# Can we predict house prices using known features of each house and a supervised learning approach?

Florence Galliers

18/10/2020

## Background:

House prices are an important part of the economy... Motivate the story... What do we already know? What lessons will the story teach us?

When purchasing a house people are more likely to look only at the number of bedrooms and number of bathrooms, rather than the square foot of the house. This may be because they do not understand how square foot is calculated, or find it difficult to visualise how this translates into the house size.

### Objectives:

- Understand which independent variables in the data set can be used to predict house price (the dependent variable)
- Minimize the differences between predicted and actual house prices by using model selection to choose the most accurate model.

## Data:

The data for this analysis contains house prices and information regarding the features of each house. It has 4600 entries. The original data set downloaded from Kaggle had 17 independent variables, however I felt that 7 of these were not relevant to this analysis and so they have been removed. The dependent variable in this analysis will be the house price (in US dollars).

This data is from houses sold in 2014 in Washington, USA. However it is an interesting data set that contains a large amount of information and number of variables. I was not able to find a similar data set from the UK, and so I am going to go ahead with this data set as I feel it will produce some interesting results.

### Data Preparation:

- There was originally a variable called sqft_basement which gave the size of the basement if present, however a lot of houses did not have basement so I felt it would be more useful to turn this variable from numeric into binary. It now shows 0 = no basement, 1 = basement.
- I also changed the variable year_renovated into a binary variable, as again, not all the houses had been renovated. 0 = not renovated, 1 = renovated.
- Check if there are any zero values for the price variable as these are not acceptable in a housing price situation, a house cannot cost nothing, so we must remove these and assume they are errors in the data set.
- Remove outliers

The remaining variables and their descriptions are shown in Table 1:

Table 1: Data Dictionary

| Variable | Description |
| --- | --- |
| price | House sale price in thousands of US dollars |
| bedrooms | Number of bedrooms |
| bathrooms | Number of bathrooms |
| sqft_living | Area of house in square feet |
| sqft_lot | Area of whole housing lot in square feet |
| floors | Number of floors in the house |
| condition | Condition of house, 1 to 5 |
| if_basement | 1 = if house has a basement, 0 = no basement |
| house_age | Year that the house was built subtracted from 2014 |
| if_renovated | 1 = if house has been renovated, 0 = if no renovation |
| city | Location of house to the nearest city in Washington, USA |

## Methods:

As house price is a continuous variable I have taken a supervised learning approach and will be using regression to look at the relationship between house price and features of each house.

Review approaches tried or considered

Summary of final approach and justification of why this approach was chosen:

## Results:

Summary of major results, graphs, diagnostic outputs

Strictly relevant to the objectives

Must include a link to the Github repository containing a fully reproducible and documented analysis Reported in scientific style.

## Conclusions:

~1 paragraph

**Literature Cited:**

## 3-5 peer reviewed references

- 
- 
- 
-