 333.000000       | 333.000000  |
| mean  | -73.947228 | 40.814401 | 5.327327       | 169.843844       | 177.894895  |
| std   | 0.004749   | 0.008532  | 8.163488       | 110.718735       | 106.076032  |
| min   | -73.957980 | 40.798910 | 1.000000       | 1.000000         | 49.000000   |
| 25%   | -73.950160 | 40.807350 | 2.000000       | 55.000000        | 120.000000  |
| 50%   | -73.946650 | 40.813330 | 3.000000       | 188.000000       | 150.000000  |
| 75%   | -73.943750 | 40.822150 | 4.000000       | 264.000000       | 200.000000  |
| max   | -73.936340 | 40.831350 | 60.000000      | 364.000000       | 1000.000000 |

# Data description

- Venue occurences
- 127 venue types



Frequency of venue category near Airbnb listing

# Methodology

- 2 types of regression: Decision tree regression and Multiple linear regression

- Fitted on 80% of data set, tested on 20%

- 131 features to consider in regression
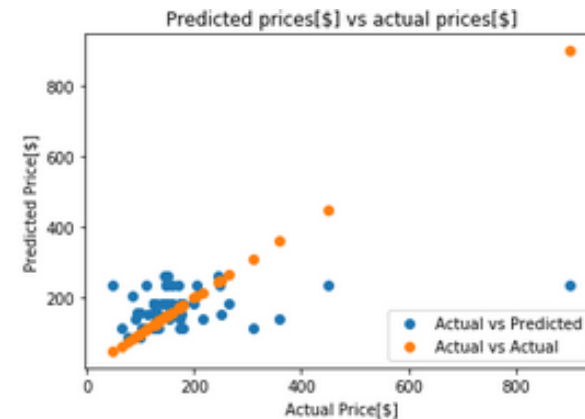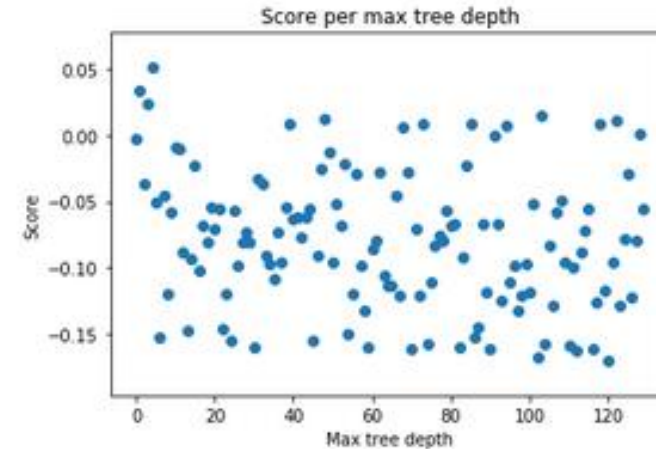
- Decision tree regression:

Multiple linear regression:

$$y = a\,x_1 + b\,x_2 + c\,x_3 \ldots$$

# Decision tree regression

- Optimal depth for decision tree:



- Best results from optimal decision tree:
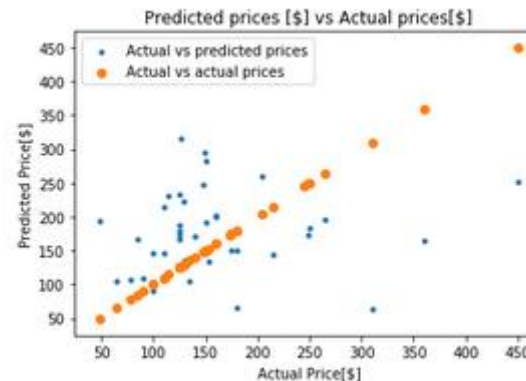  - Not very close, low accuracy and high variance

# Multiple Linear Regression

- Best fitted model has R-squared of .849 (very good)



- However, prediction accuracy still quite off:

# Results

- From multiple regression the best venues to have nearby are:
  - Gastropub
  - Food truck
  - Comfort food restaurant
- The worst venues to have nearby:
  - Chinese restaurant
  - Smoke shop
  - Hookah bar

# Discussion & conclusion

- Model should be tested with more data, as 131 features are to many for 263 observations

- However, Multiple regression performs way better and could definitely be used