



# King County House Prices

for Investors and Real estate agents

Christoph Bickle



# Contents

1. Overview of the given data
2. Parameter relations
3. Recommendations
4. Multivariate linear regression
5. Summary

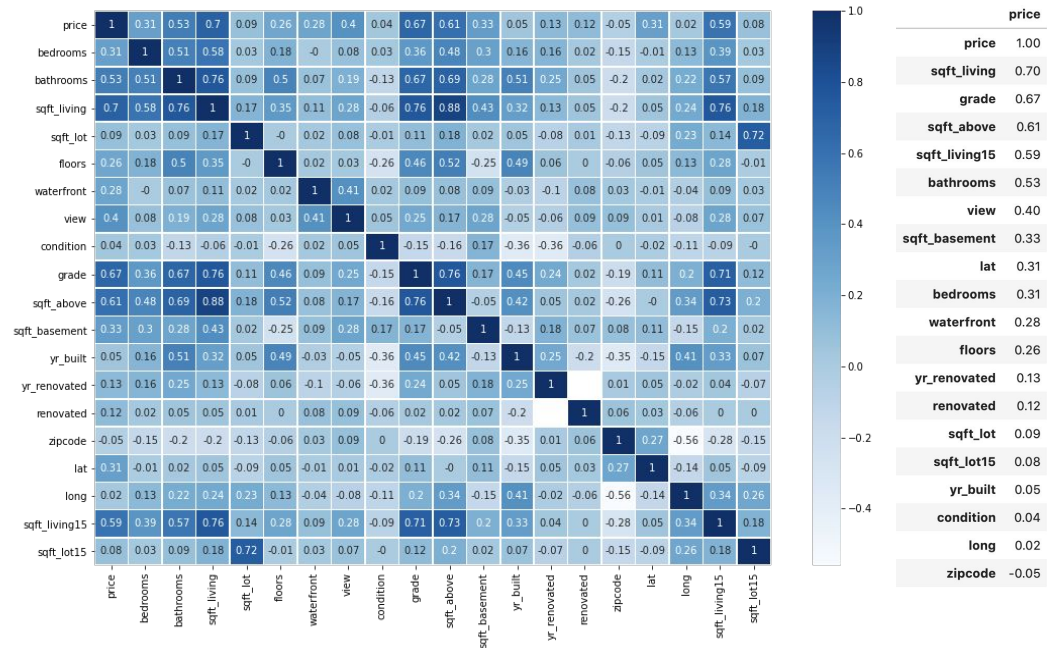


# Overview of the given data

- The data for 21597 house in King County was given
- The dataset is mostly complete with missing values in only a few variables
- Strong outliers can be found in sqft\_living, bedrooms, bathrooms and price
  - Probably luxury houses
- The outliers were not removed as these may also be interesting to high class investors and real estate agents

# Overview of the given data

- Correlation between variables
- Correlation with price
- Highest with sqft\_living and grade

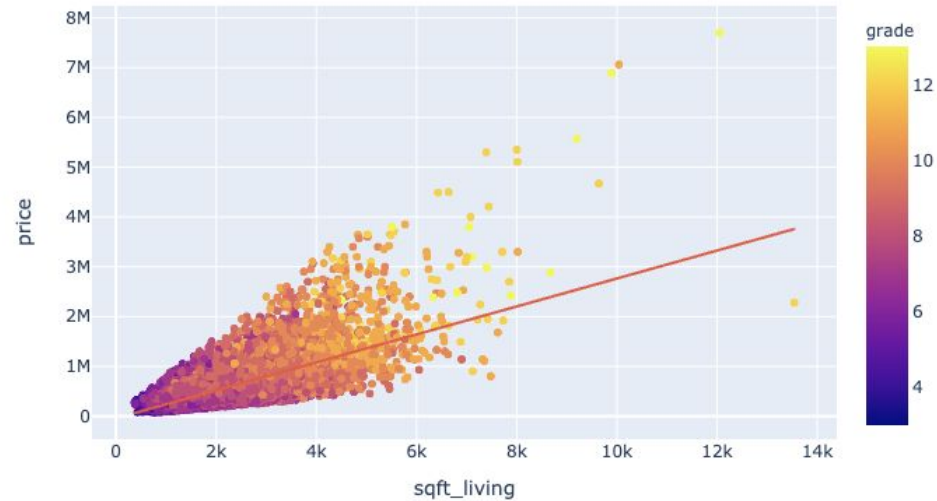


# Parameter relations

- Price is mostly influenced by sqft\_living and grade
- These go hand in hand

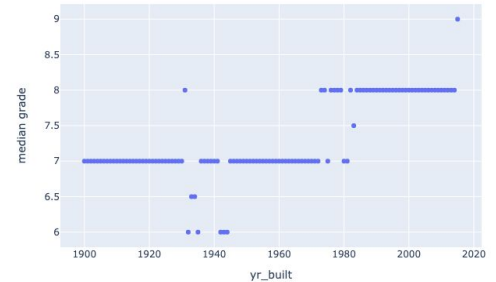
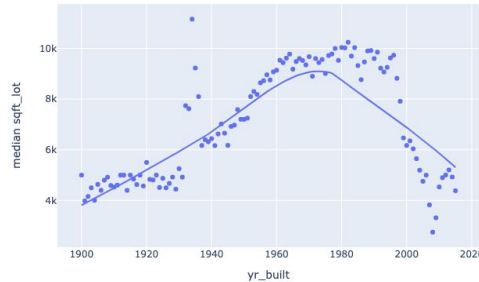
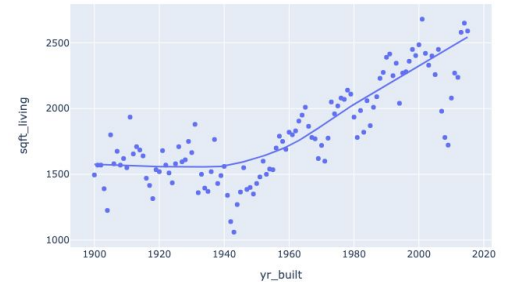
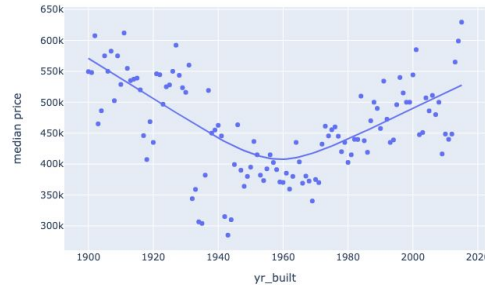
	grade
sqft_living	0.762779

- A positive linear correlation between price and sqft\_living is observable



# Parameter relations

- Older house are more expensive even though they offer less space
- Grades are somewhat constant with small dips in the 40s and an increase in the 90s



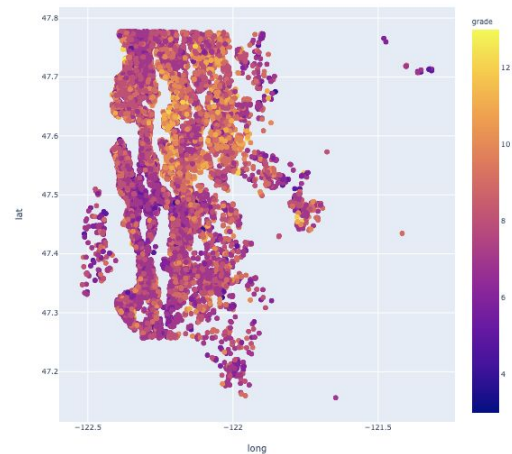
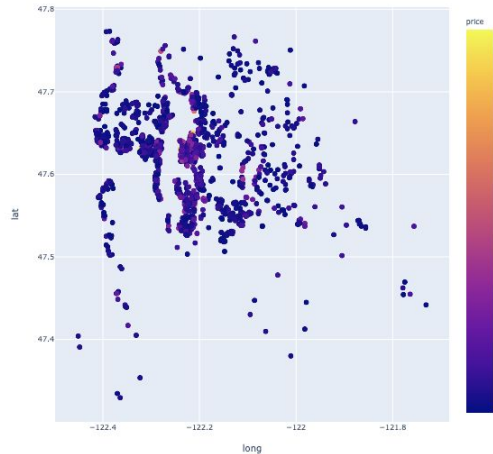
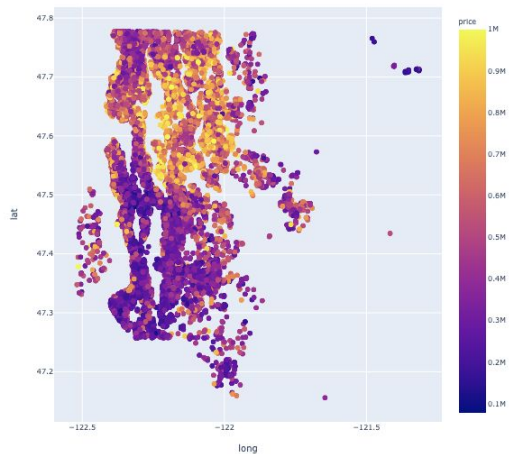
# Parameter relations

- Literature researches show that seasons influence sales prices
- The data only contains one year of data
- A small fluctuation is still visible

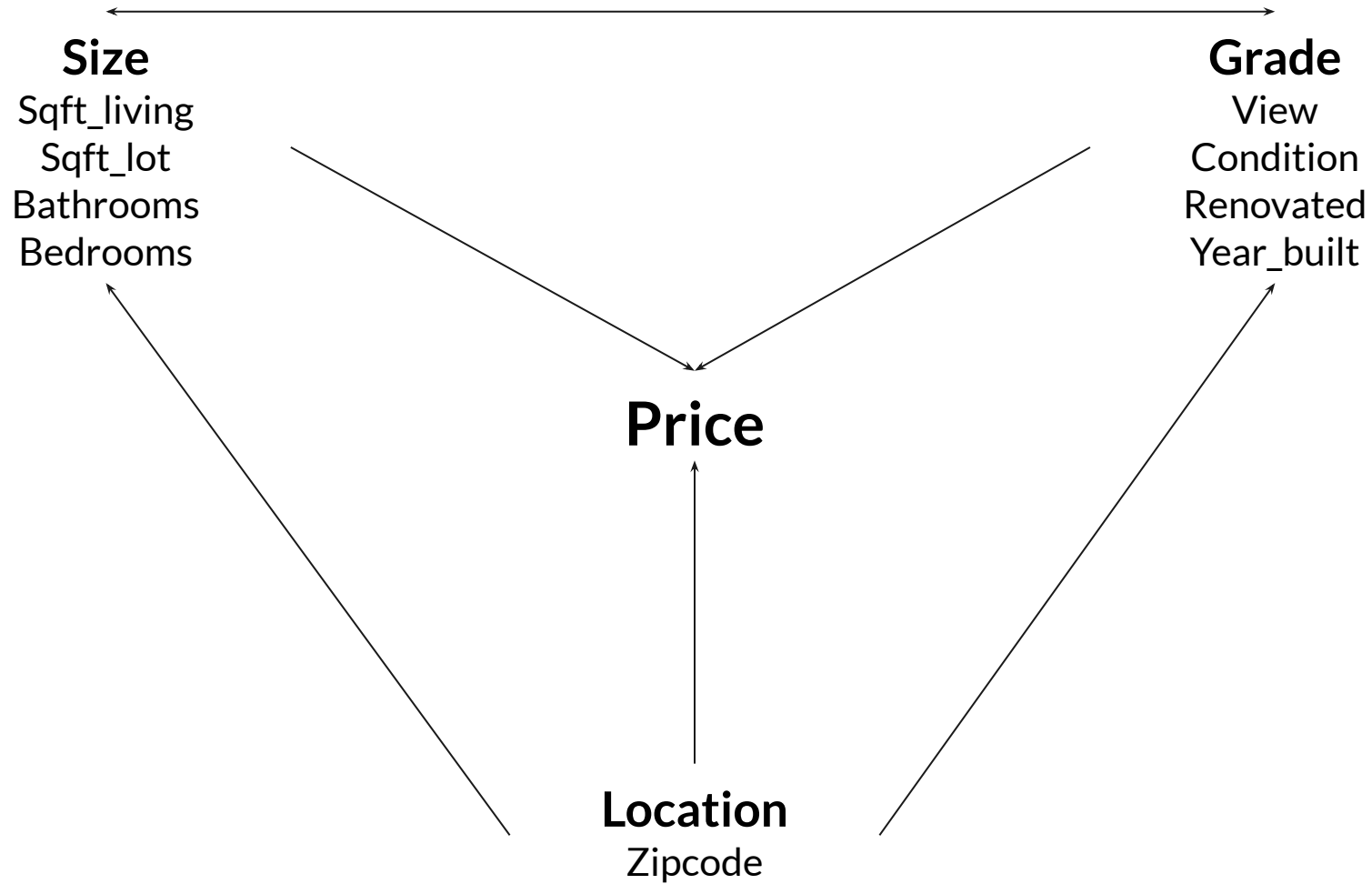


# Parameter relations

- Prices correlate with location, but this seems to be the consequence of higher grades in the area









# Recommendations

- When was the house built?
  - The newer the house the more expensive it will get to buy in the future
  - Buy houses built in the 40s to the 80s
- How big is the lot in comparison to the actual living space?
  - Buy big lots to leave room for living space expansion
  - Adding new floors is more expensive making inner city houses less desirable and hard to improve
  - Add bathrooms and bedrooms if needed, but not as sole means to improve the resale price
- Why is a grade high or low?
  - Avoid houses whose grades are high due to their location, since improvement is difficult. Thus making a buy and resell less beneficial
  - Improve grade by renovation and expansion
- Keep the seasons in mind



# Multivariate linear regression

- The algorithm uses only 3 variables:
  - sqft\_living, grade and zipcode
- Thereby being reliable and less prone to errors when using new datasets
- The 3 variables cover most other parameters, since sqft\_living also correlates with number of bathrooms, grade with view etc.
- The mean absolute percentage error at this point in time equals 18.9 %
- A better MAPE can be achieved when more variables are passed, but this could come at the cost of reliability



# Summary

- The King County house price dataset from 14/15 was used to understand influences of different parameters
- The goal was to give recommendations to real estate agents and investors
- Multivariate linear regression can be used to predict house prices for new data