

# ASC61013 Data Modelling and Machine Intelligence

## Course Work – Diamond

Ankur Singh Gulia

Registration Number: - 220181437

## *Content*

1. Introduction to Coursework
2. Domain Analysis
3. Data Pre-processing
4. Correlation
5. Regression Model and Decision Tree
6. Cross Validation
7. Learning Curve – Underfitting and Overfitting
8. Conclusion

## 1. Introduction to Coursework

The data set provided consists of 53,940 instances with 7 meta-attributes and 4 features in the raw CSV data file named diamonds\_coursework.

The 4 features presenting the data in the raw CSV datasets are Id, cut, color, and clarity.

The 7 meta-attributes that have presented themselves are carat, table, depth, x(length), y(width), z(depth), and price.

This coursework aims to present a domain analysis of the data and further associate the analysis with data cleaning and pre-processing of the datasets. With the understanding of the domain, define a correlation between the features with the target feature. After the correlations have been confirmed, with the help of a regression machine learning method to predict the price of the diamond. Also, the creation of a decision tree methodology and feature engineering to predict the diamond price based on three classes: Low, Medium, and High.

Once the predictions are confirmed, cross-validation of the machine learning pipeline is supposed to be the next step. The application of learning curves and classification evaluation metrics, prove that the pipeline is effective at preventing underfitting as well as overfitting.

## 2. Domain Analysis – Diamonds

It is said that like snowflakes, no two diamonds are alike. Each diamond has its characteristics that draw the potential of the diamond. Prominently, certain standards in features define the quality and in turn the price of the diamond.

These standards are upheld by a certain organization in the diamond industry which goes by the name of the Gemological Institute of America (GIA). According to the dataset provided, it can be classified into different features.

### 1. Id

This dataset allocates an instance of a clear address that helps in the verification of that particular instance. It does not affect the diamond in any way but it helps understand the characteristics of a particular piece in the examination.

### 2. Cut

The cut is never referred to as the diamond's shape. Instead, it provides a reference to the diamond's proportions, symmetry, and polish. The better the diamond is cut, the greater it can reflect or refract light. According to the dataset and GIA, the feature has been divided as:

- 2.1 Ideal    2.2 Premium    2.3 Very good    2.4 Good  
2.5 Fair

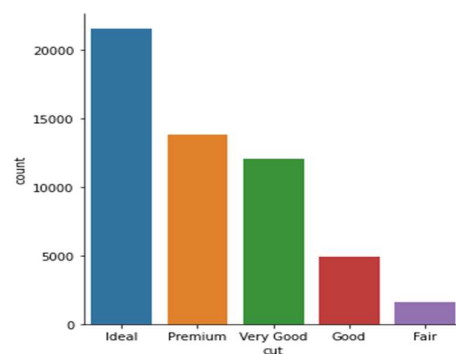


Figure (1). Cut vs Count

The graph shown in Figure (1) shows the frequency of the occurrence of the cut in the dataset provided.

### 3. Color

This feature, as the name suggests, defines the color of the diamond. Against the presumption that diamonds are clear, the fact remains that they can have subtle colors. According to the GIA standards, completely colorless diamonds are considered a rarity which makes them valuable. This feature has also been further categorized into 7 colors.

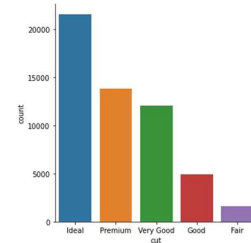


Figure (2). Color vs Count

### 4. Clarity

These diamonds are rocks formed under the ground with somewhat changes they went through which creates another feature that provides data to make up the instances. This feature entails assessing the quantity, size, relief, type, and location of the microscopic characteristics as well as their effect on the overall look of the stone.

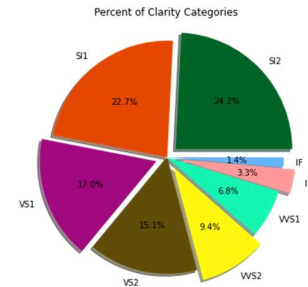


Figure (3). Clarity Category

Figure 3 clearly states the variation of clarity in the dataset.

The meta-attributes also play an essential role in defining a diamond.

According to the dataset provided above, it is clear that there are 7 attributes:

- The first three attributes which describe the dimensions of the diamond are given as X(length), Y(width), and Z(depth).

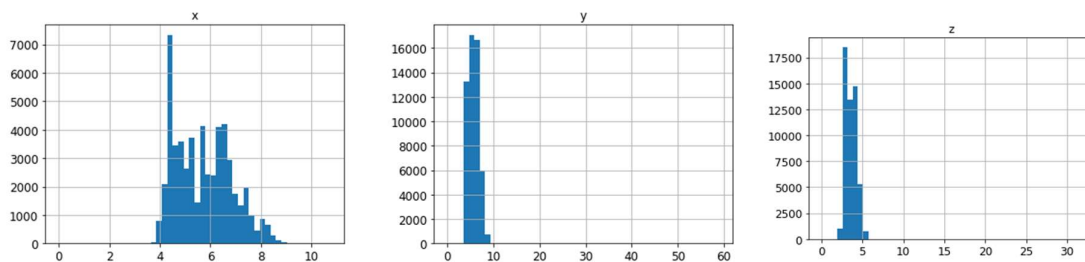


Figure (4). Dimensions Vs Count

The variation of the dimensions with count can be seen in Figure (4). For the above attributes, the range of variation, mean, and standard deviation is given as:

Description	X(length)	Y(Width)	Z(Depth)
Max Variation	10.740000	58.900000	31.800000
Mean	5.732158	5.735530	3.539362
Standard deviation	1.122585	1.143023	0.706241

Table 1

- Carat

This attribute is the most important in defining the diamond as it is the unit of diamond measurement. Carat weight is often confused with the actual weight of the system but on the contrary, it depends on various factors such as density, shape, and formulation of the jewel. Carat weight can never be observed by the naked eye. Figure (5) displays the frequency of occurrence of a variety of carats.

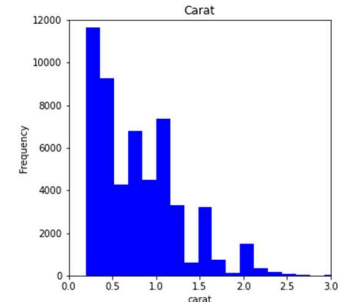


Figure (5). Carat vs Count

Description	Range	Mean	Std Deviation
Carat	5.01 – 0.20	0.798454	0.474411

Table 2

- Depth

The depth of a diamond is defined as the height or the distance from the table to the culet of the diamond. The information at the basic glance on the dataset shows:

Description	Range	Mean	Std Deviation
Depth	95.00 – 44.00	57.457981	2.233967

Table 3

- Table

A diamond's crown that extends from the top of the stone down to the girdle is known as the Table of the diamond. The dataset released the following information:

Description	Range	Mean	Std Deviation
Depth	95.00 – 44.00	57.457981	2.233967

Table 4

- Price

The last attribute but certainly not the least affecting the diamond. It is the market value that defines and differentiates between two diamonds. The instances show the variation amongst them for this attribute which can be tabulated and graphed as:

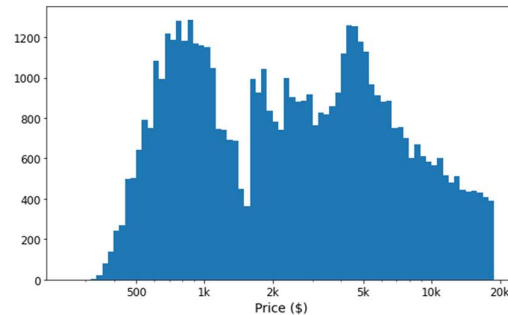


Figure (6). Price vs Count

Description	Range	Mean	Std
Depth	18823-326	3937.77	3993.28

### 3. Data Pre-processing

The dataset generally has certain unwanted or non-interactive data which does not allow the pipeline to work in accordance with the model which is required to work the dataset and features. By the book definition would be, a technique with which a raw data is converted to clean set of data is known as data Pre- processing.

For the dataset provided for this course work, it would be required to go through few steps to achieve a clean dataset. Before proceeding further, it is essential to consider the position the dataset

Name	Type	Role	Values
color	categorical	feature	??, D, E, F, G, H, I, J
clarity	categorical	feature	??, I1, IF, SI1, SI2, VS1, VS2, VVS1, VVS2
Feature 1	numeric	skip	
carat	text	meta	
depth	text	meta	
table	text	meta	
price	text	meta	
x	text	meta	
y	text	meta	
z	text	meta	

Figure (7). Meta-attributes and Features

bound	method	DataFrame.info	id	carat	cut	color	clarity	depth	table	price	x	y	z
0	1	0.23	Ideal	E	SI2	61.5	55.0	326.0	3.95				
1	2	0.21	Premium	E	SI1	59.8	61.0	326.0	3.89				
2	3	0.23	Good	E	VS1	56.9	65.0	327.0	4.05				
3	4	0.29	Premium	I	VS2	62.4	58.0	334.0	4.20				
4	5	0.31	Good	J	SI2	63.3	58.0	335.0	4.34				

Figure (8). Data Description

It is clear from a quick analysis of the dataset that there is a certain missing data which needs to be taken care of as well as unwanted data such as: '??' to be removed. Also, the 7 meta-attributes affect the diamond characteristics, hence are required to be converted to features to be included in the pipeline. For this to work in a pipeline the data set must be pre-processed first.

- Editing Domain

This is done to convert the meta-attributes into features so that they can be included in the pipeline to confirm their effect on the target feature. In orange tool, it can be done with the help of Edit domain widget. As in Figure (9), it will now be reinterpreted as numeric value instead of text.

Variable	Current Domain	New Domain
id	text	text
cut	categorical	categorical
color	categorical	categorical
clarity	categorical	categorical
carat	text	numeric (reinterpreted as numeric)
depth	text	numeric (reinterpreted as numeric)
table	text	numeric (reinterpreted as numeric)
price	text	numeric (reinterpreted as numeric)
x	text	numeric (reinterpreted as numeric)
y	text	numeric (reinterpreted as numeric)
z	text	numeric (reinterpreted as numeric)

Figure (9). Edit Domain

This step can be skipped while working for python but the next applies for both orange and python.

- Imputing

This step is considered the base for the pipeline because it consists of cleaning of data to move ahead. If slight mistake or wanted data is removed, it would necessarily affect the decision making further down the pipeline.

It is achieved by choosing the impute widget and selecting all the column to dropping the unwanted or missing data. In this case, it was necessary to remove the unknown values.

For the imputing to work, the concerned features are selected to first drop. After confirming the selection of the features, action that is needed to be taken is selected. For instance, in this case, it is

required to remove instances with unknown values and that is what is selected in the impute widget in the orange tool.

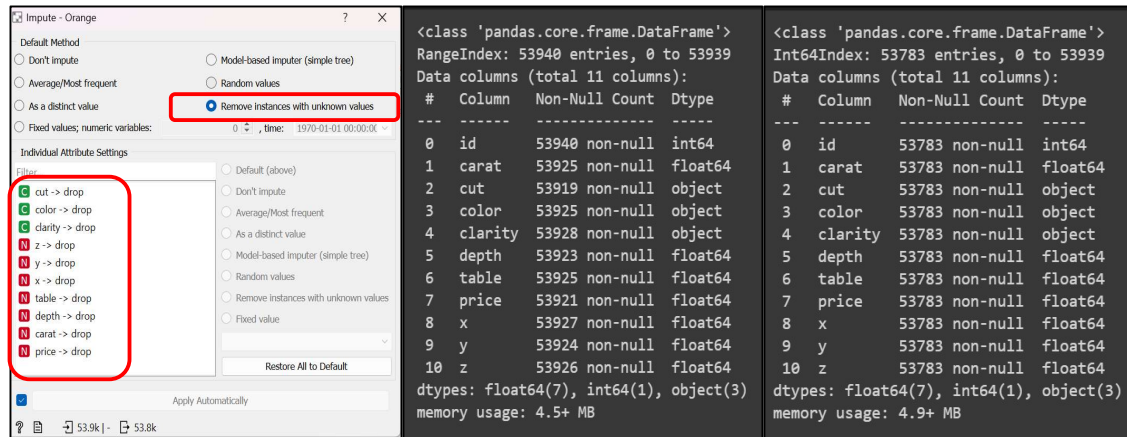


Figure (10). Imputing data, and Effect of Imputing on data

With respect to Figure (10), it clearly states that the non-null count has drop from different values to all the features having equal count of 53783.

#### 4. Correlation

Machine learning models can be classified as good or based on the data that has been provided. The selection features contribute most to the quality of resulting model. The selection of such features that helps in predicting the variable more accurately is known as Feature selection.

Feature selection is crucial task which is based on configuring the relation between the available features. With help of such relationships, one can identify the features important based on the target model.

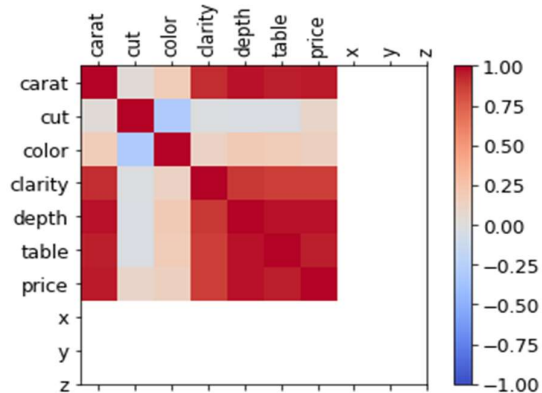


Figure (11) Correlation Plot

	carat	depth	table	price	x	y	z
carat	1.000000	0.028246	0.181732	0.921578	0.975098	0.951694	0.953359
depth	0.028246	1.000000	-0.295991	-0.010723	-0.025231	-0.029280	0.094952
table	0.181732	-0.295991	1.000000	0.127152	0.195527	0.183920	0.151075
price	0.921578	-0.010723	0.127152	1.000000	0.884465	0.865419	0.861240
x	0.975098	-0.025231	0.195527	0.884465	1.000000	0.974669	0.970746
y	0.951694	-0.029280	0.183920	0.865419	0.974669	1.000000	0.951956
z	0.953359	0.094952	0.151075	0.861240	0.970746	0.951956	1.000000

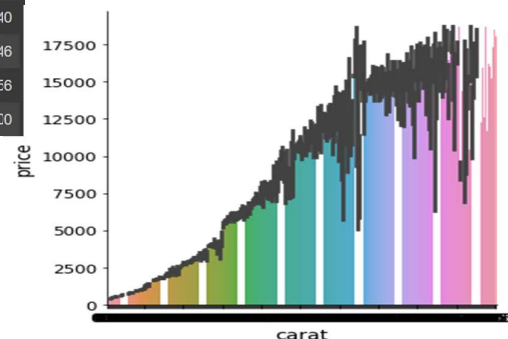
Figure (12). Correlation Table

is clearly understandable which features are having an impact where and how are they affecting the model and target

Figure (13). Price vs Carat

In case of Diamonds dataset, the correlations between features are given out in Fig (12) and Figure (11).

From the Correlation plot and the table, it







As per the domain analysis and the study of the correlation from both python and orange. It can be stated that the features that have strong relation with each other are:

### 1. X(length)-Y(Width)-Z(Depth)

As per Figure (12) and (16), the correlational values from the widget and the table suggest strong influence over each other, and also the domain analysis states that the dimension of diamond always have an impact over each other.

### 2. Price – Carat

According to the domain analysis and standard guidelines from the GIA, they have a strong co-dependency upon each other.

In accordance with the above discussion and the domain analysis, the three variables that closely correlate with the target price column are: **Carat, Cut & Clarity.**

## 5. Regression Model and Decision Tree

- Decision Tree Classifier

A Decision tree methodology is a data mining technique for establishing classification systems based on multiple covariates or for developing a prediction algorithm for a target variable. For a decision tree classifier to act on the target variable, it is first required to make the necessary changes to the data by the help of feature engineering.

### ➤ Feature Engineering

A process of creating a feature consisting the necessary requirements depending on the target to be applicable for a classifier. In this case, the target feature must be converted to classification of 3: Low, Medium and High.

It is also observed that converting the categorical features to numeric features will also affect the accuracy of the classifier as the data refines itself. Also, from the correlational analysis, it was observed that features X, Y, Z can be combined as well. Hence, the new creations are done with the help of feature constructor widget.

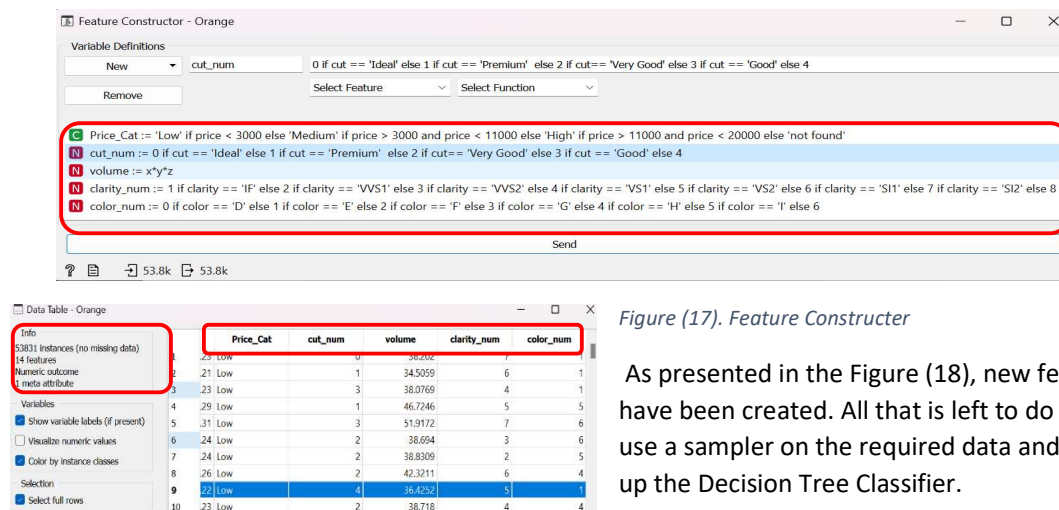


Figure (17). Feature Constructor

As presented in the Figure (18), new features have been created. All that is left to do is to use a sampler on the required data and take up the Decision Tree Classifier.

Figure (18). Data For Decision tree

Once the necessary columns are selected, a sample data and test data is created to configure the system for classifier. In this case, a ratio of 70:30 is taken for the current pipeline.

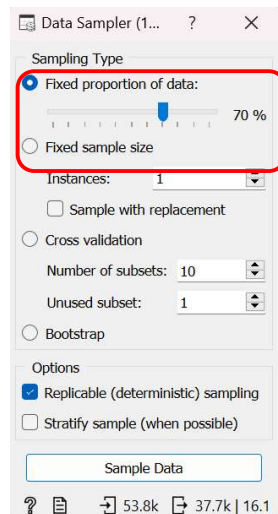


Figure (20). Data Sampler

Once the data sampler is ready for the pipeline, next step is to use this sample data for Tree Classifier and the rest remaining Data is used for testing the prediction. The pipeline for the classifier is given as:

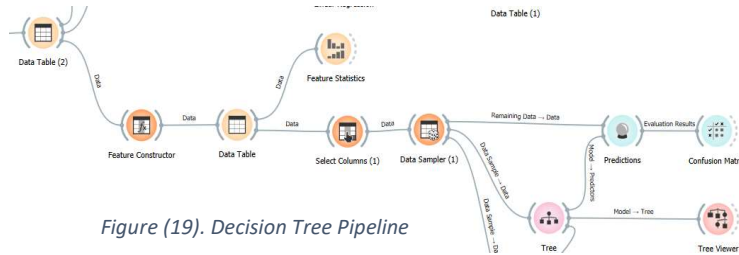


Figure (19). Decision Tree Pipeline

From the above pipeline, the predictions are passed through the confusion matrix to make the sure of underfitting and over fitting for the classification.

The confusion matrix provides an accuracy to the predictions on the system so as to validate the pipeline.

	Tree	error	Price_Cat	Selected (1)
1	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
2	0.03 : 0.03 : 0.94 : 0.00 → Medi...	0.061	Medium	No
3	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
4	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
5	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
6	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
7	0.92 : 0.00 : 0.08 : 0.00 → High	0.076	High	No
8	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
9	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
10	0.92 : 0.00 : 0.08 : 0.00 → High	0.076	High	No
11	0.08 : 0.00 : 0.92 : 0.00 → Medi...	0.080	Medium	No
12	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
13	0.03 : 0.03 : 0.94 : 0.00 → Medi...	0.061	Medium	No
14	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
15	0.03 : 0.03 : 0.94 : 0.00 → Medi...	0.061	Medium	No
16	0.03 : 0.03 : 0.94 : 0.00 → Medi...	0.061	Medium	No
17	0.98 : 0.00 : 0.02 : 0.00 → High	0.015	High	No
18	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
19	0.03 : 0.03 : 0.94 : 0.00 → Medi...	0.061	Medium	No
20	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
21	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
22	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
23	0.00 : 0.94 : 0.06 : 0.00 → Low	0.063	Low	No
24	0.03 : 0.03 : 0.94 : 0.00 → Medi...	0.061	Medium	No

Figure (21). Predictions

The Prediction shows the individual instance's predicted classification based on the pipeline into the three categories created under the feature price\_cat with presumptuous error.

Also, the precision for the classifier can also be seen in the same widget's window. As per the pipeline created, the system precision is given as 0.932.

The prediction only classifies but it is not nearly as capable to visualize the data. Hence a confusion matrix is used to give out the exact instances for the classification in terms of numbers as well as percentage.

		Predicted				Σ
		High	Low	Medium	not found	
Actual	High	90.8 %	0.0 %	4.2 %	NA	1341
	Low	0.0 %	93.5 %	2.5 %	NA	8974
	Medium	9.2 %	6.5 %	93.3 %	NA	5832
	not found	0.0 %	0.0 %	0.0 %	actual: Medium predicted: Medium	2
Σ		1223	9456	5470	0	16149

Figure (22). Confusion Matrix

The study of confusion matrix can be easily suggest out that accuracy of the pipeline and gives a better visualization of the dataset classification.

### ➤ Regression model

Looking at the domain analysis and the correlational study of the dataset, it would be better to use multiple regression and cross validate them to get an accurate result for the system/pipeline. But it is not required to use the feature engineering done on the pipeline.

For the regression model to be applicable, the pipeline would look like:

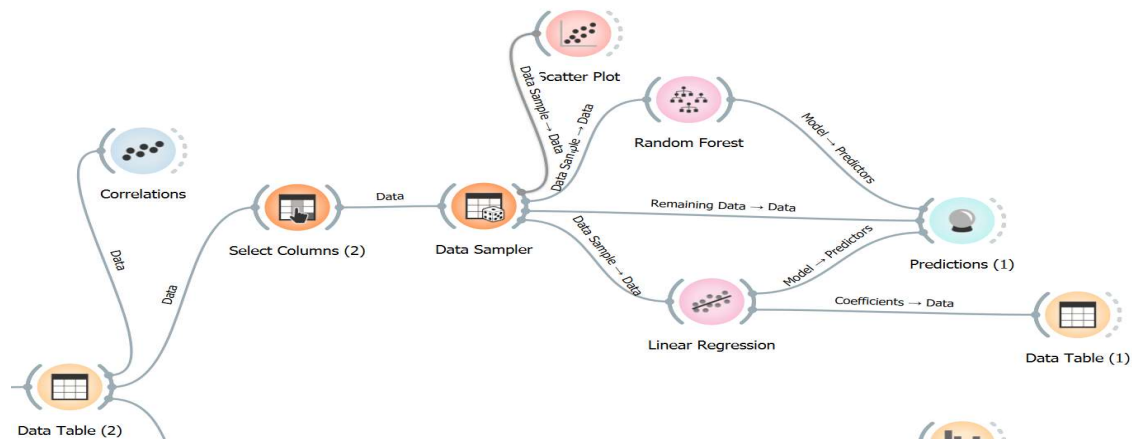


Figure (23). Regression model pipeline

As per the above pipeline, selecting the column that have an impact on the regression model based domain analysis and correlational study is the first step. Also, it could be assessed through the scatter plot between features. An example would be between price and carat.

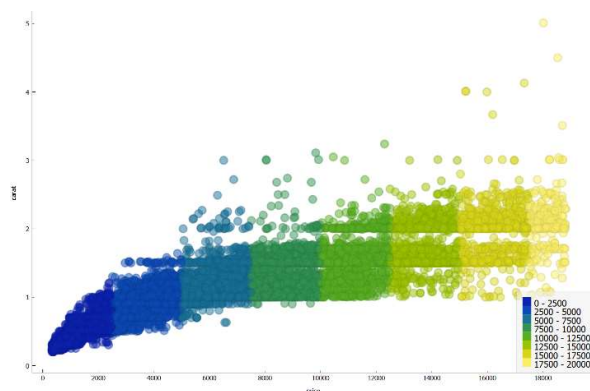


Figure (24). Scatter plot: Price vs Carat

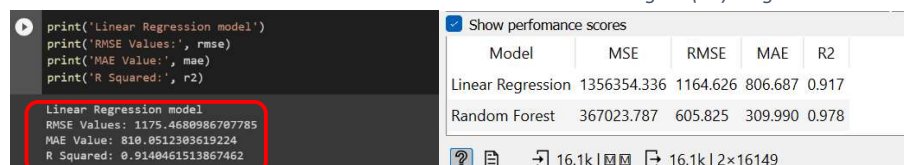
Through the scatter plot, it can be understood the variation with different carat and price.

Once that is accomplished, all that is left is to split the data into test and train dataset and to run the regression models through the prediction to figure out which one works better. The result of the predictions are visible in the Prediction widget.

The prediction stated that the random forest regression model is underfit based on the learning curve. The regression model for the linear regression seems to be a perfect fit as the validation curve and training curve work as parallel as plotted in graph. The values from orange tools are shown as:

Whereas the model for python shows the result as follows:

Figure (25). Regression Model result



From the above regression model result it is clearly visible that the Linear regression model is much more effective for the dataset based on diamonds.

## 6. Cross Validation

For a Machine learning model, a pipeline is created based on the requirement. But to actually verify the accuracy and the legitimacy of the model, cross validation is done. Essentially, it is method that evaluates and compares the learning algorithm by dividing data in two segments: one to train the model and the other one is used to test the predictive capabilities if the model such as accuracy, error, etc.

For this case, there are two methods, one is by using the orange tool's data sampler widget, wherein the subsets of the data changes to verify the pipeline.

And by changing the model, number of subset it, whether the accuracy changes or not.

The second method is by import the cross validation score from python. Which show the result for linear regression and random forest regression as:

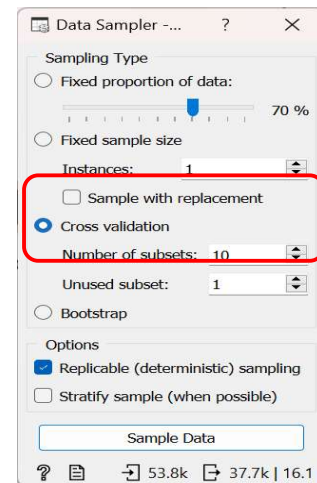


Figure (26). Validation for data Sampler

```
[45] def display_scores(scores):
      print("Scores:", scores)
      print("Mean:", scores.mean())
      print("Standard deviation:", scores.std())

[44] print('Cross validation scores for Linear regression')
      display_scores(regr_score)
```

Cross validation scores for Linear regression  
 Scores: [1138.81871203 1201.07914264 1180.57571676 1138.39881208 1102.63720905  
 1156.62652647 1257.9870992 1181.72355533 1212.50199404 1187.06669983]  
 Mean: 1175.741546744031  
 Standard deviation: 41.72825375545398

Figure (27). Cross Validation scores for Linear regression

The validation helps prevent bias and plays a crucial role in maintaining the dataset required for the system to work. It does have an impact on the learning model as the test and train dataset must always be keep apart and have the exact ratio as mentioned during the time sampling.

```
print('Cross validation scores for Random forest regression')
display_scores(forest_score)
```

Cross validation scores for Random forest regression  
 Scores: [743.1966229 750.29364409 699.78655939 720.98037864 718.25987483  
 724.71219581 806.26840526 722.02042915 718.28974275 727.73755988]  
 Mean: 733.1545412699371  
 Standard deviation: 27.693429668246637

Figure (28). Cross Validation for Random Forest regressor

## 7. Learning Curves – Underfitting and Overfitting

Learning curves are the plots that represent the performance and the impact of the model at the training data set increases. The two major causes for error in the pipeline presented are:

Bias: it describes the model so that it makes simpler assumptions such that, for the model, it is easier to approximate.

Variance: It describes the variability of the model's prediction with the change in dataset used for training the model.

What the learning curves helps in representing the exact regressor to be used for prediction as the system requires a good fit.

For this case, 3 regressors were used and trained on the same data sets. The regressors used were Linear regression, Random Forest regression and Support Vector Machine regression. The plots are shown below:

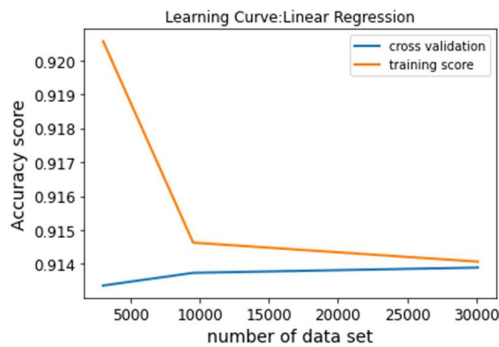


Figure (29). Learning curve for Linear regression

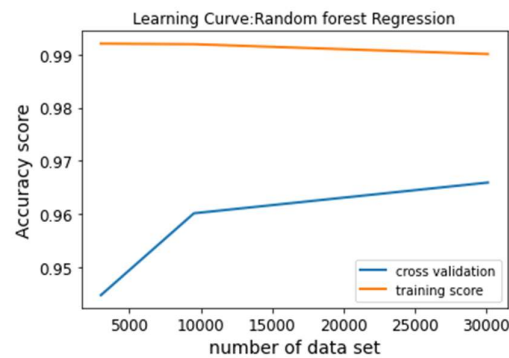


Figure (30). Learning curve for Random Forest regression

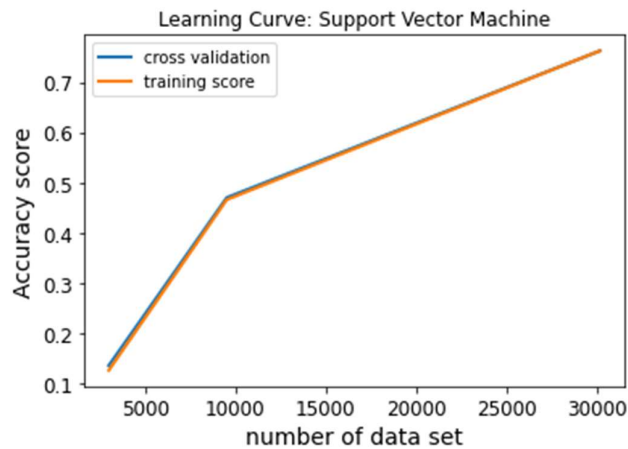


Figure (31). Learning curve for Support Vector Machine Regression

The reason for choosing the linear regression as fit for the machine learning model can be easily deduced from the above plots for all the regression models. If the training curve at any given point coincides with the validation, it is said to be an overfit.

If the training curve and the validation curve are running almost parallel to each other, it is considered as a good fit.

If the training curve and the validation curve are very far from each other it is considered to be an underfit case.

From the above plots and the knowledge based on the learning curve it can be justified that Linear regression is best suited for the Machine learning Model.

## 8. Conclusion

This course work consisted of creation of a complete pipeline for a data set based on study of 11 features of diamond. The process started with the domain analysis of the dataset and Diamond. Once the domain analysis was complete. It gave enough information to construct a pipeline based on the features for prediction of prices of diamond and classifying them in three categories: low, Medium and High.

This process started with editing the domain to make the meta-attribute into features. After the domain was complete. Next was data cleaning of unwanted or missing data instances. This was done to make the system completely based on complete datasets. With clean set of data, correlational study was done to achieve the particular features affecting the target feature price. It was noticed that X, Y, Z and depth was have been neglected as it made no impact on the target feature.

From here on, model was divided into two directions: First was creation a model that predict the price of a diamond. Second was the creation of a predictor that classified the prices of the diamonds into three categories.

For the first part of the pipeline, a sampler was taken and the data was divided into a ratio of 70:30. The bigger set was used to train the model and the smaller set to test the regression models in place. It was through cross validation and learning curve found that the Linear regression model was a good fit for the machine learning model.

For the second part, feature engineer was done to categorise the class of the price and then through the sampler, a decision tree classification was done to achieve the confusion matrix that explain the prediction with error as well. Also, the pipeline showed a ROC analysis to confirm that the tree classifier was the correct choice.

With the completion of a machine learning model, the price was predicted and the classification was done which helped in creating an effecting ML Model.