

# Predicting house prices in King County, Washington, USA

By Panagiotis Petsas

This report is part of the course “Applied Data Science Capstone” from “IBM Data Science Professional Certificate” on Coursera.

## Part 2: Data

### 1 Introduction

(in the previous document)

### 2 Data acquisition and cleaning

#### 2.1 Data sources

We obtained data on houses in King County from Kaggle [3]. This dataset contains more than 21.000 house records, providing information on features like price, year build, number of bedrooms, longitude, latitude, zip code etc. The complete list of features is listed in **Table 1**.

In order to reduce computation time, we decided to work only with a specific area. After some initial data cleaning (see below), we scraped the area phone code of each house by utilizing its zip code. We maintained records with area phone code 253. These houses are located in an area near Tacoma, a medium sized city (**Figure 1**). We used ‘www.getzips.com’ site in order to acquire the area phone codes. This procedure reduced the house records to 2766, a more manageable size for modeling.

Figure 1 Houses in King County, Washington with area phone code as 253, marked with blue dots.

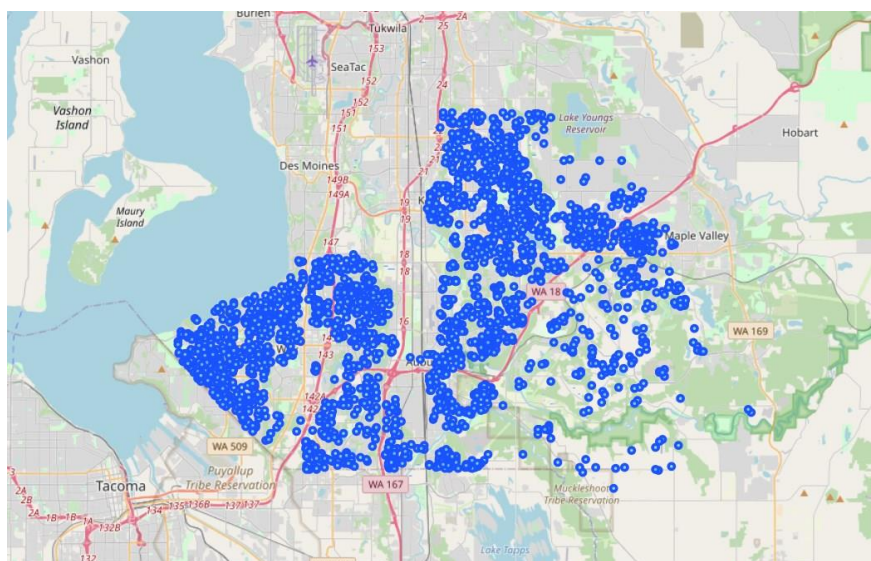


Table 1 Features in the dataset of houses in King County, Washington, along with their explanation.

Feature name	Explanation
<b>ID</b>	Unique id for each house record
<b>Date</b>	Date of home sale
<b>Price</b>	House price
<b>Bedrooms</b>	Number of Bedrooms
<b>Bathrooms</b>	Number of Bathrooms
<b>Sqft_living</b>	Apartment inter living space (in sq feet)
<b>Sqft_lot</b>	Lot space (in sq feet)
<b>Floors</b>	Number of floors
<b>Waterfront</b>	Binary (1 if the apartment overlooks the waterfront, 0 otherwise)
<b>View</b>	From 0-4 integer, a grade for house's view
<b>Condition</b>	From 1-5 integer, the condition of the house
<b>Grade</b>	From 1-13 integer, the construction and design level
<b>Sqft_above</b>	Apartment inter living space above ground (in sq feet)
<b>Sqft_basement</b>	Apartment inter living space below ground (in sq feet)
<b>Yr_built</b>	The year the house was built
<b>Yr_renovated</b>	Year of last renovation (0 if never renovated)
<b>Zipcode</b>	The zipcode of the area
<b>Lat</b>	Latitude
<b>Long</b>	Longitude
<b>Sqft_living15</b>	Apartment inter living space (in sq feet) for the 15 nearest neighboring houses
<b>Sqft_lot15</b>	Lot space (in sq feet) for the 15 nearest neighboring houses

Having the longitude and latitude of each house, we used Foursquare API [4] in order to extract the most popular venues in a proximity of 500 meters from it. We limited the API to the top 30 venues per house, in order to reduce computation time. This process derived a total of 8995 venues near the houses, providing information about their name, category and coordinates.

## 2.2 Data cleaning

The dataset did not contain any missing records. Some initial data cleaning was conducted by removing strange records (houses with zero bedrooms or bathrooms). The dataset contained a feature that indicated the year of renovation. Most houses were not renovated, thus they received a value of 0. For this reason, we created a new feature that indicated whether a house was renovated (1 if yes and 0 if not). After these adjustments, the next part of the cleaning process was to maintain the records with area phone code of 253.

The next part was to group the venue records based on their category. While Foursquare provides their own general categories, we preferred to classify the venues in our custom groups. We created 13 distinct groups (**Table 2**) to classify a total of 254 unique venue categories. The grouping was performed by searching on keywords within the venue category names. For example, venue category names including the word 'restaurant' would be put in the group of 'places to eat'. We ensured that categories that were grouped were removed from the grouping process, in order to avoid grouping them twice. There were some venue categories that could not be grouped based on keywords, which we grouped manually to the correct group.

In the next session, we summarized the amount of venues per group for each house. This process created 13 new features, one for each group, indicating the amount of venues of this group in the near proximity of each house. For example, a house with 2 places to eat and 3 venues related to services would get a value of two in the 'places to eat' feature, 3 in the 'venues related to services' features and 0 to the remaining 11 features related to venue groups.

## 2.3 Feature selection

Since our purpose is to use features to predict the house price, we must maintain features that could provide such information. We initially discarded 'id' and 'date' (of recoding the house), as they are used for reference only. Similarly we discarded the 'zipcode' and 'area phone code'

as they provide no further information. The year of renovation feature was replaced by the binary feature 'is renovated' stated above.

We performed pearson correlation analysis to identify linear patterns among the remaining features. We decided to discard four more features that caused high correlation, being 'above area in sqft' (after all, we maintained the 'total area in sqft', which is exactly equal to the summary of the 'above' and 'basement area in sqft'), 'grade' (which is a subjective value anyway), 'sqft\_living15' and 'sqft\_lot15' which indicate the total living and lot area of the 15 neighboring houses, and have high correlation with their respective features. With that procedure, we maintained 26 features for predicting the house prices.

Table 2 Groups of venue categories, along with some keywords used to identify the right venues.

Group	Some keywords used
<b>Places to eat</b>	Restaurant, food, pizza, burger, chicken, fish, breakfast, ...
<b>Places to drink</b>	Bar, wine, beer, pub, coffee, café, ...
<b>Places with sweets</b>	Candy, donut, dessert, pie, ice cream, ...
<b>Amusement</b>	Playground, arcade, casino, roller, night, ...
<b>Facilities</b>	Facility
<b>Sports</b>	Soccer, basketball, baseball, sport, gym, fit, ...
<b>Transport</b>	Airport, boat, ferry, bus, train, ...
<b>Shops</b>	Shop, market, store
<b>Services</b>	Service, assist, bank, doctor, police, ...
<b>Beauty services</b>	Spa, massage, salon, beauty, tattoo
<b>Attractions</b>	Art, history, attract, museum, ...
<b>Landmarks</b>	Mountain, lake, tree, ...
<b>Accommodation</b>	Hotel, motel, inn, rest, ...

[3] House prices in KC, Washington – Kaggle

(<https://www.kaggle.com/harlfoxem/housesalesprediction> )

[4] Foursquare API (<https://developer.foursquare.com/> )