

# **Diamonds**

Francisco Arrieta, Emily Schmidt and Lucia Camenisch

2022-12-20

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Data Exploration</b>	<b>2</b>
2.1	Dimension Summary . . . . .	2
2.2	Variable Visualisation . . . . .	7
<b>3</b>	<b>Variable Prediction and Model Performance Evaluation</b>	<b>9</b>
3.1	Linear Regression . . . . .	9
3.1.1	Data preparation . . . . .	9
3.1.2	Linear models based on complete set of predictors . . . . .	12
3.1.3	Linear models based on uncorrelated predictors . . . . .	12
3.1.4	Summary of models . . . . .	12
3.1.5	Accuracy measures of models on validation set . . . . .	13
3.2	$k$ -NN . . . . .	14
3.3	Decision Trees . . . . .	17
3.3.1	Regression Tree . . . . .	18
3.3.2	Boosted Tree . . . . .	19
3.3.3	Bagged Tree . . . . .	19
3.3.4	Random Forest . . . . .	20
3.3.5	Decision Tree Summary Table . . . . .	21
3.4	NeuralNetworks . . . . .	22
3.4.1	Basic Concept . . . . .	22
3.4.2	Data preprocessing . . . . .	22
3.4.3	Model Structures . . . . .	23
3.4.4	Neural Net Summary Table . . . . .	24
3.5	Ensembles . . . . .	25
<b>4</b>	<b>Model Performance Summary</b>	<b>25</b>
4.1	Predicted Price vs Real Price by Clarity . . . . .	27
4.2	Predicted Price vs Real Price by Color . . . . .	27
4.3	Predicted Price vs Real Price by Cut . . . . .	28
<b>5</b>	<b>Conclusions</b>	<b>28</b>
<b>6</b>	<b>References</b>	<b>28</b>

## 1 Introduction

The following report provides the result of using machine learning as a tool to estimate diamond prices for a jewelry company. Using snapshot information from their asset database, various methods were applied to train predictive models using regression analysis. The main objective was to use supervised learning methods to predict prices. These models are later compared to measure their effectiveness by using error measurements to quantify the distance between the prediction and the actual price. Finally, conclusions on the capability of each model are made and suggestions are given on which model to apply for the problem at hand.

The product being analyzed is diamonds. Being a luxury object with a long history, industry standards have been developed to serve as guidelines for estimating the value of the product. A wide variety of characteristics affect the overall price of diamonds, but this analysis will focus on their physical qualities, such as size dimensions, size ratios and color. Other factors inherent to scarce products in high demand of a capitalist economy will not be considered in our machine learning exercise.

The choice of methods shown in the report is not all inclusive and responds to the fact of only being some of the most used methods to this day. By implementing them, the analysts hope to provide a glimpse of the effects and importance in choosing the right model as well as displaying the differences between each one.

## 2 Data Exploration

In the early stages of any analysis, data exploration is a critical process aimed at understanding and analyzing the data set to gain insight and make valid decisions. The overarching goal of examining the data is to obtain intuition, identify questionable values, and strategize how to answer the problem statement. Therefore, let's summarize the findings that explain the main data characteristics, and dive into the relations between variables to understand how this analysis will be directed.

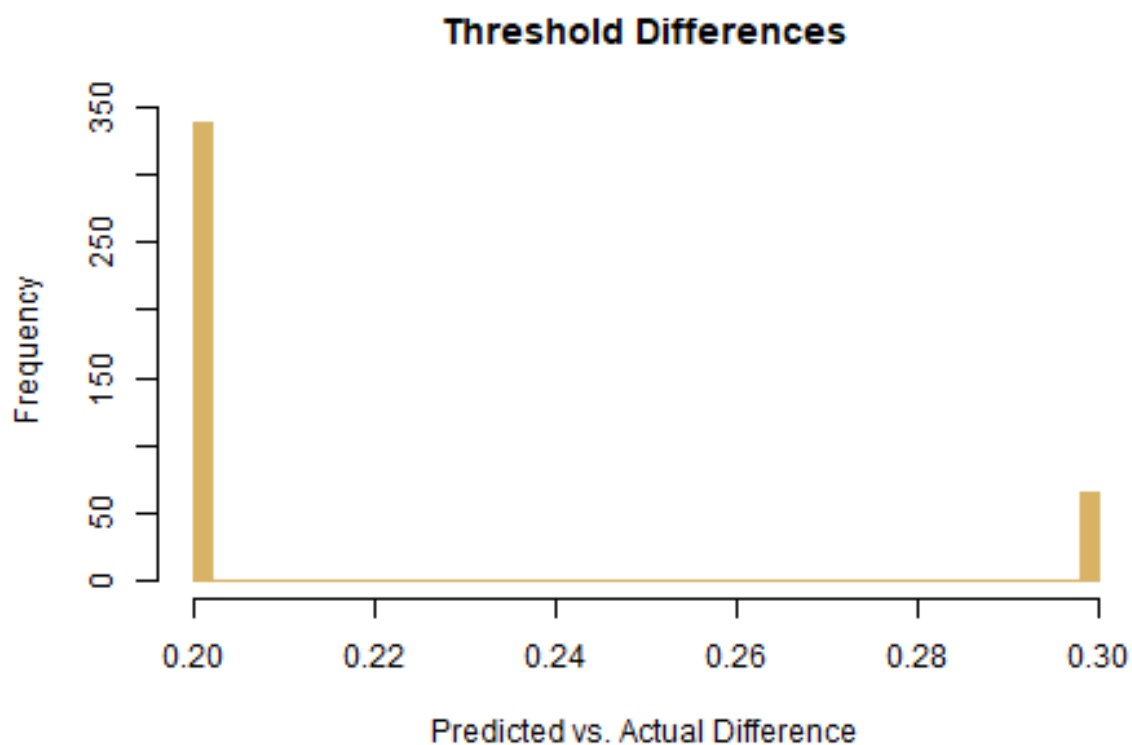
### 2.1 Dimension Summary

There is a variety of ways an analyst could approach understanding the observations in the data set. For starters, there are 53,940 records and 10 variables. The response variable (price) is an integer variable. There are an additional six numeric values and three ordinal/categorical features that have a factored structure. Those non-numeric values (cut, color, and clarity) need to be appropriately ordered, and renamed if there are any spaces within their naming convention as later this could obstruct certain code.

To dive deeper, we consider removing observations that are invalid because they add noise. For instance, price against all features were reviewed for possible impurities. In total, there are 275 (0.5%) records removed for the following reasons:

1. 20 rows had a depth of 0.0. This was considered inaccurate because a diamond needs to have this dimension specified.
2. The difference between length and width should be almost identical. If they were not, those rows were removed. In this case, only two records were not included in the data set at differences above 36.0mm.
3. The depth\_ratio is a calculation between length, width, and length. To investigate this feature further, our group computed depth\_ratio and saw there were differences between the actual and predicted values. Therefore, the executive decision was to remove any differences above a threshold of 0.3. A total of 253 rows were removed.

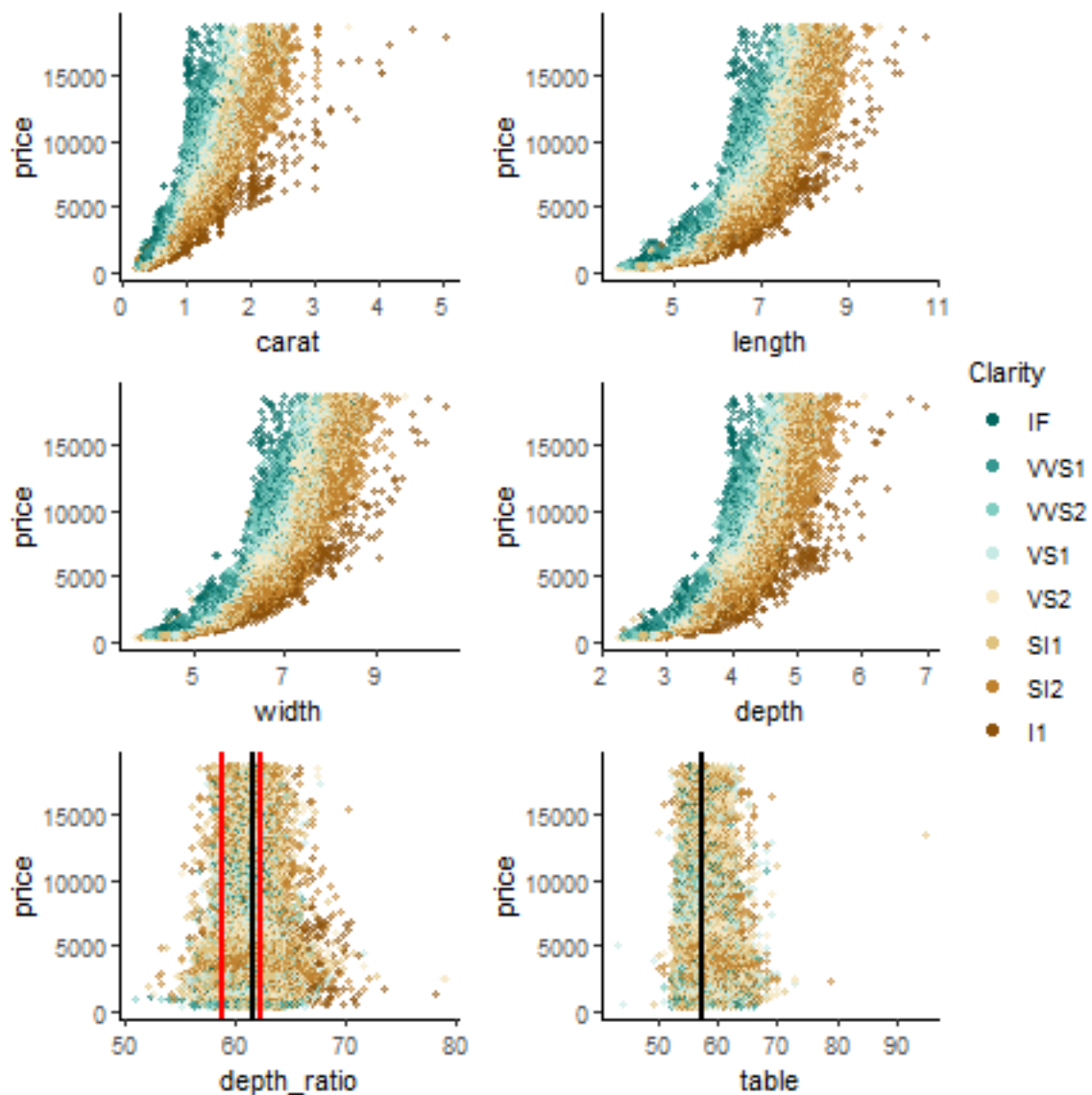
Fortunately, there were no null values that needed to be considered.



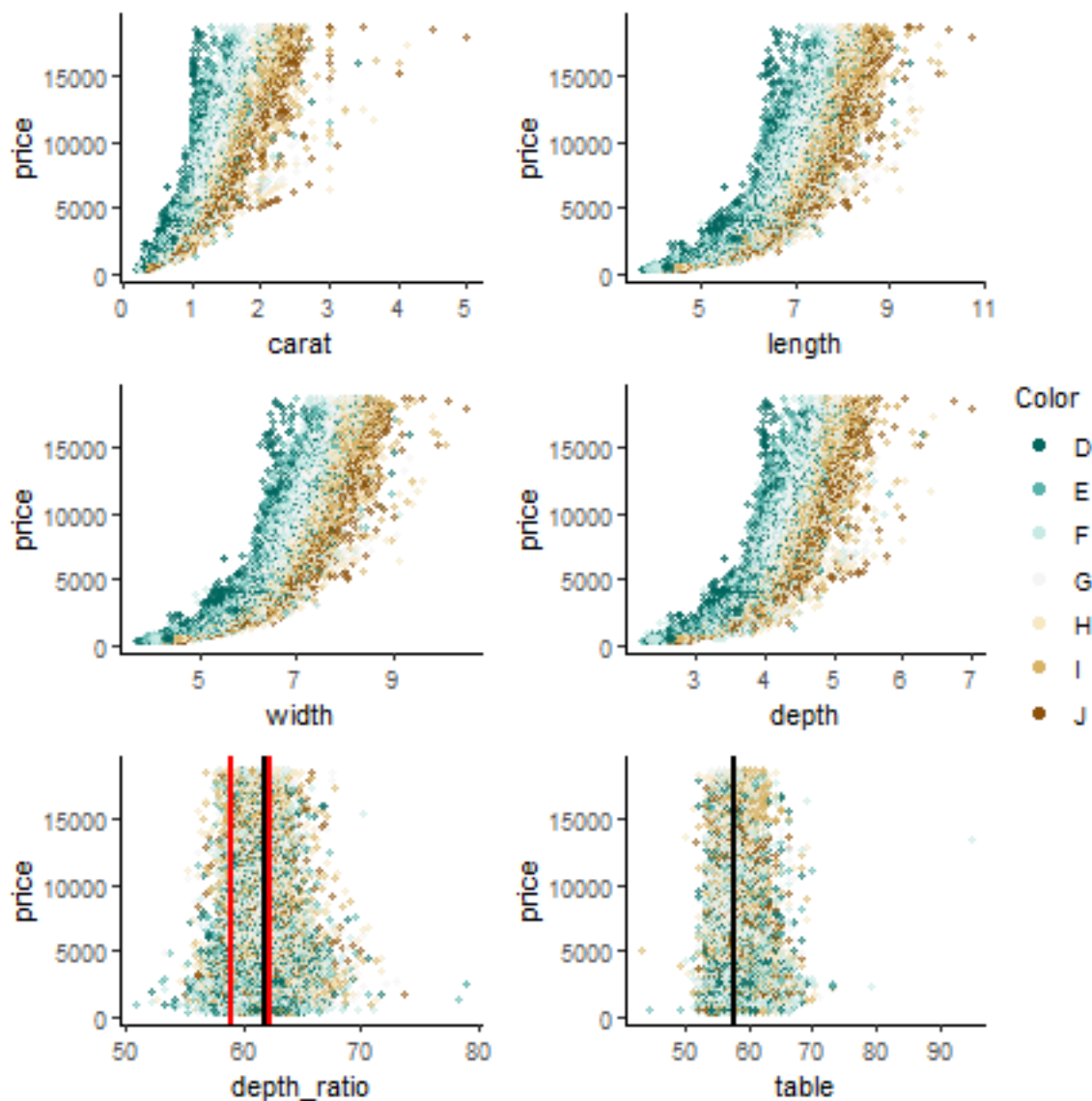
The second part of exploring the diamonds data set is to understand the relationships between vari-

ables. Prior to analyzing their correlations, we had an idea of which continuous variables would have a strong relation due to the depth\_ratio calculation. Later on, certain methods will analyze the potential issue of collinearity. By isolating price against the other continuous variables, carat, length, width, and depth, all have positively linear relations to the response feature. Some of the models used in the analysis to predict price will show these effects. In addition, the relationship between carat and price is unique because their relation is non-linear. We were also interested in showing how the categorical variables affect price. Clarity, color, and cut have distinct patterns seen between the four continuous variables mentioned above.

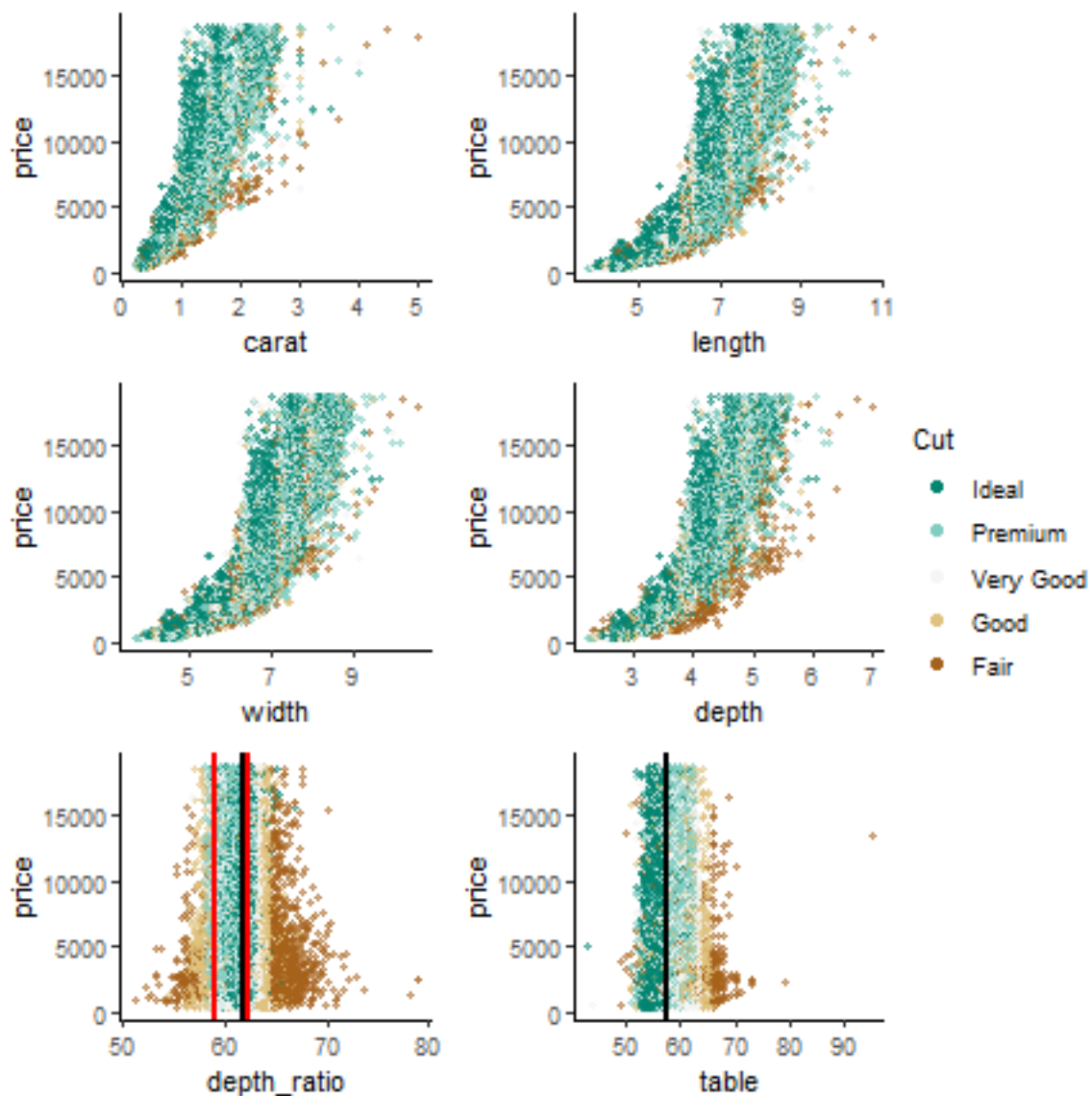
For instance, when plotting price by carat, width, length, and depth, with clarity as the color dimension, there is a clear distinction between how these points relate. Each scatter plot shows a strong, positive correlation to price. When reviewing the x-axis for each explanatory variable, the lower tiered option is typically observed more. The cheaper the gemstone, the worse quality the customer bought. Within each plot, there are outliers that can be explained by how the diamond was cut, which impacts the weight, length, and depth. By the color scale, it is noticed how there are more diamonds with a worse clarity (SI2/SI1) than there are with the best (IF/VVS1). Intuitively, this makes sense because not all customers have the available funds to afford the optimally designed diamond. On the other hand, no one wants the worst clarity either. Therefore, a majority of the diamonds lie between the second and third best, and second and third worst clarity factors. In addition, the relation between depth\_ratio and price is not distinct as those features are mixed throughout the relative ranges. Within all of the depth\_ratio plots, there are vertical lines that show the depth\_ratio mean (black) and range (red) of what is considered the prime value for depth\_ratio as long as it is above 59% but does not exceed 62.3%.



A similar relationship can be seen regarding price by the continuous variables, but as color as the color dimension. Instead of most of the observations falling closer to the best and worth clarify features, here the reader will notice how it appears that the colors E, F and G dominate (the white to light blue) the retail space.



Cut has a completely different trend than clarity and color as a majority of customers prefer and will purchase the ideal cut. This can be seen within the depth\_ratio plot especially. As stated earlier, the red lines indicate the optimal space in which a buyer would want the diamond's depth\_ratio. Most of all of the ideal cuts are within those bounds while the fair, good, and very good are outside. There is a clear distinction on how cut relates to the depth\_ratio along price.



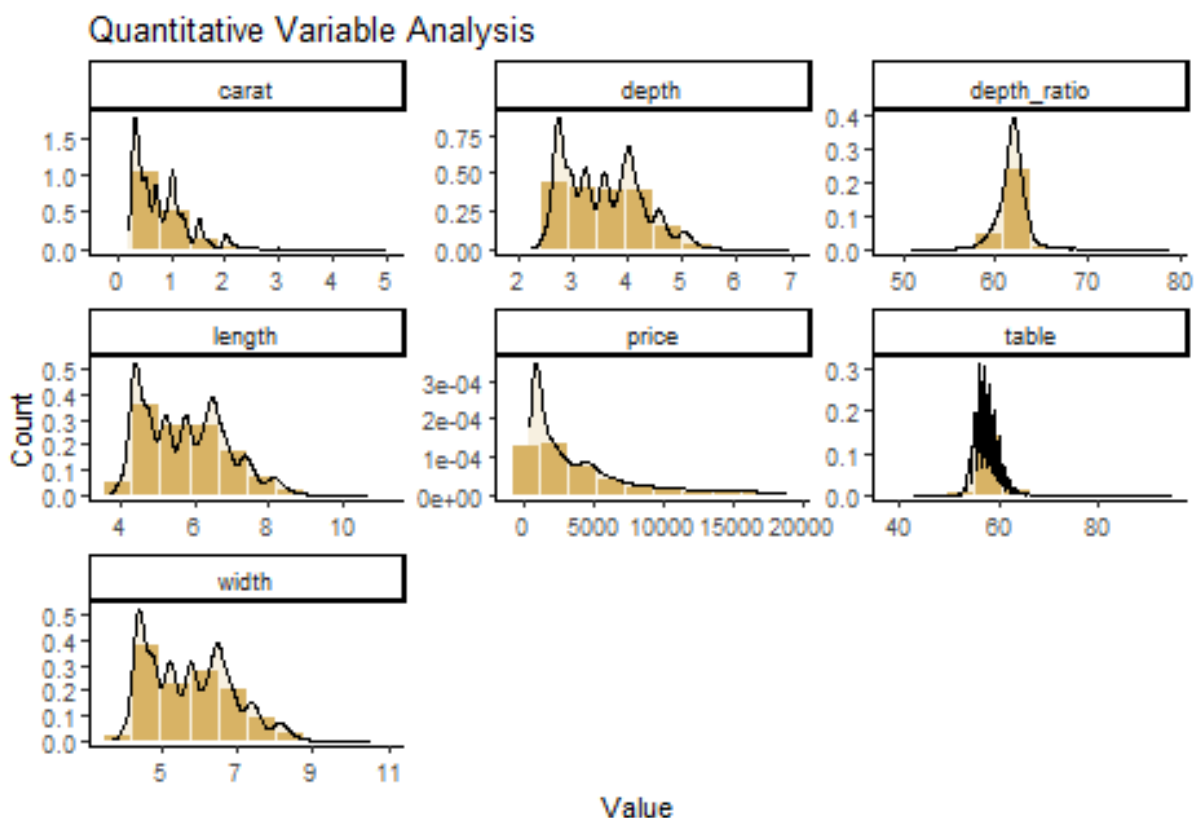
## 2.2 Variable Visualisation

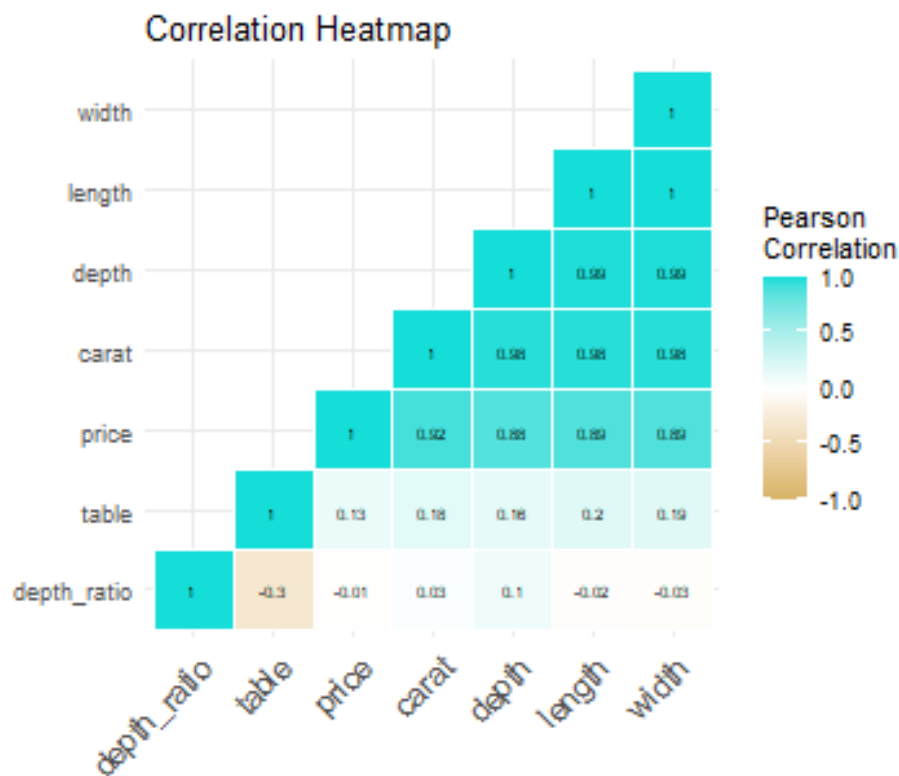
As a continuation of data exploration, variable visualization involves creating graphical representations of data to visually communicate insights and trends. It allows us to quickly and easily identify patterns that may not be immediately apparent from looking at raw data. There are three graphics that demonstrate how these continuous values are distributed. Since some models take into account transformation, those details will be described later. As a high-level overview, price and carat are positively skewed. In the scatter plot regarding price by cut for the depth\_ratio, the histogram also displays the same conclusion that the cut features fall into approximately 60 to 65 mm. The relationships in the Correlation heat map visually show how width, length, depth, and carat continue to be



highly correlated with price while depth\_ratio and table are closer to 0. Another way to visualize the relationships is through the Pearson Correlation Ellipses chart. The skinnier the ellipse, the more correlated the two values are.

Through exploration, we can gain a deeper understanding of the data and how it can be used to answer business questions or solve real-world problems, like predicting the price of a diamond.





### 3 Variable Prediction and Model Performance Evaluation

#### 3.1 Linear Regression

##### 3.1.1 Data preparation

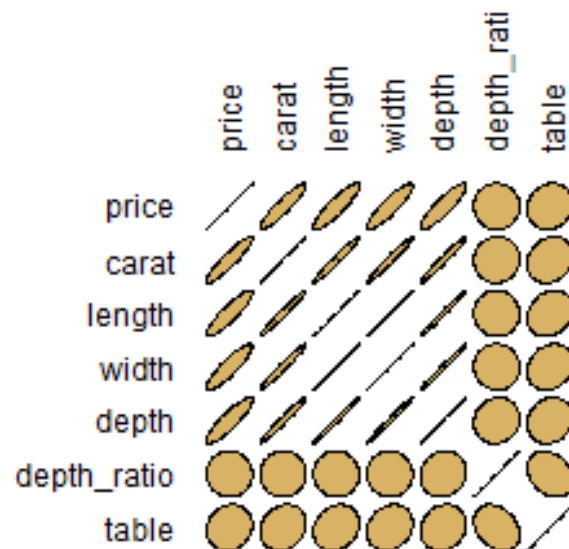
We begin by performing linear regressions on our data. As we saw during the exploration phase, some predictors (carat, length, width and depth) are highly correlated. Therefore, multicollinearity might be an issue.

We use the Generalized Variation Inflation Factors (GVIF) to measure the multicollinearity level of our data. This generalized version of the VIF allows us to take into account numerical and categorical predictors together. The GVIF clearly confirms that there is an issue, as length, width and depth all have coefficients above 1000. carat and depth\_ratio also have high values above 25, but they aren't as high extreme as the other three.

After removing length, width and depth, the GVIF coefficients of the remaining predictors are all under 2, which indicates the multicollinearity problem is solved. We display correlation ellipses of

numerical variables without length, width and depth.

### Pearson correlation ellipses for numerical variables



The only high correlation left is between carat and price, indicating that carat will surely be an important predictor for determining price.

Thus, we will perform linear regressions on two different models:

1. LM\_complete which contains all predictors
2. LM\_minus\_corr which has length, width and depth removed.

These two models will serve as basis for variable selection procedures later.

However, before starting to build our models, we also need to account for skewed variables. Linear regression might perform worse when dealing with skewed variables and it is common to use transformations such as a logarithm or a  $n$ th root to make variable more symmetrical.

We use an estimator of skewness called  $b_1$ , whose definition can be found [here](#).

The value of  $b_1$  is interpreted as follows:

- $0 \leq |b_1| < 0.5$ : variable is symmetrical;

**Table 1:** Skewness estimator for numerical variables

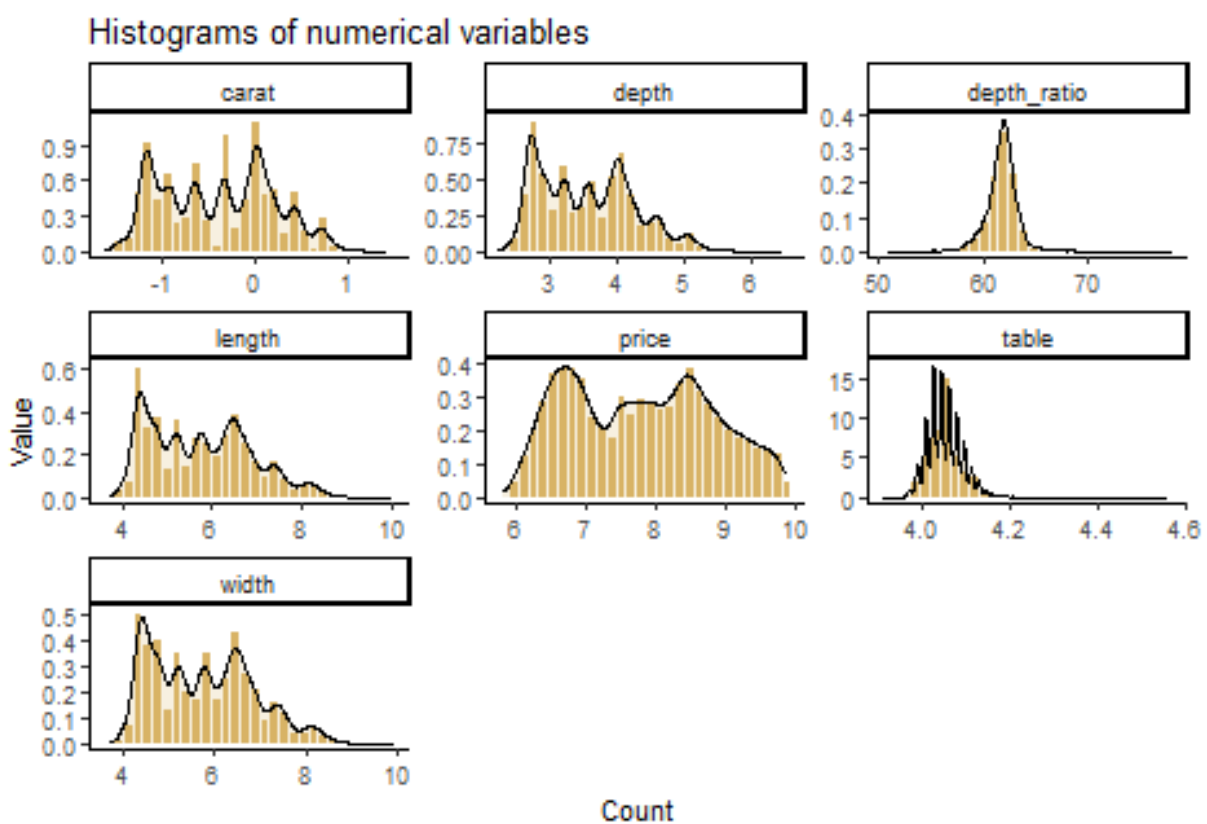
	$b_1$
price	1.6259747
carat	1.0967514
length	0.3984689
width	0.3922008
depth	0.3921093
depth_ratio	0.0182498
table	0.8728869

- $0.5 \leq |b_1| < 1$ : variable is moderately skewed;
- $|b_1| \geq 1$ : variable is highly skewed.

We compute  $b_1$  on our numerical variables and get the following results.

price and carat are highly skewed and table is moderately skewed.

We apply a logarithmic transformation to all three variables. The improvement can also be seen in the histograms, as they look more symmetrical now.



**Table 2:** Predictors used in each linear model

Model	Cut	Color	Clarity	Carat	Length	Width	Depth	Depth Ratio	Table
LM_complete	X	X	X	X	X	X	X	X	X
LM_forward_complete	X	X	X	X	X	X	X	X	X
LM_backward_complete	X	X	X	X	X		X	X	
LM_stepwise_complete	X	X	X	X	X		X	X	
LM_CpAIC_complete	X	X	X	X	X		X	X	
LM_minus_corr	X	X	X	X				X	X
LM_forward_minus_corr	X	X	X	X				X	X
LM_backward_minus_corr	X	X	X	X				X	
LM_stepwise_minus_corr	X	X	X	X				X	
LM_CpAIC_minus_corr	X	X	X	X				X	

We also standardize numerical variables by subtracting their mean and dividing by their standard deviation. This makes the comparison of  $\beta$  coefficients between variables easier.

### 3.1.2 Linear models based on complete set of predictors

The data is now ready for our linear models. For both LM\_complete and LM\_minus\_corr, we perform the following linear regressions:

1. Linear regression on the whole model
2. Forward selection on the model (iterative method)
3. Backward selection on the model (iterative method)
4. Stepwise selection on the model (iterative method)
5. Mallows's  $C_p$  and AIC selection on the model (global method)

### 3.1.3 Linear models based on uncorrelated predictors

#### 3.1.4 Summary of models

For each of these models, we summarise which variables are used as predictors in the following table.

For both basis models, forward selection doesn't discard any variables, whereas backward, stepwise and global selections all choose the same model with less variables than initially.

Thus, we have four distinct models in total. We assess the predictive performance of these four models on our validation set by computing the five accuracy measures seen during the course.

Let us recall the definitions and meaning of these measures (Data Mining for Business Analytics - Concepts, Techniques, and Applications in R, chapter 5.2, page 119). We denote the residuals by  $r = y - \hat{y}$ .

1. **ME** (Mean Error) gives an indication of whether the predictions are on average over- or under-predicting the outcome variable.

$$ME = \frac{1}{n} \sum_{i=1}^n r_i$$

2. **RMSE** (Root Mean Squared Error) is similar to the standard error of estimate in linear regression, except that it is computed on the validation data rather than on the training data.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n r_i^2}$$

3. **MAE** (Mean Absolute Error) gives the magnitude of the average absolute error.

$$MAE = \frac{1}{n} \sum_{i=1}^n |r_i|$$

4. **MPE** (Mean Percentage Error) gives the percentage score of how predictions deviate from the actual values (on average), taking into account the direction of the error.

$$MPE = 100 \cdot \frac{1}{n} \sum_{i=1}^n \frac{r_i}{y_i}$$

5. **MAPE** (Mean Absolute Percentage Error) gives a percentage score of how predictions deviate (on average) from the actual values.

$$MAPE = 100 \cdot \frac{1}{n} \sum_{i=1}^n \left| \frac{r_i}{y_i} \right|$$

### 3.1.5 Accuracy measures of models on validation set

The accuracy measures for our four models are summarized below. In order to give meaningful results, the outcome variable `price` has been rescaled to its original scale and retransformed by taking the exponential (to cancel out the logarithm). Thus, the ME, RMSE and MAE can be interpreted in the price currency, dollars.

The mean error is bigger in the models without multicollinearity issues, but their RMSE is smaller. Since the mean error is positive in all four models, we are under-predicting the price of diamonds by 36 or 50 dollars on average depending on the model.

**Table 3:** Accuracy measures of linear models

Model	ME	RMSE	MAE	MPE	MAPE
LM_complete	36.36648	846.5888	408.5364	-0.8377036	10.33643
LM_CpAIC_complete	36.32774	845.5018	408.1349	-0.8408540	10.33659
LM_minus_corr	50.38994	810.1522	405.0486	-0.8420621	10.39559
LM_CpAIC_minus_corr	50.39878	810.1610	405.0616	-0.8418121	10.39568

An increase of 14\$ in the mean error is a good trade-off to reducing the RMSE, which indicates how much the predictions will fluctuate from real values. The other three measures are quite close in all four models.

Considering the parsimony principle and the RMSE, the best choice is our fourth linear model, which contains cut, color, clarity, carat and depth\_ratio.

Overall, the RMSE for linear regression remains quite high and as we will see in the following chapters, some models will achieve much better predictive performances.

## 3.2 $k$ -NN

**3.2.0.1 Overview**  $k$ -nearest neighbors (kNN) is a simple algorithm used for classification (categorical) and regression (continuous). Since the response variable (price) is a numeric outcome, this analysis will use  $k$ -NN to predict the price of a diamond. In a regression setting, the algorithm relies on finding the most ‘similar’ records in the training data. Then, one would calculate the weighted average of the numerical target of the  $k$ -nearest neighbors for the data point.

According to the Data Mining textbook, the value of  $k$  is based on a nonparametric method since it ‘draws information from similarities between the predictor values of the records in the dataset.’ There are several ways one could choose this hyperparameter, but this analysis will focus on one method by the optimization of RMSE.  $k$  is a critical input within the  $k$ -NN function because it determines how many neighbors will be considered when making a prediction. Penn State states that a ‘larger  $k$  leads to a smoother boundary but may also introduce noise.’ On the other side, a small  $k$  ‘can increase the complexity of KNN.’

The  $k$ -nearest neighbors algorithm is very useful when predicting price because the data does not need to be transformed or predictors selected in a particular way.

**3.2.0.2 Model Structure** There are several steps that need to be completed prior to running the first  $k$ -NN model with  $k = 1$ . First, you clean the data to ensure that missing values are either removed or explained, and that all variables are essential to the response variable (see EDA section). Next, the

data is normalized so that the output remains unbiased. Normalizing simply means that the raw data is put on all the same scale ('subtracting the mean and dividing by the standard deviation'). Once the data is cleaned and normalized, it can be split into the training, validation, and test sets. Now that the foundation of the model has been built, let's start constructing the first  $k$ -NN model.

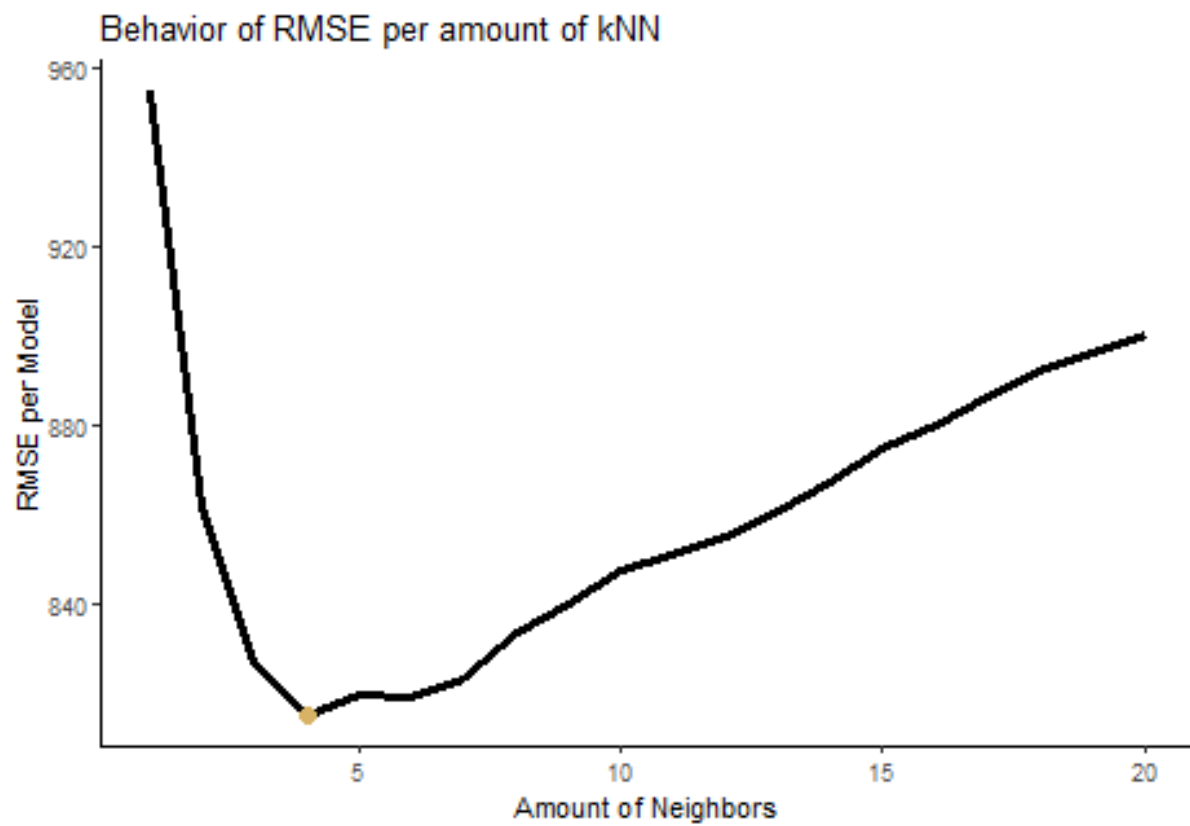
The value of  $k = 1$  is used as a starting point to see how well the model will perform. Since price is our focal point, the two error measurements that will be compared throughout all model methods are Mean Error (ME) and Root Mean Squared Error (RMSE). RMSE is of particular interest because that gives the variance of how far the predicted value is away from the actual price. Seen in the first  $k$ -NN error table, ME (\$20.60) and RMSE (\$954.60) could potentially get better if we optimized  $k$ .

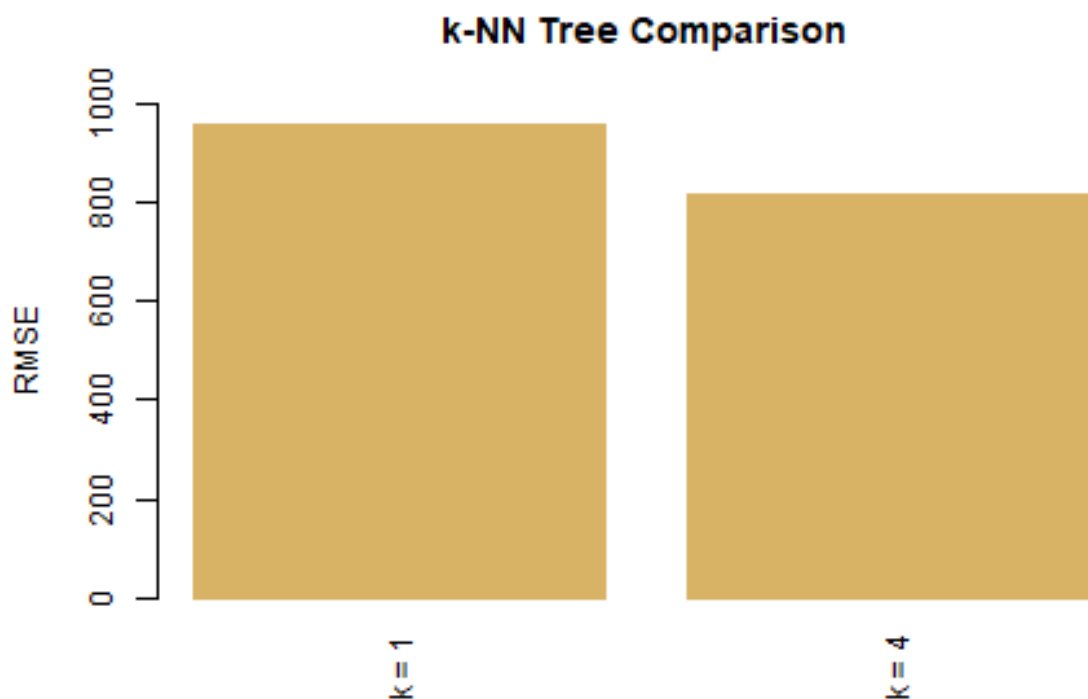
The optimal  $k$  is 4 because it produces the lowest RMSE at \$814.59. After rerunning the  $k$ -NN model, the results improved in one aspect, but not the other.

k	accuracy
1	954.5954
2	861.0887
3	827.1196
5	819.2859
6	818.5541
7	822.7554
8	833.2104
9	839.9301
10	847.3985
11	850.7010
12	854.8280
13	860.4922
14	867.3541
15	874.8586
16	879.9069
17	886.0747
18	891.9780
19	896.3521
20	900.0938



	ME	RMSE	MAE	MPE	MAPE
k = 1	20.6	954.6	487.17	-1.83	13.73
k = 4	40.86	814.59	431.29	-2.45	12.14





As stated earlier, the table above shows how the ME got worse for  $k = 4$ , but better with RMSE once the model was optimized. In this analysis, one of the main goals is to reduce the variance between predicted and real price. Therefore, the difference of \$140.00 is more important even when the ME slightly increases. Between  $k = 1$  and  $k = 4$  for the  $k$ -NN model,  $k = 4$  with an RMSE of \$814.60 will predict prices better.

### 3.3 Decision Trees

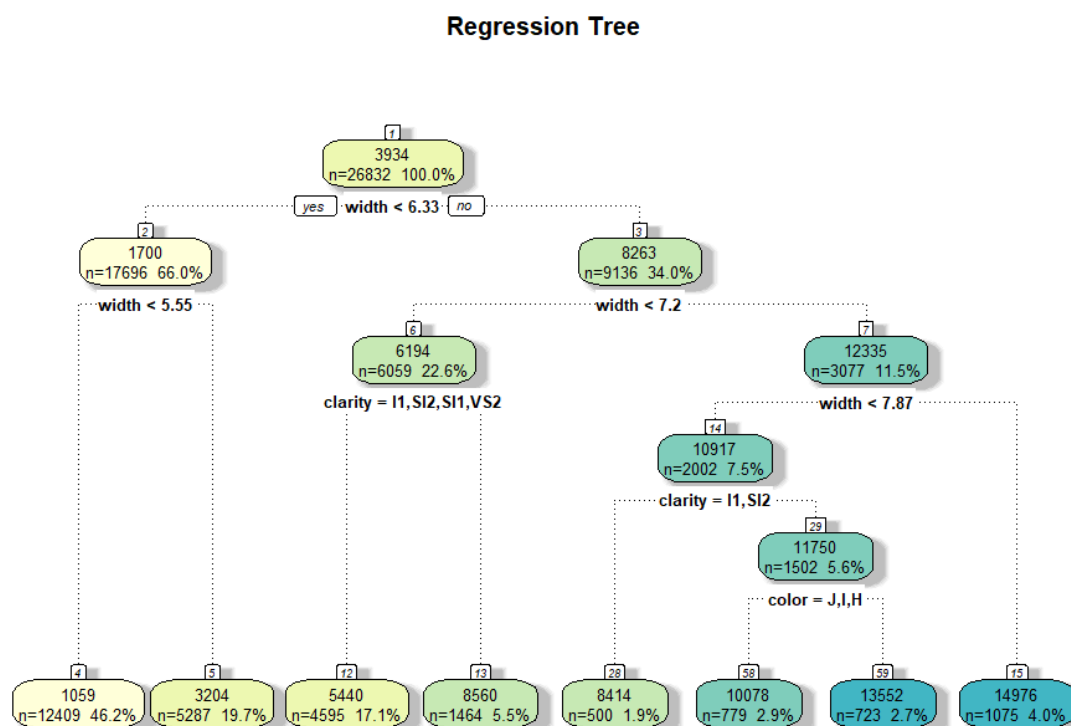
**3.3.0.1 Regression Tree Overview** There are four methods within this section that show various ways on how trees are built. They include a regression tree, boosted tree, bagged tree, and a random forest. The regression tree is most transparent and easy to interpret while the other three combine results from multiple trees.

A regression tree is a flexible data-driven method that can be used for prediction of a continuous variable. The tree separates 'records into subgroups by creating splits on predictors. These splits create logical rules that are homogeneous.' Those splits 'divide the data into subsets, that is, branches, nodes, and leaves. Like decision trees, regression trees select splits that decrease the dispersion of

target attribute values. Thus, the target attribute values can be predicted from their mean values in the leaves' which reduces the variance of the target variable.

**3.3.0.2 Model Structure and Analysis** Since the tree proactively takes into consideration the most important attributes to split on, multiple trees are created to check its validity. This idea stemmed from wanting to ensure that by removing variables due to their weak relationship to the predictor variable or by pruning, the results produced the best ME and RMSE. Within the three regression models, all have the same error rates. In addition, there are eight terminal nodes for each model and all had width, length, carat, and depth being the most important features to predict price. To exemplify this, we show the original regression tree that first splits on width < 6.33, seven splits total and uses the following for the primary splits: width < 6.325, carat < 0.985, length < 6.325 and depth < 3.935.

### 3.3.1 Regression Tree



### 3.3.2 Boosted Tree

**3.3.2.1 Overview** Boosted trees are a type of ensemble model that ‘transforms weak decision trees into strong learners. Each new tree is built considering the errors of previous trees.’ The idea behind boosting is to train a series of weak models additively, with each model attempting to correct the errors made by the previous model. Since this model is prone to overfitting, the parameters need to be carefully considered so that the lowest error rates can be achieved.

**3.3.2.2 Model Structure and Analysis** There is a wide variety of parameters that can be chosen for a boosted tree. Therefore, those options were used differently to choose the best model based off the number of predictors and how deep each tree would be allowed to interact. The main difference between each model is the number of trees ran. As the size of the tree increase, the better the model because it only considered the most important factors.

The best boosted model was the one with 100 trees and three cross-validation folds. The tree used six predictors to produce a model with a ME of \$9.00 and RMSE of \$1,170.30. What is fascinating about this tree comes from the importance variables chart. In the regression tree, clarity was one of the three least important factors, but here we see that it as the fourth which is above length.

### 3.3.3 Bagged Tree

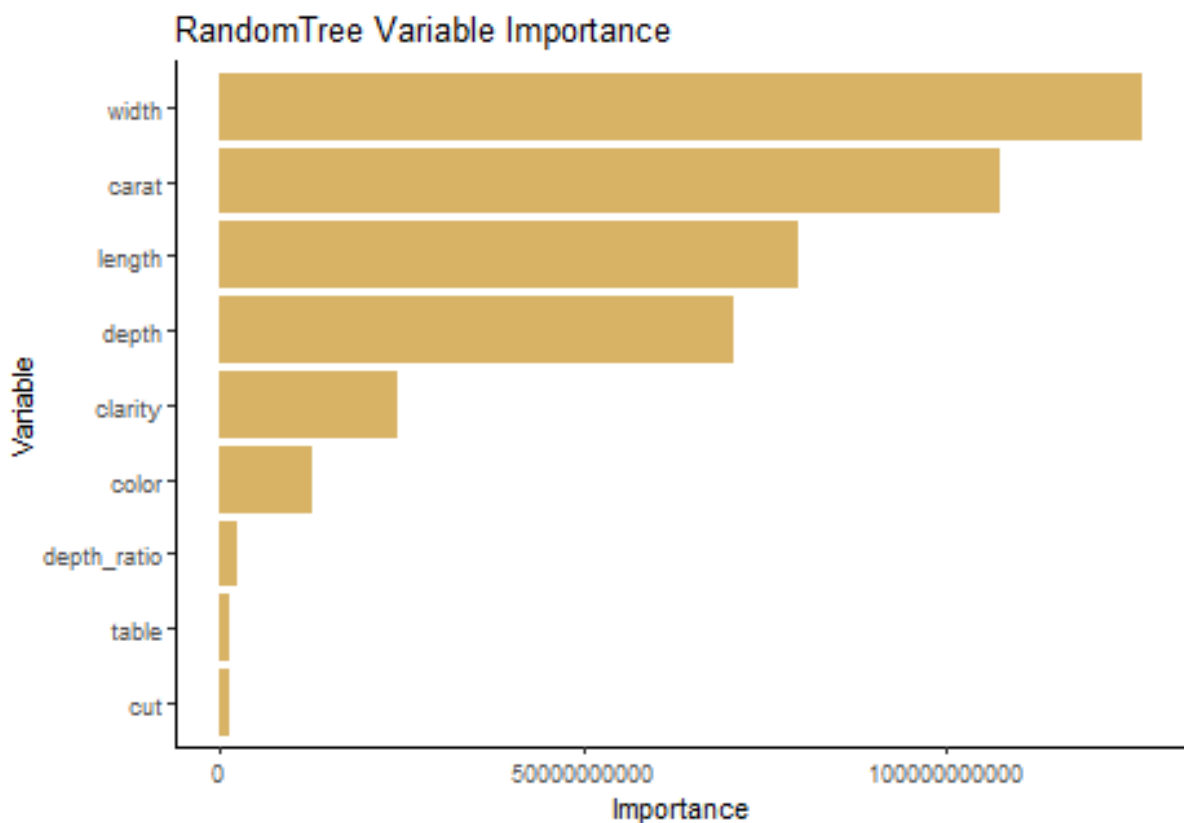
**3.3.3.1 Overview** The third method is bagged trees. This type of model ‘combines the predictions from multiple machine learning algorithms together to make more accurate predictions than any individual model.’ The objective with bagging is to train a large number of decision trees on different subsets of the training data, and then average the predictions of all the trees to make a final prediction. This model has the ability to reduce variance because it introduces randomness into the training process.

**3.3.3.2 Model Structure and Analysis** This model is pretty straight forward because the best model was built by sticking with the basics. ‘The only parameters when bagging decision trees is the number of samples and hence the number of trees to include. This can be chosen by increasing the number of trees on run after run until the accuracy begins to stop showing improvement.’ Overall, the bagging tree predicts price better than the regression tree, but worse than boosting. The ME was the lowest at \$0.93, but what is more important to consider is the RMSE which is \$1,248.77.

### 3.3.4 Random Forest

**3.3.4.1 Overview** ‘Random forests are a special case of bagging, a method for improving predictive power by combining multiple classifiers or prediction algorithms.’ They are an ensemble based on bagged trees which involves training each tree on a bootstrapped sample of the original data.

**3.3.4.2 Model Structure** Similar to the bagged trees, random forests are pretty simplistic. The main focus in this section was to check how many trees should be run within the model. There is a trade-off between computational power and RMSE. For example, what is the difference in RMSE if the model was built off of 100 trees versus 60? The model with 100 trees was the best model with an RMSE of \$575.42.

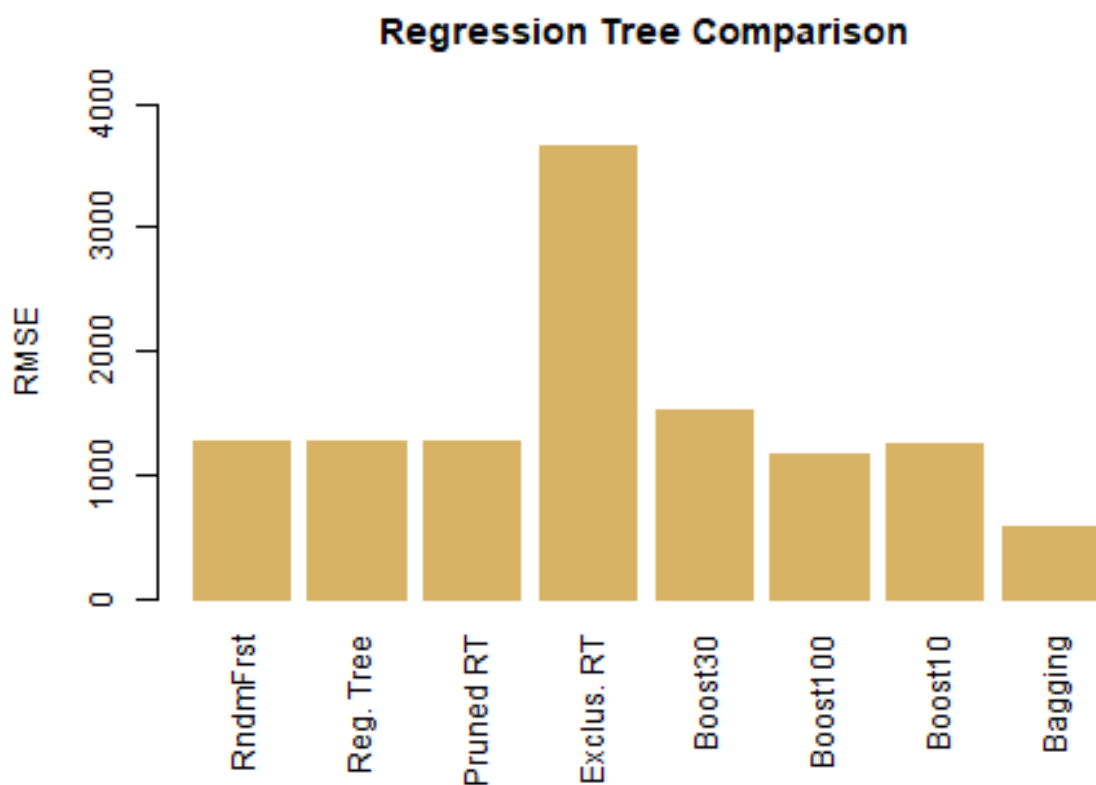


Throughout this entire section, typically the lower RMSE produced the best model. Additionally up to this point, there have been three or four variables that were most importance to the model. In the random forest though, there is much more weight on features such as clarity and color. In random forests, we proposed a trade-off with picking a model with less trees because an analyst would want

	ME	RMSE	MAE	MPE	MAPE
Reg. Tree	6.1	1267.15	846.93	-14.54	33.08
Exclus. RT	6.1	1267.15	846.93	-14.54	33.08
Pruned RT	6.1	1267.15	846.93	-14.54	33.08
Boost10	-3.9	3660.42	2765.96	-144.49	171.65
Boost30	2.54	1531.52	985.7	-35.26	46.86
Boost100	9	1170.3	680.22	-12.85	26.22
Bagging	0.93	1248.77	808.93	-14.35	31.45
RndmFrst	2.97	577.63	289.51	-1.4	7.24

to choose the parsimonious model. Therefore, the random forest with 60 trees had a slightly higher RMSE at a difference of approximately \$2 which gave a \$577.63 for the RMSE.

### 3.3.5 Decision Tree Summary Table



Throughout the four methods mentioned, the random forest model proved to be the most effective in predicting price with a RMSE of \$577.63, which is more than two times smaller than the best boosting

model. Although not all of the models results are shown, a breakdown of their error results can be seen within the table above. Overall, the regression tree, bagging and boosting were all within \$100 RMSE difference, but no match for the random forest.

### 3.4 NeuralNetworks

#### 3.4.1 Basic Concept

Artificial neural networks are a method of machine learning that receives the name due to a comparison made with actual neurons in the human brain. They consist various nodes that communicate with each other to predict an outcome based on the relationships that were determined in the process. Neural nets have proven to be very good at predicting values, through regression or classification and have been the center of much research in the past years. Though this method was once considered unuseful, with the improvement of computational power, it has acquired new popularity for prediction of images, speech recognition and non-explicit trends. (Hardesty, 2017)

The basic parts of a neural networks are: the structure, composed of layers and nodes, the weights, and the activation function. The first consists of the different components used to train a neural network, ie the explanatory variables and their connection. These independent variables are considered the *input nodes* and the outcome variable the *output node*. The layers in between these to are called the *hidden layers* as these are used to calculate the output but have no easily association of there value with regards to the outcome. This is why neural networks are considered by some like a **black box**, where the exact method for predicting is not easily explained to those who are not well informed. The weights are used to pass a certain amount of a value to the next node which after passing through the activation function will pass on to the next, and so on. This weighting can be assigned randomly or by specific methods, depending on the problem at hand and analyst discretion. Similarly, the activation function calculates the position in a curve, ie the expected value of the prediction. This as well changes per prediction problem and analyst discretion. While some are considered better for certain tasks there is no limiting factor in the way a neural net is structured.

#### 3.4.2 Data preprocessing

Due to the fact that in every node we calculate the position of the prediction in a curve, the scale of the values used affects the output of every node. Moreover, since we use all types of variables for prediction in neural networks (continuous and categorical) the difference in amplitude is very important. This is why it is standard procedure to normalize the data so they are all in the same scale.

### 3.4.3 Model Structures

There is no standard or ideal way of setting up the amount of layers in a network, but common rule of thumb is to start with the same amount of input variables and then reduce to see if this improves. (Shmueli, et al. 2018: 286) In this case, it was decided to follow the following structures:

- The first model is composed of all the variables (26 explanatory) and a single hidden layer of one node. All models have only one output node as the desired outcome is a single prediction of price. This model is considered to be the most basic and should in theory have the least predictive performance.
- The second model includes a hidden layer of 26 nodes, so it equals the input nodes.
- To see if there is an improvement, we do another model with just 13 nodes in a single hidden layer.

Of these single layer models, the best is the one with 26 nodes. We measure this by looking at the error in prediction, also called the accuracy. There are many ways of measuring this, but as mentioned before, the two we looked at are the Mean Error, as a measure of accuracy, and Root Mean Squared Error, as a measure of precision. For the 26 node model, the RMSE is \$745.16 vs \$744.04 for the 13 node net. While just looking at this configuration one could infer using less nodes is better, after multiple tests it is concluded that keeping the 26 nodes as the first hidden layer is best for reducing the RMSE going forward.

The next configuration tested is adding 2 additional hidden layers, one of 26 additional nodes, and another of 13. Just by using this configuration, the model RMSE reduced to \$627.65. As this is already a rather low RMSE, in comparison to previous neural networks, as well as other models seen above, it was decided to keep this structure as the definitive one for establishing the most accurate prediction.

The next two adjustments made are regarding the weights and the activation function. This structure was adjusted to use the “*Glorot Normal*” weight initiation method. This initiation method consists of assigning weights to the nodes using a de probability curve of a truncated normal distribution as so:

#### **Glorot Normal Initialization:**

$$w_i \sim \text{Gaussian} \left( \mu = 0, \sigma^2 = \sqrt{\frac{2}{u_{in} + u_{out}}} \right)$$

where:

$u_{in}$  = number of input nodes

$u_{out}$  = number of output nodes

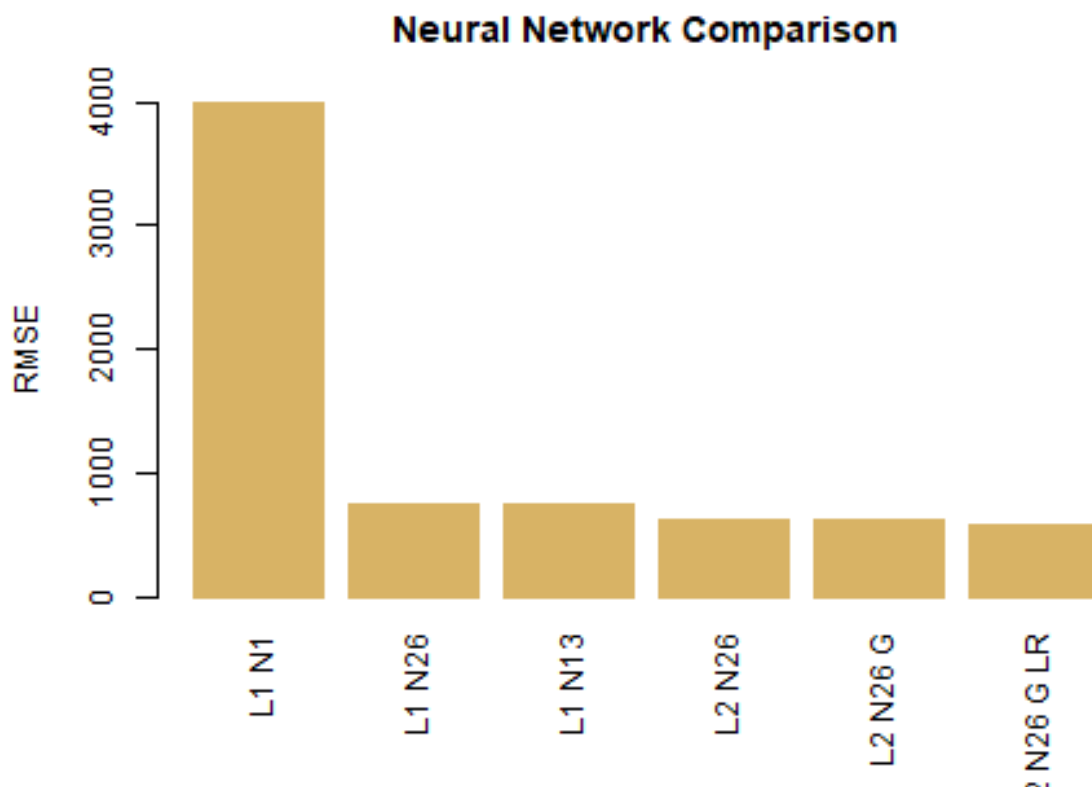


	ME	RMSE	MAE	MPE	MAPE
L1 N1	-6.62	3982.99	3025.61	-157.47	187.29
L1 N26	-83.28	745.16	450.94	-11.37	17.91
L1 N13	-86.71	744.04	454.38	-12.89	19.21
L2 N26	30.22	627.65	364.54	-5.43	13.49
L2 N26 G	-32.61	630.02	358.21	-6.62	12.72
L2 N26 G LR	8.98	577.08	322.95	-2.03	10.23

This initiation method demonstrated much better results than others such as the normal distribution, uniform distribution of the Glorot Uniform, which is why we kept it for the final model.

Finally regarding the activation function, many options are available as well. Depending on the problem at hand, one can use different functions that draw different curves and hence produce different predictions. The one used for our case was a *sigmoid* activation which by different literature is the most commonly used due to its good performance results. Though we tested others like “*relu*” or “*softplus*”, sigmoid ended being the most accurate predicting.

### 3.4.4 Neural Net Summary Table



	ME	RMSE	MAE	MPE	MAPE
Test set	25.80266	582.3353	306.1574	-1.682175	8.448477

	ME	RMSE	MAE	MPE	MAPE
Multiple Linear Regression	50.4	810.16	405.06	-0.84	10.4
Random Forest	2.97	577.63	289.51	-1.4	7.24
k-Nearest Neighbor	40.86	814.59	431.29	-2.45	12.14
Neural Network	8.98	577.08	322.95	-2.03	10.23
Ensemble	25.8	582.34	306.16	-1.68	8.45

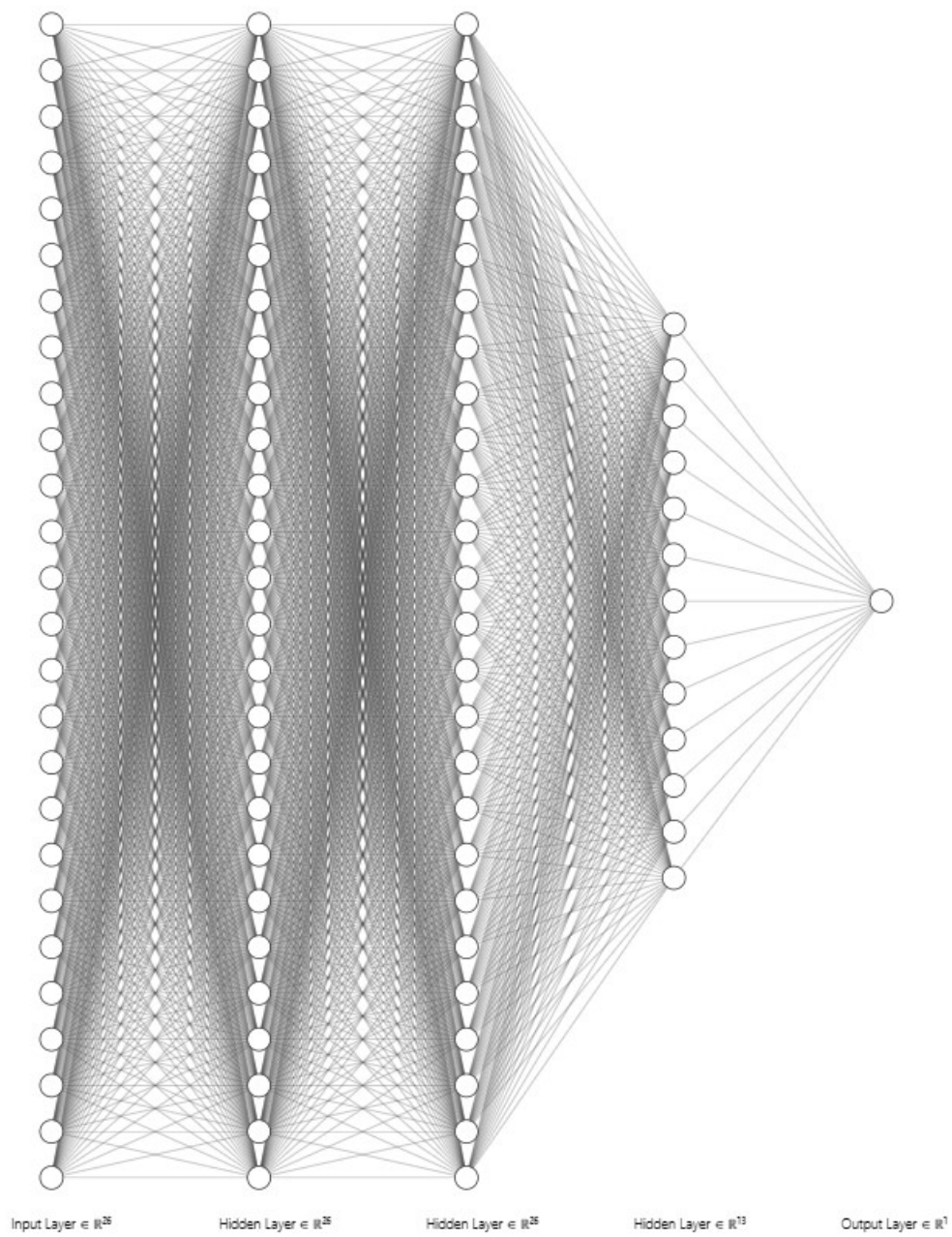
As can be seen, the predictive performance of the neural networks can vary a lot depending on the structure and parameters chosen. The negative point in this method is the requirement of human input to adjust the model to learn more precisely. Nevertheless, once the many trials have been executed, the results are very positive for a good predictive model.

### 3.5 Ensembles

**3.5.0.1 Overview, Model Structure and Analysis** ‘Ensemble methods is a machine learning technique that combines several base models in order to produce one optimal predictive model.’ The simplest approach is to combine predictions from each method above that had the best/lowest RMSE. To accurately compare models, the mean is taken over the four models to compute the average predicted price and find the error rates between the calculated price and real prices. As the table reveals, the RMSE is \$588.34 which is pretty close to random forests and the best neural network in terms of variance.

## 4 Model Performance Summary

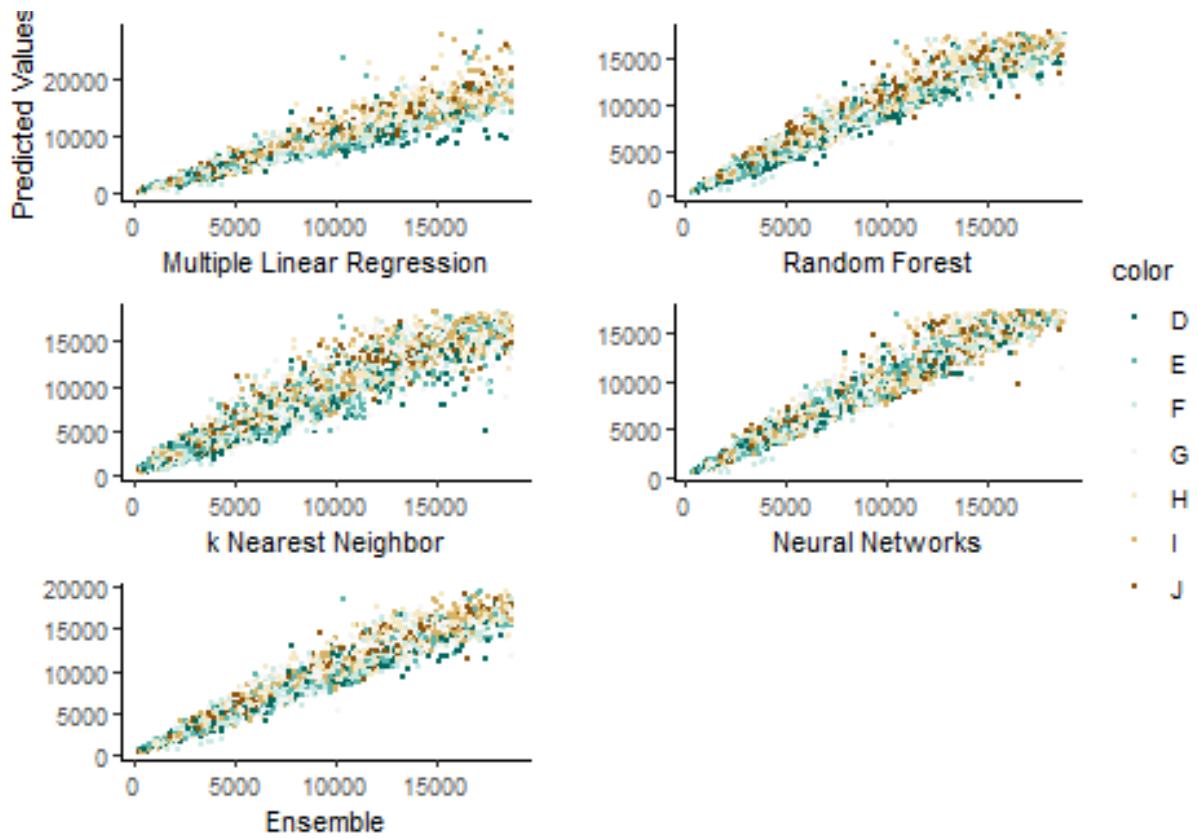
In this section, we can see how the predicted price behaves with respect to the real price, having used each of the methods and mapping per each categorical variable. This helps see how each model compares with regards to precision. The models with more precision are closer to a diagonal line. As can be seen,



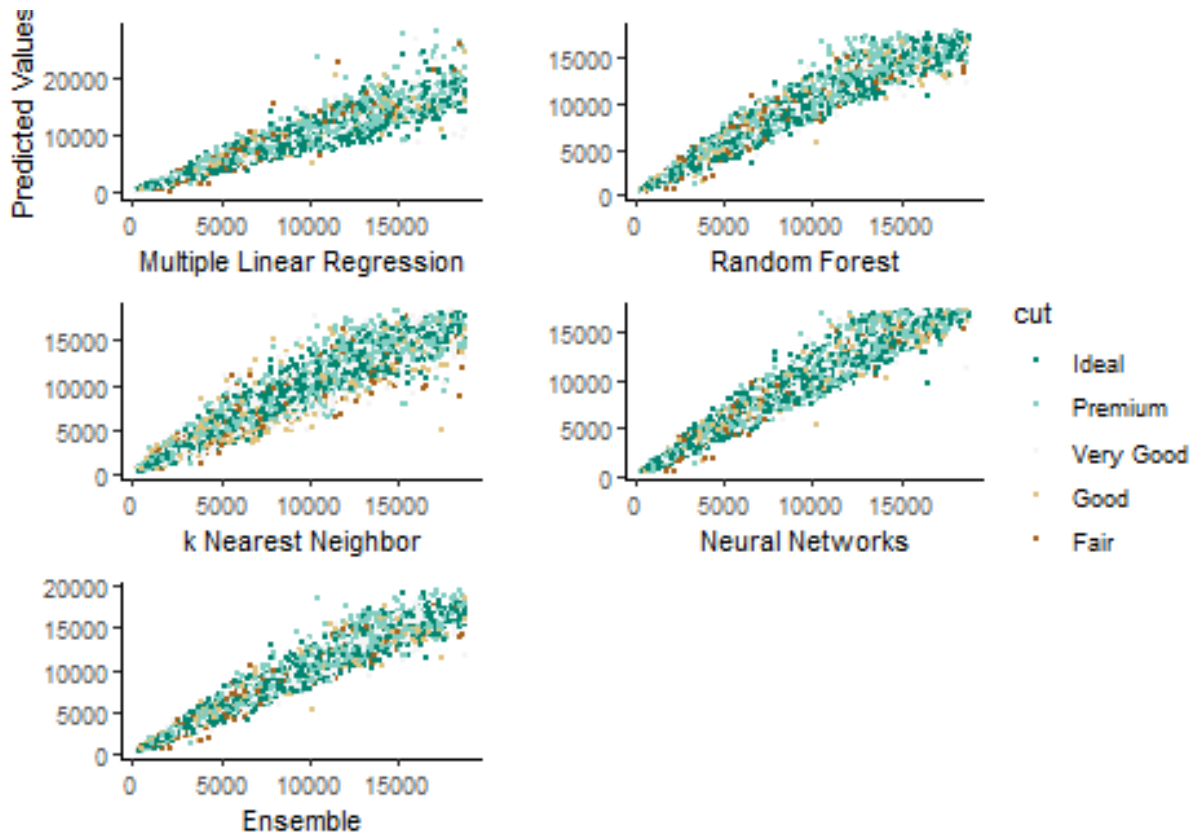
**Figure 1:** Visual description of Neural Network

#### 4.1 Predicted Price vs Real Price by Clarity

#### 4.2 Predicted Price vs Real Price by Color



### 4.3 Predicted Price vs Real Price by Cut



## 5 Conclusions

## 6 References

EDA - <https://r-graph-gallery.com/199-correlation-matrix-with-ggally.html> and <http://www.sthda.com/english/wiki/ggplot2-quick-correlation-matrix-heatmap-r-software-and-data-visualization>

k-NN - <https://online.stat.psu.edu/stat508/lesson/k#:~:text=The%20larger%20k%20is%2C%20the,Nearest%20Neigh>

- <https://www.analyticsvidhya.com/blog/2018/03/introduction-k-neighbours-algorithm-clustering/>

Regression Tree - Data Mining textbook - <https://www.ibm.com/docs/en/db2-warehouse?topic=procedures-regression-trees>

Boosted Tree - <https://towardsdatascience.com/introduction-to-boosted-trees-2692b6653b53>

- <https://leonlok.co.uk/blog/decision-trees-random-forests-gradient-boosting-whats-the-difference/>

Bagging Tree - <https://machinelearningmastery.com/bagging-and-random-forest-ensemble-algorithms-for-machine-learning/>

Random Forest - Data Mining textbook

Ensemble - <https://towardsdatascience.com/ensemble-methods-in-machine-learning-what-are-they-and-why-use-them-68ec3f9fef5f>

Hardesty, Larry. 2017. [Explained: Neural Networks] MIT News. (<https://news.mit.edu/2017/explained-neural-networks-deep-learning-0414>)