

Housing Market Analysis Project

Christian Ollen

Section 1

Question 1

Homeownership is one way millions of Americans build wealth. Homes are the foundations for millions of families, but what makes one home more valuable than another? When one first purchases a home, the first things a person considers are the size of the home, the number of bedrooms, and the number of bathrooms. But are these the most essential internal housing factors that drive up the price of a home? Our team conducted a study to determine what internal housing factors matter most within King County, Washington.

Our team analyzed the following internal housing factors: bedrooms, bathrooms, floors, condition, grade, square feet of living room, square feet of living room above ground floor, and square feet of the basement, to see what had the most significant influence on housing prices. Initially, we were able to narrow the focus of our study to just four predictor variables:

- The number of bathrooms
- The square footage of the home
- The square footage above the ground floor of a house
- The construction grade and design of a home

We further reduced our model by choosing a subset of our square footage variables to only include the overall square footage, leaving our model with three predictor variables. The number of bathrooms negatively correlated with the price, which may seem counterintuitive. From our personal experiences of buying and selling homes, more bathrooms are typically desirable, so we further reduced our model to two predictor variables.

Our study's findings reveal that the two most influential variables on housing prices in King County, Washington are the square footage of the home and the construction grade of the home. These findings are crucial for understanding the dynamics of the local housing market and can greatly assist in making informed decisions related to real estate.

Question 2

As a homeowner, gaining insight into the key factors that determine a house's value can provide you with a competitive edge. Understanding these aspects can help you make informed decisions, potentially leading to a higher sale price, improved house condition, or a faster sale.

For this next question, we are using logistic regression to estimate the odds of a house being sold above the median price. For the King County, Washington data set of homes sold is estimated to be \$450,000. We aim to find out which predictors of waterfront, condition, grade, view, year built, or year renovated contribute the most towards a high-precision model that accurately predicts when a house is worth more than the median price.

The main objective of answering this question is to assess whether a house's categorical variables also impact determining the value of a home and observe how specific characteristics and increases in quality can go a long way to improving a house's status.

Given the nature of many of the present categorical variables, our team created new indicators to determine whether a house with an above-median view, condition, or grade rating would also be sold above the median

price. We immediately observed that four out of the six variables were crucial to making good predictions, and these were

- The condition level of the house
- The view rating of the house
- The year the house was built
- The year the house was last renovated

After more careful consideration, we tested a four-variable model against a six-variable model, but the data did not support dropping the grade and waterfront variables altogether. Nevertheless, we observed how the grade level of the house was providing highly questionable results. After removing it from our model, we were left with five variables to predict results.

Our findings demonstrate the robustness of our analysis. The five-variable model we developed boasts a 73% precision, significantly outperforming random guessing. This underscores the reliability of our insights and their potential value for homeowners.

Section 2

The dataset used for this analysis is obtained from Kaggle. It contains house sale prices for King County, which includes Seattle. It includes homes sold between May 2014 and May 2015. Out of all 21 variables, we have decided to use 13 of them, and create 6 new variables based on some existing ones.

1. **price**: Sale price of the house.
2. **bedrooms**: Number of bedrooms in the house.
3. **bathrooms**: Number of bathrooms in the house.
4. **floors**: Number of floors (levels) in the house.
5. **sqft_living**: Square footage of the interior living space of the house. The bigger the living space or a house, the higher we may expect the price to be.
6. **sqft_above**: Square footage of the interior living space above ground level.
7. **waterfront**: The waterfront variable is an indicator variable, and is determining whether the house in question is located near the waterfront or not. A house near the waterfront is labeled as 1, and one which is not near the waterfront is labeled as 0. We are inclined to believe this could be a relevant predictor because houses near water are expected to be more expensive.
8. **view**: The view variable rates on a scale of 0 to 4 the quality of the view the property was. The median value for this variable is 0, which may indicate that only about half of the apartments had a somewhat favorable view.
9. **view_med**: An indicator variable pointing out whether a house's view is above the median (indicated in the previous variable's description as being 0). If the view rating is above 0, it is labeled 1; otherwise, it is labeled 0.
10. **condition**: Condition is a categorical variable, indexed from 1 to 5, which represents the overall condition of the house. Nicer apartments may be sold at a higher price than apartments which have lower ratings.
11. **condition_med**: An indicator variable of condition, pointing out whether a house's condition is above the median of 3. It is labeled 1 if the house's condition is 4 or 5, and 0 otherwise.
12. **grade**: Grade is highlighting the quality level of the building's construction and design. It is indexed from 1 to 13, and houses with ratings 11-13 are considered to be of the highest quality.
13. **grade_med**: An indicator variable which is labeled 1 if a house's grade rating is of the highest quality, that is, which has a rating of 11, 12, or 13. It is labeled 0 if its grade rating is 10 or lower.

14. **yr_built**: Year the house was built. Houses in King's County were built as early as the year 1900, and the latest were built in 2015.
15. **decade_built**: To make a more general approach towards visualizing the relationship of a house's price and the year it was built, we group houses by decade of being built. The earliest decade is 1900, while the most recent decade is of the 2010s.
16. **yr_renovated**: Year the house was last renovated. We acknowledge that houses that were renovated recently could be priced highly in comparison with those who were never renovated at all. Houses that have not been renovated are marked with the value 0.
17. **decade_renovated**: Similar to decade_built, this variable groups houses by the decades they were last renovated. This helps us make better visualizations, and to find any relationships with the results obtained from visualizations in decade_built.
18. **above_median**: A simple indicator variable, labeled 1 if a house is above median price, and 0 otherwise. Note that the median price of a house, as obtained from the training data set, was found to be \$450,000.

Section 3

Question 1 : What internal housing factors influence the price of homes in King County from May 2014 to May 2015?

- *Response Variable: price*
- Motivation: This question aims to understand the relationship between various internal features of a house (such as bedrooms, bathrooms, square footage, floors, condition, grade, and others) and its sale price. Investigating these factors can provide insights into the determinants of housing prices in the King County area during the specified time.

Question 2: Does the presence of one or more of Waterfront, Condition, View, Grade, Yr Built, and Yr Renovated cause the house to be sold above the median price?

- *Response Variable: above_median*
- Motivation: This question seeks to find whether any of the variables mentioned above impact determining the qualities that could make a house worth more than the median housing prices in the King County region.

We are using variables that identify categorical aspects of a house. Therefore, any significant results may lead homeowners to focus not only on the size of their property but also on its design and internal properties to increase their potential earnings.

Section 4

Distribution of Price



The distribution is highly right-skewed, with a large frequency of lower-priced items and very few high-priced ones. The majority of prices of the houses fall close to the lower end of the price spectrum, which means that higher prices are outliers in this particular dataset.

Distribution of Bedrooms

```
##    bedrooms    n
## 1          0    3
## 2          1  103
## 3          2 1387
## 4          3 4907
## 5          4 3437
## 6          5   806
## 7          6   133
## 8          7    22
## 9          8     5
## 10         9     2
## 11        33     1
```

The distribution of bedrooms in the dataset shows that three-bedroom houses are the most common, followed by those with four bedrooms. There is only one house with 33 bedrooms, along with those with 11 and 10 bedrooms, appear to be outliers.

(To maintain data integrity, we will exclude these outliers from further analysis.)

Distribution of Bathrooms

```
##      bathrooms      n
## 1          0.00      4
## 2          0.50      3
## 3          0.75    44
## 4          1.00  1925
## 5          1.25      5
## 6          1.50   680
## 7          1.75 1546
## 8          2.00   958
## 9          2.25 1036
## 10         2.50 2646
## 11         2.75   609
## 12         3.00   382
## 13         3.25   302
## 14         3.50   384
## 15         3.75     81
## 16         4.00     63
## 17         4.25     38
## 18         4.50     51
## 19         4.75     10
## 20         5.00     13
## 21         5.25      6
## 22         5.50      6
## 23         5.75      3
## 24         6.00      3
## 25         6.25      2
## 26         6.50      1
## 27         6.75      1
## 28         7.50      1
## 29         8.00      2
```

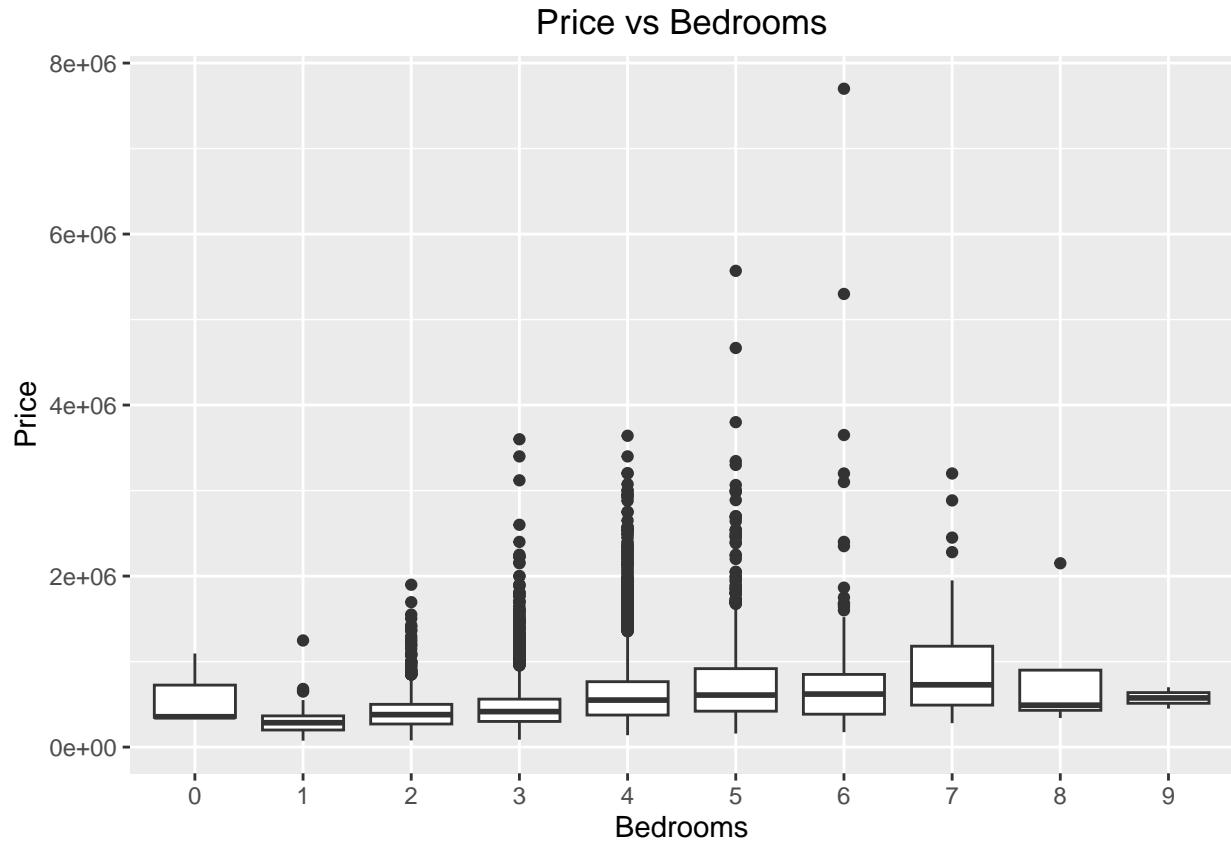
The distribution of bathrooms in the dataset shows that three-bathroom houses are the most common, followed by those with 3.5 bathrooms. The highest number of bathroom in any one house is 8, with only one house having these many bathrooms.

Price vs Square Feet Living



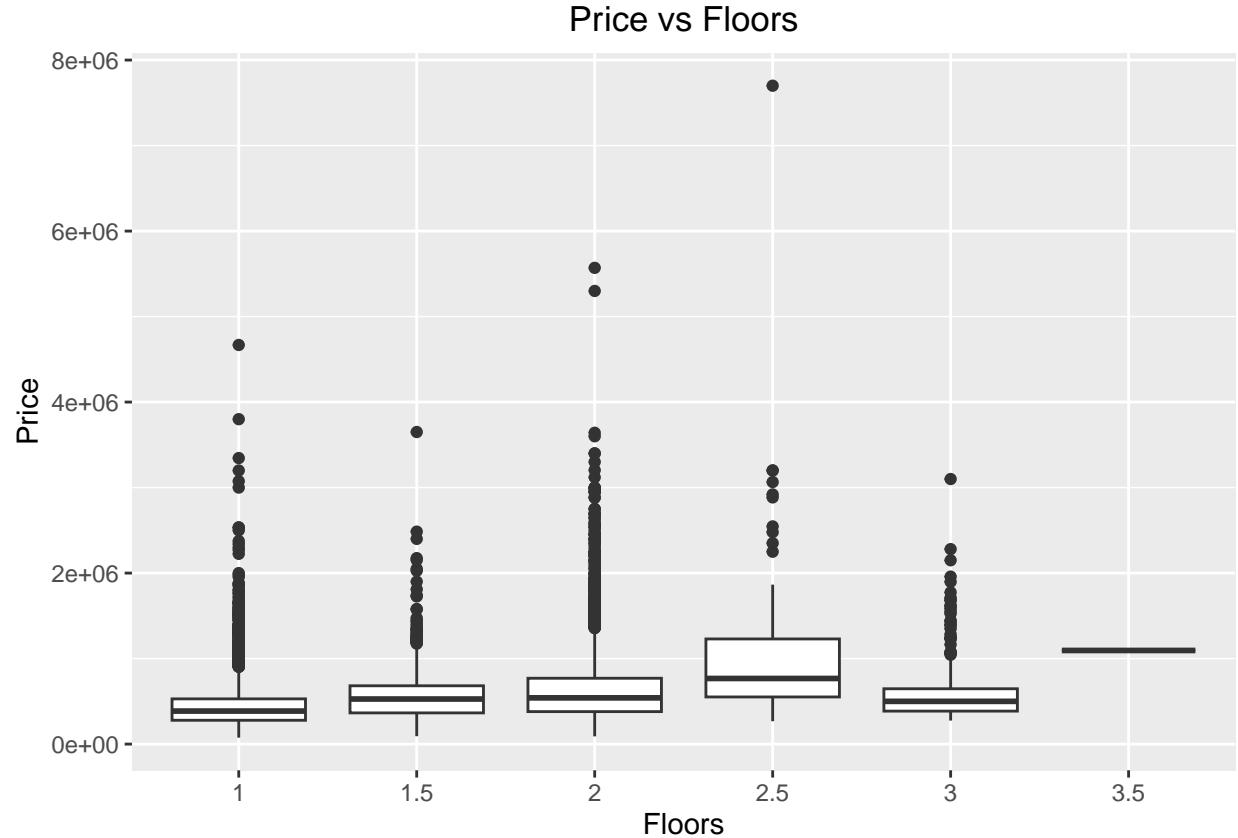
The scatter plot shows the correlation between square feet of living space and the cost of the apartments. The fact that we have an upward trend, revealing that the increase in price is proportional to the rise in the area, while it's a bit different from linear, is another proof of this cost-effectiveness. On the flip-side, there are some data points with big living area that is drastically more expensive than the rest of the listing thereby, showing these could be luxury or premium properties. These values are plotted far from the core of data and form a series of clusters that are anchored at specific points.

Price vs Bedrooms



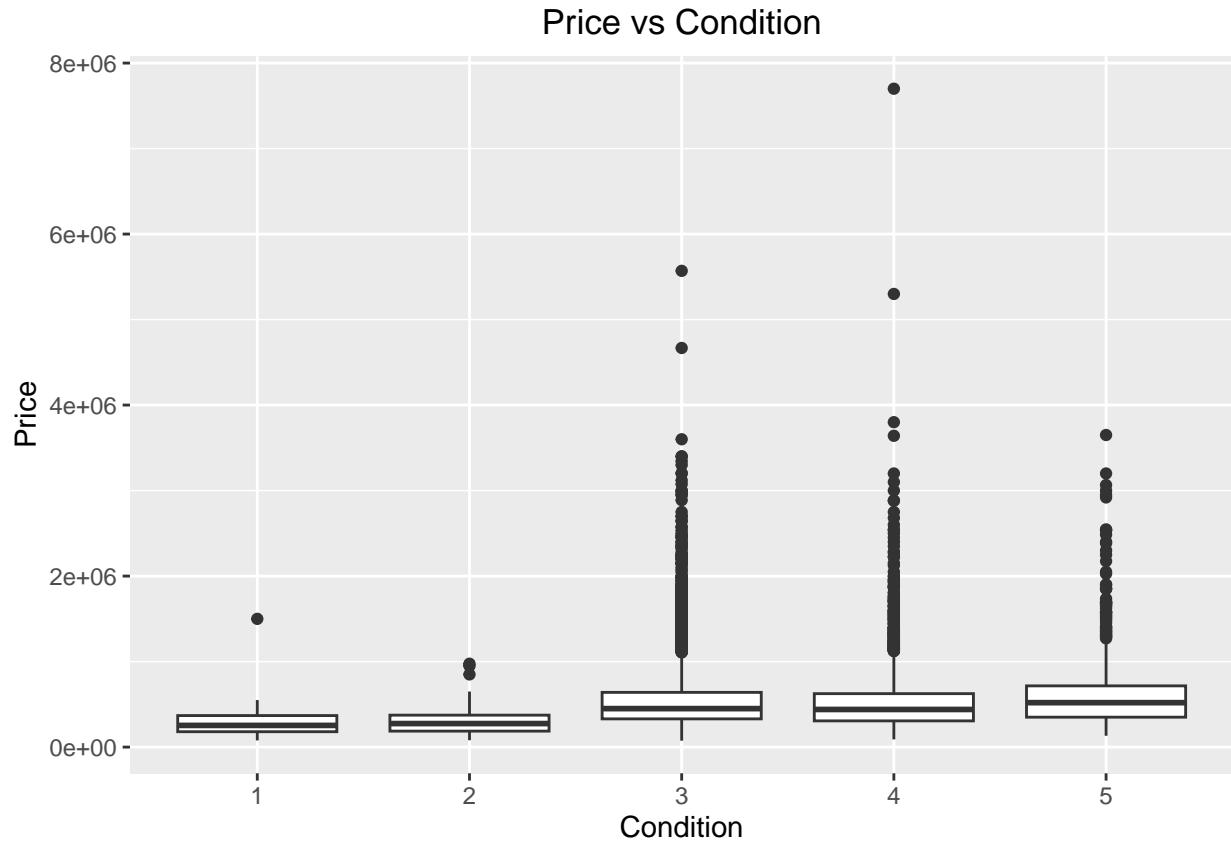
The boxplot represent the property prices distribution in line with the quantity of the bedrooms. The price continuously rises with the quality up to a number of bedrooms, after which it exactly fluctuates. The scatter plot of average prices sorted by category of bedrooms is such that the spread of each category's prices increase as the number of bedrooms increases, as shown by the longer boxes and whiskers, pointing to a bigger spread of property values. Particularly there are a great big number of outliers, especially in the upper area of rooms per property, which perhaps means that there are some properties, if many bedrooms have them, that are crushed against the mean of category, if it's in rooms.

Price vs Floors



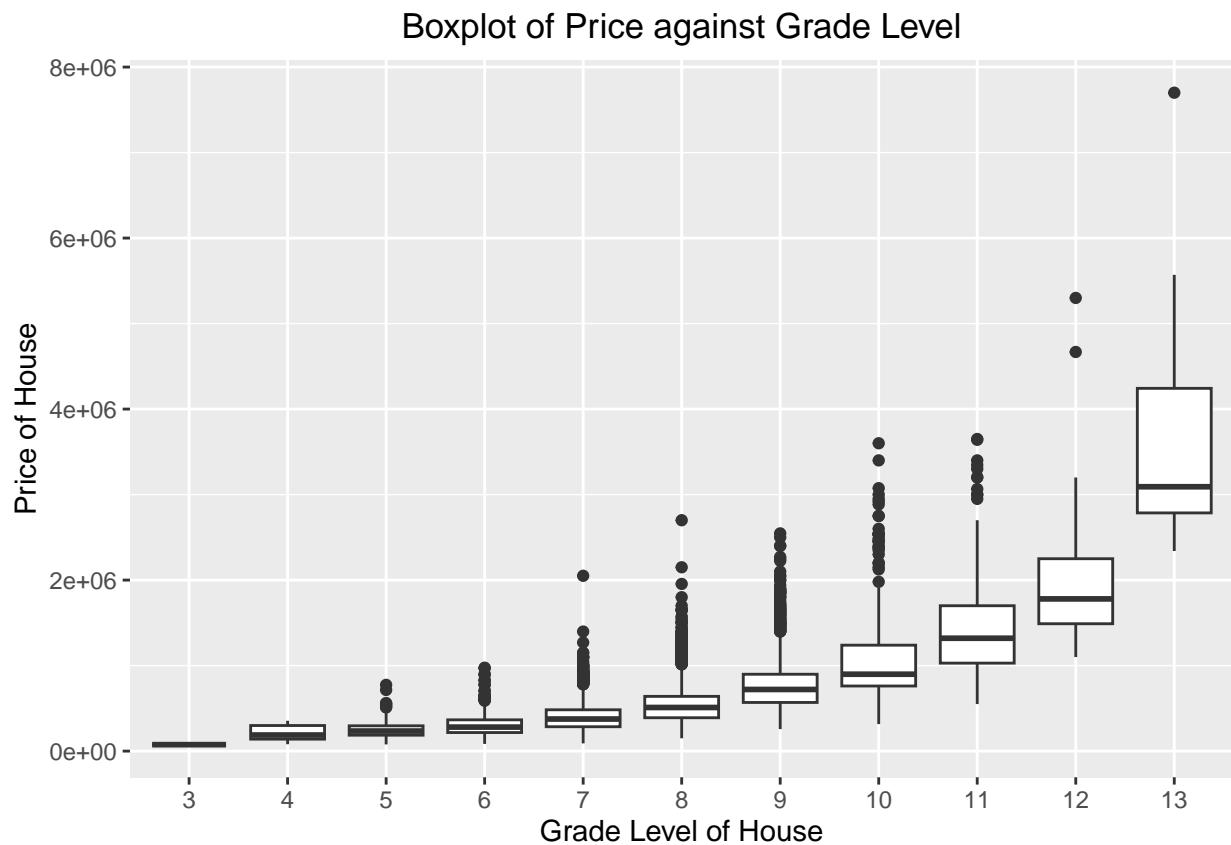
There are a significant number of community properties with 1 floor and the total price for them normally falls within a wide range. However, the median price for these types of properties is the cheapest among all the categories. Two-story boxes on the balancing beam point to the significantly larger median price, and a large part of the tale is the numerous outliers-obviously expensive homes. Residential units with one-and-half, two-and-half, and three floors have very small representation of the dataset and their price spectral is a bit uneven, but in particular, there are a few outliers among the 3-floor home residences. The scarce type of 3.5 floor homes have the highest median prices. The spread range here is narrow which implies that these house types have a pricing consistency. Outliers are present in all categories of floors, suggesting that rare or special geo-spatial attributes may be contributing to higher pricing.

Price vs Condition

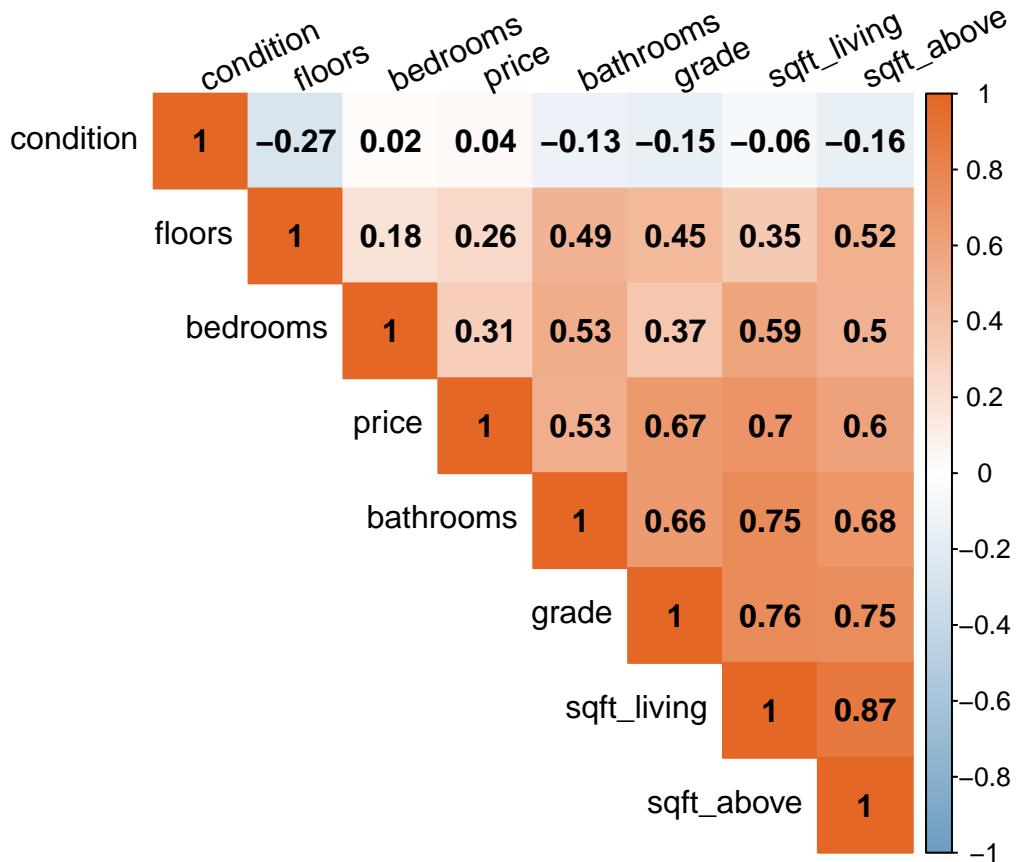


Properties in conditions rated as 1 and 2 have a narrower range of prices with lower medians, suggesting they are generally less expensive. The median prices seem to increase slightly for properties rated in condition 3, and there's a wider spread of prices, with many outliers indicating some high-priced properties. Condition 4 and 5 properties also show a higher median price compared to lower-rated conditions, with condition 5 showing the most significant spread in prices, though not necessarily the highest median price. This could imply that while good condition may contribute to a higher price, other factors like location or size might also play a significant role in determining a property's value.

Price vs Grade



It is clear that a positive linear relationship exists between grade level and price of a house. It is important to note here that most of the houses have a grade level in the range of 6 and 10, which is what could be considered an “average” rating by the Kaggle website. Houses with ratings considered to be of the highest quality are the ones rated 11, 12 and 13 on the scale. There are less than 500 houses in the data set which are given this rating, but the difference in price shows significantly.



The Correlation Matrix, a data map created from your information, shows us the relationships between various objects. The strength of the relationship between two items is indicated by each number on this map. A value around one indicates that those variables typically move in together, i.e., as one increases, the other does too. It appears as though they move in different directions when it is close to -1. We begin to see patterns as we examine the map, such as the significant correlation between cost and size—larger houses typically cost more. However, there are also subtle connections that we can miss initially. This map allows us to identify instances where two items are overly similar, which aids in selecting the most relevant data for our investigation. Comparable to a treasure map that leads us through our data, indicating which avenues to explore and which to keep clear of, enabling us to make informed decisions and reveal what's contained inside.

Section 5

For this model we are looking at how internal housing factors effect the price of the home. In this model there are multiple predictors influencing the response variable price, therefore multiple linear regression is appropriate. We start by fitting a model with as the predictor variables, and `price` as the response.

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-7.894e+05	2.533e+04	-31.172	<2e-16 ***
bedrooms	-4.143e+04	3.059e+03	-13.543	<2e-16 ***
bathrooms	-1.253e+04	5.086e+03	-2.464	0.0138 *
sqft_living	2.495e+02	6.742e+00	37.016	<2e-16 ***
floors	5.984e+03	5.763e+03	1.038	0.2992
condition	6.311e+04	3.751e+03	16.825	<2e-16 ***
grade	1.130e+05	3.270e+03	34.564	<2e-16 ***
sqft_above	-6.420e+01	6.773e+00	-9.479	<2e-16 ***

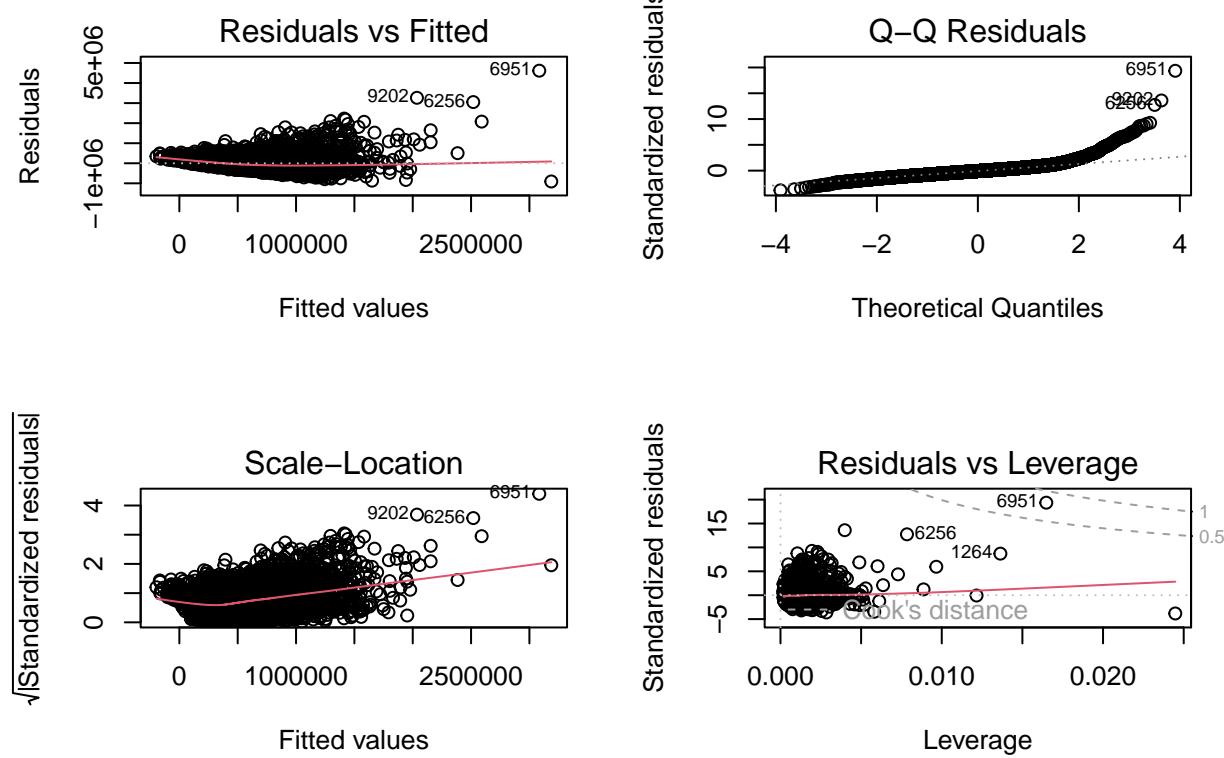
Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

For the full model, we have the following regression equation

$$\hat{y} = -789400 - 41430x_1 - 12530x_2 + 249.5x_3 + 5984x_4 + 63110x_5 + 113000x_6 - 64.2x_7$$

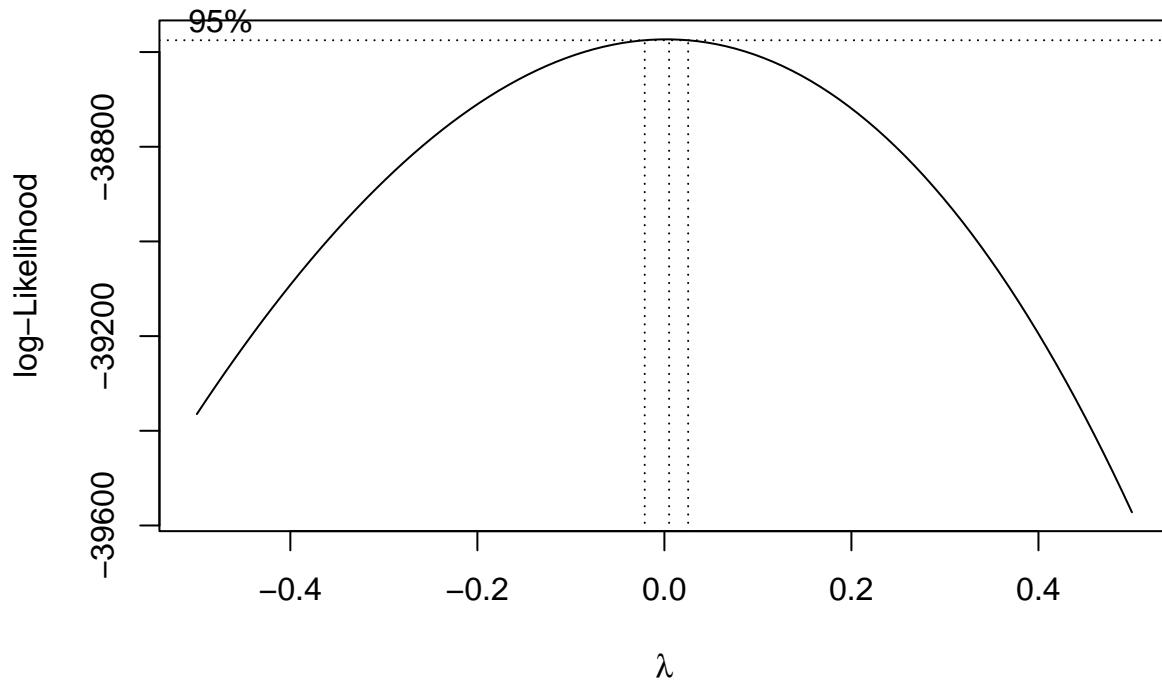
Data Transformation and Initial Model Fitting

Now that we have determined multiple linear regression model, we have to check whether our model's regression assumptions are met.

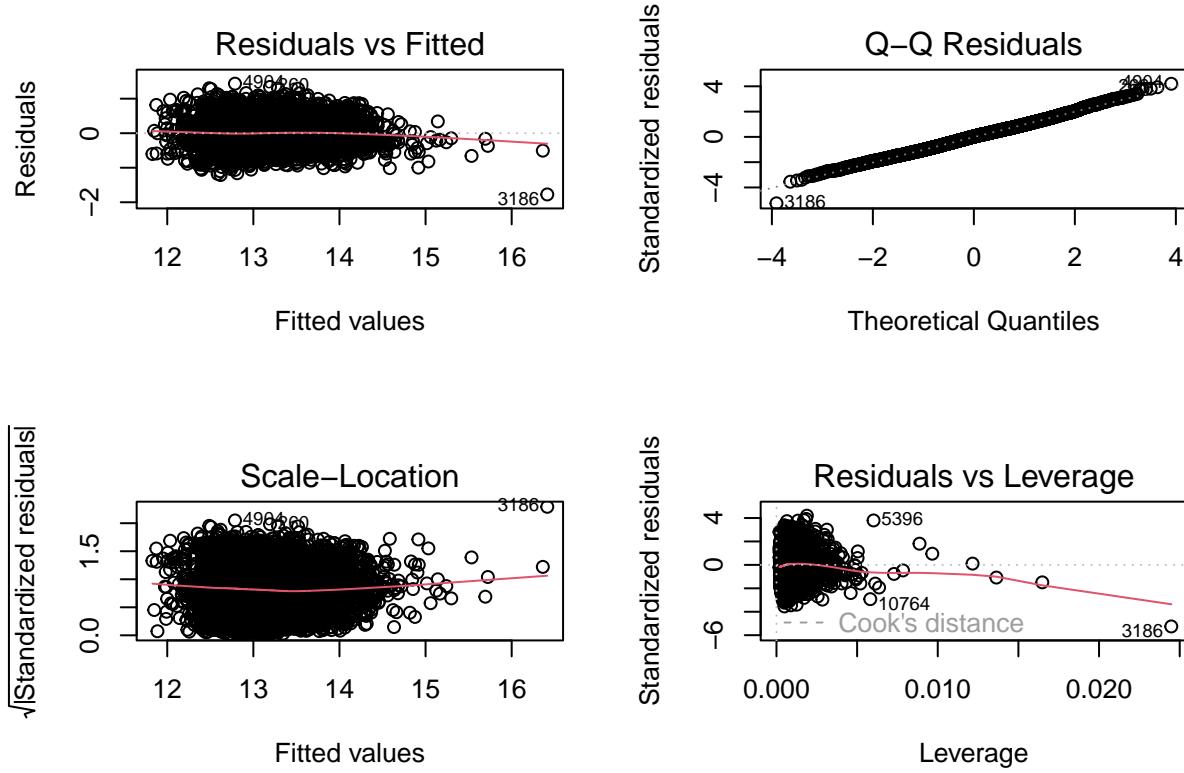


We first observe that the residuals are not evenly scattered across the horizontal axis, and also the variance has not been stabilized yet, with the variance increasing as the data points move from left to right. Thus, assumptions 1 and 2 have not been met yet. Similarly, the Q-Q residual plot points do not appear to follow the diagonal line, so assumption 4 may not be met yet.

We start by addressing assumption 2 first and stabilizing the variance. To do this, we will start by transforming the predictor variable based on the results of the Box Cox plot.



Given that zero is above the 95% confidence line on the Box Cox plot, a log transformation on the response variable should be performed, so $y^* = \log(y)$. We add this new parameter to our data set and assess the regression assumptions once again



From the above results, we arrive at the following observations:

- **Residuals vs Fitted**: The plot looks acceptable, with no clear patterns or trends. The points are evenly distributed around zero, indicating that model's first two assumptions are generally met. There's an outlier labeled "12778" that might be worth investigating.
- **Q-Q (Quantile-Quantile) Plot**: The residuals align well with the diagonal, indicating that they are approximately normally distributed. This suggests that the normality assumption holds, although the point labeled "12778" again stands out slightly, hinting at a possible outlier.
- **Scale-Location (or Spread-Location)**: The plot is relatively horizontal, suggesting consistent variance across fitted values, although there is a slight trend visible. The same outlier is visible here as well.
- **Residuals vs Leverage**: The plot shows that most points have low leverage. These might be influential points, given their positions near the Cook's distance threshold.

Overall, we can claim that all regression assumptions have been met for our model, and that they provide statistical significance based on their conclusions from the regression coefficients. After performing the transformation $y^* = \log(y)$, the results of the regression coefficients are

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.073e+01	3.596e-02	298.467	< 2e-16	***
bedrooms	-2.810e-02	4.344e-03	-6.469	1.03e-10	***
bathrooms	-1.276e-02	7.222e-03	-1.767	0.0773	.
sqft_living	3.083e-04	9.574e-06	32.201	< 2e-16	***
floors	6.787e-02	8.184e-03	8.293	< 2e-16	***
condition	1.023e-01	5.327e-03	19.204	< 2e-16	***
grade	2.048e-01	4.643e-03	44.102	< 2e-16	***

```

sqft_above -1.263e-04 9.618e-06 -13.134 < 2e-16 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

And our regression equation becomes

$$y^* = 10.73 - 0.0281x_1 - 0.01276x_2 + 0.00031x_3 + 0.06787x_4 + 0.102x_5 + 0.2048x_6 - 0.00013x_7$$

Before moving forward, let's assess whether our current model is more useful than just random guessing. We carry out an ANOVA F test, where $H_0: \beta_j = 0$, for j is an integer between 1 and 7; and let H_a : at least one of the $\beta_j \neq 0$.

```

## Analysis of Variance Table
##
## Model 1: ystar ~ 1
## Model 2: ystar ~ (bedrooms + bathrooms + sqft_living + floors + condition +
##      grade + sqft_above + price) - price
##   Res.Df   RSS Df Sum of Sq    F    Pr(>F)
## 1 10804 2989.9
## 2 10797 1261.3 7     1728.6 2114 < 2.2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

We have a very small p-value. We reject the null hypothesis, the data support that our model is doing a better job than a model with no predictors in it.

Next, we can go back and take a look at the correlations between all of the variables. We note that predictors `bathrooms`, `sqft_living`, `grade`, and `sqft_above` have moderate to high correlation with the response variable `price`. At the same time, predictors `bedrooms`, `floors`, and `condition` do not have a strong correlation to the response variable `price`.

Additionally, `sqft_living` and `sqft_above` seem to be highly correlated. Therefore, We will consider dropping `sqft_above` in the reduced model, since we would much rather choose the variable which describes the property area more generally.

We can now consider using `bathrooms`, `sqft_living` and `grade` in the reduced model.

```

##
## Call:
## lm(formula = ystar ~ bathrooms + sqft_living + grade, data = train_mlr)
##
## Residuals:
##       Min        1Q        Median        3Q        Max
## -1.644400 -0.25022  0.00162  0.23430  1.30215
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.154e+01  1.820e-02 634.187 <2e-16 ***
## bathrooms   -1.327e-02  6.785e-03 -1.955  0.0506 .
## sqft_living  2.170e-04  6.540e-06 33.177 <2e-16 ***
## grade        1.909e-01  4.484e-03 42.573 <2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3518 on 10801 degrees of freedom
## Multiple R-squared:  0.5529, Adjusted R-squared:  0.5528
## F-statistic: 4453 on 3 and 10801 DF, p-value: < 2.2e-16

```

If we check now the VIF function for the reduced model, we obtain the following

```
##   bathrooms  sqft_living      grade
##   2.428541    3.194566    2.437031
```

Whereas in the previous correlation plot we saw signs of collinearity between the predictor variables, now we observe that the VIF's for the reduced model are all below 4, suggesting there is not a huge issue with multicollinearity.

The values from the summary() function seem to indicate that, for the reduced model, the p-value for the ANOVA F test is very small. However, the `bathrooms` predictor is now just shy of being a significant variable in the model. We can drop the `bathrooms` predictor from this model since it does not pass a t-test, and now we are left with a two-predictor model.

```
##
## Call:
## lm(formula = ystar ~ sqft_living + grade, data = train_mlr)
##
## Residuals:
##       Min     1Q   Median     3Q    Max
## -1.63789 -0.25207  0.00164  0.23409  1.29588
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.154e+01  1.810e-02  637.70 <2e-16 ***
## sqft_living 2.103e-04  5.587e-06   37.65 <2e-16 ***
## grade        1.891e-01  4.386e-03   43.11 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3518 on 10802 degrees of freedom
## Multiple R-squared:  0.5528, Adjusted R-squared:  0.5527
## F-statistic:  6675 on 2 and 10802 DF,  p-value: < 2.2e-16
```

All that is left to do now is to carry a hypothesis test to assess whether we can drop five predictor variables from the full model and go with the reduced model instead. Let the null hypothesis be $H_0: \beta_1 = \beta_2 = \beta_4 = \beta_5 = \beta_7 = 0$; and let the alternative hypothesis be $H_a:$ at least one of $\beta_1, \beta_2, \beta_4, \beta_5, \beta_7 \neq 0$. We carry an ANOVA General F Test to determine the validity of the initial claim.

```
## Analysis of Variance Table
##
## Model 1: ystar ~ sqft_living + grade
## Model 2: ystar ~ (bedrooms + bathrooms + sqft_living + floors + condition +
##                   grade + sqft_above + price) - price
##             Res.Df   RSS Df Sum of Sq    F    Pr(>F)
## 1      10802 1337.2
## 2      10797 1261.3  5     75.922 129.99 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The F statistic is very large, with a small p-value. So we can reject the null. Our data suggests we can drop `bathrooms`, `bedrooms`, `floors`, `condition`, and `sqft_above` from the full model, and go with the reduced model instead. Our regression equation is given by

$$\hat{y} = 11.16 + 0.00021x_2 + 0.189x_6$$

We now interpret the results from the above regression equation. For every one percent increase in square feet

living, the predicted increase in the price variable will be 2.0895×10^{-6} , when holding grade level constant. Similarly, a one percent increase in grade level will result in a 0.0019 increase in price, when holding square footage constant.

Some facts to highlight about the model's results are the following

- Price Distribution: As seen above with the visualizations, the price distribution is right-skewed, meaning that the many houses are distributed around lower price range, while the rest of expensive houses are located in the tail of the distribution. Premium homes, as shown within this data set, are outliers.
- Bedroom Distribution: 3-bedroom houses are the most common ones, while the number of houses is decreasing as the number of bedrooms is increasing in another one. Beyond enormous houses with more 8 bedrooms are not in usual cases and this is considered as an anomaly.
- Living Space and Price Correlation: A good correlation can be drawn between the living area size and the house price, however, such a relationship doesn't perfectly fit any specific pattern before taking a log transformation.
- Floor Preferences: Of the two predominant styles of homes, either single-story or two-story are typically the most prevalent. The fewer-floor homes of 3.5 stories are almost as rare as hen's teeth, and this proves that many citizens like living in skyscrapers.
- Condition and Price: Houses in better conditions are prone to satisfy the median price, however this rule is not always the case throughout the condition hierarchy.

Model Selection

The final model's diagnostic plots and summary statistics are generated to evaluate its performance on the training data. These diagnostics help identify potential issues before moving to the testing phase. Among other criteria, we start by performing all three kinds of automated model search on our model and an intercept-only model.

```
Step:  AIC=-23180.71
ystar ~ grade + sqft_living + condition + sqft_above + floors +
bedrooms + bathrooms
```

We note that for forward selection, backward selection, and stepwise regression, all of them returned the same model with six parameters where they all recommend to drop only the `floors` predictor variable. As we move on to explain why this happens, we first call the regression subsets function from the leaps package. The results it presented are shown below

```
Selection Algorithm: exhaustive
      sqft_living  grade  condition  bedrooms  sqft_above  bathrooms
1 ( 1 )   "*"      " "       " "       " "       " "
2 ( 1 )   "*"      "*"      " "       " "       " "
3 ( 1 )   "*"
4 ( 1 )   "*"
5 ( 1 )   "*"
6 ( 1 )   "*"
7 ( 1 )   "*"
               floors
1 ( 1 )   "
2 ( 1 )   "
3 ( 1 )   "
4 ( 1 )   "
5 ( 1 )   "
6 ( 1 )   "
7 ( 1 )   "*"
```

According to the above function, the two preferred predictor variables chosen from the `regsubsets()` function were `sqft_living` and `grade`, both of which are contained in our final reduced model. These two variables seem to be the two strongest predictors out of all variables involved in this model. Given that we were able to pass an ANOVA F test using a model with three predictors against one will all predictors, this still indicates that we found the two factors which carry on the greatest contributions to the model. Similarly, these are the two which have the highest correlation to the `price` variable and which always maintained great statistical significance across all models considered.

Influential Observations, High Leverages and Outliers

This code identifies and removes influential observations that can skew the model's results. High leverage points and outliers are identified based on the "hat" values, which measure how far each observation's predictor values are from the average predictor values. Removing such influential points helps in achieving a more reliable model.

```

##      bedrooms      bathrooms      sqft_living      floors      condition
##  Min.   :1.00   Min.   :0.000   Min.   : 390   Min.   :1.000   Min.   :1.000
##  1st Qu.:4.00   1st Qu.:1.500   1st Qu.:1410   1st Qu.:1.000   1st Qu.:3.000
##  Median :4.00   Median :2.250   Median :1870   Median :1.000   Median :3.000
##  Mean   :4.33   Mean   :2.057   Mean   :1971   Mean   :1.943   Mean   :3.421
##  3rd Qu.:5.00   3rd Qu.:2.500   3rd Qu.:2440   3rd Qu.:3.000   3rd Qu.:4.000
##  Max.   :9.00   Max.   :5.250   Max.   :4090   Max.   :5.000   Max.   :5.000
##      grade      sqft_above      price      ystar
##  Min.   :4.000   Min.   : 390   Min.   : 84000   Min.   :11.34
##  1st Qu.:5.000   1st Qu.:1180   1st Qu.: 318500  1st Qu.:12.67
##  Median :5.000   Median :1530   Median : 438000  Median :12.99
##  Mean   :5.567   Mean   :1701   Mean   : 499258  Mean   :13.01
##  3rd Qu.:6.000   3rd Qu.:2100   3rd Qu.: 605000  3rd Qu.:13.31
##  Max.   :8.000   Max.   :4090   Max.   :2950000  Max.   :14.90

##
## Call:
## lm(formula = ystar ~ sqft_living + grade, data = train_no_influence)
##
## Residuals:
##      Min      1Q      Median      3Q      Max
## -1.22813 -0.25295  0.00211  0.23495  1.22756
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.151e+01  2.234e-02  515.33  <2e-16 ***
## sqft_living 2.220e-04  7.216e-06   30.77  <2e-16 ***
## grade       1.892e-01  5.478e-03   34.54  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3488 on 9911 degrees of freedom
## Multiple R-squared:  0.4566, Adjusted R-squared:  0.4564
## F-statistic: 4163 on 2 and 9911 DF,  p-value: < 2.2e-16

```

This has the following regression equation

$$\hat{y} = 11.51 + 0.00022x_2 + 0.189x_6$$

After removing the influential observations, the model is refitted to the refined dataset. This step ensures

that the model is not biased by any problematic data points.

The new regression model is very similar to the one we already had for the two-predictor model. The high leverage observations, influential observations, and the outliers did not have a significant effect on our final model.

We also run a function to find any values outside Cook's distance

```
## named numeric(0)
```

Our model contains no points found outside of a Cook's distance of 1.

Evaluate Model Performance on Test Data

The final step involves evaluating the model's performance on the testing data. The root mean squared error (RMSE) and R-squared values are calculated to gauge the model's predictive accuracy and explanatory power. These metrics provide a clear picture of how well the model generalizes to new data.

RMSE on Test Data

The Root Mean Squared Error (RMSE) calculated on the test data is printed on this line. The average discrepancy between the expected and actual values is measured by RMSE. A model with lower RMSE values performs better in terms of prediction.

R-squared

The value of R-squared, calculated using the test data, is printed on this line. R-squared shows how much of the variance in the dependent variable log_price can be accounted for by the model's independent variables (predictors). Greater R-squared values signify an improved model-data fit.

```
## RMSE on test data: NaN  
## R-squared: 0.5527659
```

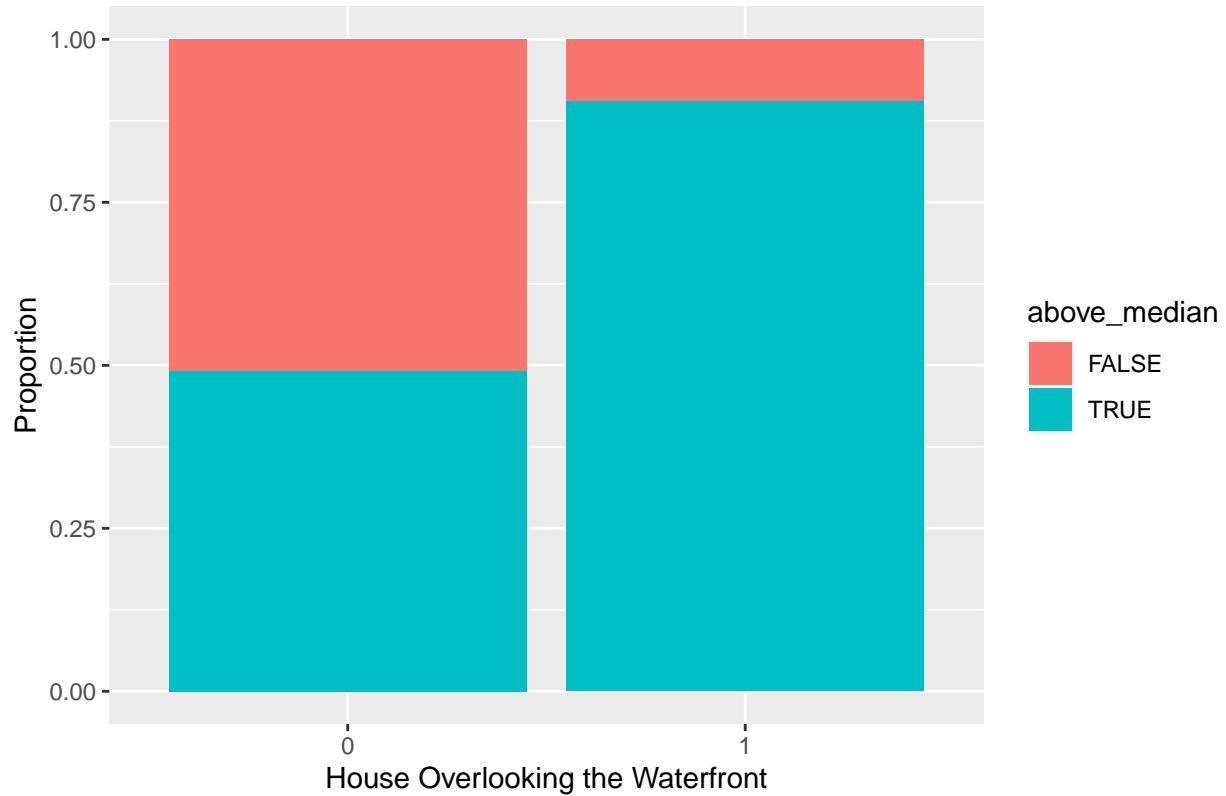
Section 6

Description of Variables

Waterfront

The waterfront variable is determining whether the house in question is located near the waterfront or not. We are inclined to believe this could be a relevant predictor because houses near water are expected to be more expensive. Consider the following visualization

Proportion of Waterfront Houses above the Median House Price



There is a considerable proportion of houses overlooking the waterfront which are also above the median house price. To have a better idea of how many such houses we have in all, we also show a two-way table with the counts of houses belonging to each category

```
## # A tibble: 2 x 3
##   above_median    0     1
##   <dbl> <int> <int>
## 1 FALSE  5455   8
## 2 TRUE  5266  76
```

Out of almost 11,000 entries of houses, it appears only 84 of them are overlooking the waterfront after all, which may lead to this variable not being too influential in the final model.

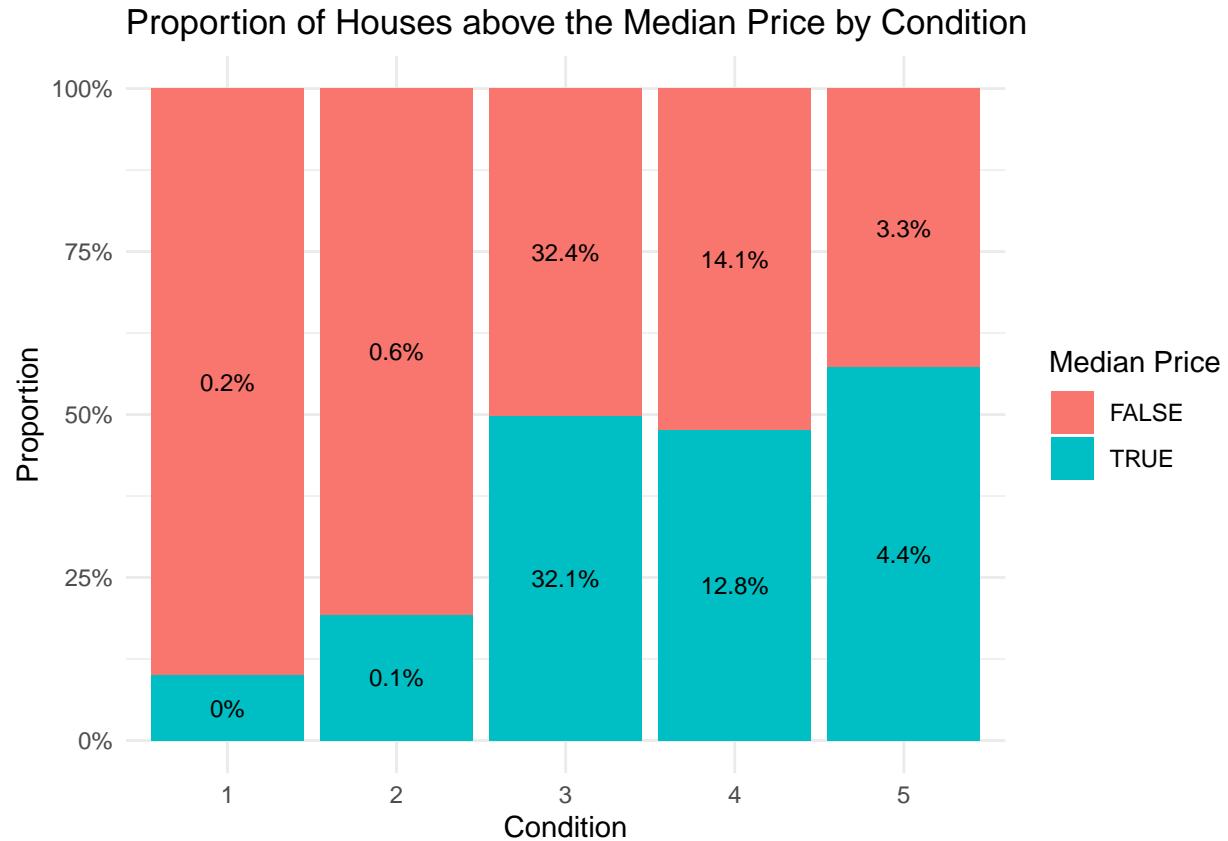
Condition

Condition is a self-explanatory categorical variable, in which it indicates the condition of an apartment in a scale of 1 to 5. Nicer apartments may be sold at a higher price than apartments which have lower ratings.

We wish to have an idea of the total number of houses for each condition rating. The following tables present a detailed summary of the number of entries for each condition level and their proportion

```
## # A tibble: 5 x 3
##   condition count percentage
##   <dbl> <int>     <dbl>
## 1 1       20      0.185
## 2 2       78      0.722
## 3 3      6970     64.5
## 4 4      2909     26.9
## 5 5      828      7.66
```

We see that very few houses have lower ratings, and most houses are in the median, 3, or above. These numbers express a strong relationship between the condition and a house being priced above the median.

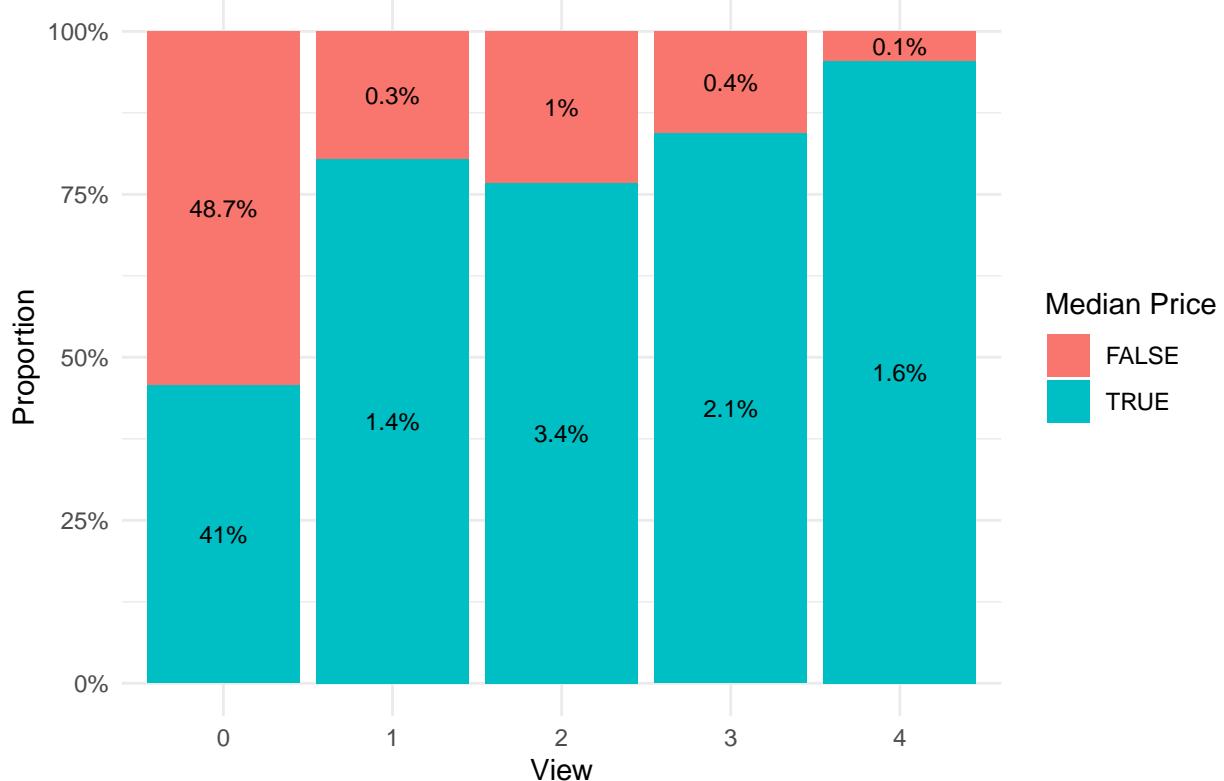


Indeed, we see a connection between the condition of a house and the house being above median price for the higher ratings. The proportion of houses which exceed the median increase as the rating increases.

View

The view variable rates on a scale of 0 to 4 how good the view of the property was. The median value for this variable is 0, which may indicate that only half of all apartments had a somewhat favorable view. We see a distribution of the proportions for each rating below

Proportion of Houses above the Median Price by View



The proportion of houses above median value appears to increase as the view gets improved, and it still shows a high percentage when having a 0 rating after all. We present the counts from the table as follows

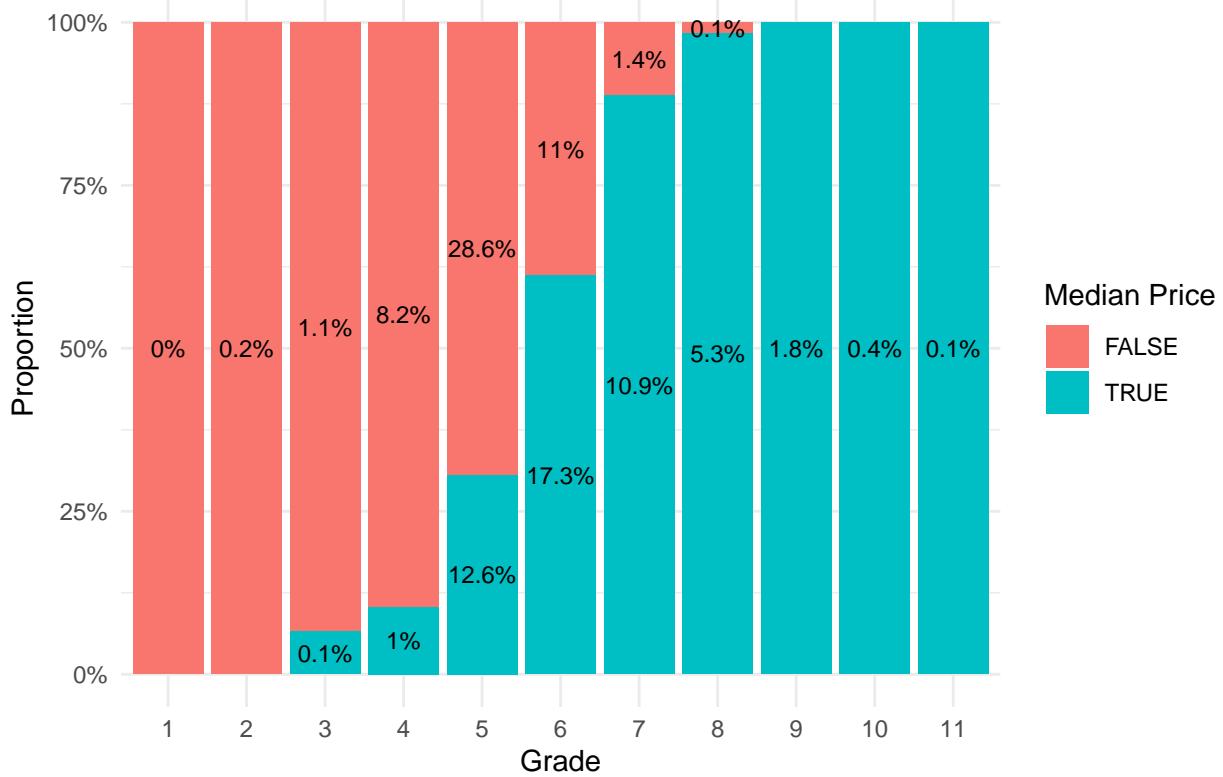
```
## # A tibble: 5 x 3
##   view count percentage
##   <int> <int>     <dbl>
## 1     0   9699     89.8
## 2     1    184     1.70
## 3     2    478     4.42
## 4     3    268     2.48
## 5     4    176     1.63
```

As mentioned before, a vast majority of houses have a view rating of 0, as we have almost 10,000 of the almost 11,000 entries having this rating. However, out of all houses with a higher view rating, over 75% of them were priced above the median across all four levels.

Grade

For the last categorical variable, grade is highlighting the quality level of the building's construction and design. It is indexed from 1 to 13, and houses with ratings 11-13 are considered to be of the highest quality. We see the proportion of houses being above or below the median price by grade.

Proportion of Houses above the Median Price by Grade



We observe that there are no houses with grade levels 1 or 2. And that beginning from grade level 8, all subsequent levels have over 50% of the houses being above the median price. More importantly, this bar chart appears to highlight a positive relationship between the grade and being above the median price for houses.

Again, we present a two-way table to highlight the total amounts of houses over each grade level.

```
## # A tibble: 11 x 3
##   grade count percentage
##   <dbl> <int>     <dbl>
## 1     1     1 0.00925
## 2     2     21 0.194
## 3     3    122 1.13
## 4     4    994 9.20
## 5     5   4447 41.2
## 6     6   3062 28.3
## 7     7   1326 12.3
## 8     8    587 5.43
## 9     9    193 1.79
## 10    10    44 0.407
## 11    11     8 0.0740
```

We see that the vast majority of houses (it's over 9,000) have been given ratings between 6 and 9, which means that they have an average level of construction and design. Very few of them were given the highest quality level ratings.

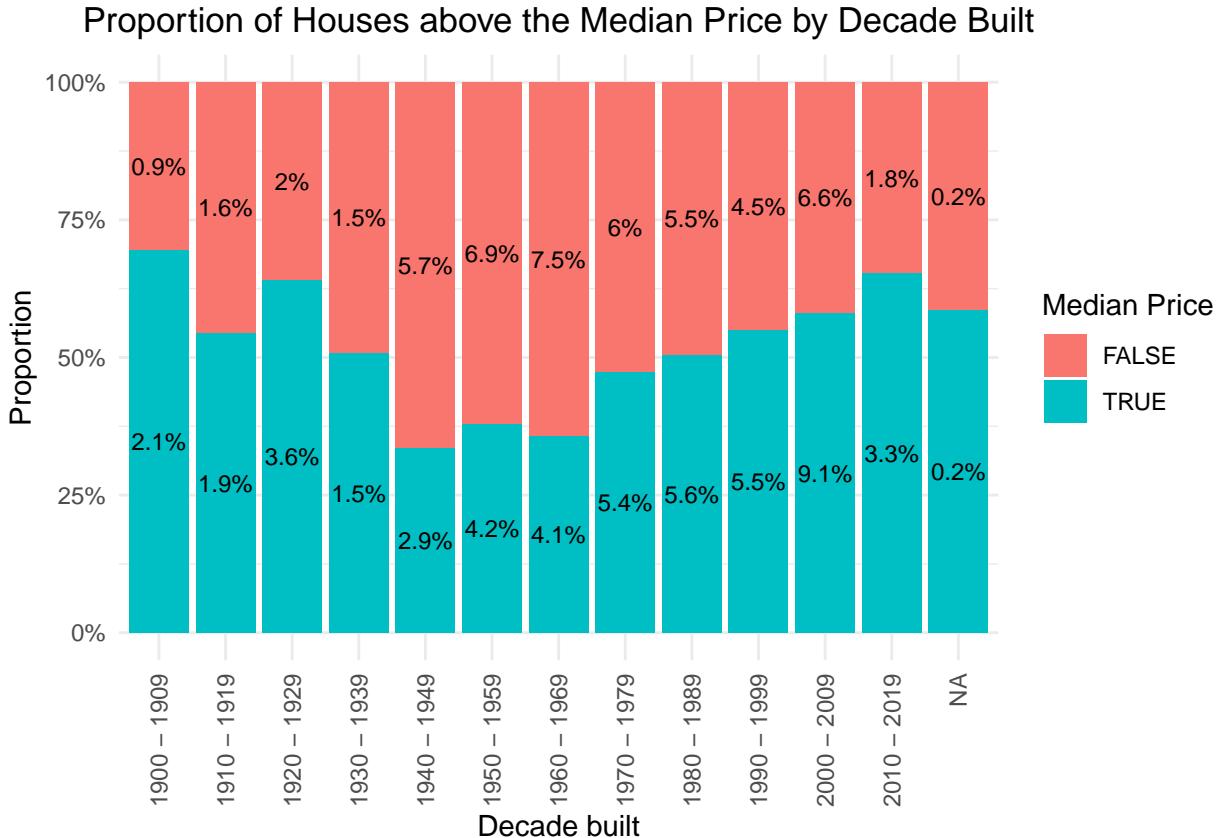
Year Built

We now present the first of two numerical predictor variables, and that is the year the house was first built. Given that this data set includes houses built from the early 1900s, we will split the houses into sets of

decades, starting with the 1900s, then the 1910s, and so on.

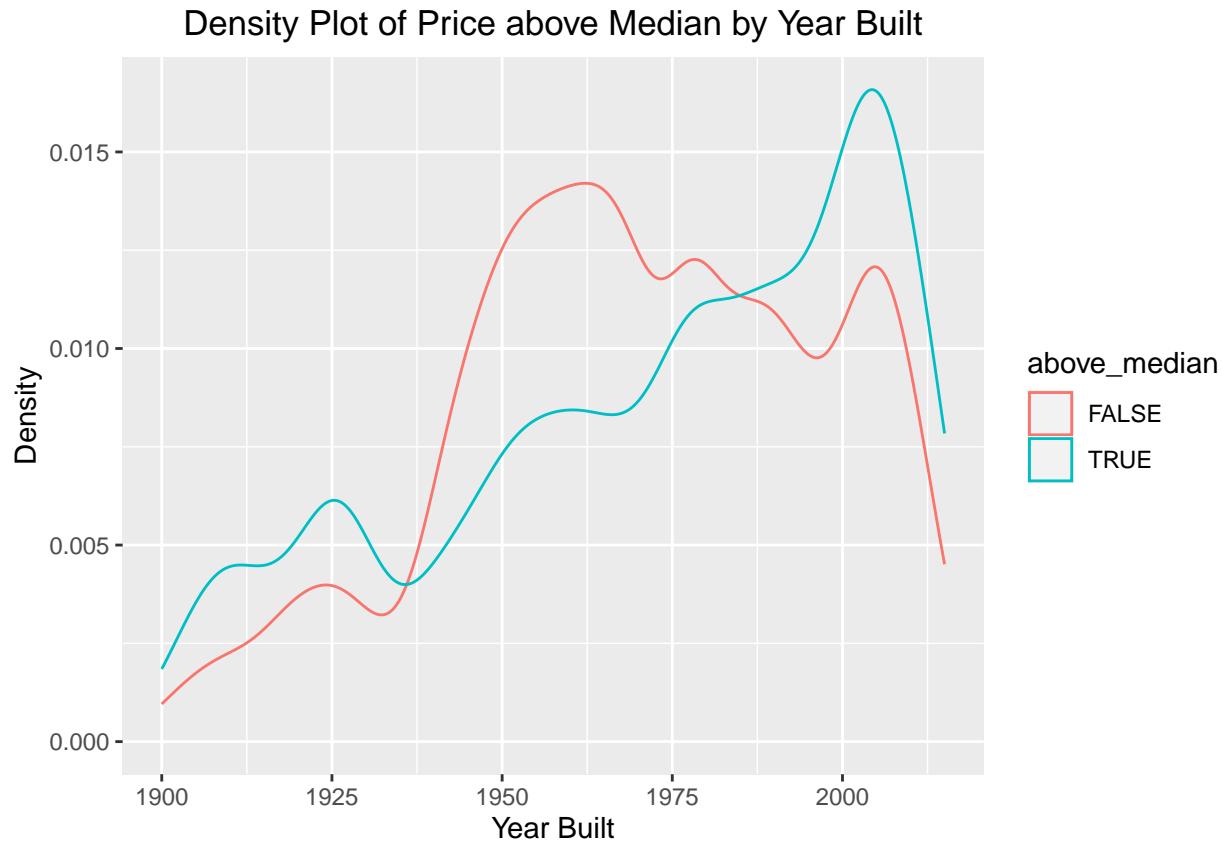
```
## # A tibble: 13 x 3
##   decade_built count percentage
##   <fct>        <int>     <dbl>
## 1 1900 - 1909    330     3.05
## 2 1910 - 1919    368     3.41
## 3 1920 - 1929    607     5.62
## 4 1930 - 1939    323     2.99
## 5 1940 - 1949    924     8.55
## 6 1950 - 1959   1201    11.1
## 7 1960 - 1969   1253    11.6
## 8 1970 - 1979   1239    11.5
## 9 1980 - 1989   1193    11.0
## 10 1990 - 1999  1088    10.1
## 11 2000 - 2009  1686    15.6
## 12 2010 - 2019   552     5.11
## 13 <NA>           41     0.379
```

From the above table, we can observe that a majority of the houses were built between the decades of 1950 and 2000, but otherwise, the decades the houses were built are very evenly spread, consisting of a fair proportion from every decade.



We can see a good proportion of house prices remain steadily above median until we enter the 1950s decade, where we experience a sudden drop on the total houses above the median from that era. Beginning from the 1950s, more and more house prices have been slowly making their way up above the median, but it is still unclear whether a linear relationship exists at all.

Finally, we will consider a density plot that considers the houses being above median price by year built.



From the above graph, we see that most of the houses which are below median price were built between the 1940s and the 1970s. For houses which are above the median, they are found mostly around the 1920s decade, as well as the present day, with a majority found from 1980 to 2010.

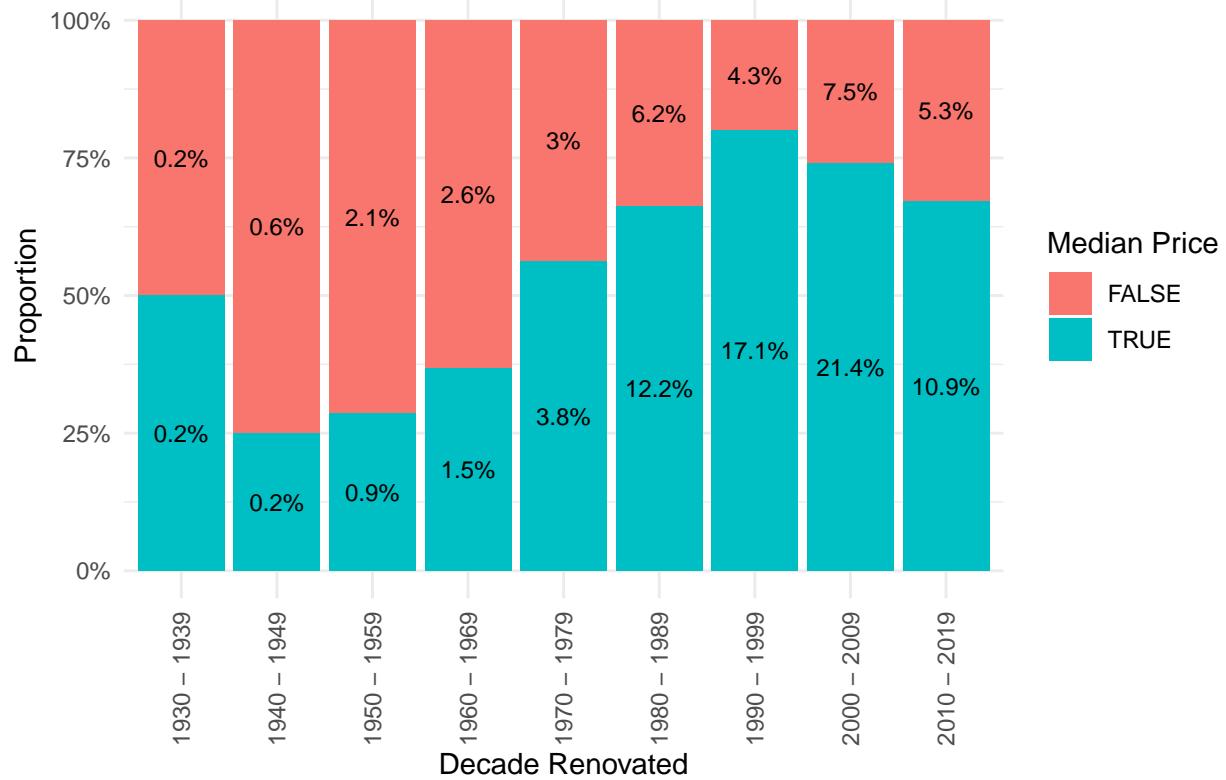
Year Renovated

```
## # A tibble: 10 x 3
##   decade_renovated count percentage
##   <fct>           <int>     <dbl>
## 1 1930 - 1939      2     0.0185
## 2 1940 - 1949      4     0.0370
## 3 1950 - 1959     14     0.130
## 4 1960 - 1969     19     0.176
## 5 1970 - 1979     32     0.296
## 6 1980 - 1989     86     0.796
## 7 1990 - 1999    100     0.925
## 8 2000 - 2009    135     1.25
## 9 2010 - 2019     76     0.703
## 10 <NA>          10337    95.7
```

The table indicates that a staggering 95% of all houses in this data set have not been renovated yet, with less than 500 out of 10,806 houses having undergone this process. We are yet to determine whether this will prove to be a strong enough indicator during our analysis.

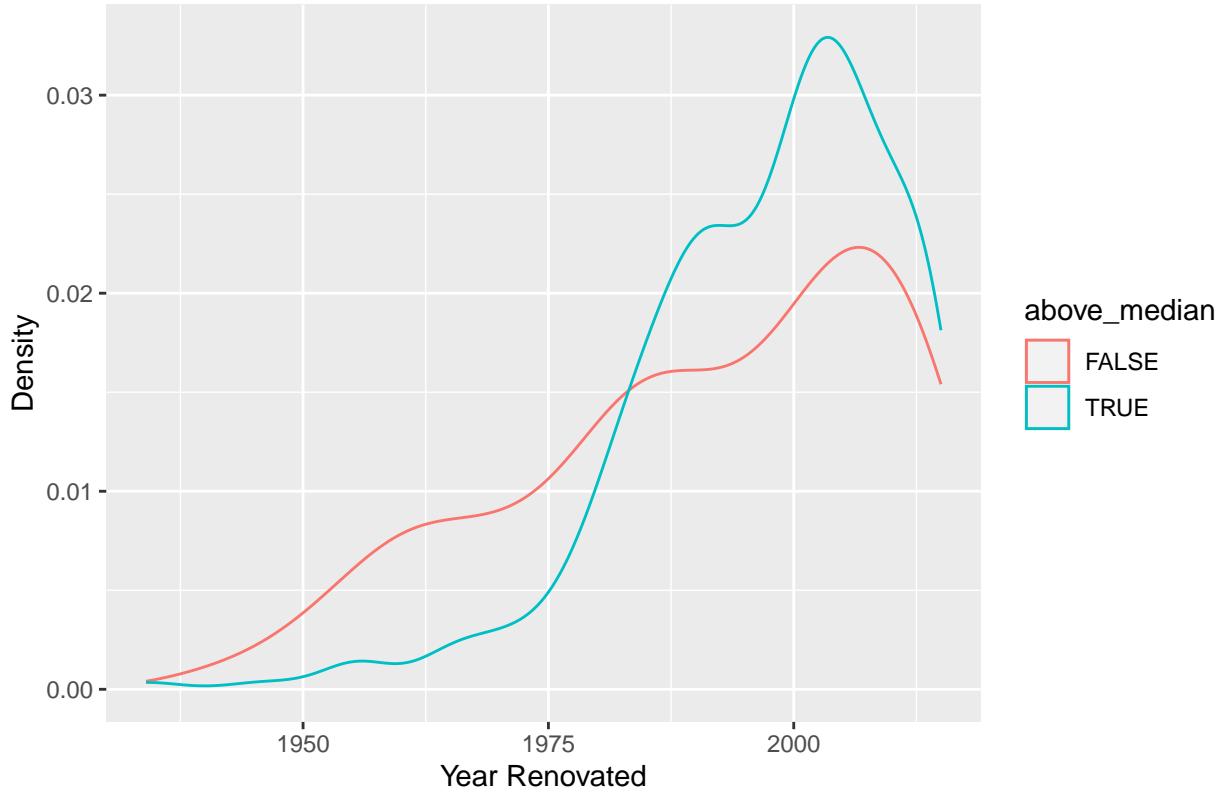
Also, out of the homes that were actually renovated, a good chunk of them were renovated from 1980 onwards.

Proportion of Houses above the Median Price by Decade Renovated



After removing the values which are zero (indicating the house has never been renovated), we notice that, unlike the bar chart for year built, this bar chart provides a better look at how a house being renovated can positively impact its chances of being priced above median.

Density Plot of Price above Median by Year Renovated



As mentioned earlier, we see that the majority of houses above the median price were renovated after the 1980s. Houses which are below the median follow a steady growth along the decades, with some of them peaking after the 2000s decade.

Section 7

Before we can start to work in the regression formula, we will define the response variable as a binary indicator. Specifically, our response's name will be `above_median`, and it will determine whether a house's price is above the median or not. The houses that are above the median will be labeled as $\hat{y} = 1$, and those which are not will be labeled $\hat{y} = 0$.

We want to assess our response variable with both numerical and categorical predictor variables. The numerical variables for this model are year built and year renovated, while the categorical variables we will use include waterfront, condition, view, and grade. Waterfront is already a binary indicator variable, but the latter three of these are labeled using indices ranging from 1 to 5, 0 to 4, and 1 to 13, respectively.

Given the nature of the categorical variables, there is no quantifiable way to measure the change from, say, a 9 to a 10 for the grade variable. Thus, instead of using the numbers from their given range, we will transform all these variables into indicator variables, where we determine whether the values of the data set are above or below the set value.

First, we find the median value of the condition variable. We know that it is 3, so variables who have a condition of 4 or 5 are going to be labeled as 1, while those which have a condition level of 1, 2, or 3 will be labeled as 0.

We will also compare the view variable, and it has a given range of 0 to 4. However, the median of this variable is 0, and as we saw before, over 89% of all houses were given a view rating of 0. So again the value of the indicator variable will be 1 if the view is greater than 0, and 0 otherwise.

Finally, we will transform the grade variable. According to the Kaggle website where we draw our data, the indices 11-13 suggest a high quality level for the design of a house. Taking this as the standard level, we reserve the value of 1 to those houses with grade level in this range, with 0 if they fail to attain it.

Our converted model will therefore contain four indicator variables, as well as two numerical variables, in order to determine whether a house is above or below the median price.

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.137e+01	1.536e+00	-7.400	1.36e-13 ***
yr_built	5.619e-03	7.752e-04	7.248	4.22e-13 ***
yr_renovated	4.331e-04	5.530e-05	7.831	4.86e-15 ***
waterfront1	5.763e-01	3.865e-01	1.491	0.135894
condition_medTRUE	1.670e-01	4.624e-02	3.612	0.000304 ***
view_medTRUE	1.564e+00	8.470e-02	18.461	< 2e-16 ***
grade_medTRUE	1.606e+01	1.394e+02	0.115	0.908288
<hr/>				
Signif. codes:	0 ‘***’	0.001 ‘**’	0.01 ‘*’	0.05 ‘.’
	0.1 ‘ ’	1		

Regression Equation is as follows

$$\log \left(\frac{\hat{\pi}}{1 - \hat{\pi}} \right) = -11.36 + 0.005616x_1 + 0.000433x_2 + 0.5762I_1 + 0.1676I_2 + 1.563I_3 + 16.06I_4$$

We can interpret this logistic equation formula by considering each predictor variable's regression coefficient

(i) For every single year increase, the estimated probability that the house is above the median price is multiplied by $e^{0.005616} = 1.005632$, when controlling for all the other variables. This means that the newer the house is, the more likely it is to be above the median price.

(ii) For each additional year a house was last renovated, the estimated probability that the house will be above median price is multiplied by $e^{0.000433} = 1.000433$, when controlling for all the other variables. The more recent a house was renovated, the more likely it will be to be worth more.

(iii) The estimated probability that a house overlooking the waterfront is above the median price is $e^{0.5762} = 1.779264$ times the probability for houses that are not near the waterfront, when controlling for all the other variables. There may exist a link after all between houses near the waterfront and them having a higher price.

- We can now assess the coefficient for Waterfront using Wald test :

- H_0 is $\beta_1=0$. H_a is $\beta_1 \neq 0$
- Test statistic is $Z = \frac{\hat{\beta}_1 - 0}{se(\hat{\beta}_1)} = \frac{0.2296}{0.4139} = 0.5547$
- P value is 0.5791 which is more than test statistic and also much more than 0.05. This suggests we could drop the Waterfront predictor from our model, in presence of condition, grade, view, year built and year renovated.

(iv) The estimated probability that a house with above-median condition is above the median price is $e^{0.1676} = 1.182464$ times the probability for houses that are below the median in condition level, when controlling for all the other variables. The inside condition of a house determines a role in its price, as expected.

(v) The estimated probability that a house with above-median view is above the median price is $e^{1.563} = 4.773119$ times the probability for houses that are below the median in view level, when controlling for all the other variables. The view index of a house determines a role in its price.

(vi) Finally, according to the regression model, the estimated probability that a house with great grade level is above the median price is $e^{16.06} = 9435597$ times more likely than houses which are below the median, when controlling for all the other variables. This number appears to be very questionable.

We also notice that, since the standard error is very high, the predictor is not reliable and it might be recommended to drop the predictor in presence of all the others.

Model Assessment using Likelihood Ratio test

We now compare our model to determine whether it is useful in predicting whether a house's price is above the median better than random sampling. Let the null hypothesis be that $\beta_j = 0$ for $j = 1, 2, 3, 4, 5, 6$, and the alternative hypothesis be that at least one of the $\beta_j \neq 0$.

Consider a model with no predictors. Then we will compute the test statistic by finding the difference of the deviances between our full model and the no-predictor model.

$$\Delta G^2 = D(R) - D(F) = 14978.96 - 14054.81 = 924.1545$$

The p value from this test statistic is exactly 0. So we reject the null. The data supports our model over an intercept-only model.

Next, we want to check the individual predictor variables, and whether we can remove some of them if they are not contributing enough to the full model. We will start by considering whether we can remove the waterfront and the grade variables.

Let the null hypothesis be that $\beta_3 = \beta_6 = 0$; that is, that both predictors are not contributing any additional information in presence of the other variables. The alternative hypothesis is that at least one of β_3 or β_6 are not equal to 0. By carrying out the likelihood ratio test, we find that

$$\Delta G^2 = D(R) - D(F) = 14300.03 - 14054.81 = 245.2246$$

This has also a corresponding p-value of 0. So we reject the null. This means we cannot drop both waterfront and grade from our model. However, we can try to drop only the grade predictor, since it is the one that poses the more questionable results in the regression formula. We do this by using the Wald test.

Let $H_0: \beta_6 = 0; H_a: \beta_6 \neq 0$.

$$Z = \frac{\hat{\beta}_6 - 0}{se(\hat{\beta}_6)} = \frac{16.06}{139.4} = 0.115$$

The corresponding p-value is 0.9084451. This means that we fail to reject the null hypothesis, and the data supports removing the grade predictor from the model. Our new logistic regression equation will be given by

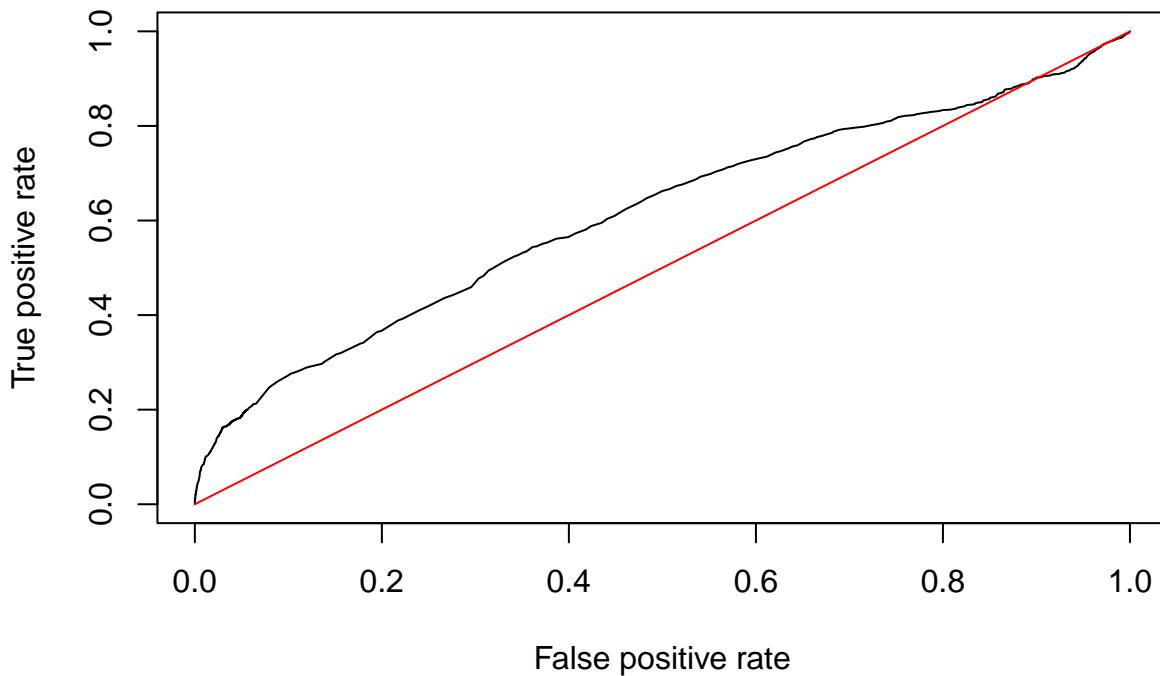
$$\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right) = -13.15 + 0.006537x_1 + 0.000441x_2 + 0.1585I_1 + 1.653I_2 + 0.6518I_3$$

We check as well for any multicollinearity in our model by using the Variance Inflation Factors (VIF).

```
##          yr_builtin      yr_renovated condition_medTRUE      view_medTRUE
##            5.485094           5.428401           5.189692           6.948286
##      waterfront1
##            12.178093
```

When we run the tests, we find that all the remaining predictor variables in our model are well below 5, with all of them being between 1 and 1.3. This indicates that there is little correlation between all predictors present.

ROC Curve for Reduced Model



The ROC curve above is above the diagonal for almost all values of the false positive rate, so it does better job than random guessing. We can also see that the point which maximizes the True Positive rate while minimizing the False Positive rate is around 0.3 for the FPR indicator.

Next, we take a look at the confusion matrix.

```
##          FALSE  TRUE
##  FALSE  4865  536
##  TRUE   3937 1469
```

The sample size of our data is $n = 10807$. From the above table, we can find the following values:

The **error rate** is $\frac{536+3937}{10807} = 0.4138984$

The **accuracy** is $\frac{4865+1469}{10807} = 0.5861016$

The **false positive rate** is $\frac{536}{536+4865} = 0.09924088$

The **false negative rate** is $\frac{3937}{3937+1469} = 0.7282649$

The **true positive rate** is $\frac{1469}{3937+1469} = 0.2717351$

The **true positive rate** is $\frac{4865}{4865+536} = 0.9007591$

The **precision** is $\frac{1469}{1469+536} = 0.7326683$

We will now compute the Area Under the Curve (AUC), which also assesses whether our model is better at predicting the value response than just random sampling.

```
## [[1]]
## [1] 0.6138072
```

The AUC of our ROC curve is 0.6138072, which means our logistic regression does better than random guessing.

We can now carry model diagnostic procedures by considering the three main approaches to model selection: forward selection, backward selection, and stepwise regression.

After calling the relevant functions, we found that all three selection methods returned the same type of model: a five-predictor model containing all our variables so far for the reduced model: `view_med`, `yr_renovated`, `yr_built`, `condition_med`, and `waterfront`.

Our R output was as follows:

```
Step:  AIC=14308.64
compare ~ view_med + yr_renovated + yr_built + condition_med +
    waterfront
```