

```

---
title: "A Statistician's Priority List for Boosting Home Value"
theme: simplex
output:
  html_document: default
  pdf_document: default
  word_document: default
---

```

Introduction

For many people, a house is not only a residence but also a place where they have been investing throughout their stay. Hence, how to maintain and boost home value during their stay has been a question for many house owners. Generally speaking, there are many commonly known factors that would help increase home values, however, for most people with a limited budget, it is hard to take everything into consideration when they want to boost their home values. Therefore, it is of great significance to learn what should be prioritized during the home improvement with a purpose of value boost. Although home owners are unable to obtain everything they want with a tight budget, they can do the things that really matter and bring up the sale prices of the house by prioritizing the controllable things. To obtain a priority list for house improvement and home value bringing up, statistical methods like linear regression and random forest would be utilized in this project to analyze a Kaggle dataset containing house sale prices of King County, Washington from May 2014 to May 2015. Statistical models would be constructed to find out the most significant house attributes that are related to home prices.

Purpose

The primary objective of this project is to find out a priority list for home improvement that can be helpful for boosting home values from a statistician's view. Commonly used house attributes would be analyzed and the project would be planned to figure out the most relevant features of a house regarding sale prices. Hopefully, this project could offer some suggestions on house improvement and home value boosting for home investors during their stay.

Background

As the most populous county in Washington with the largest city of the state, Seattle, sitting in the west, King County embraces quite a few large corporations (e.g., Boeing, Microsoft, Amazon). The map below highlights the exact location of King County.

```

```{r load packages, include = FALSE}
Load required packages
library(caret)
library(ggplot2)
library(GGally)
library(dplyr)
library(randomForest)
library(leaps)
library(knitr)
library(kableExtra)
```

```{r map, echo=FALSE, fig.cap = "Map of King County, Washington", fig.align = "center"}
include_graphics("king_county_map.png")
```

```

The strong local economy has been driving the increase in the number of households in this area, which has been creating the demand from the real estate market. The graphic below displays the right-skewed distribution of house sale prices in the time window we analyzed. The house price in King County ranges from $\$75,000$ to $\$7,700,000$, and the average price is slightly above $\$540,000$.

Data

A real-world dataset that contains house sale price information and the corresponding house features

of King County, Washington from 2014 to 2015 will be used. It is originated from Kaggle, and can be imported to R from `mlr3data` package.

Basically, there are 21,613 observations along with 19 house features such as the number of bathrooms, bedrooms, floors, and square footage of the house in the original data. The code that help us load the data and print the first few lines of the original data is shown as below

```
```{r include = TRUE,warning=FALSE}
Load the required dataset
library(mlr3data)
data("kc_housing")
head(kc_housing)
```
```

```
```{r include = FALSE,warning=FALSE}
Print out the summary statistics
summary(kc_housing)
```
```

```
```{r include = FALSE,warning=FALSE}
Print out the data dimensions
dim(kc_housing)
```
```

Variables

The full variable dictionary is summarized as below:

****id:**** unique ID of the house

****date:**** the sale date of the house

****price:**** the final sale price of the house

****bedrooms:**** count of bedrooms in the house

****bathrooms:**** count of bathrooms in the house

****sqft_living:**** square footage of the living area in the house

****sqft_lot:**** square footage of the lot for the house

****floors:**** total levels in the house

****waterfront:**** whether the house has a waterfront view. If yes, the value is 1. Otherwise, the value is 0.

****view:**** how many times the house has been viewed

****condition:**** the overall condition of the house

****grade:**** the overall grade given to the housing unit by King County grading system. According to [King County Assessor's webpage](<http://info.kingcounty.gov/assessor/esales/Glossary.aspx?type=r>), this represents the construction quality of improvements. Grades run from grade 1 to 13.

****sqft_basement:**** square footage of the basement

****sqft_above:**** square footage of the house apart from the basement

****yr_built:**** which year the house was built

****yr_renovated:**** which year the house was renovated. If no renovation has been done, the value is 0

****zipcode:**** the zip code for the house address

```

**lat:** latitude coordinate of the house location

**long:** longitude coordinate of the house location

**sqft_living15:** square footage of the living area in the house measured in 2015

**sqft_lot15:** square footage of the lot for the house measured in 2015

**renovated:** whether the house has been renovated. If yes, the value is 1. Otherwise, the value is 0.

**basemt:** whether the house has basement. If yes, the value is 1. Otherwise, the value is 0.

```

Packages

The first package that will be used in this project is ``mlr3data``, which offers the dataset that we are going to analyze. Besides, we will use the ``ggplot2`` and ``lattice`` packages for the purpose of data visualization and ``randomForest`` for random forest models that can be helpful to analyze the effects of the house factors on the house price. Also, we will utilize the ``stargazer`` package to offer neat and more readable model results of linear regressions. Another important package that can be useful in this project is ``GGally``, in which the ``ggcorr`` function can help us obtain the correlation matrix. We also use the ``dplyr`` package to manipulate and modify data frames.

And after receiving the feedback on the first draft of the project, I decide to add new packages of ``caret``, which is used to tuning random forest models so that the model performs best regarding cross-validation root mean square errors (CV RMSEs) can be found.

Also, the packages of ``knitr`` and ``kableExtra`` are also considered to help us generate a neat and well-formatted final document via R markdown.

```

```{r message=FALSE, warning=FALSE,include=FALSE}
Load the required packages
library(ggplot2)
library(lattice)
library(randomForest)
library(stargazer)
library(GGally)
library(dplyr)
library(caret)
library(knitr)
library(kableExtra)
```

```

Exploratory Data Analysis (EDA)

We start our data exploration with the variable of interest ``price``.

```

```{r EDA - Price, fig.align = "center",message=FALSE, warning=FALSE,include=TRUE,}
Check the distribution of house sale price
qplot(x = price, data = kc_housing, bins = 50,
 main = "King County House Sale Prices")
```

```

```

```{r summary,message=FALSE, warning=FALSE,include=TRUE,}
5-point summary of price
summary(kc_housing$price)
```

```

It is clear that the distribution of the home price is positively skewed with a long right tail, which implies that some houses are expected to have higher values than others.

Next, we explore the relationship between the features of the house and the home prices by plotting ``price`` with each feature. We find that some features like ``bathrooms``, ``sqft_living``, ``grade`` and ``waterfront`` have relatively stronger relationships with ``price`` than others.

```

```{r graphs, message=FALSE, warning=FALSE,include=TRUE}
Create scatterplot for price and bathrooms
ggplot(kc_housing, aes(x = bathrooms, y = price)) +
 geom_point(shape = 1) +
 geom_smooth(method = lm, color = "orange", se = FALSE) +
 ggtitle("Scatterplot of Price vs. Number of Bathrooms")
Create scatterplot for price and sqft_living
ggplot(kc_housing, aes(x = sqft_living, y = price)) +
 geom_point(shape = 1) +
 geom_smooth(method = lm, color = "orange", se = FALSE) +
 ggtitle("Scatterplot of Price vs. Square Footage of Living Area")
Create scatterplot for price and grade
ggplot(kc_housing, aes(x = grade, y = price)) +
 geom_point(shape = 1) +
 geom_smooth(method = lm, color = "orange", se = FALSE) +
 ggtitle("Scatterplot of Price vs. Grade")
Create scatterplot for price and waterfront
ggplot(kc_housing, aes(x = waterfront, y = price)) +
 geom_point(shape = 1) +
 geom_smooth(method = lm, color = "orange", se = FALSE) +
 ggtitle("Scatterplot of Price vs. Waterfront")
```

```

We also include the correlation matrix that reflects how variables are correlated with each other.

```

```{r correlation matrix, fig.height = 10, fig.width = 15, message = FALSE, warning =
FALSE,include=TRUE, fig.align = "center"}
Correlation matrix for numeric variables
ggcorr(kc_housing[, -c(1:2)], method = c("everything","pearson"))
```

```

We noticed that `sqft_living` and `sqft_above` are highly correlated with a correlation of `r cor(kc_housing\$sqft_living, kc_housing\$sqft_above)`. This makes a lot of sense because most of living area is usually above the basement. The univariate correlation between `sqft_living` and `price` (`r cor(kc_housing\$sqft_living, kc_housing\$price)`) is higher than that between `sqft_above` and `price` (`r cor(kc_housing\$sqft_above, kc_housing\$price)`). Similarly, `sqft_living` and `sqft_living15` are highly correlated with a correlation of `r cor(kc_housing\$sqft_living, kc_housing\$sqft_living15)`. The univariate correlation between `sqft_living` and `price` (`r cor(kc_housing\$sqft_living, kc_housing\$price)`) is higher than that between `sqft_living15` and `price` (`r cor(kc_housing\$sqft_living15, kc_housing\$price)`).

We will consider to remove the variables that are not likely to affect the house price later.

Data Modification

Based on above findings, we will modify our dataset by introducing two new binary features:

- renovated: Equal to 1 if the house have been renovated and 0 otherwise
- basemt: Equal to 1 if a house has basement and 0 otherwise.

We also exclude the useless information from the original set and remain the variables that are most relevant to the house price. Moreover, we would drop the variables of sqft_above and sqft_living15 in the further analysis to avoid the issue of multicollinearity. The new modified version of the dataset is named as house, and it's summary statistics are printed.

```

```{r message=FALSE, warning=FALSE, include=TRUE}
Create new variable "renovated" based on the existing variable "yr_renovated"
kc_housing$renovated = as.factor(ifelse(kc_housing$yr_renovated > 0, "1", "0"))
Create new variable "basemt" based on the existing variable "sqft_basement"
kc_housing$basemt = as.factor(ifelse(kc_housing$sqft_basement > 0, "1", "0"))
Create boxplot for price and basemt
ggplot(kc_housing,
 aes(x = basemt, y = price, fill = basemt)) +

```

```
geom_boxplot(alpha = 0.5) +
stat_summary(fun.y = mean,
 geom = "point",
 shape = 20,
 size = 4,
 color= "orange",
 fill= "orange") +
theme(legend.position = "none") +
ggtitle("Boxplot of Price between Houses with and without Basements")
```

```

However, as we noticed in the boxplot of `basemt`, this new variable does not likely to have significant influence on the home prices. So this is also a needless variable to make predictive model later.

Combined with the variable modification investigation with the previous EDA part, we then create a new dataset that contains the necessary predictors only, and summarize the final dataset as below.

```
```{r}
Create a new dataset for further analysis
house = subset(
 kc_housing,
 select = -c(
 date,
 sqft_basement,
 sqft_living15,
 sqft_above,
 sqft_lot,
 sqft_lot15,
 yr_built,
 yr_renovated,
 zipcode
)
)
house$waterfront = as.factor(house$waterfront)
house = na.omit(house)
summary(house)
```

```

Modeling

Before fitting the models, we partitioned the dataset in which 70% are used as train sets and 30% are used as test sets. And in order to make our results consistent every time we run the data, I set a seed of 913 to remove the sampling randomness.

```
```{r test-train split,message=FALSE, warning=FALSE,include=FALSE}
Test-Train split
set.seed(913)
house_idx = createDataPartition(house$price, p = 0.7, list = FALSE)
house_trn = house[house_idx,]
house_tst = house[-house_idx,]
```

```

Then we use the training set `house_trn` to tuning the models and reporting the cross-validated errors measured as RMSEs. Here, we will use a 5-fold cross validation.

Next we construct the functions that we may need during the modeling procedure.

- Create a utility function for calculating RMSEs

```
```{r RMSE function}
Create a utility function for calculating RMSE later
rmse = function(actual, predicted) {
 sqrt(mean((actual - predicted) ^ 2))
}
```

```

Random Forest Model

First, we trained random forest models using all the predictors in the `house` dataset with `price` as response variable. The default tuning parameters chosen by the `caret` package would be used.

```
```{r train rf, eval=FALSE, include=FALSE}
Tune a random forest model without normalizing predictors
set.seed(913)
rf_unscale_mod = train(
 price ~ ., data = house_trn,
 trControl = trainControl(method = "cv", number = 5),
 method = "rf"
)
```
```

The best random forest model with scaling uses 7 for `mtry`, meaning that 7 features will be randomly chosen every time a tree is grown. The number of trees is 500 by default. We find that about 88% of the variation in prices can be explained by this model.

```
```{r message=FALSE, warning=FALSE }
Fit the best random forest model
rf_mod = randomForest(price ~ .,
 data = house_trn,
 mtry = 7,
 importance = T)
importance = importance(rf_mod)
VarImportance = data.frame(variables = row.names(importance),
 importance = round(importance[, '%IncMSE'], 2))
```

```
Rank variables by importance
rank = VarImportance %>% mutate(rank = paste0('#', dense_rank(desc(importance))))
```

```
ggplot(rank, aes(
 x = reorder(variables, importance),
 y = importance,
 fill = importance
)) +
 geom_bar(stat = 'identity') +
 geom_text(
 aes(x = variables, y = 0.5, label = rank),
 hjust = 0,
 vjust = 0.6,
 size = 4,
 color = 'white'
) +
 labs(x = 'Variables', y = 'Importance') +
 ggtitle("Variable Importance from Selected Random Forest Model") +
 coord_flip()
```
```

It turns out that the `grade` (ranging from 1 to 13), which represents the construction quality, is the most important variable regarding house sale price. For houses only meet the minimum building standards, their grades are low and ranged from 1 to 3, and their prices are expected to be the lowest on average. For houses that have achieved average performance in terms of construction and design, they are graded as 7. And their averaged prices are expected to be moderate among all houses. As for the houses graded over 12, they are thought to have excellent designs and use the best materials while construction. As a result, their prices are expected to be highest too.

`sqft_living`, the spacing of the living rooms, turns out to be the next most important feature that is related to home values. And the further next significant factor that is highly correlated with the home prices is the location, consistent with the `latitude` and `longitude` ranked as third and fourth important variables.

Additive Linear Regression Model

Although, we learned the effects of location on house price are likely to be significant, we are not going to interpret them in the multiple linear regression model because location is not a factor that could be changed for home owners after the house being purchased. Hence, in this linear model, only the variables of `grade`, `sqft_living`, `view`, `waterfront` and `bathrooms` will be included.

```
```{r message=TRUE, warning=TRUE,include=TRUE}
Fit the linear regression model
lm.mod = lm(log(price) ~ grade + sqft_living + view + bathrooms +
 waterfront, data = house_trn)
lm.mod = step(lm.mod, trace = FALSE)
Print the model results
stargazer(lm.mod, type = "text")
```
```

The estimated model is

```
$$
\hat{\log\{price\}} = 0.249grade + 0.0001sqft_living + 0.038view + 0.091bathrooms + 0.183waterfront +
10.976
$$
```

The R-squared of this model is 0.728, which means this linear model containing the most relevant predictors could explain about 72.8% of the total variation in the house prices. And F statistic of this model is 0.314 ($p = 0.000$), which means the overall model is statistically significant. Also, at 5% level of significance, we notice all predictors are statistically significant individually.

As we notice, here we used the log-transformation on the dependent variable because the variable is right skewed based on previous analysis. Hence, the model results imply that one additional point in the house grade is expected to increase the house price by 24.9% on average when other factors are assumed to be the same. And the house with a waterfront is expected to be 18.3% higher in price than the house without a waterfront when other conditions are the same. And an additional bathroom in a house is expected to bring up the home value by 9.1%, and with a one-sqft increase in the living room is expected to enhanced the house value by 0.01% on average conditional to all other factors respectively. As for the time of the house being viewed, it is shown that each time the house is viewed, the house value is expected to be boost by 3.8% on average.

Model comparisons

Finally we compare the predictive performance of the two models.

```
```{r}
rf.prediction = predict(rf_mod, house_tst)
rf.rmse = rmse(rf.prediction, house_tst$price)
rf.rmse
lm.prediction = exp(predict(lm.mod, house_tst))
lm.rmse = rmse(lm.prediction, house_tst$price)
lm.rmse
```

```{r message=FALSE, warning=FALSE,include=TRUE}
library(testthat)

test_that("matrix",{
 expect_false(
 isTRUE(
 all.equal(
 lm.rmse,
 rf.rmse
)
)
)
})
```
```

The reported test rmses for random forest model is 322215.6 and for the linear regression model is 351018.4. That means the best random forest model we found by `caret` performs better in predicting than the linear model. So we decide this as a better model that can help us predict the house values.

Limitations

The first limitation of this project comes from the fact that some unchangable external factors such as the location of the house and the real estate market such as unexpected economic crisis could impact the house price unwantedly, even if the priority list is completed by the owner.

Besides, it is quite time consuming to tuning the parameters in the random forest model, which limited my attempts of a wider range of possible values in the parameters should be tuned. That means, the random forest model may still be improved if we can take more time and tuning more models.

Thirdly, we only considered two methods in predicting the house price so there may be better models using different methods.

Lastly, we should notice that the omitted variable biases may still exist, and there are other potential factors that are not included in the data may also affect the home values significantly.

Conclusion

In conclusion, from a view of statistician, a priority list for boosting home value should contain the following two things:

1. Try to maintain the house in a good condition and add more custom design so that it can be graded higher.
2. Try to expand the space of living room in the house

And if we want to use the easily collected house features to predict the house prices, a random forest model with 7 features being selected per tree and considers 500 trees at one time is recommended. Although it may not be precise, it can offer a general idea of how valuable the homes would be for the house owners.