# Modelling property prices in King County

Georg F.K. Höhn

Ironhack

6 June 2025

# Overview

# Introduction

# Introduction

- dataset from property sales in King County, Seattle, WA, USA
- timeframe: between May 2014 and May 2015
- source: https://www.kaggle.com/datasets/minasameh55/king-country-houses-aa

## Aims

- predict sales prices based on dataset
- (closer look at subset of data with price $\geq$ \$650k)

EDA

# Overview

- ▶ 21 613 datarows
- ▶ 21 columns
- ▶ descriptions in meta-data, but some clarifications
  - ▶ grade up to 13 (instead of 11)
  - ▶ views 0–4 (instead of binary)
  - ▶ id identifies properties not transactions
- ▶ target variable `price`
  - ▶ range: 75k–7700k
  - ▶ mean: 540088,14
  - ▶ median: 450k
  - ▶ right-skewed

# Pre-processing and feature engineering

- ▶ no missing values`
- ▶ duplicates for `id` explained by multiple sales of same property
  - ▶ created new column `prev_sale_within_year`
- ▶ extracted month and day (as ints) from `date`
- ▶ later:
  - ▶ drop zipcode
  - ▶ drop day

# Transforming data

- transformed following data into categorial:
    - `yr_renovated`
    - `yr_built`
    - `sqft_basement`
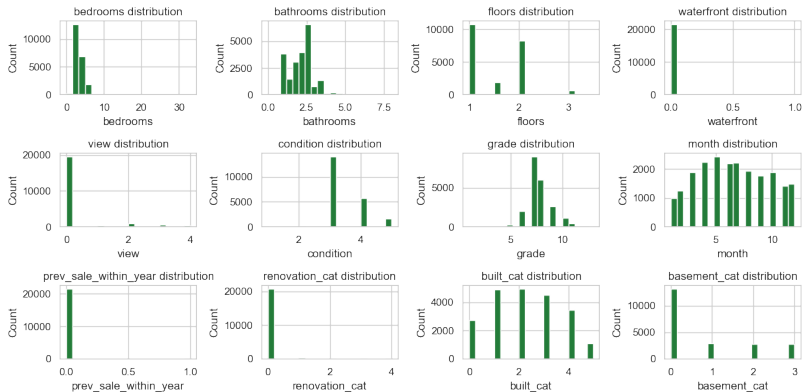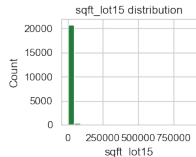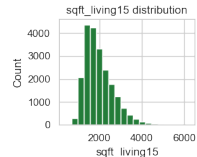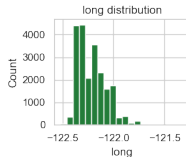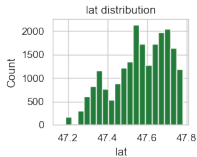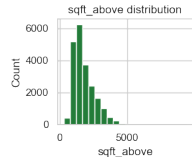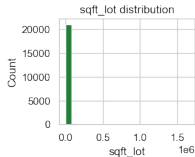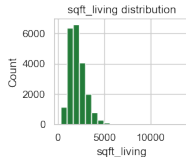- two variants:
    - labels
    - one-hot-encoding

# Scaling

- kept a version of unscaled data
- scaling using `sklearn.compose.ColumnTransformer`
  - apply StandardScaler to continuous columns
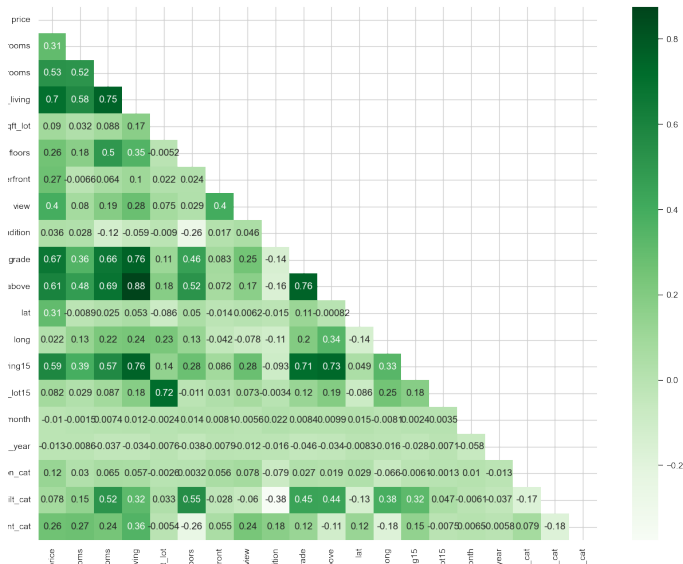  - pass through categorial variables unchanged

# Distributions

# Barplots for categorial data

# Histograms for continuous data

# Correlation matrix

| | price | ooms | ooms | ving | t_lot | oors | front | view | ition | rade | bove | lat | long | ng15 | ot15 | onth | year | _cat | _cat | _cat |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| price | | | | | | | | | | | | | | | | | | | | |
| rooms | 0.31 | | | | | | | | | | | | | | | | | | | |
| rooms | 0.53 | 0.52 | | | | | | | | | | | | | | | | | | |
| _living | 0.7 | 0.58 | 0.75 | | | | | | | | | | | | | | | | | |
| qft_lot | 0.09 | 0.032 | 0.088 | 0.17 | | | | | | | | | | | | | | | | |
| floors | 0.26 | 0.18 | 0.5 | 0.35 | -0.0052 | | | | | | | | | | | | | | | |
| rfront | 0.27 | -0.0066 | 0.064 | 0.1 | 0.022 | 0.024 | | | | | | | | | | | | | | |
| view | 0.4 | 0.08 | 0.19 | 0.28 | 0.075 | 0.029 | 0.4 | | | | | | | | | | | | | |
| dition | 0.036 | 0.028 | -0.12 | -0.059 | -0.009 | -0.26 | 0.017 | 0.046 | | | | | | | | | | | | |
| grade | 0.67 | 0.36 | 0.66 | 0.76 | 0.11 | 0.46 | 0.083 | 0.25 | -0.14 | | | | | | | | | | | |
| above | 0.61 | 0.48 | 0.69 | 0.88 | 0.18 | 0.52 | 0.072 | 0.17 | -0.16 | 0.76 | | | | | | | | | | |
| lat | 0.31 | -0.0089 | 0.025 | 0.053 | -0.086 | 0.05 | -0.014 | 0.0062 | -0.015 | 0.11 | -0.0082 | | | | | | | | | |
| long | 0.022 | 0.13 | 0.22 | 0.24 | 0.23 | 0.13 | -0.042 | -0.078 | -0.11 | 0.2 | 0.34 | -0.14 | | | | | | | | |
| ng15 | 0.59 | 0.39 | 0.57 | 0.76 | 0.14 | 0.28 | 0.086 | 0.28 | -0.093 | 0.71 | 0.73 | 0.049 | 0.33 | | | | | | | |
| _lot15 | 0.082 | 0.029 | 0.087 | 0.18 | 0.72 | -0.011 | 0.031 | 0.073 | -0.0034 | 0.12 | 0.19 | -0.086 | 0.25 | 0.18 | | | | | | |
| month | -0.01 | -0.0015 | 0.0074 | 0.012 | -0.0024 | 0.014 | 0.0081 | 0.0056 | 0.022 | 0.0084 | 0.0099 | 0.015 | -0.008 | 0.0024 | 0.0035 | | | | | |
| _year | -0.013 | -0.0086 | -0.037 | -0.034 | -0.0076 | -0.038 | -0.0079 | -0.012 | -0.016 | -0.046 | -0.034 | -0.0083 | -0.016 | -0.028 | -0.0071 | -0.058 | | | | |
| n_cat | 0.12 | 0.03 | 0.065 | 0.057 | -0.0026 | 0.0032 | 0.056 | 0.078 | -0.079 | 0.027 | 0.019 | 0.029 | -0.0066 | 0.0061 | 0.0013 | 0.01 | -0.013 | | | |
| ilt_cat | 0.078 | 0.15 | 0.52 | 0.32 | 0.033 | 0.55 | -0.028 | -0.06 | -0.38 | 0.45 | 0.44 | -0.13 | 0.38 | 0.32 | 0.047 | -0.0061 | -0.037 | -0.17 | | |
| nt_cat | 0.26 | 0.27 | 0.24 | 0.36 | -0.0054 | -0.26 | 0.055 | 0.24 | 0.18 | 0.12 | -0.11 | 0.12 | -0.18 | 0.15 | -0.0075 | 0.0065 | 0.0058 | 0.079 | -0.18 | |

# Implementation aspects

# Helper functions

- `mk_histplots`: generate histplots for a (subset of a) dataframe, optionally saving output
- `mk_barplots`: generate countplots for a (subset of a) dataframe, optionally saving output
- `mk_boxplots`: generate barplots for a (subset of a) dataframe, optionally saving output
- `plot_corrmatrix`: plot a correlation matrix for a dataframe, optionally saving output
- evaluation function
- function for plotting coefficients in linear models

# Model comparison

# Overall approach

- ▶ running different models including some hyper-feature tuning
  - ▶ Linear Regression Models
  - ▶ KNN
  - ▶ GradientBoost
  - ▶ RandomForests
  - ▶ XGBRegressor
- ▶ evaluation function printing for training and test sets:
  - ▶ $r^2$-score
  - ▶ adjusted $r^2$-score (penalising models with more features)
  - ▶ mean-squared-errors
  - ▶ mean-absolute-errors
  - ▶ ratio of MAE to maximum target value
- ▶ evaluation function returns adj. $r^2$ and MAE for training and test for collection in dataframes

# Data collection

- ▶ lists for each model family transformed into DataFrame
- ▶ sorted by adj. r² and MAE for test sets
- ▶ top 3 for each family included in an overall list

## Results

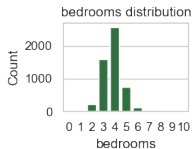| type | scaled | $r^2$-a-train | $r^2$-a-test | mae-train | mae-test |
|------|--------|---------------|--------------|-----------|----------|
| XGBoost | False | 0.975 | 0.913 | 40600.868 | 65516.939 |
| XGBoost | False | 0.975 | 0.913 | 40600.871 | 65516.938 |
| XGBoost | False | 0.975 | 0.913 | 40600.864 | 65516.941 |
| GradientBoost | True | 0.998 | 0.908 | 3511.192 | 64412.386 |
| GradientBoost | True | 0.996 | 0.906 | 5168.707 | 64658.581 |
| GradientBoost | True | 0.993 | 0.905 | 10322.906 | 63306.014 |
| RandomForests | True | 0.980 | 0.879 | 30229.928 | 71517.825 |
| RandomForests | True | 0.980 | 0.879 | 30229.928 | 71517.825 |
| RandomForests | True | 0.980 | 0.877 | 30653.158 | 72256.111 |
| knn | True | 0.818 | 0.787 | 86150.484 | 95555.056 |
| knn | True | 0.829 | 0.783 | 84326.308 | 97208.734 |
| knn | True | 0.838 | 0.782 | 82691.138 | 97932.769 |
| lm-lasso | True | 0.704 | 0.684 | 125477.227 | 126751.772 |
| lm-lasso | True | 0.704 | 0.684 | 125477.227 | 126751.772 |
| lm-lasso | True | 0.704 | 0.684 | 125477.226 | 126751.891 |

# Hyperparameters of best performing model

**XGBoost**

- ▶ estimators: 1.0
- ▶ booster: gbtree
- ▶ subsample: 1
- ▶ dart_nominalized_type: TREE
- ▶ min_child_weight: 1.0
- ▶ max_depth: 6.0
- ▶ reg_alpha: 0.3
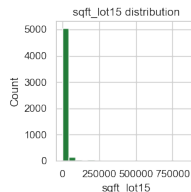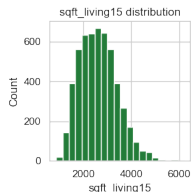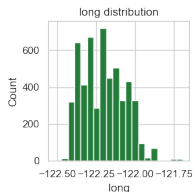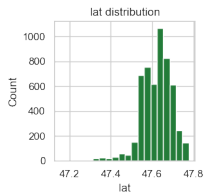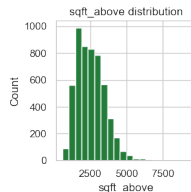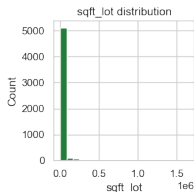- ▶ reg_lambda: 1.0
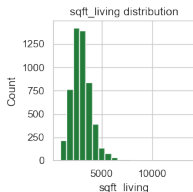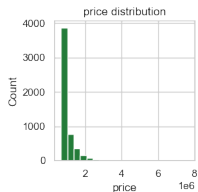- ▶ learning rate: 0.3
- ▶ max_iterations: 20.0

Over $650k

- much smaller dataset (5324 datapoints)
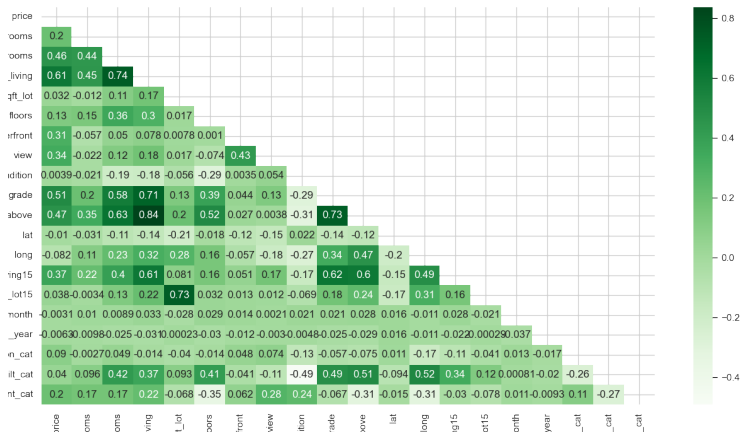
# Barplots for categorial data

# Histograms for continuous data

# Correlation matrix

# Model fitting

- ▶ fit XGBoost on unscaled data using label encoding with a range of hyperparameters
- ▶ best result:
  - ▶ $r^2$-adjusted (train): 0.990268
  - ▶ $r^2$-adjusted (test): 0.761220
  - ▶ MAE (train): 33374.973380
  - ▶ MAE (test): 132851.792547
- ▶ strong overfitting

# Conclusion/outlook

- ▶ more cleanup
- ▶ visualising results would be nice
- ▶ using RFE (?) to consider rank of features
- ▶ GridSearchCV instead of manual looping?
- ▶ idea for lat/long:
    - ▶ calculate distance matrix
    - ▶ dimensional scaling to 1 or 2 dimensions
    - ▶ provides small metric of relative distance for each datapoint (might help to model that house price is also influenced by spatially defined neighbourhoods)

Thanks for your attention!