

# UNIVERSITI MALAYA

**Master of Data Science (Semester 1 – 2022/2023)**

**Faculty of Computer Science & Information Technology**

**WQD7005 Data Mining**

**Group Project**

**Project Title: House Price Prediction**

**Instructor: Prof Dr Teh Ying Wah**

Name	Matric Number
Jasmeen Kah Ying Bong	S2142739
Hii Yew Han	S2037987
Dinesh V M Ramachandran	S2119167
Aw Yeong Fung Mun	17197465
Lee Ziteng Nicholas	S2132376

## Table of Contents

1	Introduction.....	1
2	Our Dataset .....	2
2.1	Table Information.....	2
2.2	Column Metadata .....	4
3	Business Understanding.....	6
3.1	Analysis Goal .....	6
3.2	Analysis Data .....	6
4	Methodology.....	7
4.1	SEMMA Process .....	7
4.2	SEMMA Description.....	8
5	Results.....	9
5.1	SAMPLE – Accessing and Assaying Prepared Data .....	9
5.2	EXPLORE – Exploring Data Science.....	13
5.2.1	Univariate Analysis.....	13
5.2.2	Bivariate Analysis.....	26
5.2.3	Multivariate Analysis.....	31
5.2.4	Interesting Visualizations.....	33
5.3	MODIFY – Data Modification.....	40
5.3.1	Modifying and Correcting Source Data.....	40
5.3.2	Examining Exported Data.....	48
5.3.3	Creating Training and Validation Data.....	50
5.4	MODEL – Data Modelling .....	52
5.4.1	Constructing a Decision Tree Predictive Model.....	52
5.4.2	Constructing a Gradient Boosting Predictive Model .....	56
5.4.3	Constructing a Logistics Regression Predictive Model .....	57
5.4.4	Constructing a Neural Network Predictive Model.....	59

5.5 ASSESS – Data Assessment .....	63
6 Conclusion .....	65
7 Appendix.....	67
7.1 SAMPLE .....	67
7.1.1 Create a SAS Enterprise Miner Project .....	67
7.1.2 Create a SAS Enterprise Miner Diagram.....	70
7.1.3 Create a SAS Enterprise Miner Library.....	78
7.2 EXPLORE.....	85
7.2.1 Create Histogram, Pie Chart and Boxplot.....	85
7.2.2 Add Missing Bin and Changing Graph Properties - Number of Bin in Histogram	
86	
7.2.3 Display Variable Association .....	88
7.2.4 Create Summary Statistics .....	88
7.2.5 Perform Variable Clustering .....	89
7.2.6 Perform Variable Selection.....	90
7.2.7 Perform Variable Correlation .....	92
7.2.8 Interesting Visualizations.....	94
7.2.9 Correlation Table .....	96
7.3 MODIFY .....	101
7.3.1 Modify Inconsistent Data.....	101
7.3.2 Modify Noisy Data and Incomplete Data .....	102
7.3.3 Data Partition .....	104
7.4 MODEL.....	105
7.4.1 Decision Tree Model.....	105
7.4.2 Gradient Boosting Model.....	106
7.4.3 Logistics Regression Model.....	107
7.4.4 Neural Network Model .....	109

7.5 ASSESS.....	111
7.5.1 Compare between Decision Tree, Logistic Regression, and Gradient Boosting 111	
7.5.2 Compare Neural Network.....	111
7.5.3 Compare the Best Model .....	112

## 1 Introduction

The fluctuation of house prices is a contentious issue because it can have a significant impact on the economy as a whole. An increase in house prices means an increase in non-financial assets, which in turn increases personal wealth and stimulates household consumption, boosting the economy. However, a decrease in house prices limits an individual's borrowing capacity, causing investments to decline due to the decrease in the value of collateral. The 2008 housing bubble is a perfect example of the importance of stable, predictable house prices. Unexpected changes in house prices can lead to an increase in real long-term interest rates, bankruptcies in financial institutions, and global economic depression. While it may be difficult to control house prices, it is possible to predict them.

Therefore, in this report, we used a dataset from Kaggle covering house sales from May 2014 through May 2015 at King County, which is in the U.S. state of Washington. According to the data gathered by the 2020 United States census, King County has the highest estimated population of 2,269,675 among all counties in Washington and is also the 13th-most populous county in the United States. With a higher population, King County has more potential house buyers and higher house demands; thus, house price data in King County will be more complete and more precise. In this dataset, it includes categorical data and numerical data, we can perform multi-faceted analysis and prediction on the dataset.

## 2 Our Dataset

### 2.1 Table Information

Data Source	<a href="https://www.kaggle.com/datasets/harlfoxem/housesalesprediction">https://www.kaggle.com/datasets/harlfoxem/housesalesprediction</a>		
Data Name	kc_house_data.csv		
File Size	2.52 MB		
Year	2016		
Dimension	21613 rows, 21 columns		
Columns	Attributes	Types	Description
	id	Numerical	A notation for a house
	date	Numerical	Date house was sold
	price	Numerical	Price of each home sold
	bedrooms	Numerical	Number of bedrooms
	bathrooms	Numerical	Number of bathrooms
	sqft_living	Numerical	Square footage of the apartments interior living space
	sqft_lot	Numerical	Square footage of the land space
	floors	Numerical	Number of floors
	waterfront	Categorical	Overlooking the waterfront or not: • 1 = Yes • 0 = No
	view	Numerical	An index from 0 to 4 of how good the view of the property was
	condition	Numerical	An index from 1 to 5 on the condition of the property
	grade	Numerical	An index from 1 to 13 based on King County grading system
	sqft_above	Numerical	The square footage of the interior housing space that is above ground level
	sqft_basement	Numerical	The square footage of the interior housing space that is below ground level

	yr_built	Numerical	Year when the house was initially built
	yr_renovated	Numerical	Year when house was renovated
	zipcode	Numerical	5-digit zip code area the house is in
	lat	Numerical	Latitude
	long	Numerical	Longitude
	sqft_living15	Numerical	The square footage of interior housing living space for the nearest 15 neighbours
	sqft_lot15	Numerical	The square footage of the land lots of the nearest 15 neighbours

## 2.2 Column Metadata

This is our metadata; we have a total of 21 variables. They are divided into two types of data, namely numeric data, and character data. The following is a detailed display of the metadata:

Column	Model Role	Measurement Level	Description
bathrooms	Input	Interval	Number of bathrooms
bedrooms	Input	Interval	Number of bedrooms
condition	Input	Ordinal	Condition of the property
date	Time ID	Nominal	Date house was sold
floors	Input	Interval	Number of floors
grade	Input	Ordinal	Grade level of the property
id	ID	Interval	A notation for a house
lat	Input	Interval	Latitude
long	Input	Interval	Longitude
price	Target	Interval	Price of each home sold
sqft_above	Input	Interval	The square footage of the interior housing space that is above ground level
sqft_basement	Input	Interval	The square footage of the interior housing space that is below ground level
sqft_living	Input	Interval	Square footage of the apartments interior living space
sqft_living15	Input	Interval	The average square footage of interior housing living space for the nearest 15 neighbours
sqft_lot	Input	Interval	Square footage of the land space
sqft_lot15	Input	Interval	The average square footage of the land lots of the nearest 15 neighbours
view	Input	Ordinal	View index of the property
waterfront	Input	Binary	House which has a view to a waterfront
yr_built	Input	Interval	Year when the house was initially built
yr_renovated	Input	Interval	Year when house was renovated
zipcode	Input	Interval	5-digit zip code area the house is in

Data Source Wizard -- Step 5 of 8 Column Metadata

Name	Role	Level	Report	Drop	Type	Format	Informat	Length	Format Type
bathrooms	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
bedrooms	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
condition	Input	Ordinal	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
date	Time ID	Nominal	No	No	Character	\$15.	\$15.	15	CATEGORY
floors	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
grade	Input	Ordinal	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
id	ID	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
lat	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
long	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
price	Target	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
sqft_above	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
sqft_basement	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
sqft_living	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
sqft_living15	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
sqft_lot	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
sqft_lot15	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
view	Input	Ordinal	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
waterfront	Input	Binary	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
yr_builtin	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
yr_renovated	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY
zipcode	Input	Interval	No	No	Numeric	BEST12.0	BEST32.0	8	QUANTITY

Show code   Explore   Refresh Summary   < Back   Next >   Cancel

Figure 2.1 Column Metadata

## 3 Business Understanding

### 3.1 Analysis Goal

We wish to predict house prices. Therefore, our objectives are:

1. To find features that are relevant in predicting house prices.
2. To build a predictive model for house prices.

### 3.2 Analysis Data

This dataset contains house sales sold between May 2014 and May 2015 for King County, which includes Seattle. It contains essential attributes of a house including but not limited to bedrooms, bathrooms, and square foot size. This dataset is available publicly on Kaggle. The target variable is the price of the house sold.

## 4 Methodology

The methodology we adopted for this report is SEMMA, which stands for Sampling, Exploring, Modifying, Modelling, and Assessing. Developed by the SAS Institute as a process for data mining large amounts of data to uncover previously unknown patterns which can be utilized as a business advantage. In this report, the first two items are mainly discussed, namely Sample and Explore.

### 4.1 SEMMA Process



Figure 4.1 SEMMA Process

## 4.2 SEMMA Description

### Sample

- In this step, a subset of the relevant volume dataset is selected from a sizable dataset provided for the model's building. Identification of variables or factors (both dependent and independent) impacting the process is the aim of the first stage of the process. The gathered data is then divided into categories for preparation and validation.

### Explore

- Univariate, bivariate and multivariate analysis are carried out in this step, in order to investigate interrelated relationships between variables and to find data gaps. The univariate analysis examines each factor separately to comprehend its role in the broader scheme, whereas the bivariate and multivariate analysis explores the link between variables. With a focus on data visualisation, all the influencing factors that might have an impact on the study's conclusion are examined.

### Modify

- In this step, business logic is used to derive the lessons discovered during the exploration phase from the data gathered during the sample phase. In other words, the data is analysed to determine whether it needs to be refined or transformed before moving on to the modelling stage.

### Model

- After the variables have been clarified and the data have been cleaned, the modelling step employs a variety of data mining techniques to create a projected model of how this data achieves the process's final, desired result.

### Assess

- In the SEMMA process' last stage, the model's applicability and dependability to the subject under study are assessed. Now that the data has been tested, it can be used to determine how effectively it performs.

## 5 Results

### 5.1 SAMPLE – Accessing and Assaying Prepared Data

The dataset we acquired is a comma-separated values (CSV) file, in order to make use of it in SAS Enterprise Miner, we need a compatible SAS file, which is an .sas7bdat file. This is done by running the File Import tool under the Sample Tab of the SEMMA tools palette and connecting it to the Save Data tool under the Utility Tab. Subsequently, running Save Data will produce an object file in a SAS compatible format that can be imported in the data source.

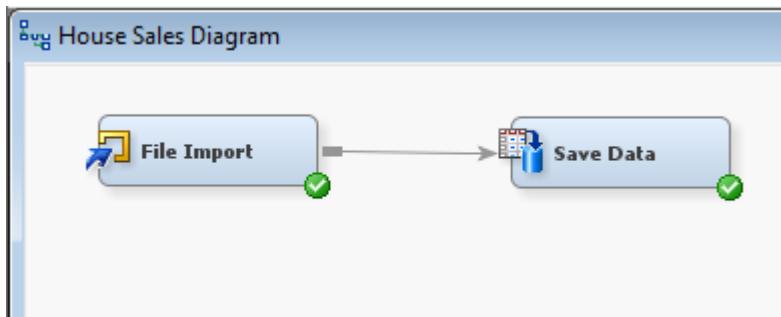


Figure 5.1 Convert Data Format

The dataset we have chosen is the data on house sales in King County, USA. In the following data, the variables related to the house are the ID, price of the house, number of bedrooms, number of bathrooms, square foot of living space, square foot of whole lot, number of floors, if there is a waterfront, view, condition, grade, square foot above and square foot basement. The variables related to time are dates of the house being sold, built year, and renovated year. Whereas the variables related to the location are zip code, latitude, and longitude of the house. In addition, additional information is provided such as, the average square footage of living space and whole lot for the nearest 15 neighbours.

In this house sales dataset, we have a total of 21 variables. They are divided into two types of data, namely numeric data, and character data. The following is a detailed display of the data:

Name	Role	Level	Report	Order	Drop	Lower	Upper	Type	Format	Informat	Length
bathrooms	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
bedrooms	Input	Nominal	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
condition	Input	Nominal	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
date	Rejected	Nominal	No		No	.	.	Character	\$15.	\$15.	15
floors	Input	Nominal	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
grade	Input	Nominal	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
id	ID	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
lat	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
long	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
price	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
sqft_above	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
sqft_basement	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
sqft_living	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
sqft_living15	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
sqft_lot	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
sqft_lot15	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
view	Input	Nominal	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
waterfront	Input	Binary	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
yr_builtin	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
yr_renovated	Input	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8
zipcode	Rejected	Interval	No		No	.	.	Numeric	BEST12.0	BEST32.0	8

Figure 5.2 Details of Data

After the data set is imported into SAS Enterprise Miner, we can see that the system will automatically set the measurement level for the metadata according to the attributes of the variables. By default, numeric data with only 1's and 0's is set to binary level, other numeric data is set to the interval level, and character data is set to the nominal level. In addition, there are some attributes which roles have been set rejected because SAS Enterprise Miner does not automatically recognize the format.

Metadata Completed.		
<b>Library:</b>	GROUP5A	
<b>Data Source:</b>	EM_SAVE_TRAIN	
<b>Role:</b>	Raw	
<b>Role</b>	<b>Level</b>	<b>Count</b>
ID	Interval	1
Input	Binary	1
Input	Interval	12
Input	Nominal	5
Rejected	Interval	1
Rejected	Nominal	1

Figure 5.3 Result of Basic Setting

Our assessment shows that the default metadata set by SAS Enterprise Miner is not meeting our requirements. Therefore, we need to make manual changes to the roles and levels of the attributes. For levels, we have changed bedrooms, condition, floors, grade, and view to levels that reflect the data such as intervals and ordinals. For roles, we have changed date and zipcode from rejected to TimeID and Input respectively. Finally, price has been changed from input to target since it is a target variable. The following is a comparison of the changes that we have made:



Name	Role	Level
bathrooms	Input	Interval
bedrooms	Input	Nominal
condition	Input	Nominal
date	Rejected	Nominal
floors	Input	Nominal
grade	Input	Nominal
id	ID	Interval
lat	Input	Interval
long	Input	Interval
price	Input	Interval
sqft_above	Input	Interval
sqft_basement	Input	Interval
sqft_living	Input	Interval
sqft_living15	Input	Interval
sqft_lot	Input	Interval
sqft_lot15	Input	Interval
view	Input	Nominal
waterfront	Input	Binary
yr_builtin	Input	Interval
yr_renovated	Input	Interval
zipcode	Rejected	Interval

Name	Role	Level
bathrooms	Input	Interval
bedrooms	Input	Interval
condition	Input	Ordinal
date	TimeID	Nominal
floors	Input	Interval
grade	Input	Ordinal
id	ID	Interval
lat	Input	Interval
long	Input	Interval
price	Target	Interval
sqft_above	Input	Interval
sqft_basement	Input	Interval
sqft_living	Input	Interval
sqft_living15	Input	Interval
sqft_lot	Input	Interval
sqft_lot15	Input	Interval
view	Input	Ordinal
waterfront	Input	Binary
yr_builtin	Input	Interval
yr_renovated	Input	Interval
zipcode	Input	Interval

Figure 5.4 Setting Comparison

In these metadata, we have four measurement levels, they are interval, binary, ordinal and nominal attributes. The following is the updated content based on our understanding of the data sources, and it is also the sample result that we will use in the end:

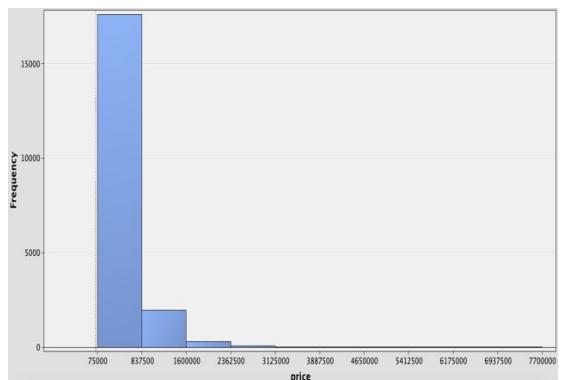
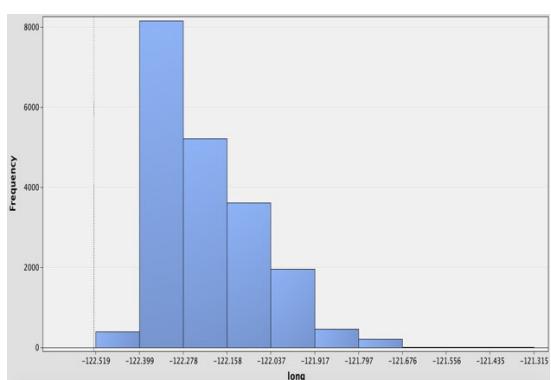
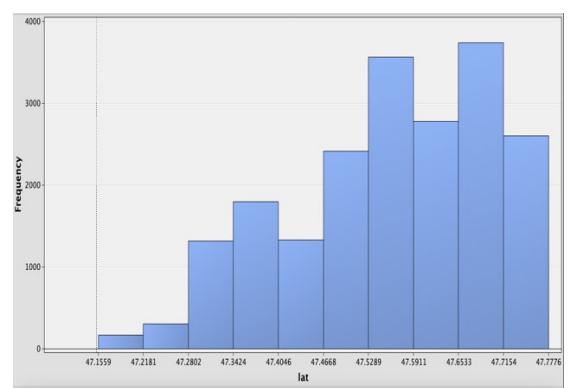
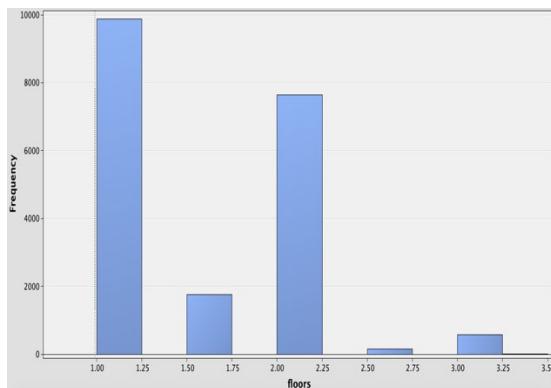
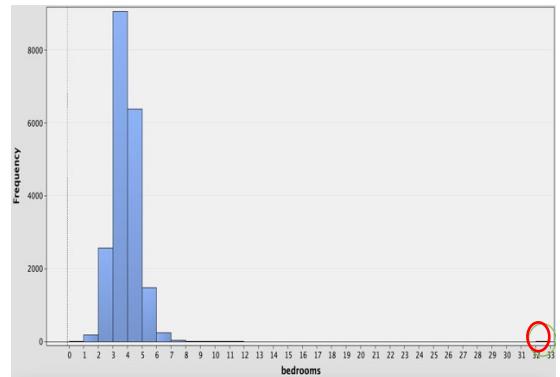
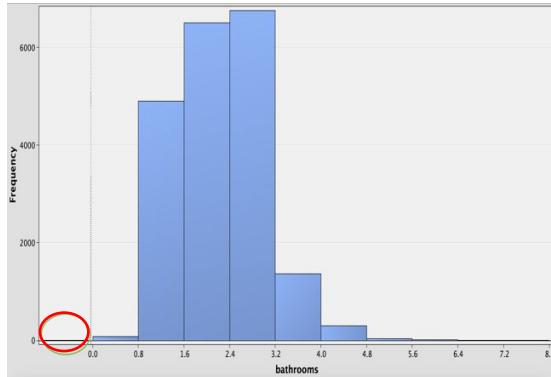
Metadata Completed.		
Library:	GROUP5A	
Data Source:	EM_SAVE_TRAIN	
Role:	Raw	
Role	Level	Count
ID	Interval	1
Input	Binary	1
Input	Interval	14
Input	Ordinal	3
Target	Interval	1
Time ID	Nominal	1

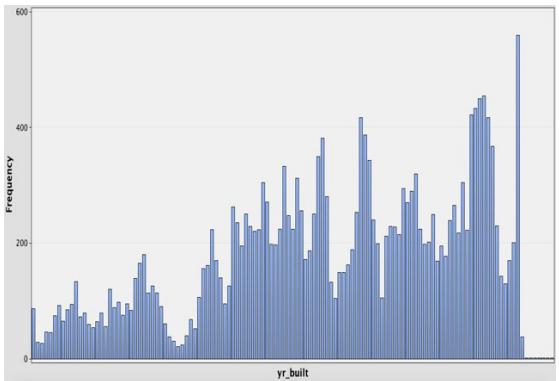
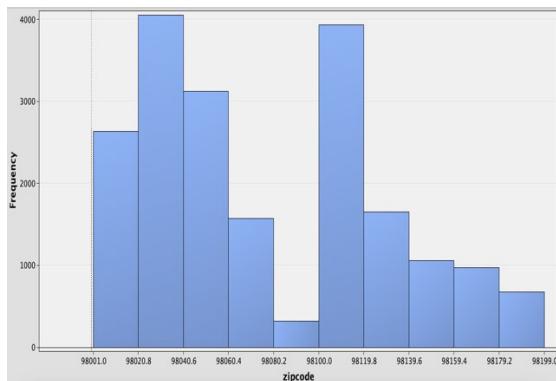
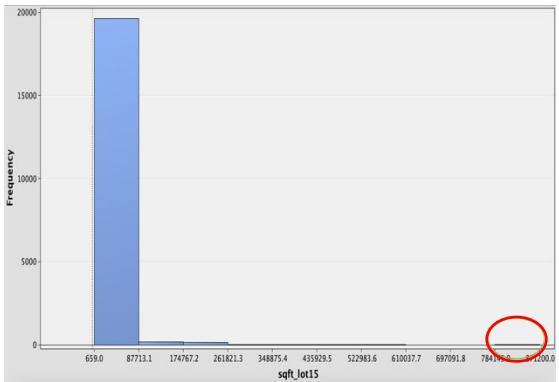
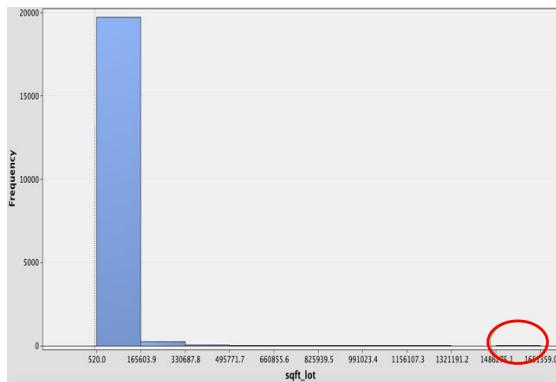
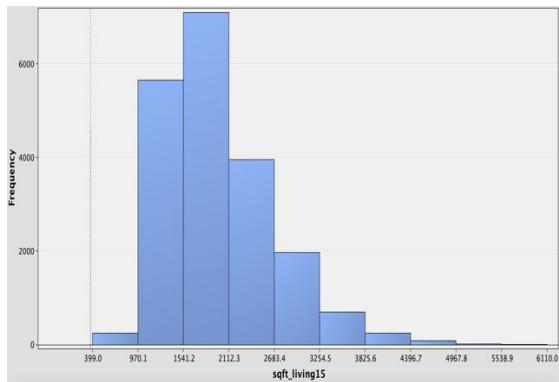
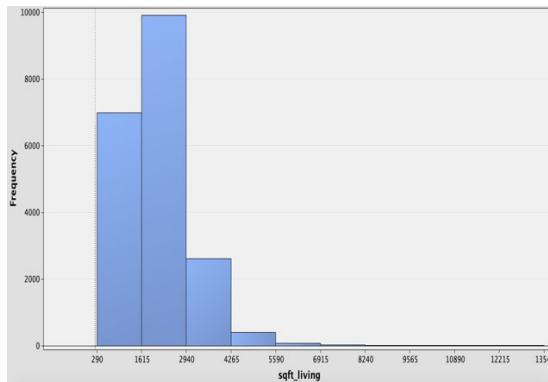
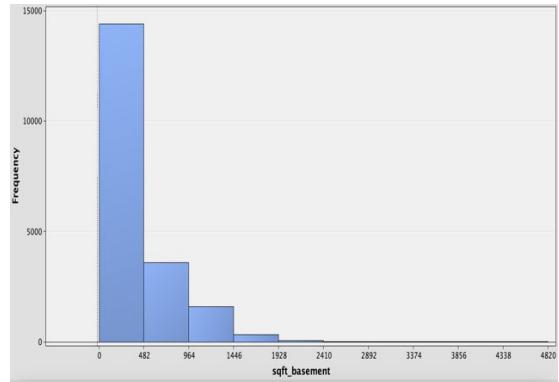
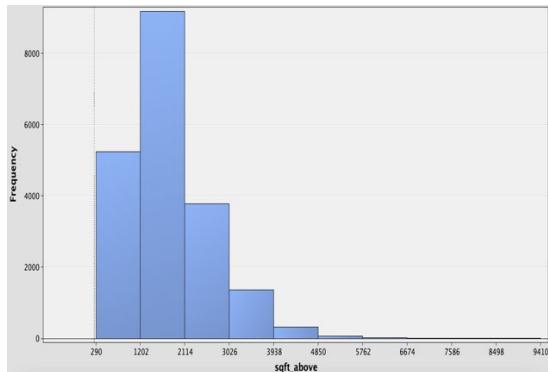
Figure 5.5 Result of Advanced Setting

## 5.2 EXPLORE – Exploring Data Science

### 5.2.1 Univariate Analysis

A descriptive analysis that uses just one variable is called a univariate analysis. This analysis's objectives are to provide us with the initial impression of each variable and to look for incomplete, noisy, and inconsistent data.





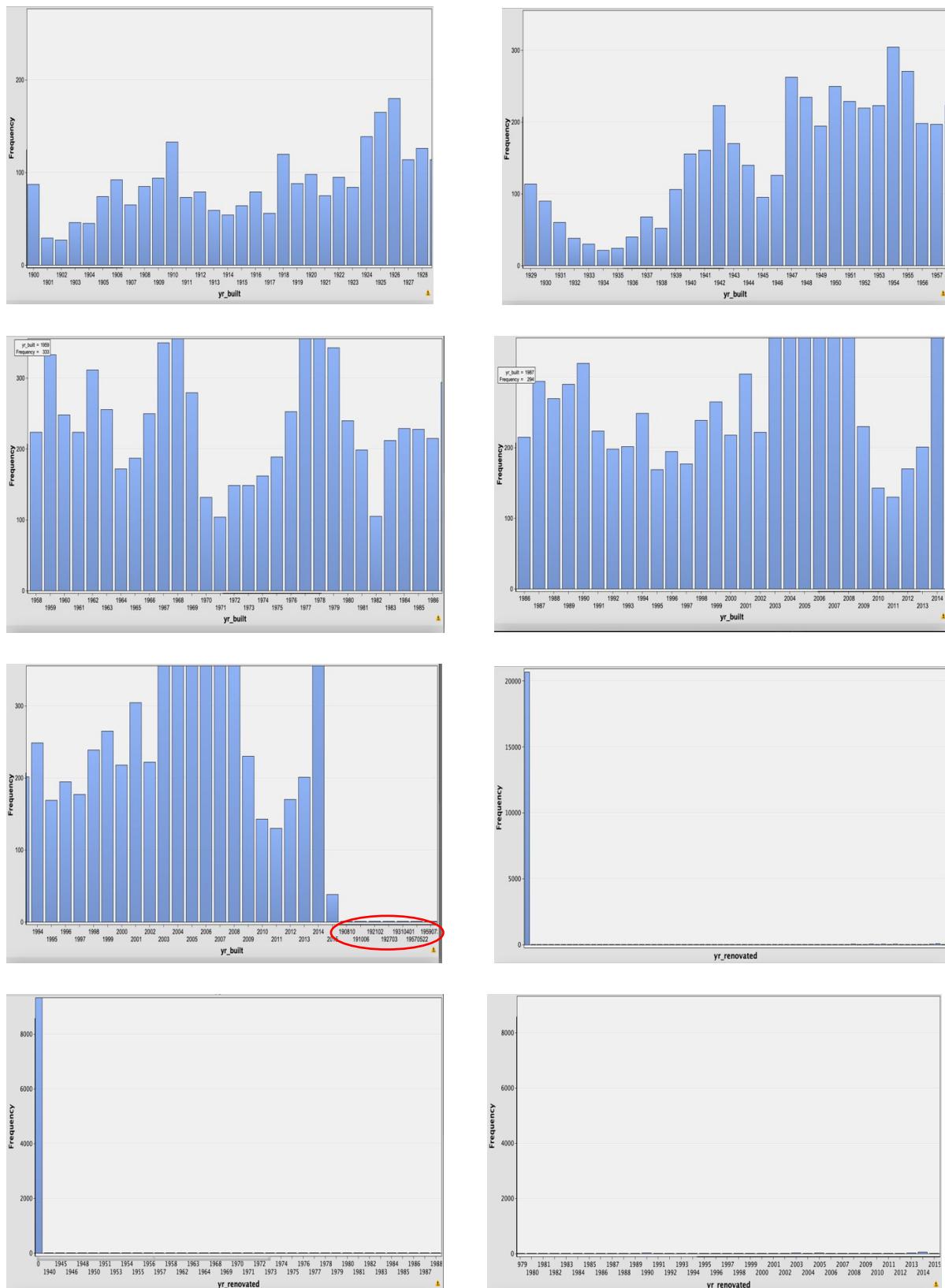


Figure 5.6 Histogram

We use histogram to show the distribution of interval variables and nominal variables. The above shows the histogram of different interval variable. We observed that the bathrooms variable has min value of 0 and max value of 8. The mode is the 4<sup>th</sup> bin ranges between 2.4 to 3.2 with around 6000 observations in this range. We found an error from the histogram. There is an incomplete data error in the bathrooms variable as highlighted in the red circle showing the missing bin in the bathrooms histogram. This may be due to failure recording the value due to human error. We observed that bedrooms variable has min value of 0 and max value of 33. The mode is the 4<sup>th</sup> bin ranges between 3 to 4 with around 8000 observations in this range. We found an error from the histogram. There is an outlier in the bedrooms variable as highlighted in the red circle. We observed that floors variable has min value of 1 to max value of 3.5. The mode is the 1<sup>st</sup> bin ranges between 1 to 1.25 with around 10000 observations. We observed that lat variable has min value of 47.1559 to max value of 47.7776. The mode is the 9<sup>th</sup> bin ranges between 47.6533 to 47.7154 with around 4000 observations. We observed that long variable has min value of -122.519 and max value of -121.315. The mode is the 2<sup>nd</sup> bin ranges between -122.399 to -122.278 with around 8000 observations. We observed that price variable has min value of 75000 and max value of 7700000. The mode is the 1st bin ranges between 75000 to 837500 with around 15000 observations. We observed that sqft\_above variable has min value of 290 and max value of 9410. The mode is the 2<sup>nd</sup> bin ranges between 1202 to 2114 with around 8000 observations. We observed that sqft\_basement variable has min value of 0 and max value of 4820. The mode is the 1<sup>st</sup> bin ranges between 0 to 482 with around 15000 observations. We observed that sqft\_living variable has min value of 290 and max value of 13540. The mode is the 2<sup>nd</sup> bin ranges between 1615 to 2940 with around 10000 observations. We observed that sqft\_living15 variable has min value of 399 and max value of 6210. The mode is 3<sup>rd</sup> bin ranges between 1541.2 to 2112.3 with around 6000 observations. We observed that sqft\_lot has min value of 520 and max value of 1651359. The mode is 1<sup>st</sup> bin ranges between 520 to 165603.9 with around 20000 observations. We found an error from the histogram. There is an outlier in the sqft\_lot variable as highlighted in the red circle. We observed that sqft\_lot15 variable has min value of 651 and max value of 871200. The mode is 1<sup>st</sup> bin ranges between 659 to 87713.3 with around 20000 observations. We found an error from the histogram. There is an outlier in the sqft\_lot15 variable as highlighted in the red circle. We observed that zipcode variable has min value of 98001 to max value of 98199. The mode is 2<sup>nd</sup> bin ranges between 98020.8 to 98040.6 with around 4000 observations. We observed that yr\_built has a min value of 1900 and max value of 2015. The mode is 2014 which has the highest frequency of 559 with around 600 observations. We found an error from the histogram.

There is an inconsistency data error in the yr\_built variable as highlighted in the red circle, there are eight values, which are 190810, 191006, 192102, 192703, 19310401, 19570522 and 19590731 which are supposed to be record in years but are recorded together with months and dates due to human error. We observed that yr\_renovated has a min value of 0 and max value of 2015. The mode is 0 which has the highest frequency of 20699 with around 20000 observations.

Ordered Inputs	Data Role	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Abs C.V.	Coefficient of Variation	Sign
1TRAIN	soft_lot	7617	0	21613	520	1651359	15106.97	74420.51	13.06002	285.0778	INPUT	soft_lot	2.741815	2.41815+		
2TRAIN	soft_lot15	7620	0	21613	651	871200	12768.46	527304.189	506743	150.7631	INPUT	soft_lot15	2.138409	1.138409		
3TRAIN	soft_basement	0	0	21613	0	4520	291.509	442.575	2.715574	719.4579	INPUT	soft_bas...	1.518221	1.518221+		
4TRAIN	price	450000	0	21613	75000	7700000	540182.2	367362.24	0.21716	34.52244	INPUT	price	0.680071	0.680071+		
5TRAIN	soft_above	1560	0	21613	290	9410	1788.391	828.091	1.446664	3.402304	INPUT	soft_ab...	0.463037	0.463037+		
6TRAIN	soft_living	1910	0	21613	290	13540	918.4409	1.471555	5.243093	715.5243	INPUT	soft_living	0.441579	0.441579+		
7TRAIN	bathrooms	2.25	5	21608	0	8.2	114807.0	0.770164	0.511142	1.280413	INPUT	bathroo...	0.364177	0.364177+		
8TRAIN	floors	1.5	0	21613	1	3.5	1494309	0.539989	0.616177	-0.48472	INPUT	floors	0.361363	0.361364+		
9TRAIN	soft_living15	1840	0	21613	398	624	1986.532	687.393	1.179743	1.19743	48.5265	INPUT	soft_liv...	0.451515	0.451515+	
10TRAIN	bathrooms	3	0	21613	0	33	12768.46	12768.46	1.093662	1.093662	1.093662	INPUT	bathrooms	0.24913	0.24913	
11TRAIN	lat	47.5718	0	21613	47.1559	47.7776	47.56005	0.138564	-0.48527	-0.6769	INPUT	lat	0.002913	0.002913+		
12TRAIN	long	-122.231	0	21613	-122.519	-121.315	-122.214	0.140828	0.885027	0.959911	INPUT	long	0.001152	-0.001152		
13TRAIN	zipcode	98065	0	21613	98001	98199	98077.94	53.50503	0.405497	-0.33233	INPUT	zipcode	.0005455	.0005455+		

Figure 5.7 Summary Statistics for Interval Variables

The above summary statistics for interval variables are shown to support the visualization.

From the summary statistics, there is 5 missing values in bathrooms as highlighted in the red circle.

Data Role	Variable Name	Level	CODE	Frequency Count	Type	Percent	Level Index	Role	Label	Plot
TRAIN	yr_built	2014	34	559N		2.586406	115	INPUT	yr_built	1
TRAIN	yr_built	2006	45	45N		2.106288	109	INPUT	yr_built	1
TRAIN	yr_built	2007	26	450N		2.08208	106	INPUT	yr_built	1
TRAIN	yr_built	2004	73	433N		2.003424	105	INPUT	yr_built	1
TRAIN	yr_built	2003	9	422N		1.952529	104	INPUT	yr_built	1
TRAIN	yr_built	1977	12	417N		1.929394	78	INPUT	yr_built	1
TRAIN	yr_built	2012	75	416N		1.894494	109	INPUT	yr_built	1
TRAIN	yr_built	1978	64	587N		1.790589	79	INPUT	yr_built	1
TRAIN	yr_built	1968	20	381N		1.762282	69	INPUT	yr_built	1
TRAIN	yr_built	2008	40	367N		1.698052	109	INPUT	yr_built	1
TRAIN	yr_built	1967	59	350N		1.619396	68	INPUT	yr_built	1
TRAIN	yr_built	1979	14	343N		1.588108	80	INPUT	yr_built	1
TRAIN	yr_built	1959	36	333N		1.40739	60	INPUT	yr_built	1
TRAIN	yr_built	1990	70	320N		1.48059	91	INPUT	yr_built	1
TRAIN	yr_built	1962	56	312N		1.443576	63	INPUT	yr_built	1
TRAIN	yr_built	1954	42	305N		1.411188	55	INPUT	yr_built	1
TRAIN	yr_built	2010	5	294N		1.360292	109	INPUT	yr_built	1
TRAIN	yr_built	1987	4	294N		1.360292	88	INPUT	yr_built	1
TRAIN	yr_built	1989	45	290N		1.341785	90	INPUT	yr_built	1
TRAIN	yr_built	1969	18	280N		1.295517	70	INPUT	yr_built	1
TRAIN	yr_built	1955	0	272N		1.253875	56	INPUT	yr_built	1
TRAIN	yr_built	1978	55	270N		1.242648	89	INPUT	yr_built	1
TRAIN	yr_built	1999	78	265N		1.226114	100	INPUT	yr_built	1
TRAIN	yr_built	1947	19	263N		1.21686	48	INPUT	yr_built	1
TRAIN	yr_built	1963	7	256N		1.184472	64	INPUT	yr_built	1
TRAIN	yr_built	1976	76	253N		1.170592	77	INPUT	yr_built	1
TRAIN	yr_built	1960	39	250N		1.156711	51	INPUT	yr_built	1
TRAIN	yr_built	1966	37	250N		1.156711	67	INPUT	yr_built	1
TRAIN	yr_built	1994	15	249N		1.152084	95	INPUT	yr_built	1
TRAIN	yr_built	1960	8	248N		1.147458	61	INPUT	yr_built	1
TRAIN	yr_built	1980	61	244N		1.147443	81	INPUT	yr_built	1
TRAIN	yr_built	1978	95	239N		1.105116	99	INPUT	yr_built	1
TRAIN	yr_built	1948	25	235N		1.087309	49	INPUT	yr_built	1
TRAIN	yr_built	2009	107	230N		1.064174	110	INPUT	yr_built	1
TRAIN	yr_built	1951	33	229N		1.059547	52	INPUT	yr_built	1
TRAIN	yr_built	1984	21	228N		1.059547	85	INPUT	yr_built	1
TRAIN	yr_built	1953	21	228N		1.059547	86	INPUT	yr_built	1
TRAIN	yr_built	1958	99	224N		1.036413	59	INPUT	yr_built	1
TRAIN	yr_built	1961	53	224N		1.036413	62	INPUT	yr_built	1
TRAIN	yr_built	1991	41	224N		1.036413	92	INPUT	yr_built	1
TRAIN	yr_built	1942	10	223N		1.031786	43	INPUT	yr_built	1
TRAIN	yr_built	1953	38	223N		1.031786	54	INPUT	yr_built	1

Data Role	Variable Name	Level	CODE	Frequency Count	Type	Percent	Level Index	Role	Label	Plot
TRAIN	yr_built	2002	49	222N		1.02716	103	INPUT	yr_built	1
TRAIN	yr_built	1952	52	201N		1.01626	55	INPUT	yr_built	1
TRAIN	yr_built	2000	32	218N		1.008652	101	INPUT	yr_built	1
TRAIN	yr_built	1986	47	215N		0.994772	87	INPUT	yr_built	1
TRAIN	yr_built	1983	63	212N		0.980891	84	INPUT	yr_built	1
TRAIN	yr_built	1973	80	203N		0.97253	94	INPUT	yr_built	1
TRAIN	yr_built	2013	97	201N		0.929996	114	INPUT	yr_built	1
TRAIN	yr_built	1981	28	199N		0.920742	82	INPUT	yr_built	1
TRAIN	yr_built	1956	48	198N		0.916115	57	INPUT	yr_built	1
TRAIN	yr_built	1953	50	198N		0.911515	91	INPUT	yr_built	1
TRAIN	yr_built	1957	84	197N		0.911488	58	INPUT	yr_built	1
TRAIN	yr_built	1949	77	195N		0.902235	50	INPUT	yr_built	1
TRAIN	yr_built	1996	31	195N		0.902235	97	INPUT	yr_built	1
TRAIN	yr_built	1975	60	189N		0.892744	76	INPUT	yr_built	1
TRAIN	yr_built	1955	3	188N		0.865232	66	INPUT	yr_built	1
TRAIN	yr_built	1926	72	180N		0.832832	27	INPUT	yr_built	1
TRAIN	yr_built	1997	82	177N		0.818952	98	INPUT	yr_built	1
TRAIN	yr_built	1964	51	172N		0.795817	65	INPUT	yr_built	1
TRAIN	yr_built	1974	83	170N		0.765644	44	INPUT	yr_built	1
TRAIN	yr_built	2012	100	170N		0.786564	113	INPUT	yr_built	1
TRAIN	yr_built	1995	6	169N		0.781937	96	INPUT	yr_built	1
TRAIN	yr_built	1925	44	165N		0.763429	26	INPUT	yr_built	1
TRAIN	yr_built	1974	89	162N		0.747149	75	INPUT	yr_built	1
TRAIN	yr_built	1941	22	161N		0.744922	42	INPUT	yr_built	1
TRAIN	yr_built	1940	85	156N		0.721788	41	INPUT	yr_built	1
TRAIN	yr_built	1972	46	149N		0.68494	73	INPUT	yr_built	1
TRAIN	yr_built	1973	43	149N		0.68494	74	INPUT	yr_built	1
TRAIN	yr_built	2010	67	143N		0.661639	111	INPUT	yr_built	1
TRAIN	yr_built	1944	105	140N		0.647758	45	INPUT	yr_built	1
TRAIN	yr_built	1924	69	139N		0.643131	25	INPUT	yr_built	1
TRAIN	yr_built	1971	62	133N		0.610747	110	INPUT	yr_built	1
TRAIN	yr_built	1970	109	133N		0.610744	71	INPUT	yr_built	1
TRAIN	yr_built	2011	102	130N		0.60149	112	INPUT	yr_built	1
TRAIN	yr_built	1928	87	126N		0.582982	29	INPUT	yr_built	1
TRAIN	yr_built	1946	58	126N		0.582982	47	INPUT	yr_built	1
TRAIN	yr_built	1918	80	126N		0.552121	19	INPUT	yr_built	1
TRAIN	yr_built	1927	11	114N		0.52746	28	INPUT	yr_built	1
TRAIN	yr_built	1929	27	114N		0.52746	30	INPUT	yr_built	1
TRAIN	yr_built	1939	57	106N		0.490446	40	INPUT	yr_built	1
TRAIN	yr_built	1952	92	105N		0.481319	83	INPUT	yr_built	1
TRAIN	yr_built	1971	66	104N		0.48119	72	INPUT	yr_built	1
TRAIN	yr_built	1920	81	98N		0.453431	21	INPUT	yr_built	1

Data Role	Variable Name ▼	Level	CODE	Frequency Count	Type	Percent	Level Index	Role	Label	Plot
TRAIN	yr_built	1922		35	95N	0.43955	23INPUT	yr_built		1
TRAIN	yr_built	1945		68	95N	0.43955	46INPUT	yr_built		1
TRAIN	yr_built	1909		24	94N	0.434923	10INPUT	yr_built		1
TRAIN	yr_built	1906		114	92N	0.42567	7INPUT	yr_built		1
TRAIN	yr_built	1930		29	90N	0.416416	31INPUT	yr_built		1
TRAIN	yr_built	1919		113	88N	0.407162	20INPUT	yr_built		1
TRAIN	yr_built	1900		13	87N	0.402536	1INPUT	yr_built		1
TRAIN	yr_built	1908		93	85N	0.393282	9INPUT	yr_built		1
TRAIN	yr_built	1923		74	84N	0.388655	24INPUT	yr_built		1
TRAIN	yr_built	1912		101	79N	0.365521	13INPUT	yr_built		1
TRAIN	yr_built	1916		16	79N	0.365521	17INPUT	yr_built		1
TRAIN	yr_built	1921		17	75N	0.347013	22INPUT	yr_built		1
TRAIN	yr_built	1905		65	74N	0.342387	6INPUT	yr_built		1
TRAIN	yr_built	1911		89	73N	0.33776	12INPUT	yr_built		1
TRAIN	yr_built	1937		91	68N	0.314625	38INPUT	yr_built		1
TRAIN	yr_built	1907		98	65N	0.300745	8INPUT	yr_built		1
TRAIN	yr_built	1915		23	64N	0.296118	16INPUT	yr_built		1
TRAIN	yr_built	1931		94	60N	0.277611	32INPUT	yr_built		1
TRAIN	yr_built	1913		96	59N	0.272984	14INPUT	yr_built		1
TRAIN	yr_built	1917		103	56N	0.259103	18INPUT	yr_built		1
TRAIN	yr_built	1914		71	54N	0.24985	15INPUT	yr_built		1
TRAIN	yr_built	1938		112	52N	0.240596	39INPUT	yr_built		1
TRAIN	yr_built	1903		108	46N	0.212835	4INPUT	yr_built		1
TRAIN	yr_built	1904		30	45N	0.088208	5INPUT	yr_built		1
TRAIN	yr_built	1936		90	40N	0.185074	37INPUT	yr_built		1
TRAIN	yr_built	192		104	38N	0.17582	33INPUT	yr_built		1
TRAIN	yr_built	2015		110	38N	0.17582	116INPUT	yr_built		1
TRAIN	yr_built	1933		2	30N	0.038895	34INPUT	yr_built		1
TRAIN	yr_built	1901		79	29N	0.134179	2INPUT	yr_built		1
TRAIN	yr_built	1902		106	27N	0.124925	3INPUT	yr_built		1
TRAIN	yr_built	1935		115	24N	0.111044	36INPUT	yr_built		1
TRAIN	yr_built	1934		111	21N	0.097164	35INPUT	yr_built		1
TRAIN	yr_built	190810		118	1N	0.004627	117INPUT	yr_built		1
TRAIN	yr_built	191006		122	1N	0.004627	118INPUT	yr_built		1
TRAIN	yr_built	192102		116	1N	0.004627	119INPUT	yr_built		1
TRAIN	yr_built	192703		120	1N	0.004627	120INPUT	yr_built		1
TRAIN	yr_built	19310401		119	1N	0.004627	121INPUT	yr_built		1
TRAIN	yr_built	19570522		117	1N	0.004627	122INPUT	yr_built		1
TRAIN	yr_built	19590731		121	1N	0.004627	123INPUT	yr_built		1

Data Role	Variable Name ▲	Level	CODE	Frequency Count	Type	Percent	Level Index	Role	Label	Plot
TRAIN	yr_renovated	0		0	20699N	0.9577106	61INPUT	yr_renovated		1
TRAIN	yr_renovated	2014		14	91N	0.421043	69INPUT	yr_renovated		1
TRAIN	yr_renovated	2003		6	37N	0.101393	68INPUT	yr_renovated		1
TRAIN	yr_renovated	2000		11	36N	0.166566	58INPUT	yr_renovated		1
TRAIN	yr_renovated	2005		25	35N	0.16194	55INPUT	yr_renovated		1
TRAIN	yr_renovated	2007		9	35N	0.16194	60INPUT	yr_renovated		1
TRAIN	yr_renovated	2004		32	35N	0.16194	62INPUT	yr_renovated		1
TRAIN	yr_renovated	1990		29	26N	0.120298	59INPUT	yr_renovated		1
TRAIN	yr_renovated	2006		19	25N	0.115671	45INPUT	yr_renovated		1
TRAIN	yr_renovated	1989		35	24N	0.111044	61INPUT	yr_renovated		1
TRAIN	yr_renovated	1987		28	22N	0.101791	44INPUT	yr_renovated		1
TRAIN	yr_renovated	2002		2	22N	0.101791	57INPUT	yr_renovated		1
TRAIN	yr_renovated	2009		31	22N	0.101791	64INPUT	yr_renovated		1
TRAIN	yr_renovated	1993		1	20N	0.092337	46INPUT	yr_renovated		1
TRAIN	yr_renovated	1993		50	19N	0.08791	48INPUT	yr_renovated		1
TRAIN	yr_renovated	1994		7	19N	0.08791	49INPUT	yr_renovated		1
TRAIN	yr_renovated	1998		26	19N	0.08791	53INPUT	yr_renovated		1
TRAIN	yr_renovated	2001		37	19N	0.08791	56INPUT	yr_renovated		1
TRAIN	yr_renovated	1983		17	18N	0.083283	38INPUT	yr_renovated		1
TRAIN	yr_renovated	1984		12	18N	0.083283	39INPUT	yr_renovated		1
TRAIN	yr_renovated	1987		33	18N	0.083283	42INPUT	yr_renovated		1
TRAIN	yr_renovated	2008		10	18N	0.083283	63INPUT	yr_renovated		1
TRAIN	yr_renovated	2010		3	18N	0.083283	65INPUT	yr_renovated		1
TRAIN	yr_renovated	1985		36	17N	0.078656	40INPUT	yr_renovated		1
TRAIN	yr_renovated	1986		30	17N	0.078656	41INPUT	yr_renovated		1
TRAIN	yr_renovated	1992		5	17N	0.078656	47INPUT	yr_renovated		1
TRAIN	yr_renovated	1999		4	17N	0.078656	54INPUT	yr_renovated		1
TRAIN	yr_renovated	1995		24	16N	0.07403	50INPUT	yr_renovated		1
TRAIN	yr_renovated	2015		45	16N	0.07403	70INPUT	yr_renovated		1
TRAIN	yr_renovated	1988		20	15N	0.069403	43INPUT	yr_renovated		1
TRAIN	yr_renovated	1996		52	15N	0.069403	51INPUT	yr_renovated		1
TRAIN	yr_renovated	1997		41	15N	0.069403	52INPUT	yr_renovated		1
TRAIN	yr_renovated	2011		15	13N	0.060149	66INPUT	yr_renovated		1
TRAIN	yr_renovated	1980		38	11N	0.050895	35INPUT	yr_renovated		1
TRAIN	yr_renovated	1982		56	11N	0.050895	37INPUT	yr_renovated		1
TRAIN	yr_renovated	1985		47	11N	0.050895	57INPUT	yr_renovated		1
TRAIN	yr_renovated	1979		40	10N	0.046268	34INPUT	yr_renovated		1
TRAIN	yr_renovated	1970		27	9N	0.041642	25INPUT	yr_renovated		1
TRAIN	yr_renovated	1968		46	8N	0.037015	23INPUT	yr_renovated		1
TRAIN	yr_renovated	1977		22	8N	0.037015	32INPUT	yr_renovated		1
TRAIN	yr_renovated	1975		61	6N	0.027761	30INPUT	yr_renovated		1
TRAIN	yr_renovated	1978		8	6N	0.027761	33INPUT	yr_renovated		1

Data Role	Variable Name ▲	Level	CODE	Frequency Count	Type	Percent	Level Index	Role	Label	Plot
TRAIN	yr_renovated	1958		62	5N	0.023134	15INPUT	yr_renovated		1
TRAIN	yr_renovated	1964		63	5N	0.023134	20INPUT	yr_renovated		1
TRAIN	yr_renovated	1965		67	5N	0.023134	21INPUT	yr_renovated		1
TRAIN	yr_renovated	1973		34	5N	0.023134	28INPUT	yr_renovated		1
TRAIN	yr_renovated	1981		23	5N	0.023134	36INPUT	yr_renovated		1
TRAIN	yr_renovated	1960		65	4N	0.018507	17INPUT	yr_renovated		1
TRAIN	yr_renovated	1963		48	4N	0.018507	19INPUT	yr_renovated		1
TRAIN	yr_renovated	1969		43	4N	0.018507	24INPUT	yr_renovated		1
TRAIN	yr_renovated	1972		53	4N	0.018507	27INPUT	yr_renovated		1
TRAIN	yr_renovated	1945		18	3N	0.013881	5INPUT	yr_renovated		1
TRAIN	yr_renovated	1953		54	3N	0.013881	10INPUT	yr_renovated		1
TRAIN	yr_renovated	1955		55	3N	0.013881	12INPUT	yr_renovated		1
TRAIN	yr_renovated	1956		57	3N	0.013881	13INPUT	yr_renovated		1
TRAIN	yr_renovated	1957		21	3N	0.013881	14INPUT	yr_renovated		1
TRAIN	yr_renovated	1974		16	3N	0.013881	29INPUT	yr_renovated		1
TRAIN	yr_renovated	1976		59	3N	0.013881	31INPUT	yr_renovated		1
TRAIN	yr_renovated	1940		58	2N	0.009254	3INPUT	yr_renovated		1
TRAIN	yr_renovated	1946		60	2N	0.009254	6INPUT	yr_renovated		1
TRAIN	yr_renovated	1950		42	2N	0.009254	8INPUT	yr_renovated		1
TRAIN	yr_renovated	1962		51	2N	0.009254	18INPUT	yr_renovated		1
TRAIN	yr_renovated	1967		66	2N	0.009254	22INPUT	yr_renovated		1
TRAIN	yr_renovated	1971		39	2N	0.009254	26INPUT	yr_renovated		1
TRAIN	yr_renovated	1934		68	1N	0.004627	2INPUT	yr_renovated		1
TRAIN	yr_renovated	1944		69	1N	0.004627	4INPUT	yr_renovated		1
TRAIN	yr_renovated	1948		44	1N	0.004627	7INPUT	yr_renovated		1
TRAIN	yr_renovated	1951		49	1N	0.004627	9INPUT	yr_renovated		1
TRAIN	yr_renovated	1954		13	1N	0.004627	11INPUT	yr_renovated		1
TRAIN	yr_renovated	1959		64	1N	0.004627	16INPUT	yr_renovated		1

Figure 5.8 Summary Statistics for Nominal Variables

The above summary statistics for nominal variables are shown to support the visualization.

There is an inconsistency data error in the yr\_built variable as highlighted in the red circle, which are 190810, 191006, 192102, 192703, 19310401, 19570522 and 19590731.

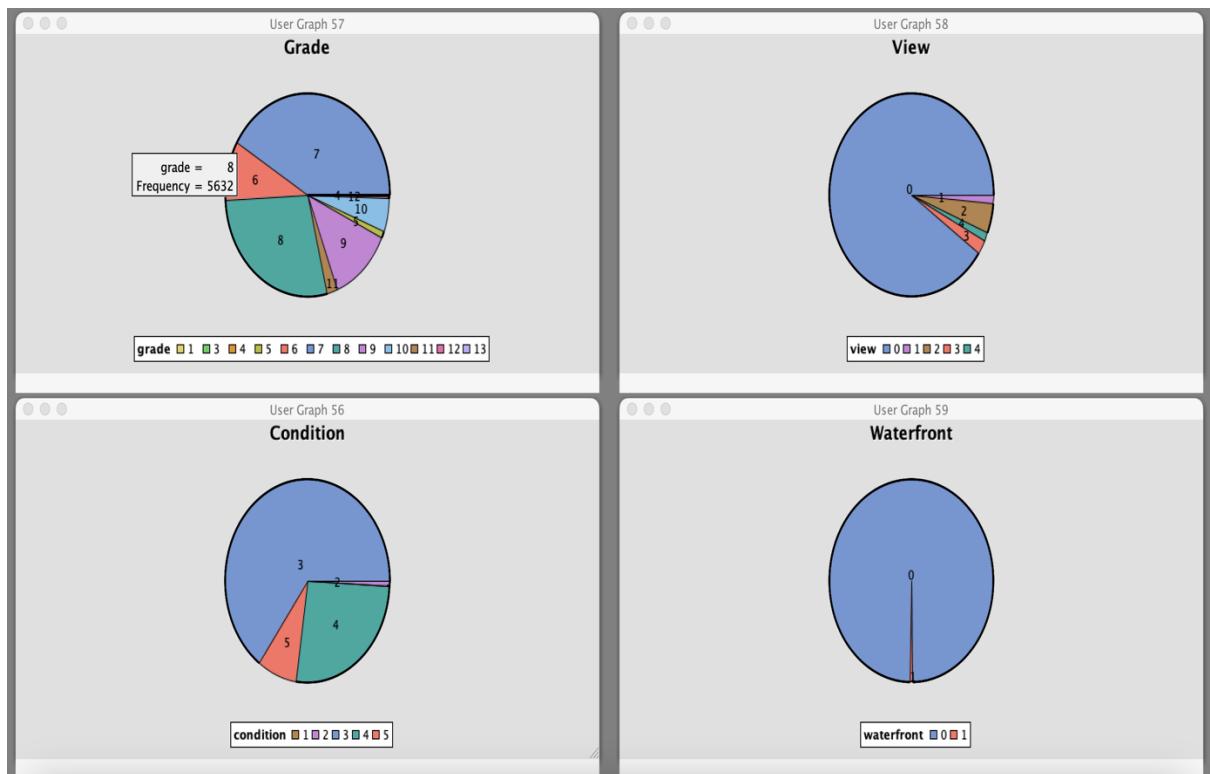


Figure 5.9 Pie Chart

We use Pie Charts to check the distribution of categorical variables. We observed that all our categorical variables are imbalance in term of the distribution.

Data Role	Variable Name	Level	CODE	Frequency Count	Type ▾	Percent	Level Index	Role	Label	Plot
TRAIN	condition	3		0	14031N	64.91926	3INPUT	condition	1	
TRAIN	condition	4		2	5679N	26.27585	4INPUT	condition	1	
TRAIN	condition	5		1	1701N	7.870263	5INPUT	condition	1	
TRAIN	condition	2		4	172N	0.795817	2INPUT	condition	1	
TRAIN	condition	1		3	30N	0.138805	1INPUT	condition	1	
TRAIN	grade	7		0	9881N	41.55369	6INPUT	grade	1	
TRAIN	grade	8		2	6068N	28.0757	7INPUT	grade	1	
TRAIN	grade	9		4	2615N	12.02092	8INPUT	grade	1	
TRAIN	grade	6		1	2018N	9.42811	5INPUT	grade	1	
TRAIN	grade	10		6	1134N	5.246842	9INPUT	grade	1	
TRAIN	grade	11		3	399N	1.846111	10INPUT	grade	1	
TRAIN	grade	5		5	242N	1.119696	4INPUT	grade	1	
TRAIN	grade	12		7	90N	0.416416	11INPUT	grade	1	
TRAIN	grade	4		8	29N	0.134179	3INPUT	grade	1	
TRAIN	grade	13		10	13N	0.060149	12INPUT	grade	1	
TRAIN	grade	3		9	3N	0.013881	2INPUT	grade	1	
TRAIN	grade	1		11	1N	0.004627	1INPUT	grade	1	
TRAIN	view	0		0	19489N	90.17258	1INPUT	view	1	
TRAIN	view	2		963N	4.455652	3INPUT	view	1		
TRAIN	view	3		1	510N	2.359691	4INPUT	view	1	
TRAIN	view	1		4	332N	1.536113	2INPUT	view	1	
TRAIN	view	4		2	319N	1.475964	5INPUT	view	1	
TRAIN	waterfront	0		0	21450N	99.24582	1INPUT	waterfront	1	
TRAIN	waterfront	1		1	163N	0.754176	2INPUT	waterfront	1	

Figure 5.10 Summary Statistics for Class Variables

The above summary statistics for class variables are shown to support the visualization.

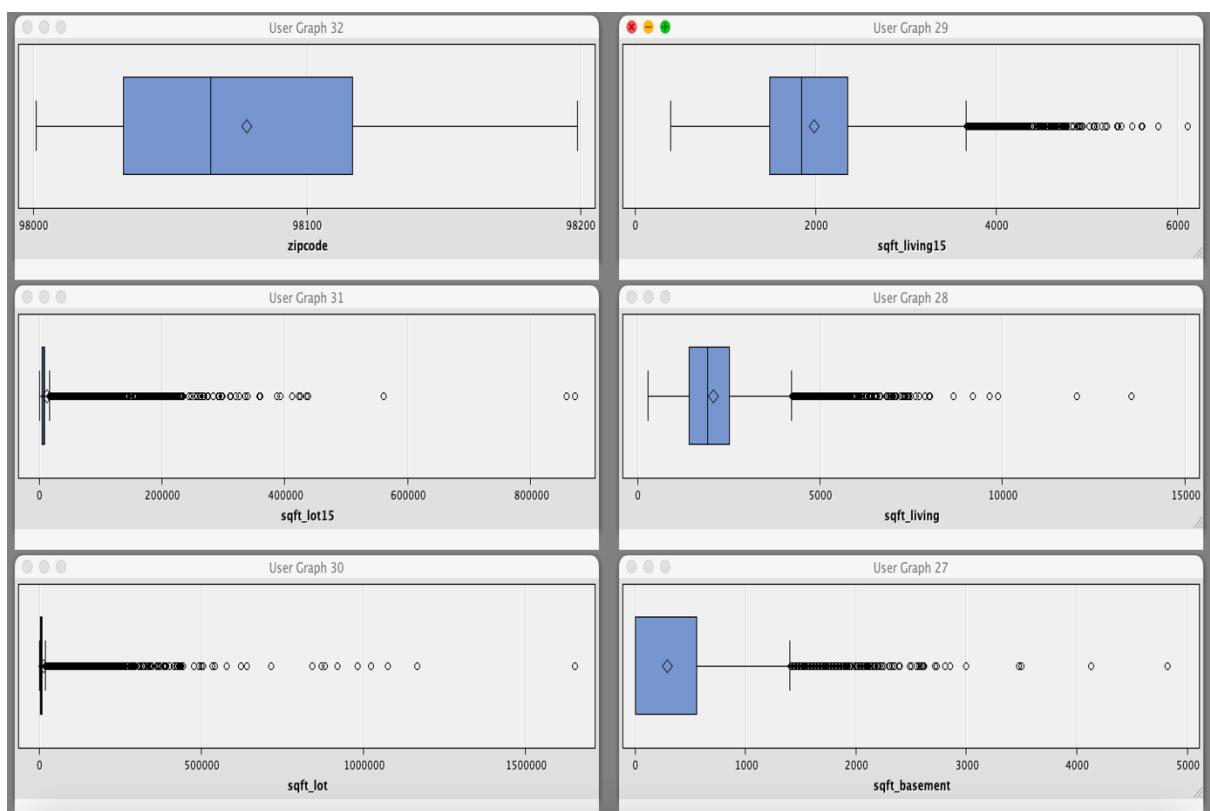
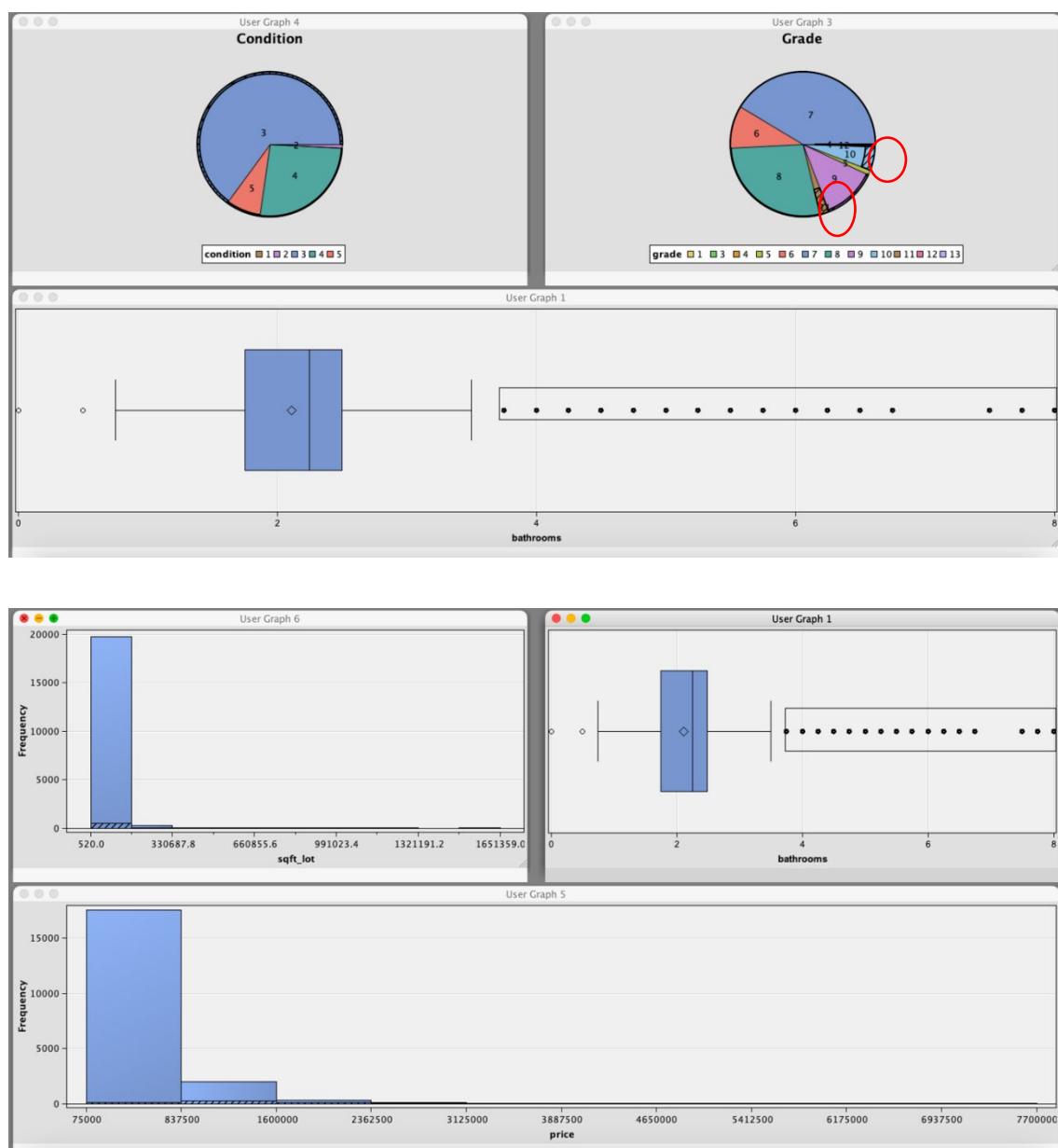


Figure 5.11 Boxplot

We use boxplot to check for outliers in an interval variable. The figure above shows the boxplot of different interval variables. There are outliers in bathrooms variable, bedrooms variable, sqft\_above variable, price variable, lat variable, long variable, sqft\_living15 variable, sqft\_lot15 variable, sqft\_living variable, sqft\_lot variable and sqft\_basement variable. If the outliers are due to error, it should be removed. For sqft\_above variables, price variables, lat variables, long variables, sqft\_living15 variables, sqft\_lot15 variables, sqft\_living variables, sqft\_lot variables and sqft\_basement variables, we observed that outliers have mostly the same distribution as the non-outliers, because there are many points outside the fence, the outliers are unlikely due to some error and should be kept.



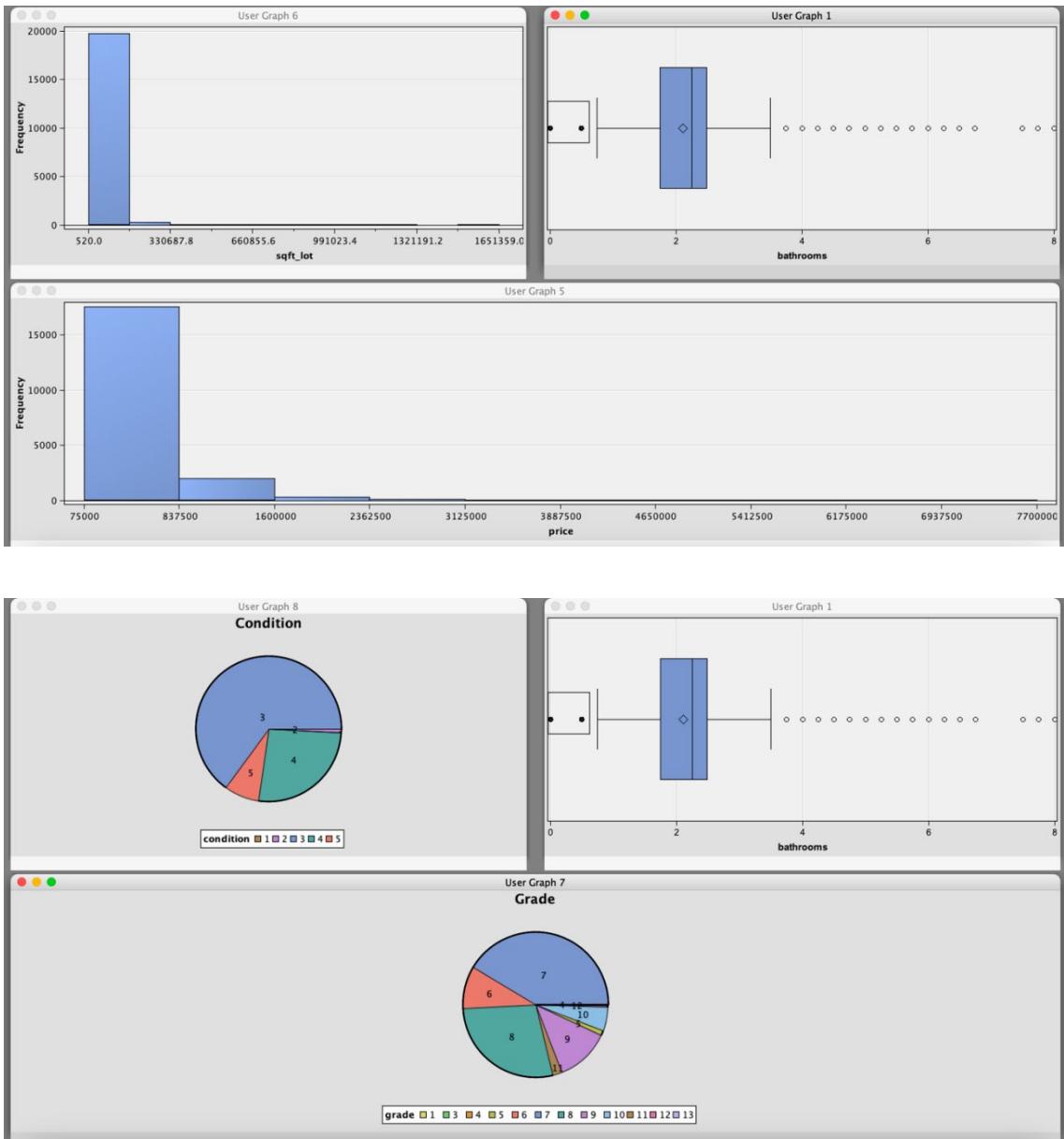
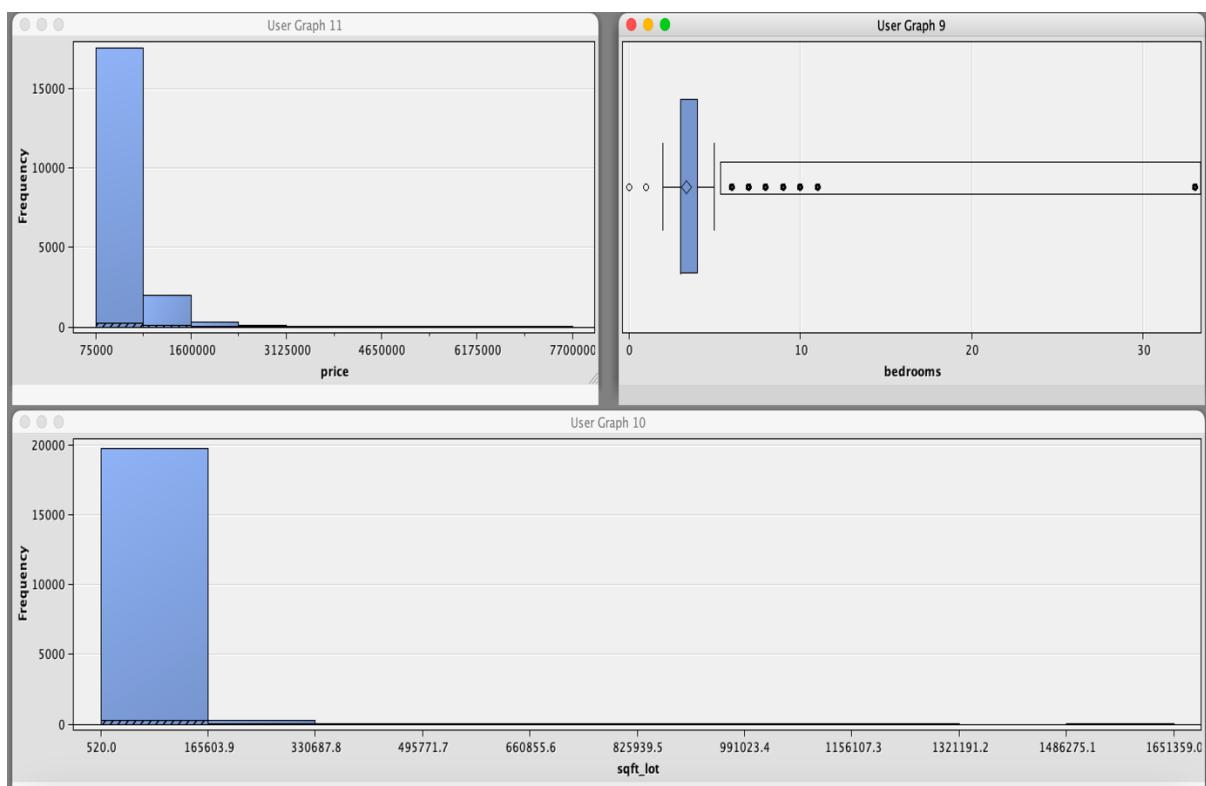
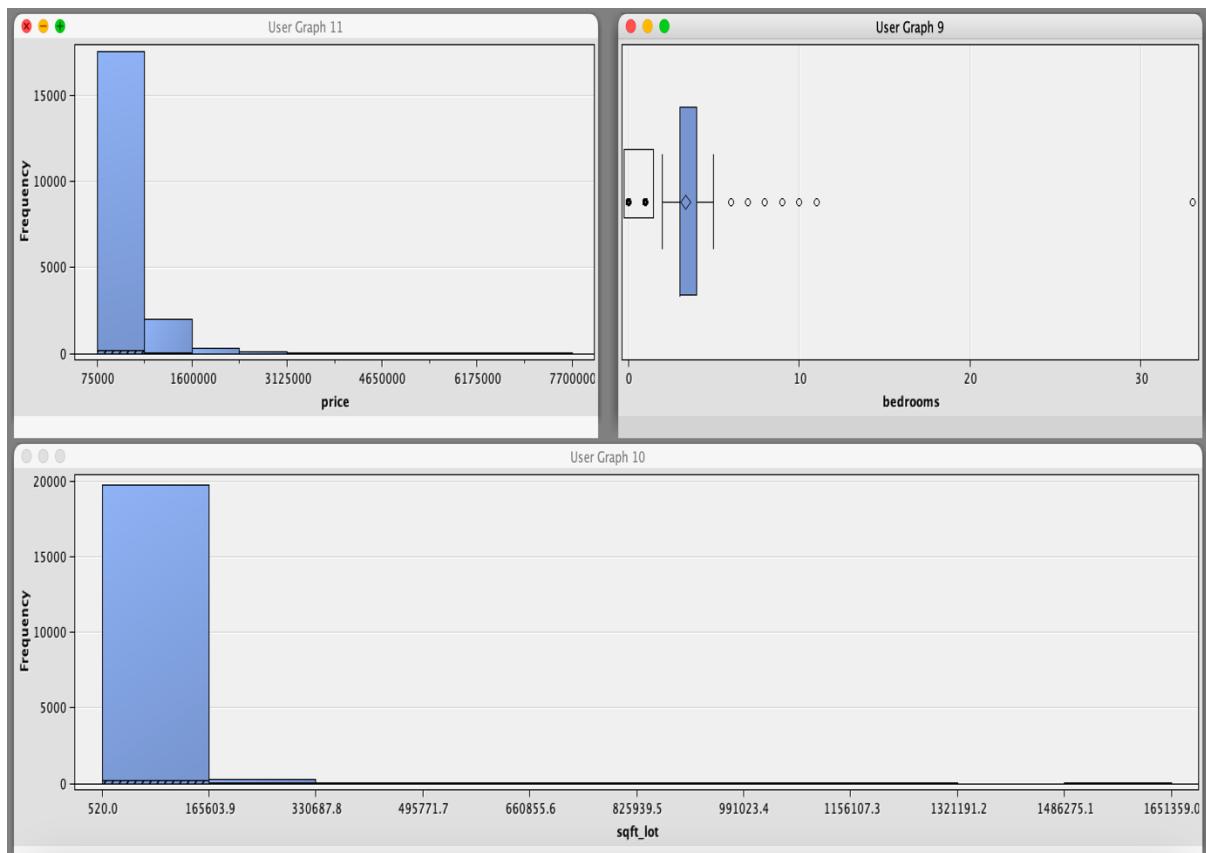


Figure 5.12 Variable Association of Outliers in Bathrooms and Other Variables

We investigate whether the outlier should be removed by examining the variable association between the outlier and other variables. First is the bathrooms variable. We observe that the bathrooms variable has an outlier of 0 bathrooms. There are ten observations with 0 bathrooms. Since it is not conventional to have houses without bathroom, these outliers should be removed. For other outliers, we observed that house with more bathrooms tend be higher grade as highlighted in the red circle. These characteristics are consistent with price. As the outliers have mostly the same distribution as the non-outliers, so the outliers are unlikely to be due to some error and should be kept.



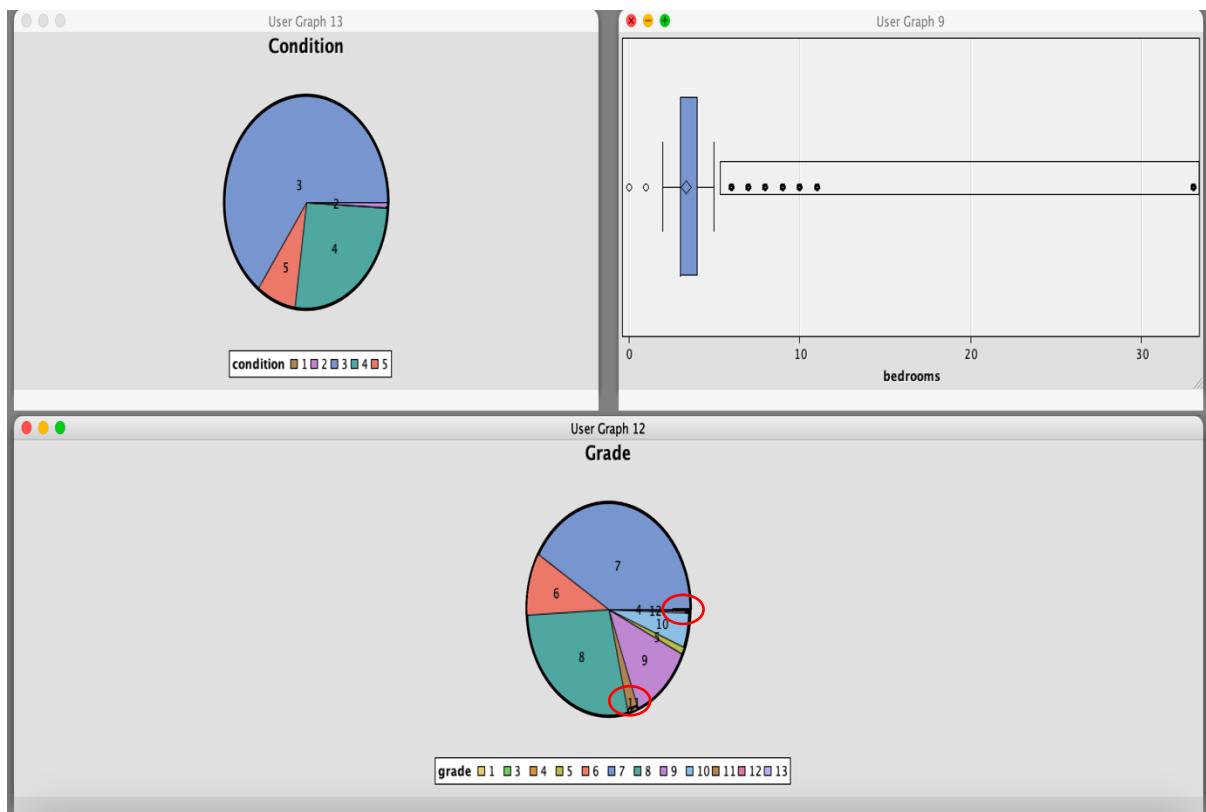
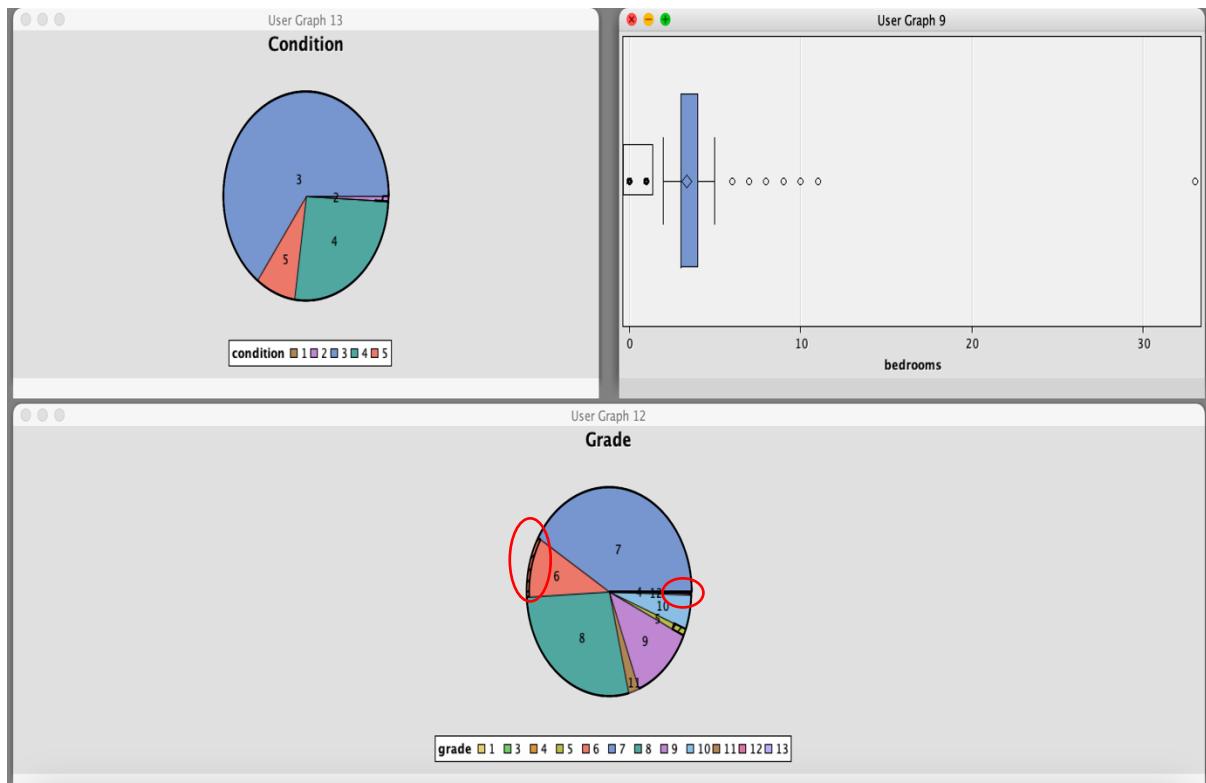


Figure 5.13 Variable Association of Outliers in Bedrooms and Other Variables

Second is the bedrooms variable. We observe that there is a suspicious 33 room entry which is totally unrealizable. Therefore, it is likely that human error causes the outlier, and the point should be removed. We observe that the bedrooms variable has an outlier of 0 bedrooms. There are ten observations with 0 bathrooms. Although we noticed that the house with 0 bathrooms has grade of 5 to 6 as highlighted in the red circle, but since it is not conventional to have house without bathrooms, these outliers should be removed. For other outliers, we observed that house with more bedrooms tend be higher grade as highlighted in the red circle. These characteristics are consistent with price. As the outliers have mostly the same distribution as the non-outliers, so the outliers are unlikely to be due to some error and should be kept.

As a conclusion for univariate analysis, we learned that 1) There is incomplete error in the bathrooms variable. 2) There are inconsistency errors in the yr\_built variable. 3) There are noisy errors in the bathrooms variable, bedrooms variable, sqft\_above variable, price variable, lot variable, long variable, sqft\_living15 variable, sqft\_lot15 variable, sqft\_living variable, sqft\_lot variable and sqft\_basement variable.

## 5.2.2 Bivariate Analysis

Bivariate analysis is an analysis conducted with 2 different variables mainly **to understand if there is a relationship among the chosen variables**. At the same time, this **analysis is also used to identify the redundant variables** and to **examine for relevant variables**.

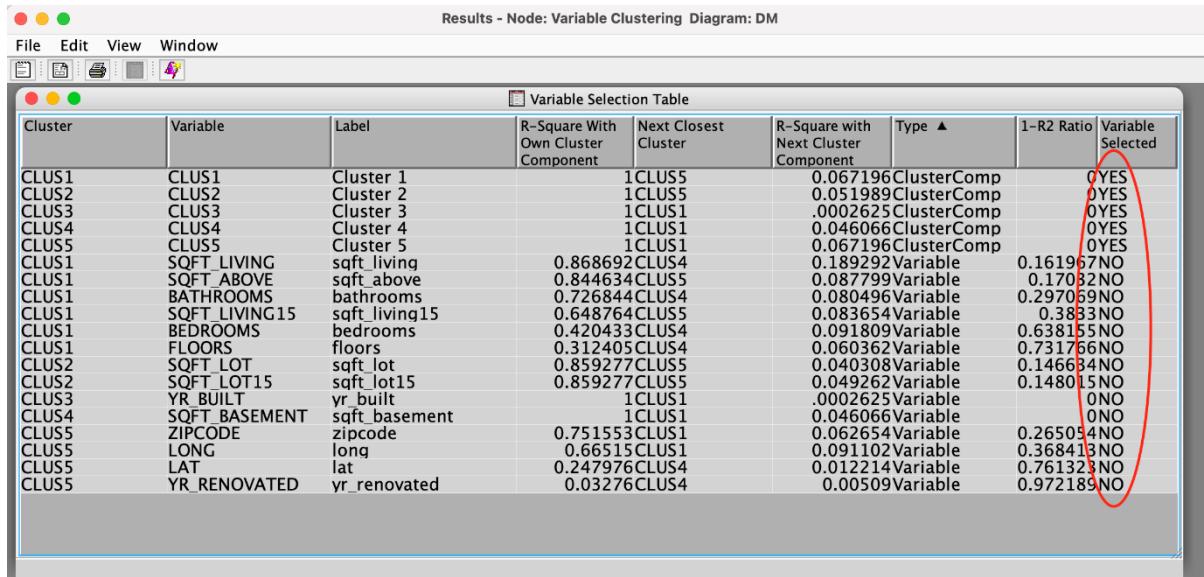


Figure 5.14 Variable Clustering

First, we will be looking at the redundant variable as shown in Figure 5.14. SAS Enterprise Miner tool, namely variable clustering was used to spot for redundant variables. In references to the figure above, no variable belongs to a cluster which indicates that the variables are not redundant and should not be removed from the dataset for next phase.

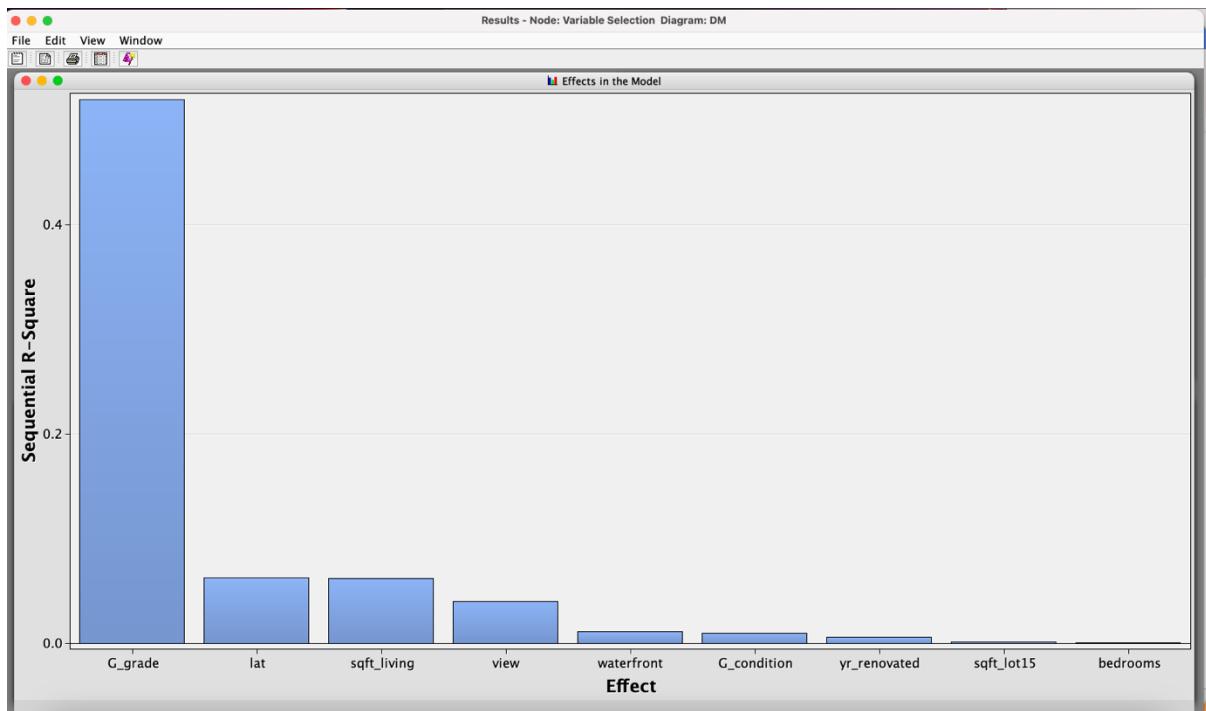


Figure 5.15 Variable Selection

Next, SAS Enterprise Miner was used to locate the relevant variables that contributes to the target variable house sales price (price) using the variable selection tool. Relevant variables identified are the grade, latitude, sqft\_living, view, waterfront, condition, yr\_renovated, sqft\_lot15 and bedrooms. The other variable not listed will be dropped from the dataset.

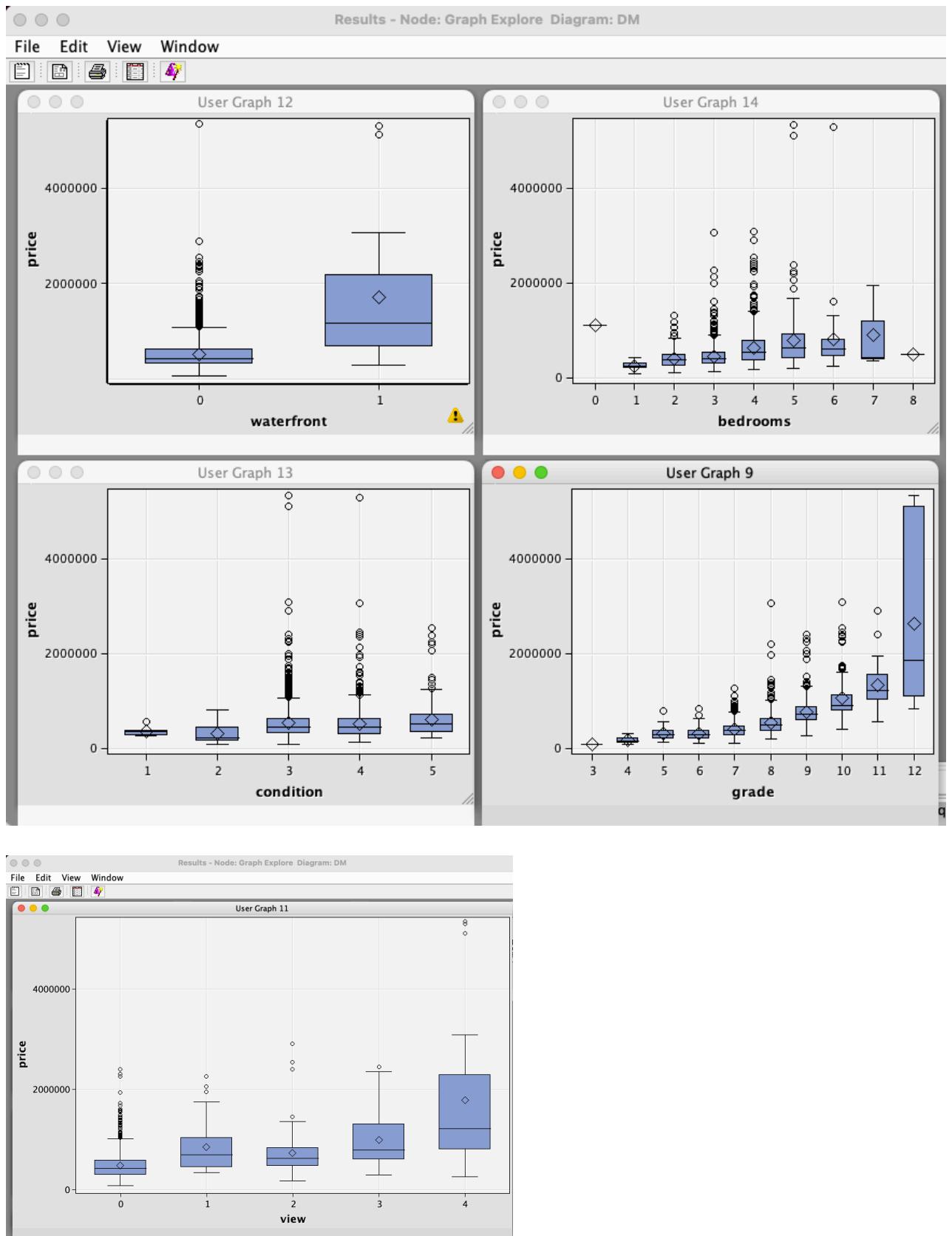


Figure 5.16 Selected Categorical Variable with Target

To further verify the selected findings by visualization, box plot and scatter plot are plotted. Figure 5.16 above shows the bivariate relationship between the target variable, house price (price) and the relevant variables as highlighted earlier.

From the figure above, we noticed that the waterfront (0-without waterfront, 1-with waterfront), bedrooms, condition, grade, view variables increase, the house price variable increases with those variables.

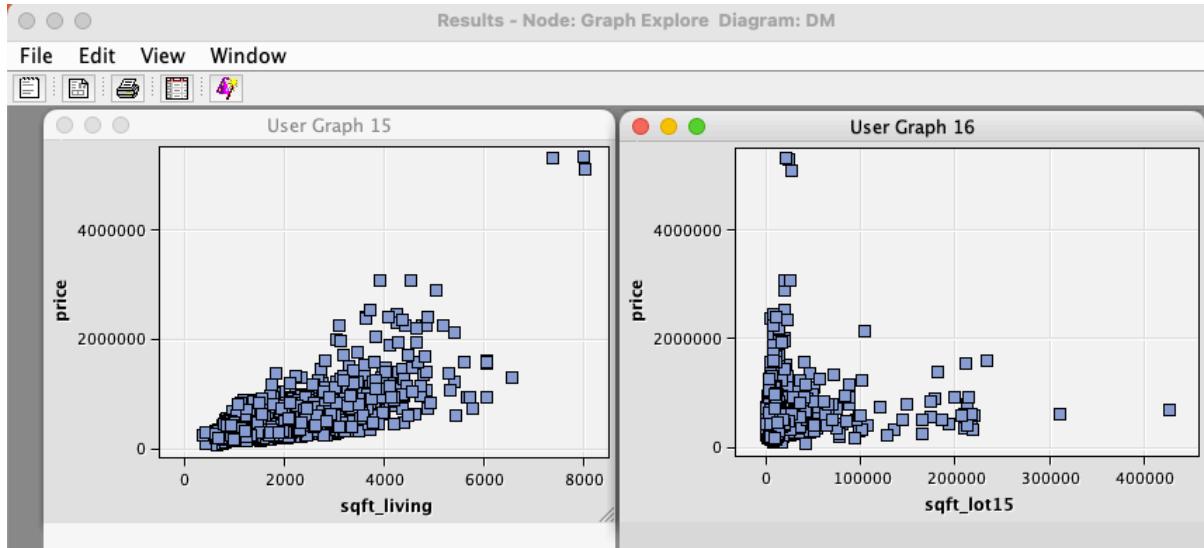


Figure 5.17 Selected Continuous Variable with Target

Referring to Figure 5.17, the scatter plots were plotted for the price variable and both sqft\_living and sqft\_lot15 respectively. In the figure above, we can see that as the sqft\_living increases the price variable increases proportionally. For the sqft\_lot15 variable we can see that the higher the increase in the sqft\_lot15 the increase in price were in a complex mix or trend given that the outlier positioned at 0 are not taken into consideration.

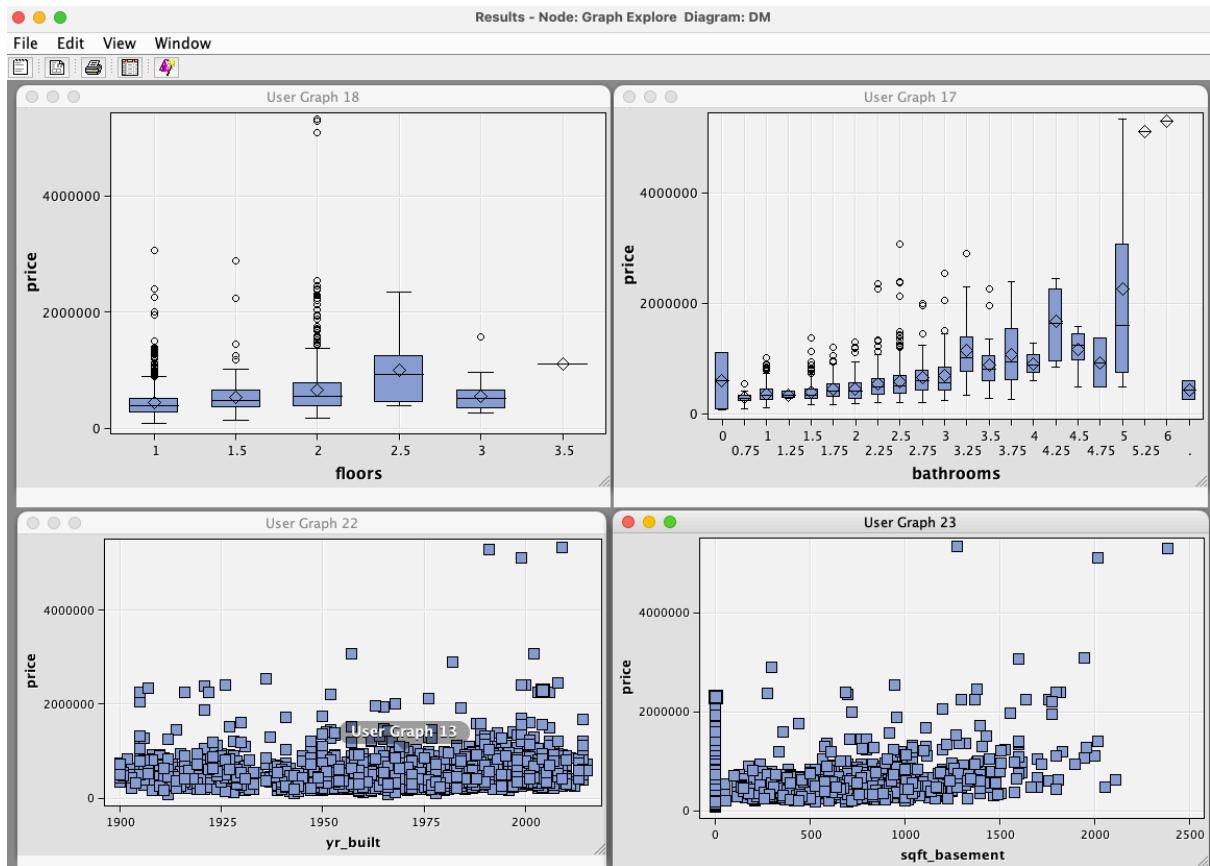


Figure 5.18 Dropped Variables

The figure above shows the several scatter plot and box plot for sqft\_basement, yr\_built, floor, bathrooms variable with house price (price) variable. From our observations the variable to be dropped are plotted in a constant points instead of following the price trend which proves that these variables do not contribute towards the house the prices and should be dropped as per the results of the previous charts.

In conclusion, the visualization for the relevant variable selection namely grade, latitude, sqft\_living, view, waterfront, condition, yr\_renovated, sqft\_lot15 and bedrooms should be retained and to drop the other variables.

### 5.2.3 Multivariate Analysis

Multivariate analysis is a type of statistical method that analyses the relationship between several variables. With multivariate analysis, we can observe the patterns and make comparisons between multiple factors at once. It is more useful as we always compare how several factors affects the target variable in real life situations.

Figure 5.19 shows a correlation matrix that shows the correlation between the interval variables. If the color is red, the variables are in positive relationship, blue shows negative relationship and grey shows weaker relationship. The square footage of the apartments is mostly in red, showing positive relationships with several variables based on the correlation matrix.

Besides, we can validate the relationship between the variables using correlation table. Figure 5.20 shows part of the correlation table where the full result is listed in Appendix. We observed that the strongest positive relationship is between the square footage of the apartments interior living space (sqft\_living) and the square footage of the interior housing space that is above ground level (sqft\_above) with a correlation coefficient of 0.8766. Meanwhile, the strongest negative relationship is between zipcode and longitude with correlation coefficient of -0.5640.

Thus, we can conclude that square footage of the apartments interior living space (sqft\_living) and the square footage of the interior housing space that is above ground level (sqft\_above) are highly correlated and have a strong positive relationship.

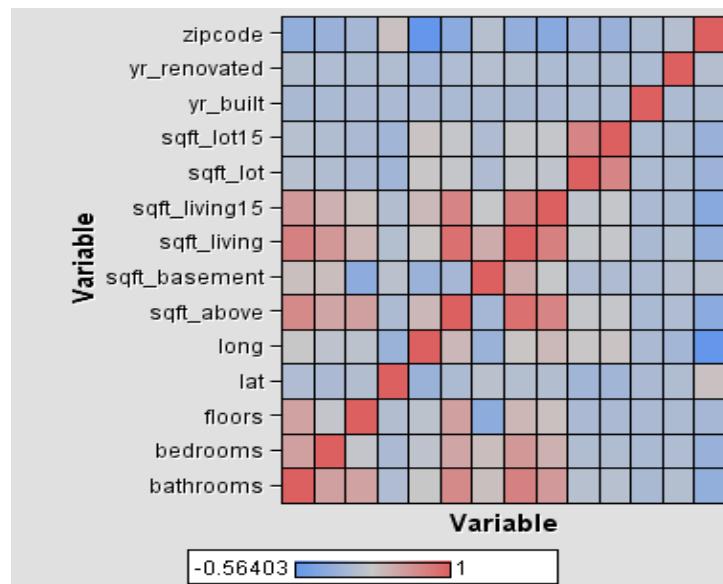


Figure 5.19 Correlation Matrix

Variable	Variable	Correlation
sqft_lot	lat	-0.08568
sqft_lot15	lat	-0.08641
yr_built	lat	-0.0106
yr_renovated	lat	0.029408
zipcode	lat	0.26703
bathrooms	long	0.223076
bedrooms	long	0.12954
floors	long	0.125339
lat	long	-0.13547
long	long	1
sqft_above	long	0.343752
sqft_basement	long	-0.14468
sqft_living	long	0.240213
sqft_living15	long	0.334574
sqft_lot	long	0.229509
sqft_lot15	long	0.254436
yr_built	long	-0.00924
yr_renovated	long	-0.0684
zipcode	long	-0.56403
bathrooms	sqft_above	0.685412
bedrooms	sqft_above	0.477615
floors	sqft_above	0.523794
lat	sqft_above	-0.000777
long	sqft_above	0.343752
sqft_above	sqft_above	1
sqft_basement	sqft_above	-0.0519
sqft_living	sqft_above	0.876601
sqft_living15	sqft_above	0.731873
sqft_lot	sqft_above	0.183487
sqft_lot15	sqft_above	0.194019
yr_built	sqft_above	-0.01339
yr_renovated	sqft_above	0.023243
zipcode	sqft_above	-0.26119
bathrooms	sqft_basement	0.283718
bedrooms	sqft_basement	0.302999
floors	sqft_basement	-0.24569
lat	sqft_basement	0.110518
long	sqft_basement	-0.14468
sqft_above	sqft_basement	-0.0519
sqft_basement	sqft_basement	1
sqft_living	sqft_basement	0.435077
sqft_living15	sqft_basement	0.200411
sqft_lot	sqft_basement	0.01531
sqft_lot15	sqft_basement	0.017304
yr_built	sqft_basement	-0.00426
yr_renovated	sqft_basement	0.071345
zipcode	sqft_basement	0.074778
bathrooms	sqft_living	0.754688
bedrooms	sqft_living	0.576627
floors	sqft_living	0.353868
lat	sqft_living	0.052554
long	sqft_living	0.240213
sqft_above	sqft_living	0.876601
sqft_basement	sqft_living	0.435077
sqft_living	sqft_living	1

Figure 5.20 Part of Correlation Table

#### 5.2.4 Interesting Visualizations

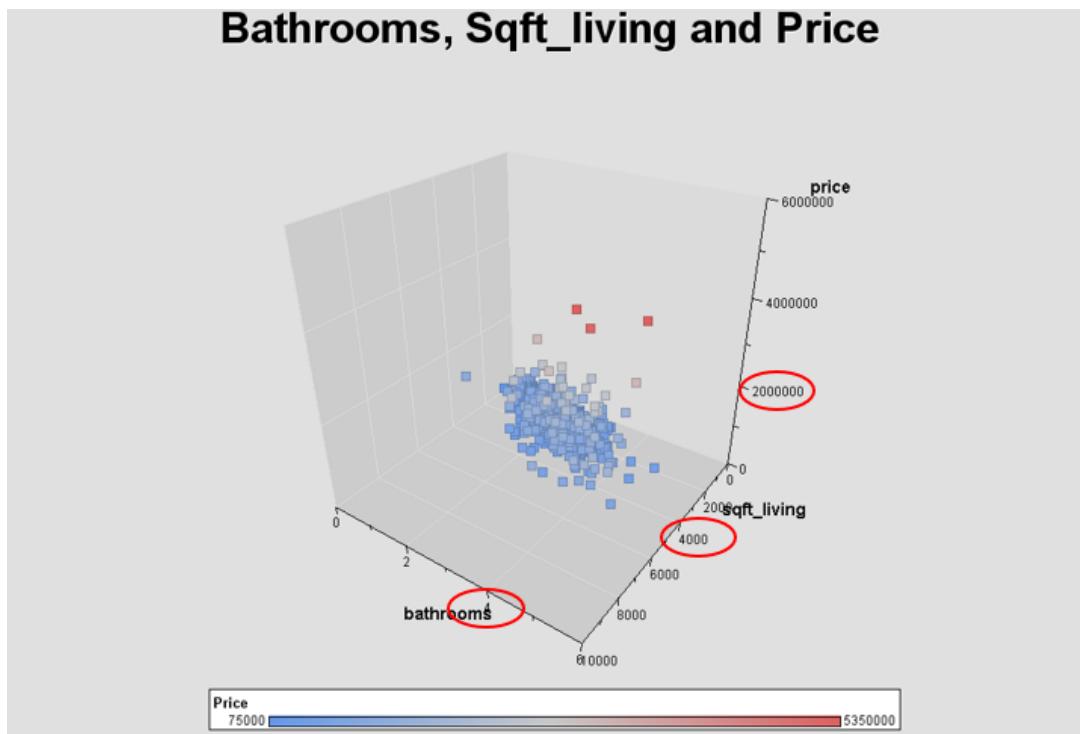


Figure 5.21 Bathrooms, Sqft\_living and Price

(1) What does the distribution look like for Bathrooms and Sqft\_living by Price?

Figure 5.21 shows a scatterplot of the number of bathrooms (bathrooms), the square footage of the apartment's interior living space (stationing), and the price of each house sold. Based on the scatterplot, we can observe that most of the data points are concentrated in the center, showing most of the price of each house sold is less than \$2,000,000, acquiring less than 4 bathrooms and less than 4000 square feet of interior living space. It is also clearly shown that there are some outliers at extreme house prices.

## Bedrooms, Sqft\_living and Price

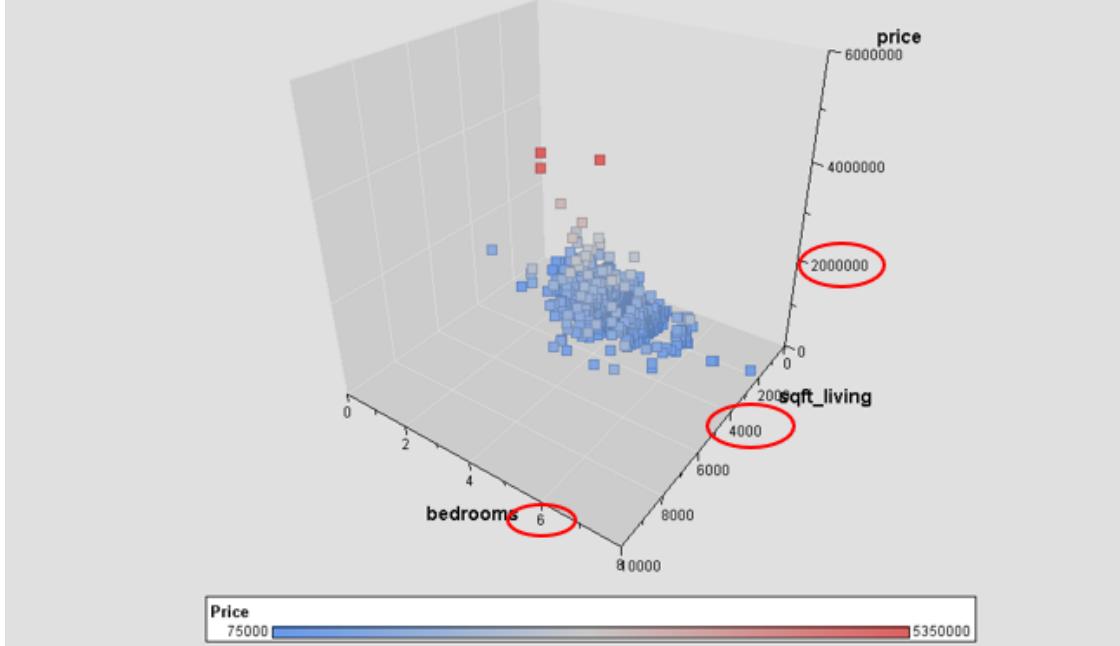


Figure 5.22 Bedrooms, Sqft\_living and Price

(2) What does the distribution look like for Bedrooms and Sqft\_living by Price?

Figure 5.22 shows a scatterplot of the number of bedrooms (bedrooms), the square footage of the apartment's interior living space (sqft\_living), and the price of each house sold. Based on the scatterplot, we can observe that most of the data points are concentrated in the center same as Figure 5.21. This clearly shows that most of the price of each house sold is also less than \$2,000,000, acquiring less than 6 bedrooms and less than 4000 square feet of interior living space. It is also clearly shown that there are some outliers at extreme house prices and number of bedrooms.

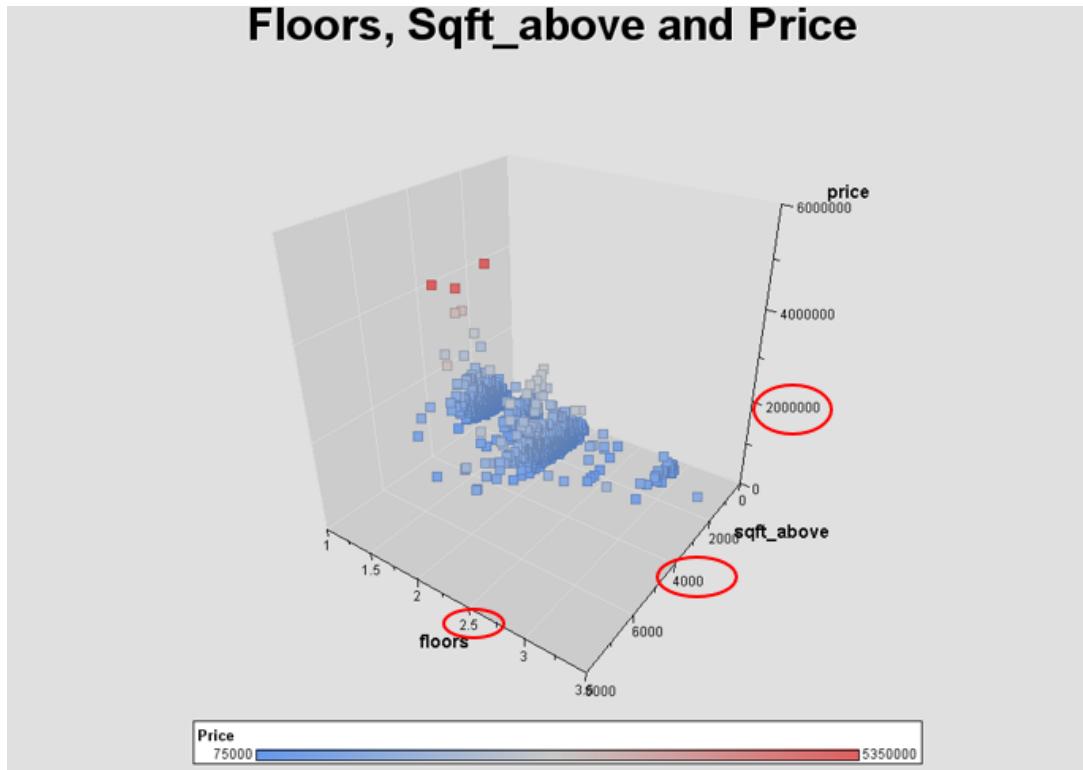


Figure 5.23 Floors, Sqft\_above and Price

(3) What does the distribution look like for Floors and Sqft\_above by Price?

Figure 5.23 shows a scatterplot of the number of bedrooms (bedrooms), the square footage of the interior housing space that is above ground level (sqft\_above), and the price of each house sold. Based on the scatterplot, we can observe that most of the data points are too concentrated in the center with each house mostly sold equal or less than \$2,000,000, acquiring less than 2.5 floors and less than 4000 square feet of interior house space. It is also clearly shown that there are some outliers at extreme house prices (\$4,000,000 and above) having only less than 1.5 floors.

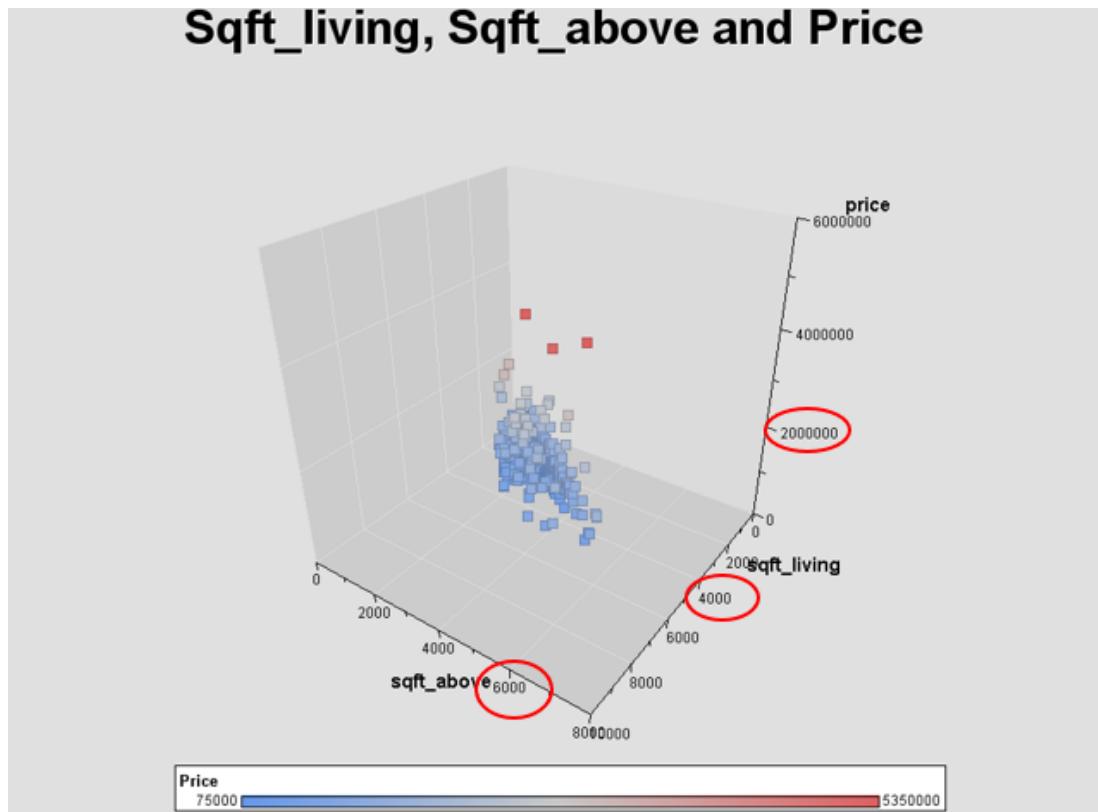


Figure 5.24 Sqft\_living, Sqft\_above and Price

(4) What does the distribution look like for Sqft\_living and Sqft\_above by Price?

Figure 5.24 shows a scatterplot of the square footage of the apartment's interior living space (sqft\_living), the square footage of the interior housing space that is above ground level (sqft\_above), and the price of each house sold. Based on the scatterplot, we can observe that most of the data points are too concentrated in the center with each house mostly sold around \$2,000,000, acquiring less than 4000 square feet of interior living space and less than 6000 square feet of interior house space above ground level. Outliers are also identified at extreme house prices (\$4,000,000).

## Sqft\_living, Sqft basement and Price

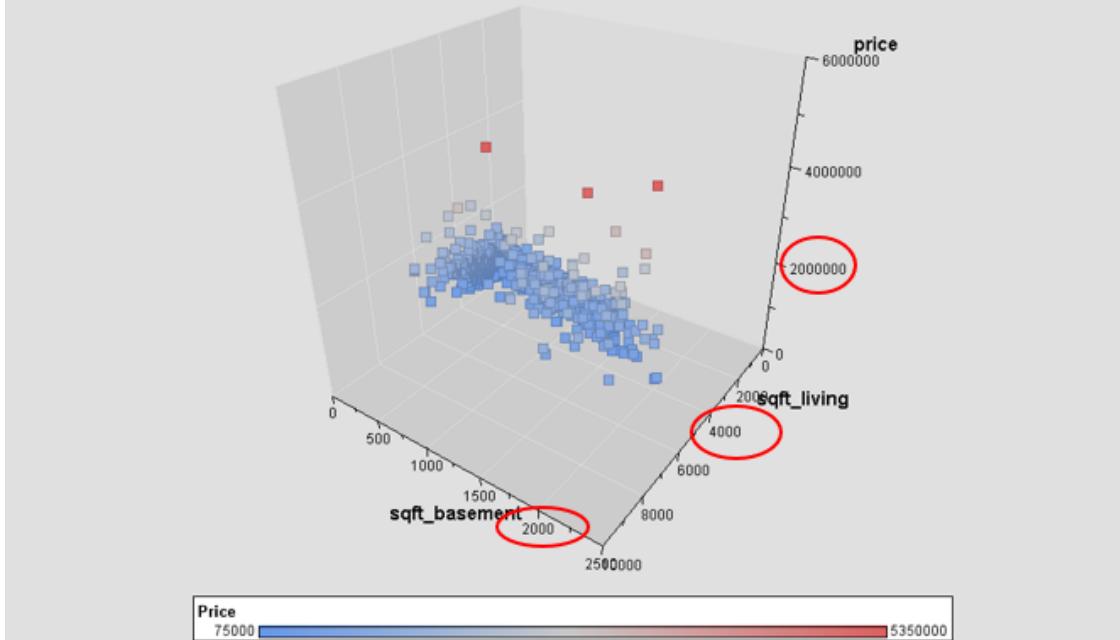


Figure 5.25 Sqft\_living, Sqft\_basement and Price

(5) What does the distribution look like for Sqft\_living and Sqft\_basement by Price?

Figure 5.25 shows a scatterplot of the square footage of the apartment's interior living space (sqft\_living), the square footage of the interior housing space that is below ground level (sqft\_basement), and the price of each house sold. Based on the scatterplot, we can observe that most of the data points are evenly distributed in contrast with previous figures. Each house mostly sold around \$2,000,000, acquiring less than 4000 square feet of interior living space and 2000 square feet of interior house basement.

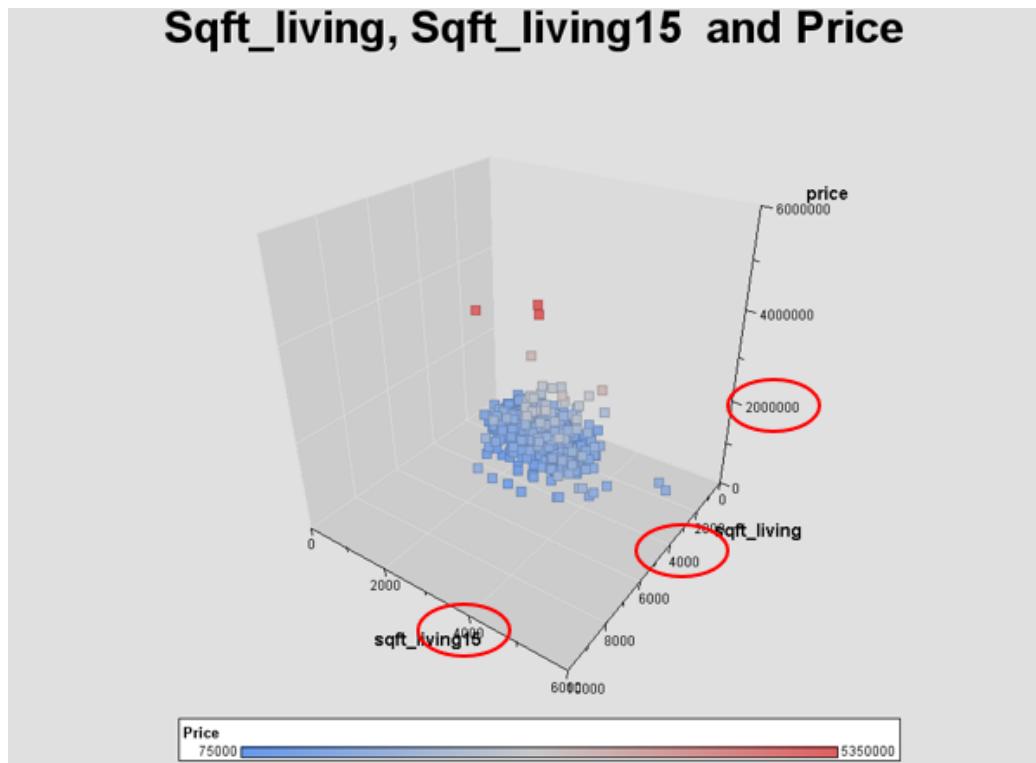


Figure 5.26 Sqft\_living, Sqft\_living15 and Price

(6) What does the distribution look like for Sqft\_living and Sqft\_living15 by Price?

Figure 5.26 shows a scatterplot of the square footage of the apartment's interior living space (sqft\_living), the square footage of the interior housing living space for the nearest 15 neighbours and the price of each house sold. Based on the scatterplot, we can observe that most customers prefer purchasing houses for around \$2,000,000, acquiring less than 4000 square feet of interior living space and 4000 square feet for the nearest 15 neighbours. Outliers are also detected like the previous figures.

## Sqft\_lot, Sqft\_lot15 and Price

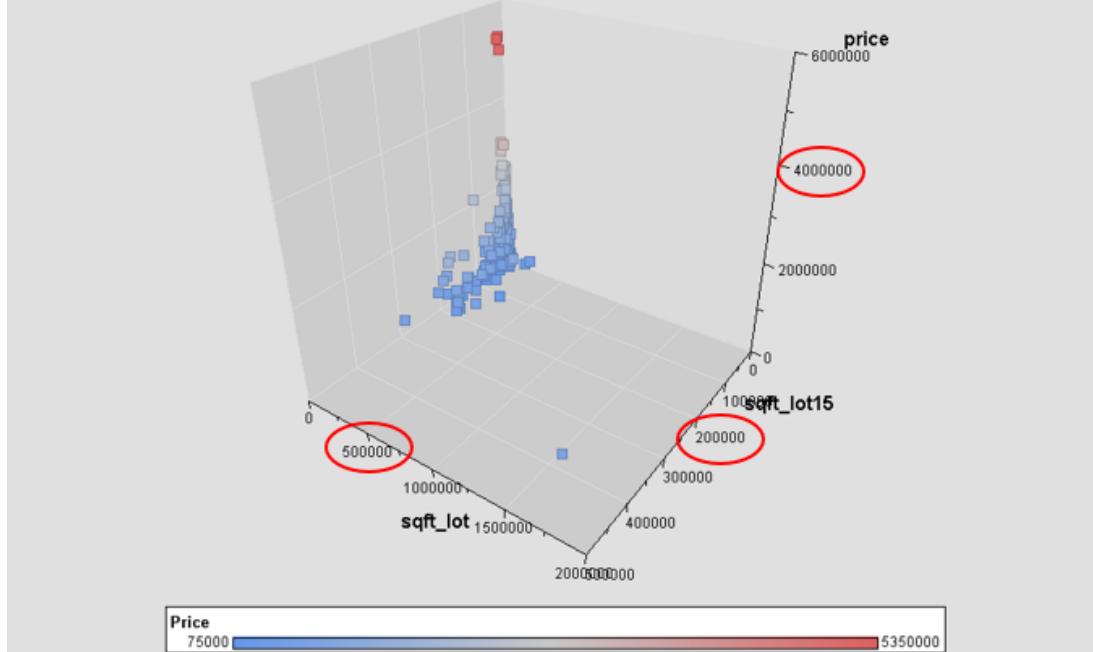


Figure 5.27 Sqft\_lot, Sqft\_lot15 and Price

(7) What does the distribution look like for Sqft\_lot and Sqft\_lot15 by Price?

Figure 5.27 shows a scatterplot of the square footage of the land space (sqft\_lot), the square footage of the land lots of the nearest 15 neighbours, and the price of each house sold. Based on the scatterplot, we can observe that most customers prefer purchasing houses less than \$4,000,000, acquiring less than 500000 square feet of land and 200000 square feet for the nearest 15 neighbors. There are still a few customers who still prefer to buy \$6,000,000 during their stay in King County.

## 5.3 MODIFY – Data Modification

In this phase, we will alter, clean, reduce, and transform the data after we have explored it. This stage is crucial for further data modelling and has a direct impact on the precision of the prediction model.

### 5.3.1 Modifying and Correcting Source Data

To perform classification, we have added an attribute based on the price named price\_range. We have taken the median of the price attributes as our baseline, if the value of the price is over the median the price range is 1, if lower or equal to the median the value of price\_range is 0. And we have set the price\_range attributes as target.

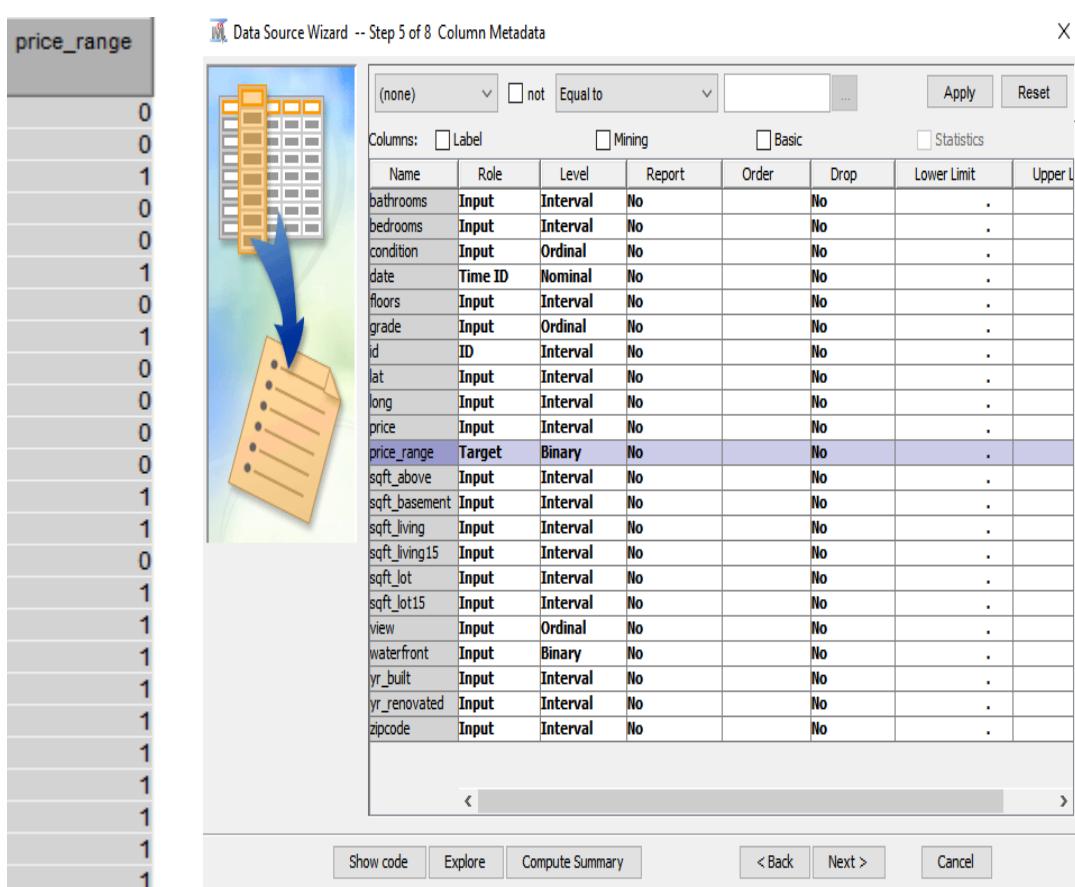


Figure 5.28: Create price\_range variable

During data exploration, we found incomplete data, noisy data, and inconsistent data. These variables are as follows:

<b>Variable</b>	<b>Error Type</b>
bathrooms	Incomplete
yr_built	Inconsistent
bathrooms bedrooms sqft_above price lat long sqft_living15 sqft_lot15 sqft_living sqft_lot sqft_basement	Noisy

Next, we will modify these error type. To modify inconsistent data, we will use Talend, which is open-source data integration platform which provides various software and services for data integration, data management, enterprise application integration, data quality, cloud storage and Big Data. In the data exploration, we found the following inconsistent errors in our yr\_built attributes. And we replace them with the correct values.

<b>Inconsistent yr_built value</b>	<b>After replacement</b>
192102	1921
19570522	1957
190810	1908
19310401	1931
192703	1927
19590731	1959
191006	1910

yr_built	yr_built_substring
integer	integer
1970	1970
1948	1948
2003	2003
192102	1921
2007	2007
1958	1958
1916	1916
1950	1950
2007	2007

yr_built	yr_built_substring
integer	integer
1968	1968
1978	1978
1994	1994
1996	1996
1912	1912
191006	1910
1976	1976

yr_built	yr_built_substring
integer	integer
1952	1952
1954	1954
2005	2005
1983	1983
1976	1976
19590731	1959
1908	1908
1955	1955

yr_built	yr_built_substring
integer	integer
1968	1968
1918	1918
192703	1927
1960	1960
1959	1959
1928	1928
1997	1997

yr_built	yr_built_substring
integer	integer
1964	1964
1945	1945
1996	1996
19310401	1931
1968	1968
1995	1995
1981	1981

yr_built	yr_built_substring
integer	integer
1910	1910
2004	2004
1976	1976
190810	1908
2005	2005
1966	1966
1904	1904
1926	1926

yr_built	yr_built_substring
integer	integer
2002	2002
1981	1981
1990	1990
19570522	1957
1980	1980
2005	2005
1952	1952
1967	1967

Figure 5.29 Replace inconsistent data with correct values

Next, we modify the noise data and missing data. To modify the noise data and missing data, we will use SAS Enterprise Miner. Limited methods are followed, and we replace these values with missing values.

Variable	Error Type	Limits method
bathrooms bedrooms sqft_above price lat long sqft_living15 sqft_lot15 sqft_living sqft_lot sqft_basement	Noisy	Extreme Percentiles
bathrooms	Incomplete	Mean

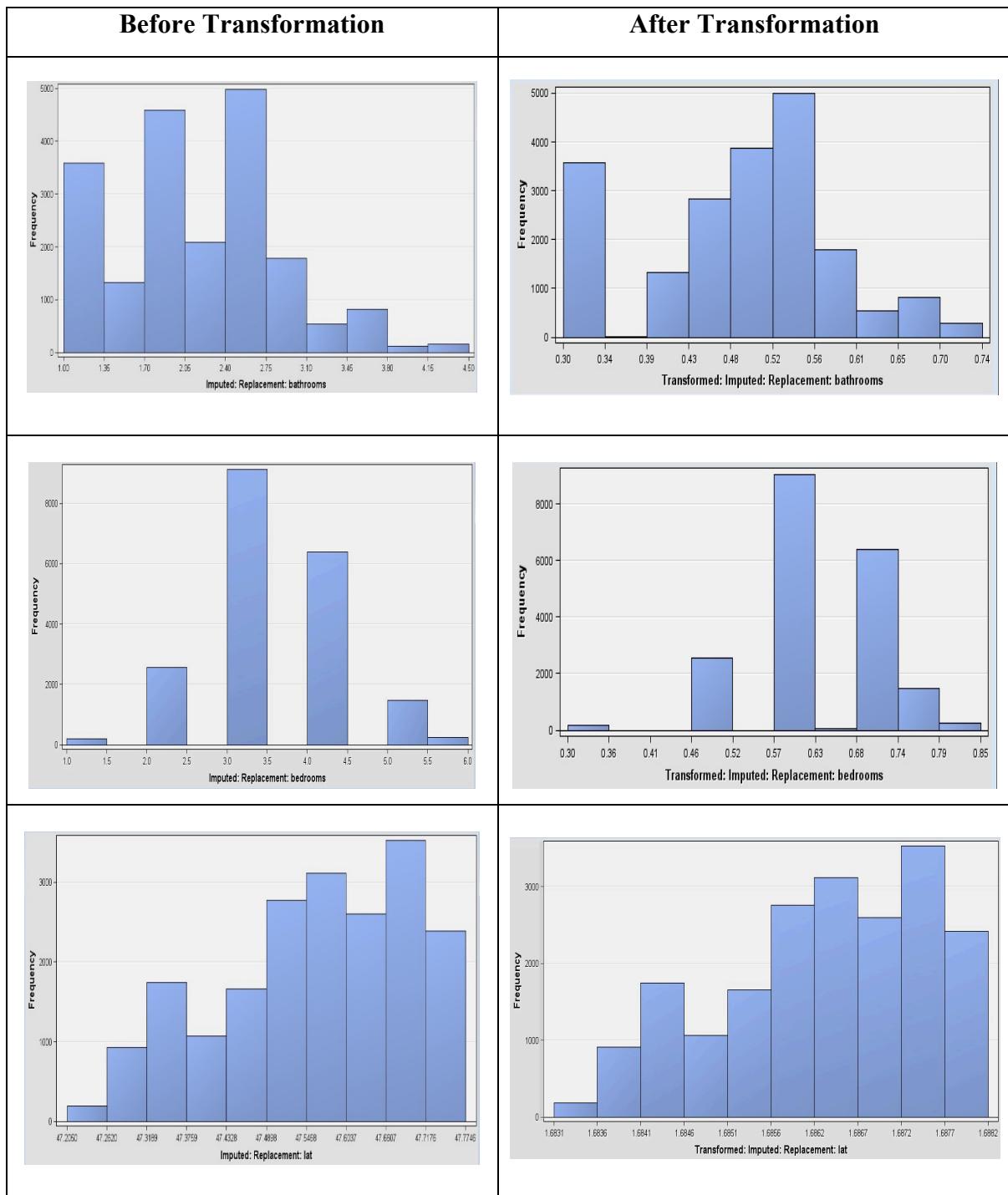
We padded incomplete and missing data with mean. The result is follows:

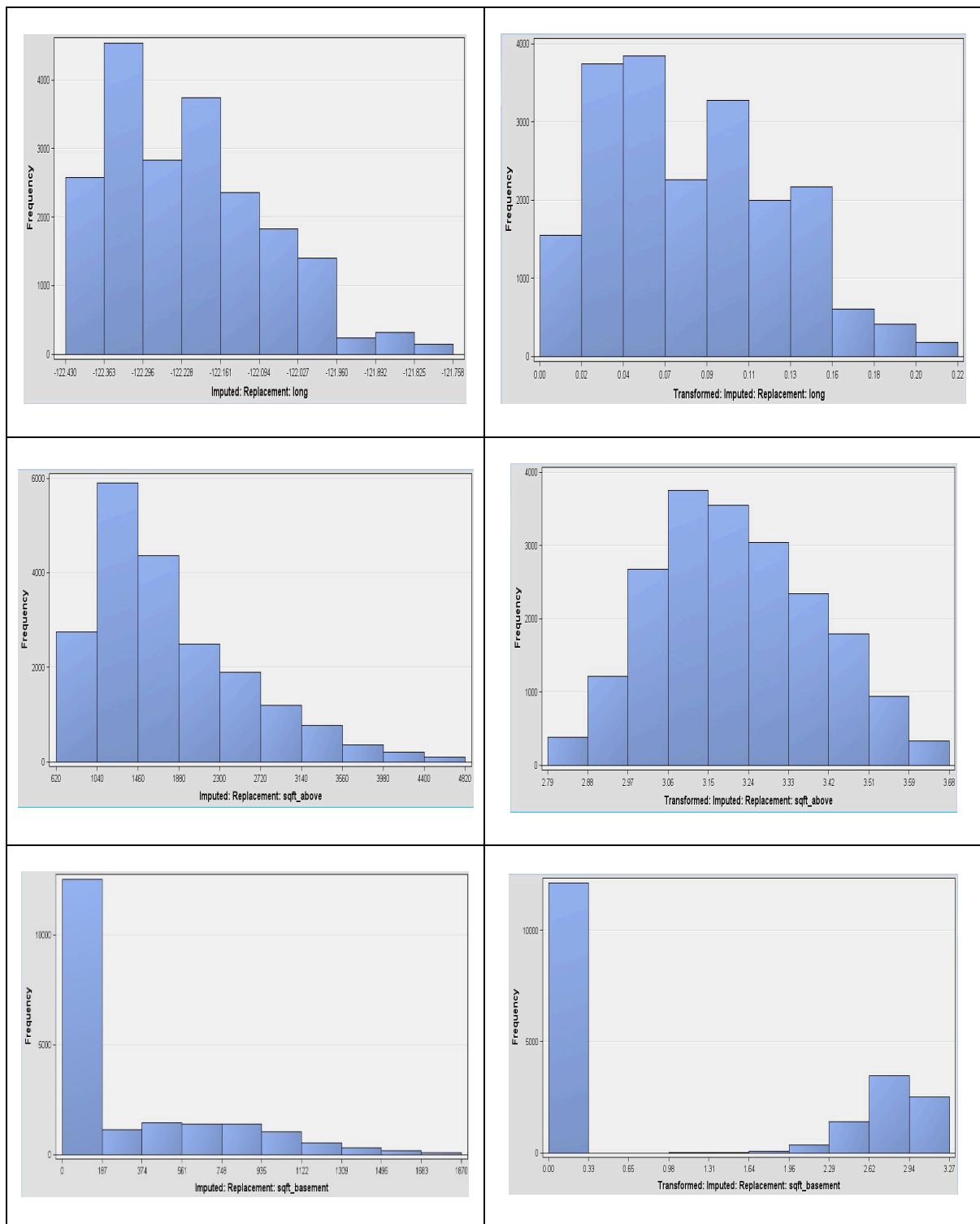
Variable Name	Impute Method	Imputed Variable	Impute Value	Role	Measurement Level	Label	Number of Missing for TRAIN
REP_bathrooms	MEAN	IMP REP_bathrooms	2.107488INPUT	INTERVAL		Replacement: bath...	178
REP_bedrooms	MEAN	IMP REP_bedrooms	3.359458INPUT	INTERVAL		Replacement: bedro...	75
REP_lat	MEAN	IMP REP_lat	47.560779INPUT	INTERVAL		Replacement: lat	215
REP_long	MEAN	IMP REP_long	-122.215INPUT	INTERVAL		Replacement: long	214
REP_price	MEAN	IMP REP_price	528856.8INPUT	INTERVAL		Replacement: price	216
REP_sqft_above	MEAN	IMP REP_sqft_above	1774.22INPUT	INTERVAL		Replacement: sqft...	201
REP_sqft_basement	MEAN	IMP REP_sqft_bas...	282.1509INPUT	INTERVAL		Replacement: sqft...	103
REP_sqft_living	MEAN	IMP REP_sqft_living	2063.344INPUT	INTERVAL		Replacement: sqft_li...	200
REP_sqft_living15	MEAN	IMP REP_sqft_livin...	1978.331INPUT	INTERVAL		Replacement: sqft_li...	213
REP_sqft_lot	MEAN	IMP REP_sqft_lot	13092.99INPUT	INTERVAL		Replacement: sqft_lot	214
REP_sqft_lot15	MEAN	IMP REP_sqft_lot15	11534.16INPUT	INTERVAL		Replacement: sqft_li...	213

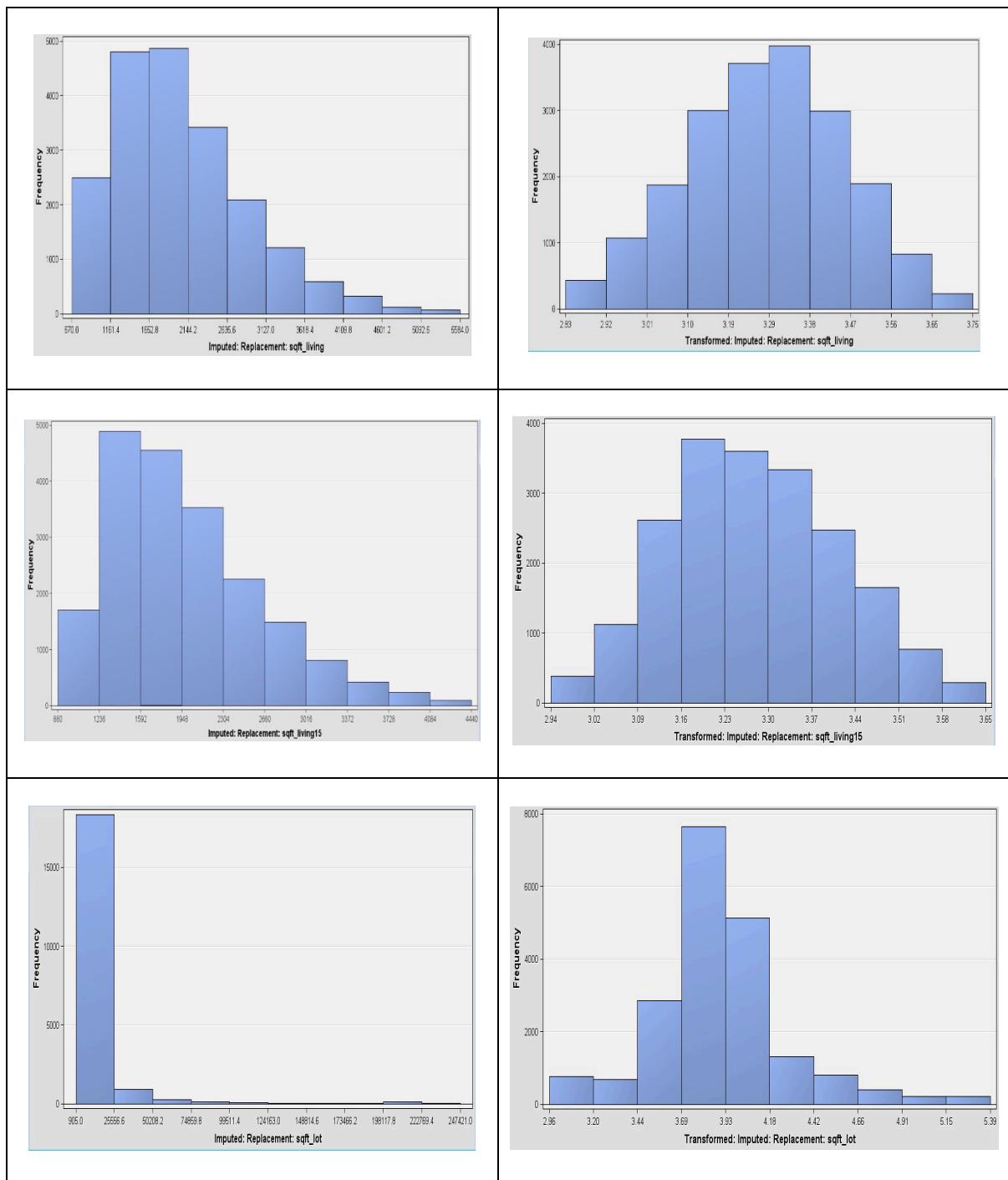
Figure 5.30 Imputation Summary

We will modify the variables and apply normalisation after substituting the data. The advantage of normalisation is that it lessens skewness. This step is advantageous for machine learning algorithms that assume the feature variable has a normal distribution, lowering the level of measurement while keeping the ratio constant. This is another advantage of normalisation because it enhances the model's training efficiency.

All the variables that we used for the log 10 transformation normalisation approach are listed in the table below.







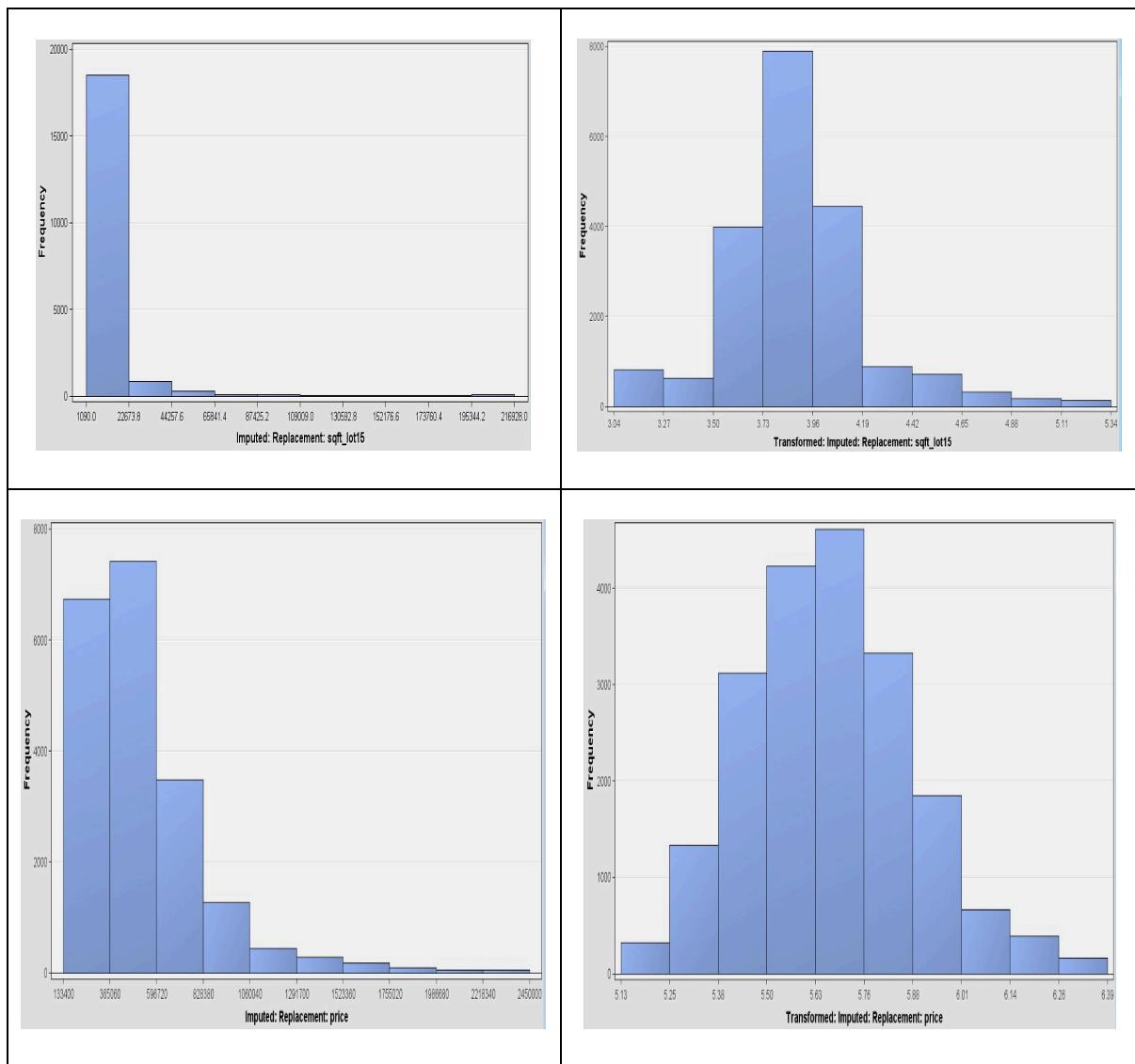


Figure 5.31 Transformation Data

### 5.3.2 Examining Exported Data

In the previous section, we identify that there is an incomplete data error in the bathrooms variable. The following is the histogram of the bathrooms variable. We can see that the missing bin is 0. Therefore, the issue is fixed.

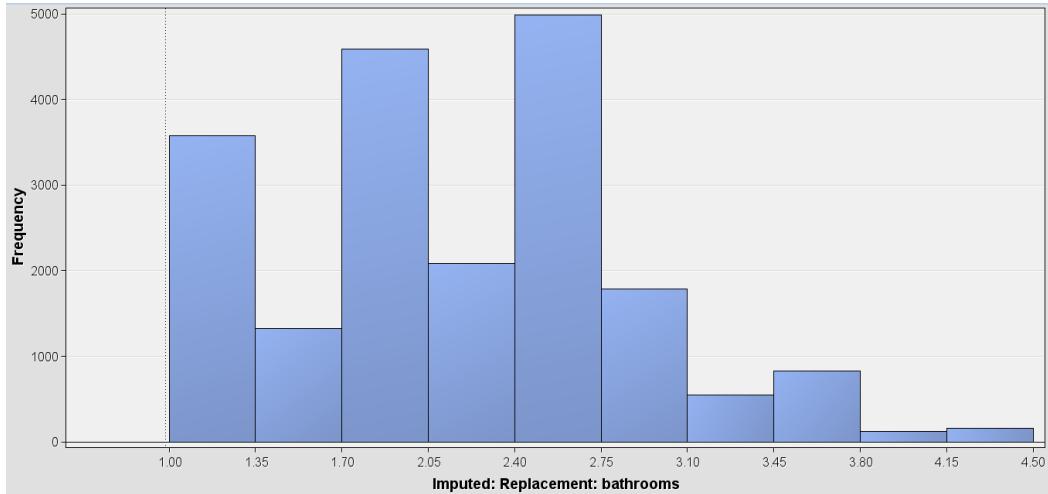


Figure 5.32 Histogram of bathrooms

In the previous section, we identify there is an inconsistent data in the yr\_built variable. The following is the histogram of the yr\_built variable. We can see that now there is no value that has month and date. Therefore, the issue is fixed.

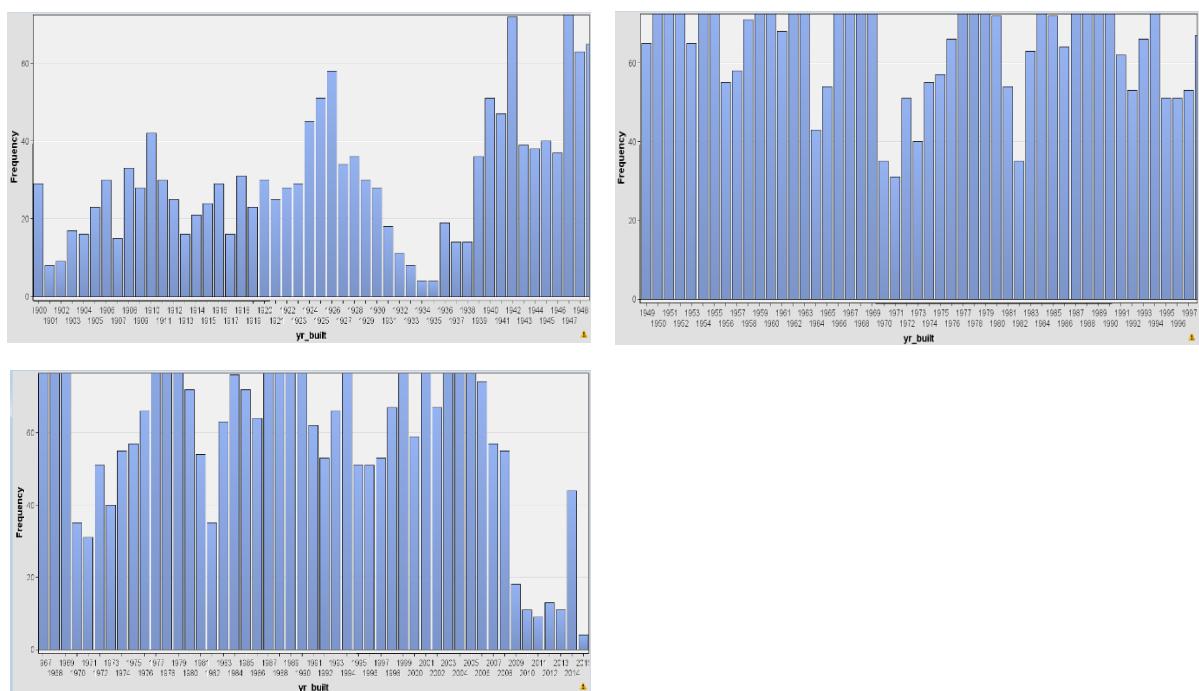


Figure 5.33 Histogram of yr\_built

In the previous section, we identify that there are noisy outliers error in bathrooms variable, bedrooms variable, sqft\_above variable, price variable, lat variable, long variable, sqft\_living15 variable, sqft\_lot15 variable, sqft\_living variable, sqft\_lot variable and sqft\_basement variable. The following is the boxplot of the abovementioned variable. We observed that after replacement and imputation there are still some outliers in the box plot. To maintain the originality of the dataset to prevent overfitting, we decide to keep the remaining outliers.

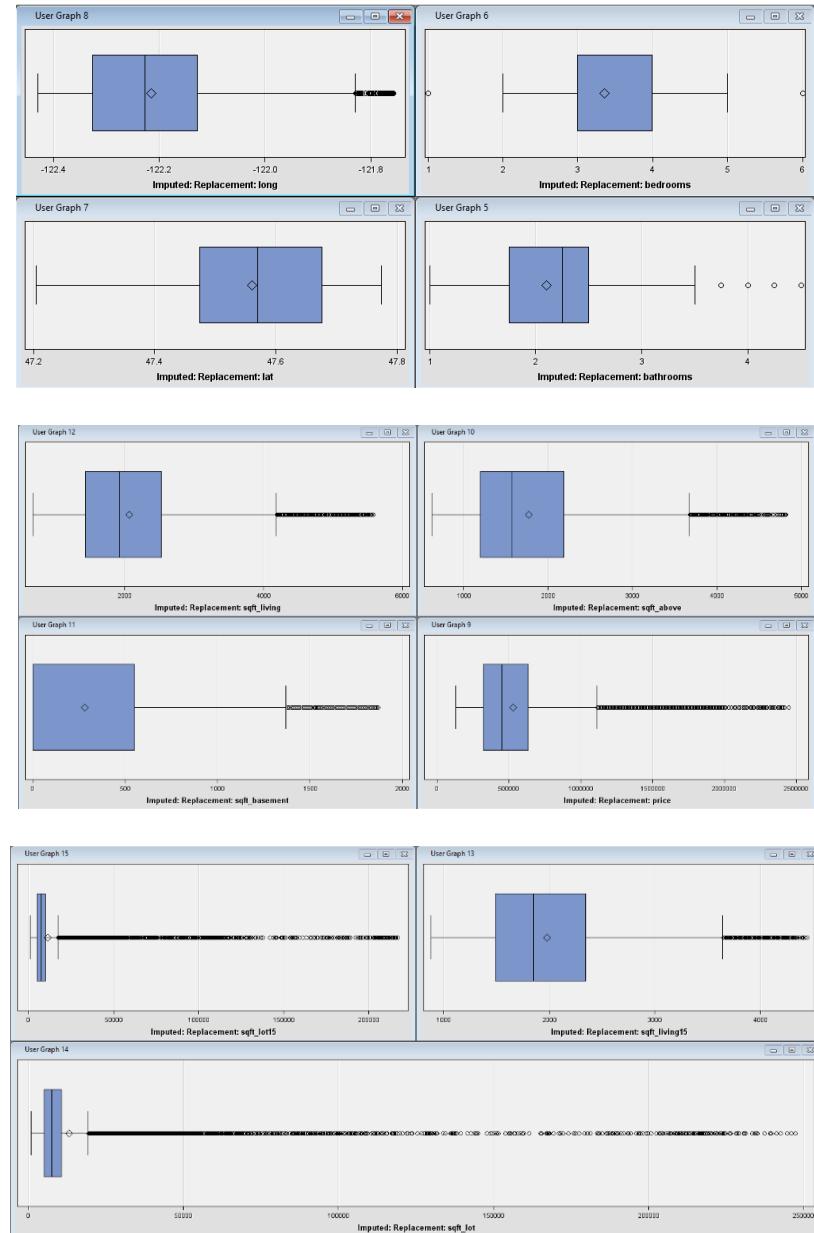


Figure 5.34 Boxplot

### 5.3.3 Creating Training and Validation Data

By using the Data Partition node, we are able to split the dataset to create Training and Validation sets. The split ratio for this project is 50:50 on Training and Validation as shown on the figure below:

Data Set Allocations	
Training	50.0
Validation	50.0
Test	0.0

Figure 5.35 Data Set Allocations

Output						
39	*	-----*				
40	*	-----*				
41	* Report Output					
42	*	-----*				
43						
44						
45						
46						
47	Summary Statistics for Class Targets					
48						
49	Data=DATA					
50						
51		Numeric	Formatted	Frequency		
52	Variable	Value	Value	Count	Percent	Label
53						
54	price_range	0	0	10864	50.2660	
55	price_range	1	1	10749	49.7340	
56						
57						
58	Data=TRAIN					
59						
60		Numeric	Formatted	Frequency		
61	Variable	Value	Value	Count	Percent	Label
62						
63	price_range	0	0	5431	50.2591	
64	price_range	1	1	5375	49.7409	
65						
66						
67	Data=VALIDATE					
68						
69		Numeric	Formatted	Frequency		
70	Variable	Value	Value	Count	Percent	Label
71						
72	price_range	0	0	5433	50.2730	
73	price_range	1	1	5374	49.7270	
74						

Figure 5.36 Report Output for Data Partition node

Based on our analysis previously, we would then focus on specific variables and drop variables that are insignificant to our modelling. We would keep LG10\_IMP\_REP\_bedrooms, LG10\_IMP\_REP\_lat, LG10\_IMP\_REP\_living, LG10\_IMP\_REP\_lot15, condition, grade, price\_range, view, waterfront and yr\_renovated.

Name	Drop	Role	Level
LG10_IMP_REP_bathrooms	Yes	Input	Interval
LG10_IMP_REP_bedrooms	No	Input	Interval
LG10_IMP_REP_lat	No	Input	Interval
LG10_IMP_REP_long	Yes	Input	Interval
LG10_IMP_REP_price	Yes	Input	Interval
LG10_IMP_REP_sqft_above	Yes	Input	Interval
LG10_IMP_REP_sqft_basement	Yes	Input	Interval
LG10_IMP_REP_sqft_living	No	Input	Interval
LG10_IMP_REP_sqft_living15	Yes	Input	Interval
LG10_IMP_REP_sqft_lot	Yes	Input	Interval
LG10_IMP_REP_sqft_lot15	No	Input	Interval
condition	No	Input	Ordinal
date	Yes	Time ID	Nominal
floors	Yes	Input	Interval
grade	No	Input	Ordinal
id	Yes	ID	Interval
price_range	No	Target	Binary
view	No	Input	Ordinal
waterfront	No	Input	Binary
yr_built	Yes	Input	Interval
yr_renovated	No	Input	Interval
zipcode	Yes	Input	Interval

Figure 5.3.7 Drop Variables

## 5.4 MODEL – Data Modelling

After the data are split into training and validation set with a ratio of 50:50, the training data is used to build several models. In this study, we build a decision tree, gradient boosting, logistic regression, and neural network models to predict house prices.

### 5.4.1 Constructing a Decision Tree Predictive Model

Firstly, a decision tree predictive model is built. A decision tree is a supervised model, and it has a tree-based structure that consists of a root node, internal nodes, branches, and leaf nodes. In each node, decision rules are generated to divide the targets to solve the classification task. The decision tree model is preferred as it is simple, quick, and able to manage a high volume of data as we have 21,613 rows in this study.

The decision tree is built as shown below.

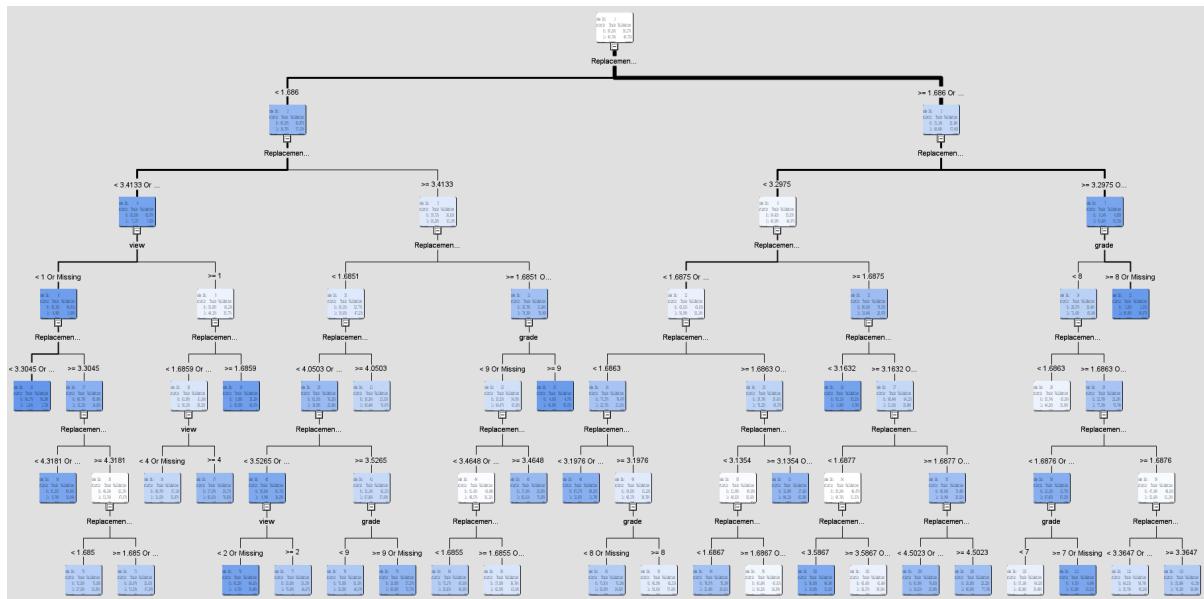


Figure 5.4.1.1: Decision Tree

Based on the decision tree developed, we found several interesting points where (1) latitude has a high level of information gain as it is selected to be the root node and appears in internal nodes several times as it is able to split the data into classes effectively, (2) the shortest split is after 3 layers of rules and (3) at most situations when the grade is high, the price range is high. The list of decision rules is as below:

- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} < 1$  Or Missing and  $\text{sqft\_living} < 3.3045$  Or Missing, then  $\text{price\_range} = \text{low}$ .

- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} < 1$  Or Missing,  $\text{sqft\_living} \geq 3.3045$  and  $\text{sqft\_lot15} < 4.3181$  Or Missing, then  $\text{price\_range} = \text{low}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} < 1$  Or Missing,  $\text{sqft\_living} \geq 3.3045$ ,  $\text{sqft\_lot15} \geq 4.3181$  and  $\text{lat} < 1.685$ , then  $\text{price\_range} = \text{low}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} < 1$  Or Missing,  $\text{sqft\_living} \geq 3.3045$ ,  $\text{sqft\_lot15} \geq 4.3181$  and  $\text{lat} \geq 1.685$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} \geq 1$  and  $\text{lat} \geq 1.6859$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} \geq 1$ ,  $\text{lat} < 1.6859$  Or Missing and  $\text{view} < 4$  Or Missing, then  $\text{price\_range} = \text{low}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} < 3.4133$  Or Missing,  $\text{view} \geq 1$ ,  $\text{lat} < 1.6859$  Or Missing and  $\text{view} \geq 4$  Or Missing, then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} < 1.6851$ ,  $\text{sqft\_lot15} < 4.0503$  Or Missing,  $\text{sqft\_living} < 3.5265$  Or Missing and  $\text{view} < 2$  Or Missing, then  $\text{price\_range} = \text{low}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} < 1.6851$ ,  $\text{sqft\_lot15} < 4.0503$  Or Missing,  $\text{sqft\_living} < 3.5265$  Or Missing and  $\text{view} \geq 2$  Or Missing, then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} < 1.6851$  and  $\text{sqft\_lot15} \geq 4.0503$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} < 1.6851$ ,  $\text{sqft\_lot15} < 4.0503$  Or Missing,  $\text{sqft\_living} \geq 3.5265$  and  $\text{grade} < 9$ , then  $\text{price\_range} = \text{low}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} < 1.6851$ ,  $\text{sqft\_lot15} < 4.0503$  Or Missing,  $\text{sqft\_living} \geq 3.5265$  and  $\text{grade} \geq 9$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} \geq 1.6851$  Or Missing and  $\text{grade} \geq 9$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} \geq 1.6851$  Or Missing,  $\text{grade} < 9$  Or Missing and  $\text{sqft\_living} \geq 3.4648$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} \geq 1.6851$  Or Missing,  $\text{grade} < 9$  Or Missing,  $\text{sqft\_living} < 3.4648$  Or Missing and  $\text{lat} < 1.6855$ , then  $\text{price\_range} = \text{low}$ .
- If  $\text{lat} < 1.686$ ,  $\text{sqft\_living} \geq 3.4133$ ,  $\text{lat} \geq 1.6851$  Or Missing,  $\text{grade} < 9$  Or Missing,  $\text{sqft\_living} < 3.4648$  Or Missing and  $\text{lat} \geq 1.6855$ , then  $\text{price\_range} = \text{high}$ .
- If  $\text{lat} \geq 1.686$  Or Missing,  $\text{sqft\_living} < 3.2975$ ,  $\text{lat} < 1.6875$  Or Missing,  $\text{lat} < 1.6863$  and  $\text{sqft\_living} < 3.1976$  Or Missing, then  $\text{price\_range} = \text{low}$ .

- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $< 1.6875$  Or Missing, lat  $< 1.6863$ , sqft\_living  $\geq 3.1976$  and grade  $< 8$  Or Missing, then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $< 1.6875$  Or Missing, lat  $< 1.6863$ , sqft\_living  $\geq 3.1976$  and grade  $\geq 8$  Or Missing, then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $< 1.6875$  Or Missing, lat  $\geq 1.6863$  Or Missing and sqft\_living  $\geq 3.1354$  Or Missing, then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $< 1.6875$  Or Missing, lat  $\geq 1.6863$  Or Missing, sqft\_living  $< 3.1354$  and lat  $< 1.6867$ , then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $< 1.6875$  Or Missing, lat  $\geq 1.6863$  Or Missing, sqft\_living  $< 3.1354$  and lat  $\geq 1.6867$ , then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $\geq 1.6875$  and sqft\_living  $< 3.1632$ , then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $\geq 1.6875$ , sqft\_living  $\geq 3.1632$  Or Missing, lat  $< 1.6877$  and sqft\_lot15  $< 3.5867$ , then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $\geq 1.6875$ , sqft\_living  $\geq 3.1632$  Or Missing, lat  $< 1.6877$  and sqft\_lot15  $\geq 3.5867$ , then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $\geq 1.6875$ , sqft\_living  $\geq 3.1632$  Or Missing, lat  $\geq 1.6877$  Or Missing and sqft\_lot15  $< 4.5023$  Or Missing, then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $< 3.2975$ , lat  $\geq 1.6875$ , sqft\_living  $\geq 3.1632$  Or Missing, lat  $\geq 1.6877$  Or Missing and sqft\_lot15  $\geq 4.5023$  Or Missing, then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $\geq 3.2975$  and grade  $\geq 8$  Or Missing, then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $\geq 3.2975$  Or Missing, grade  $< 8$  and lat  $< 1.6863$ , then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $\geq 3.2975$  Or Missing, grade  $< 8$ , lat  $\geq 1.6863$ , lat  $< 1.6876$  Or Missing and grade  $< 7$ , then price\_range =low.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $\geq 3.2975$ , grade  $< 8$ , lat  $\geq 1.6863$ , lat  $< 1.6876$  Or Missing and grade  $\geq 7$ , then price\_range =high.
- If lat  $\geq 1.686$  Or Missing, sqft\_living  $\geq 3.2975$  Or Missing, grade  $< 8$ , lat  $\geq 1.6863$  Or Missing, lat  $\geq 1.6876$ , sqft\_living  $< 3.3647$  Or Missing, then price\_range =low.

- If lat  $\geq 1.686$  Or Missing, sqft\_living  $\geq 3.2975$  Or Missing, grade  $< 8$ , lat  $\geq 1.6863$  Or Missing, lat  $\geq 1.6876$ , sqft\_living  $\geq 3.3647$  Or Missing, then price\_range =high.

The decision tree predictive model was built with the train data and the statistical results are obtained as shown below. The result shows that the misclassification rate of the decision tree model is 0.12751.

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price_range	_NOBS_	Sum of Frequencies		10806	10807	
price_range	_MISC_	Misclassification Rate		0.117897	0.12751	
price_range	_MAX_	Maximum Absolute Error		0.981655	0.981655	
price_range	_SSE_	Sum of Squared Errors		1884.061	2050.38	
price_range	_ASE_	Average Squared Error		0.087177	0.094863	
price_range	_RASE_	Root Average Squared Error		0.295257	0.307999	
price_range	_DIV_	Divisor for ASE		21612	21614	
price_range	_DFT_	Total Degrees of Freedom		10806	.	

Figure 5.4.1.2: Statistical output

The variable importance for decision tree is acquired as below. It is found that the three variables that have highest importance are latitude (lat), square footage of the apartments interior living space (sqft\_living) and grade with importance 1.000, 0.8081 and 0.2468 respectively.

Variable Name	Label	Number of Splitting Rules	Importance	Validation Importance	Ratio of Validation to Training Importance
REP_LG10_IMP_REP_lat	Replacement: Transformed: Imputed: Replac...	11	1.0000	1.0000	1.0000
REP_LG10_IMP_REP_sqft_living	Replacement: Transformed: Imputed: Replac...	9	0.8081	0.8390	1.0383
grade		5	0.2468	0.2648	1.0733
REP_LG10_IMP_REP_sqft_lot15	Replacement: Transformed: Imputed: Replac...	4	0.2153	0.2150	0.9984
view		3	0.1921	0.1637	0.8525
REP_LG10_IMP_REP_bedrooms	Replacement: Transformed: Imputed: Replac...	0	0.0000	0.0000	.
REP_yr_renovated	Replacement: yr_renovated	0	0.0000	0.0000	.
condition		0	0.0000	0.0000	.
waterfront		0	0.0000	0.0000	.

Figure 5.4.1.3: Variable importance

### 5.4.2 Constructing a Gradient Boosting Predictive Model

Gradient boosting is a machine learning method used for both regression and classification tasks. In this project, it will be applied specifically for classification tasks. Unlike a single decision tree, gradient boosting is less prone to overfitting the data. The accuracy of a model is crucial, and the gradient boosting algorithm can aid in reducing the bias error of the model.

The Gradient Boosting Node of SAS Enterprise Miner can be utilized to connect and analyze previously divided data. The statistics resulting from running this node are illustrated in the figure below. The results indicate that the misclassification rate of the gradient boosting model produced is 0.122421.

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price_range	_NOBS_	Sum of Frequencies		10806	10807	.
price_range	SUMW	Sum of Case Weig...		21612	21614	.
price_range	_MISC_	Misclassification R...		0.1191	0.122421	.
price_range	_MAX_	Maximum Absolute...		0.962634	0.971617	.
price_range	_SSE_	Sum of Squared Er...		1924.193	1989.676	.
price_range	_ASE_	Average Squared ...		0.089034	0.092055	.
price_range	_RASE_	Root Average Squ...		0.298385	0.303406	.
price_range	_DIV_	Divisor for ASE		21612	21614	.
price_range	_DFT_	Total Degrees of F...		10806	.	.

Figure 5.4.2.1: Fit Statistics

The table below shows the importance of variables in the Gradient Boosting model. The three key factors that have the greatest impact on the results of this model are LG10\_IMP\_REP\_lat, LG10\_IMP\_REP\_sqft\_living, and grade.

Variable Name	Label	Number of Splitting Rules	Importance	Validation Importance	Ratio of Validation to Training Importance
LG10_IMP_REP_lat	Transformed: Imputed...	82	1	1	1
LG10_IMP_REP_sqft_living	Transformed: Imputed...	35	0.751581	0.785669	1.045356
grade		14	0.601595	0.668711	1.111564
view		7	0.154703	0.151371	0.978463
LG10_IMP_REP_sqft_lot15	Transformed: Imputed...	12	0.141388	0.149133	1.054778
LG10_IMP_REP_bedrooms	Transformed: Imputed...	0	0	0	.
yr_renovated		0	0	0	.
waterfront		0	0	0	.
condition		0	0	0	.

Figure 5.4.2.2: Variable Importance

### 5.4.3 Constructing a Logistics Regression Predictive Model

Logistic regression is a statistical model used for classification and predictive analytics. It is used to compute the probability of an event taking place by having the log-odds for the event be a linear combination of one or more independent variables. Logistic regression uses a logistic function to run a regression analysis where the dependent variable is binary ranging from 0 and 1. At the core of logistic regression analysis is the task of estimating the log odds of an event. Logistic regression is used to explain the relationship between of all data type (nominal, interval, ordinal, ratio) of independent variables with binary dependant variables.

The SAS Enterprise miner model number regression was run through the stepwise method.

The regression node was used to link with the data partition node. Logistic regression was selected as the regression type. The figure shows a misclassification rate of. 0.202461 This gives us an accuracy of 79.75%. The result of the statistics test is shown as figure below.

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price_range		_AIC_	Akaike's Information ...	9272.608		
price_range		_ASE_	Average Squared Error	0.138747	0.138657	
price_range		_AVERR_	Average Error Function	0.425995	0.427032	
price_range		_DFE_	Degrees of Freedom ...	10773		
price_range		_DFM_	Model Degrees of Fr...	33		
price_range		_DFT_	Total Degrees of Fre...	10806		
price_range		_DIV_	Divisor for ASE	21612	21614	
price_range		_ERR_	Error Function	9206.608	9229.877	
price_range		_FPE_	Final Prediction Error	0.139597		
price_range		_MAX_	Maximum Absolute E...	0.99586	0.999712	
price_range		_MSE_	Mean Square Error	0.139172	0.138657	
price_range		_NOBS_	Sum of Frequencies	10806	10807	
price_range		_NW_	Number of Estimate ...	33		
price_range		_RASE_	Root Average Sum of...	0.372487	0.372367	
price_range		_RFPE_	Root Final Prediction ...	0.373627		
price_range		_RMSE_	Root Mean Squared ...	0.373057	0.372367	
price_range		_SBC_	Schwarz's Bayesian ...	9513.107		
price_range		_SSE_	Sum of Squared Errors	2998.597	2996.932	
price_range		_SUMW_	Sum of Case Weight...	21612	21614	
price_range		_MISC_	Misclassification Rate	0.201832	0.202461	

Figure 5.4.3.1 Statistical results for Logistic Regression

The Odds Ratio Estimates shows on the strength an event is associated with exposure where in this scenario would be on the price range.

Odds Ratio Estimates		
Effect	price_range	Point Estimate
LG10_IMP REP_bathrooms	1	4.283
LG10_IMP REP_bedrooms	1	0.103
LG10_IMP REP_lat	1	.
LG10_IMP REP_long	1	29.239
LG10_IMP REP_sqft_above	1	9.406
LG10_IMP REP_sqft_basement	1	1.464
LG10_IMP REP_sqft_living	1	3.435
LG10_IMP REP_sqft_living15	1	57.098
LG10_IMP REP_sqft_lot	1	0.873
LG10_IMP REP_sqft_lot15	1	0.452
condition	1 vs 5	0.412
condition	2 vs 5	0.357
condition	3 vs 5	0.583
condition	4 vs 5	0.629
floors	1	1.962

Figure 5.4.3.2 Odd Ratio Estimates for Logistic Regression

Based on the output of Odds Ratio Estimates from we found that:

- 1) Sqft\_living15 has 57 times the odds of having higher price range than lower price range level.
- 2) Long(Longitude) has 29 times the odds of having higher price range than lower price range level.
- 3) Sqft\_above has 9 times the odd of having higher price range than lower price range level.

#### 5.4.4 Constructing a Neural Network Predictive Model

A neural network is a type of machine learning process (deep learning) and a collection of algorithms that employ linked neurons or nodes in a layered framework to mimic the human brain. Therefore, it develops an adaptive system that continually learns from previous errors and improvements. With this adaptation to the changes of input, the network can achieve success without having to change the output criterion.

The existing neural network node of SAS Enterprise Miner can be utilized to connect and analyze previously divided data. Different hidden layers with optimized hidden units of neural network models are drawn and shown below. The statistics resulting from running neural network node are also illustrated below in figures.

From the statistics below, the three models with the lowest misclassification rate are the number of hidden units of 5, 3 and 4. The misclassification rate of a neural network model with several hidden units of 5 is 0.111872. The misclassification rate of a neural network model with several hidden units of 3 is 0.11363. The misclassification rate of a neural network model with several hidden units of 4 is 0.113723.

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10724	.	.
price range		DFM	Model Degrees of ...	82	.	.
price range		NW	Number of Estima...	82	.	.
price range		AIC	Akaike's Informati...	5785.875	.	.
price range		SBC	Schwarz's Bayesi...	6383.479	.	.
price range		ASE	Average Squared ...	0.079959	0.081487	.
price range		MAX	Maximum Absolut...	0.999171	0.99956	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.282771	0.285459	.
price range		SSE	Sum of Squared E...	1728.079	1761.256	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.081182	.	.
price range		MSE	Mean Squared Error	0.080571	0.081487	.
price range		RFPE	Root Final Predicti...	0.284925	.	.
price range		RMSE	Root Mean Squar...	0.28385	0.285459	.
price range		AVERR	Average Error Fun...	0.260127	0.265749	.
price range		ERR	Error Function	5621.875	5743.899	.
price range		MISC	Misclassification ...	0.110124	0.11363	.
price range		WRONG	Number of Wrong ...	1190	1228	.

Figure 5.4.4.1: Statistical results for 3 layers neural network (Misclassification rate= 0.11363)

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10697	.	.
price range		DFM	Model Degrees of ...	109	.	.
price range		NW	Number of Estima...	109	.	.
price range		AIC	Akaike's Informati...	5906.991	.	.
price range		SBC	Schwarz's Bayesi...	6701.368	.	.
price range		ASE	Average Squared ...	0.080438	0.081345	.
price range		MAX	Maximum Absolut...	0.999294	0.999642	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.283616	0.285211	.
price range		SSE	Sum of Squared E...	1738.43	1758.198	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.082077	.	.
price range		MSE	Mean Squared Error	0.081258	0.081345	.
price range		RFPE	Root Final Predicti...	0.286492	.	.
price range		RMSE	Root Mean Squar...	0.285058	0.285211	.
price range		AVERR	Average Error Fun...	0.263233	0.266123	.
price range		ERR	Error Function	5688.991	5751.972	.
price range		MISC	Misclassification ...	0.109754	0.113723	.
price range		WRONG	Number of Wrong ...	1186	1229	.

Figure 5.4.4.2: Statistical results for 4 layers neural network (Misclassification rate= 0.113723)

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10670	.	.
price range		DFM	Model Degrees of ...	136	.	.
price range		NW	Number of Estima...	136	.	.
price range		AIC	Akaike's Informati...	6033.828	.	.
price range		SBC	Schwarz's Bayesi...	7024.977	.	.
price range		ASE	Average Squared ...	0.081782	0.082381	.
price range		MAX	Maximum Absolut...	0.999582	0.999534	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.285976	0.28702	.
price range		SSE	Sum of Squared E...	1767.474	1780.575	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.083867	.	.
price range		MSE	Mean Squared Error	0.082824	0.082381	.
price range		RFPE	Root Final Predicti...	0.289598	.	.
price range		RMSE	Root Mean Squar...	0.287792	0.28702	.
price range		AVERR	Average Error Fun...	0.266603	0.270603	.
price range		ERR	Error Function	5761.828	5848.813	.
price range		MISC	Misclassification ...	0.112345	0.111972	.
price range		WRONG	Number of Wrong ...	1214	1209	.

Figure 5.4.4.3: Statistical results for 5 layers neural network (Misclassification rate= 0.111872)

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10535	.	.
price range		DFM	Model Degrees of ...	271	.	.
price range		NW	Number of Estima...	271	.	.
price range		AIC	Akaike's Informati...	6580.888	.	.
price range		SBC	Schwarz's Bayesi...	8555.897	.	.
price range		ASE	Average Squared ...	0.087004	0.08806	.
price range		MAX	Maximum Absolut...	0.998638	0.998207	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.294964	0.296749	.
price range		SSE	Sum of Squared E...	1880.332	1903.322	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.09148	.	.
price range		MSE	Mean Squared Error	0.089242	0.08806	.
price range		RFPE	Root Final Predicti...	0.302457	.	.
price range		RMSE	Root Mean Squar...	0.298734	0.296749	.
price range		AVERR	Average Error Fun...	0.279423	0.28473	.
price range		ERR	Error Function	6038.888	6154.147	.
price range		MISC	Misclassification ...	0.123357	0.122791	.
price range		WRONG	Number of Wrong ...	1333	1327	.

Figure 5.4.4.4: Statistical results for 10 layers neural network (Misclassification rate= 0.122791)

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10400	.	.
price range		DFM	Model Degrees of ...	406	.	.
price range		NW	Number of Estima...	406	.	.
price range		AIC	Akaike's Informati...	6588.019	.	.
price range		SBC	Schwarz's Bayesi...	9546.889	.	.
price range		ASE	Average Squared ...	0.082742	0.08381	.
price range		MAX	Maximum Absolut...	0.99859	0.999119	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.287648	0.2895	.
price range		SSE	Sum of Squared E...	1788.212	1811.476	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.089202	.	.
price range		MSE	Mean Squared Error	0.085972	0.08381	.
price range		RFPE	Root Final Predicti...	0.298867	.	.
price range		RMSE	Root Mean Squar...	0.293209	0.2895	.
price range		AVERR	Average Error Fun...	0.26726	0.271869	.
price range		ERR	Error Function	5776.019	5876.174	.
price range		MISC	Misclassification ...	0.113363	0.115758	.
price range		WRONG	Number of Wrong ...	1225	1251	.

Figure 5.4.4.5: Statistical results for 15 layers neural network (Misclassification rate= 0.115758)

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10265	.	.
price range		DFM	Model Degrees of ...	541	.	.
price range		NW	Number of Estima...	541	.	.
price range		AIC	Akaike's Informati...	7541.683	.	.
price range		SBC	Schwarz's Bayesi...	11484.41	.	.
price range		ASE	Average Squared ...	0.094492	0.096462	.
price range		MAX	Maximum Absolut...	0.99644	0.998993	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.307396	0.310583	.
price range		SSE	Sum of Squared E...	2042.162	2084.925	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.104452	.	.
price range		MSE	Mean Squared Error	0.099472	0.096462	.
price range		RFPE	Root Final Predicti...	0.323191	.	.
price range		RMSE	Root Mean Squar...	0.315392	0.310583	.
price range		AVERR	Average Error Fun...	0.298893	0.305257	.
price range		ERR	Error Function	6459.683	6597.816	.
price range		MISC	Misclassification ...	0.13548	0.138244	.
price range		WRONG	Number of Wrong ...	1464	1494	.

Figure 5.4.4.6: Statistical results for 20 layers neural network (Misclassification rate= 0.138244)

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	10130	.	.
price range		DFM	Model Degrees of ...	676	.	.
price range		NW	Number of Estima...	676	.	.
price range		AIC	Akaike's Informati...	8316.96	.	.
price range		SBC	Schwarz's Bayesi...	13243.55	.	.
price range		ASE	Average Squared ...	0.102997	0.105123	.
price range		MAX	Maximum Absolut...	0.997256	0.997939	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.320931	0.324227	.
price range		SSE	Sum of Squared E...	2225.963	2272.138	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.116743	.	.
price range		MSE	Mean Squared Error	0.10987	0.105123	.
price range		RFPE	Root Final Predicti...	0.341677	.	.
price range		RMSE	Root Mean Squar...	0.331466	0.324227	.
price range		AVERR	Average Error Fun...	0.322273	0.32943	.
price range		ERR	Error Function	6964.96	7120.309	.
price range		MISC	Misclassification ...	0.147881	0.152679	.
price range		WRONG	Number of Wrong ...	1598	1650	.

Figure 5.4.4.7: Statistical results for 25 layers neural network (Misclassification rate= 0.152679)

Fit Statistics						
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	9995	.	.
price range		DFM	Model Degrees of ...	811	.	.
price range		NW	Number of Estima...	811	.	.
price range		AIC	Akaike's Informati...	8162.731	.	.
price range		SBC	Schwarz's Bayesi...	14073.18	.	.
price range		ASE	Average Squared ...	0.095985	0.097755	.
price range		MAX	Maximum Absolut...	0.992798	0.998342	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.309815	0.312658	.
price range		SSE	Sum of Squared E...	2074.438	2112.874	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.111562	.	.
price range		MSE	Mean Squared Error	0.103774	0.097755	.
price range		RFPE	Root Final Predicti...	0.334009	.	.
price range		RMSE	Root Mean Squar...	0.322139	0.312658	.
price range		AVERR	Average Error Fun...	0.302643	0.308684	.
price range		ERR	Error Function	6540.731	6671.903	.
price range		MISC	Misclassification ...	0.136868	0.143888	.
price range		WRONG	Number of Wrong ...	1479	1555	.

Figure 5.4.4.8: Statistical results for 30 layers neural network (Misclassification rate= 0.143888)

Fit Statistics						
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	9860	.	.
price range		DFM	Model Degrees of ...	946	.	.
price range		NW	Number of Estima...	946	.	.
price range		AIC	Akaike's Informati...	8255.297	.	.
price range		SBC	Schwarz's Bayesi...	15149.61	.	.
price range		ASE	Average Squared ...	0.093227	0.094524	.
price range		MAX	Maximum Absolut...	0.992583	0.99905	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.305331	0.307447	.
price range		SSE	Sum of Squared E...	2014.821	2043.038	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.111116	.	.
price range		MSE	Mean Squared Error	0.102171	0.094524	.
price range		RFPE	Root Final Predicti...	0.333341	.	.
price range		RMSE	Root Mean Squar...	0.319643	0.307447	.
price range		AVERR	Average Error Fun...	0.294434	0.30045	.
price range		ERR	Error Function	6363.297	6493.917	.
price range		MISC	Misclassification ...	0.134277	0.13445	.
price range		WRONG	Number of Wrong ...	1451	1453	.

Figure 5.4.4.9: Statistical results for 35 layers neural network (Misclassification rate= 0.13445)

Fit Statistics						
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
price range		DFT	Total Degrees of ...	10806	.	.
price range		DFE	Degrees of Freed...	9725	.	.
price range		DFM	Model Degrees of ...	1081	.	.
price range		NW	Number of Estima...	1081	.	.
price range		AIC	Akaike's Informati...	9016.2	.	.
price range		SBC	Schwarz's Bayesi...	16894.37	.	.
price range		ASE	Average Squared ...	0.101709	0.103186	.
price range		MAX	Maximum Absolut...	0.993489	0.997395	.
price range		DIV	Divisor for ASE	21612	21614	.
price range		NOBS	Sum of Frequencies	10806	10807	.
price range		RASE	Root Average Squ...	0.318918	0.321226	.
price range		SSE	Sum of Squared E...	2198.128	2230.262	.
price range		SUMW	Sum of Case Wei...	21612	21614	.
price range		FPE	Final Prediction Er...	0.12432	.	.
price range		MSE	Mean Squared Error	0.113014	0.103186	.
price range		RFPE	Root Final Predicti...	0.35259	.	.
price range		RMSE	Root Mean Squar...	0.336176	0.321226	.
price range		AVERR	Average Error Fun...	0.317148	0.322847	.
price range		ERR	Error Function	6854.2	6978.019	.
price range		MISC	Misclassification ...	0.147418	0.150828	.
price range		WRONG	Number of Wrong ...	1593	1630	.

Figure 5.4.4.10: Statistical results for 40 layers neural network (Misclassification rate= 0.150828)

## 5.5 ASSESS – Data Assessment

At first, we compare the performance between Decision Tree, Logistic Regression, and Gradient Boosting. We observed that the Gradient Boosting model has the lowest misclassification rate of 0.122421, resulting with the highest accuracy of 88%, a very close margin with Decision Tree with 0.12751 misclassification rate, resulting in 87% of accuracy.

Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Train: Sum of Frequencies	Train: Misclassification Rate	Selection Criterion: Valid: Misclassification Rate
Y	Boost	Boost	Gradient Bo...	price_range	10806	0.1191	0.122421
	Tree	Tree	Decision Tr...	price_range	10806	0.117897	0.12751
	Reg	Reg	Regression	price_range	10806	0.201832	0.202461

Figure 5.5.1: Performance between Decision Tree, Logistic Regression, Gradient Boosting

Second, we compare the performance between the neural networks. The 5 layers of neural network has the lowest misclassification rate of 0.111872 and therefore achieve an accuracy of 88.81%

Model Node	Model Description	Target Variable	Target Label	Selection Criterion: Valid: Misclassification Rate
Neural3	5 layer N...	price ran...		0.111872
Neural	3 layer N...	price ran...		0.11363
Neural2	4 layer N...	price ran...		0.113723
Neural5	15 layer ...	price ran...		0.115758
Neural9	10 layer ...	price ran...		0.122791
Neural7	35 layer ...	price ran...		0.13445
Neural10	20 layer ...	price ran...		0.138244
Neural8	30 layer ...	price ran...		0.143888
Neural6	40 layer ...	price ran...		0.150828
Neural4	25 layer ...	price ran...		0.152679

Figure 5.5.2 Performance between the neural networks

Finally, we compare the two best models from the above which is the Gradient Boosting Model and the 5-layer Neural Network. We observe that 5-layer Neural Network has the best classification accuracy with 89%.

Model Node	Model Description	Target Variable	Target Label	Selection Criterion: Valid: Misclassification Rate
Neural11	5 layer N...	price ran...		0.111872
Boost	Gradient ...	price ran...		0.122421

Figure 5.5.3: Performance between 5-layer Neural Network and Gradient Boosting

## 6 Conclusion

After we revised the data, there are 21 variables in the house sales dataset including 16 interval variables, 1 binary variable, 3 ordinal variables, and 1 nominal variable they became as follows:

Role	Variable	Count
ID	Interval	1
Input	Binary	1
Input	Interval	14
Input	Ordinal	3
Target	Interval	1
Time ID	Nominal	1

During data exploration, we found incomplete data, noisy data, inconsistent data, and intentional data. These variables are as follows:

Variable	Error Type
bathrooms	Incomplete
yr_built	Inconsistent
bathrooms bedrooms sqft_above price lat long sqft_living15 sqft_lot15 sqft_living sqft_lot sqft_basement	Noisy

SAS Enterprise Miner was used to locate the relevant variables as listed in the table below using the variable selection tool.

Number	Variable
1	grade
2	lat
3	sqft_living
4	view
5	waterfront
6	condition
7	yr_renovated
8	sqft_lot15
9	bedrooms

Through the Modelling and Assessment Phase of SEMMA, we are able to build several predictive models such as Decision Tree, Gradient Boosting, Logistic Regression, and Neural Networks and subsequently, our results find that 5-layer Neural Network produces the best accuracy.

Overall, through SEMMA, we find some interesting patterns that stands out.

<b>Interesting Patterns</b>	<b>Found In</b>
Higher grade level shows a high positive relationship with price. Longer boxplot body length and higher price were observed as the grade increased	Exploration
Sqft_living shows a positive relationship with price	Exploration
Sqft_living and sqft_above show a strong correlation coefficient, 0.8766.	Exploration
Latitude is selected as the root node as it has a high level of information gain	Modelling Phase – Decision Tree
When the grade is high, the price range is more likely to be high.	Modelling Phase – Decision Tree

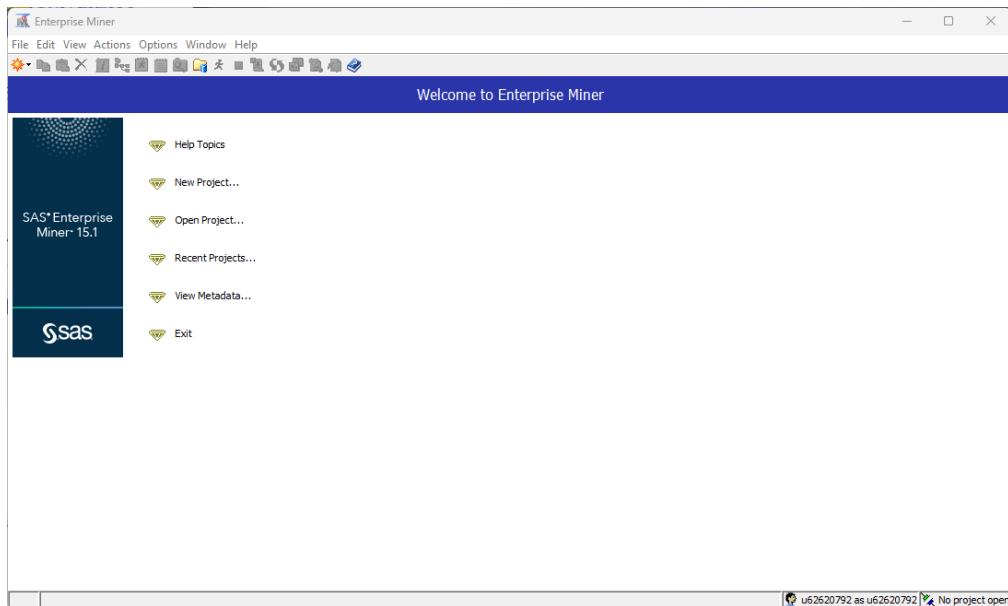
## 7 Appendix

In 7.1 and 7.2, there are steps consisting of the first 2 characters of the SAS SEMMA for creating a SAS Enterprise Miner project. For 7.3, it shows the full correlation table.

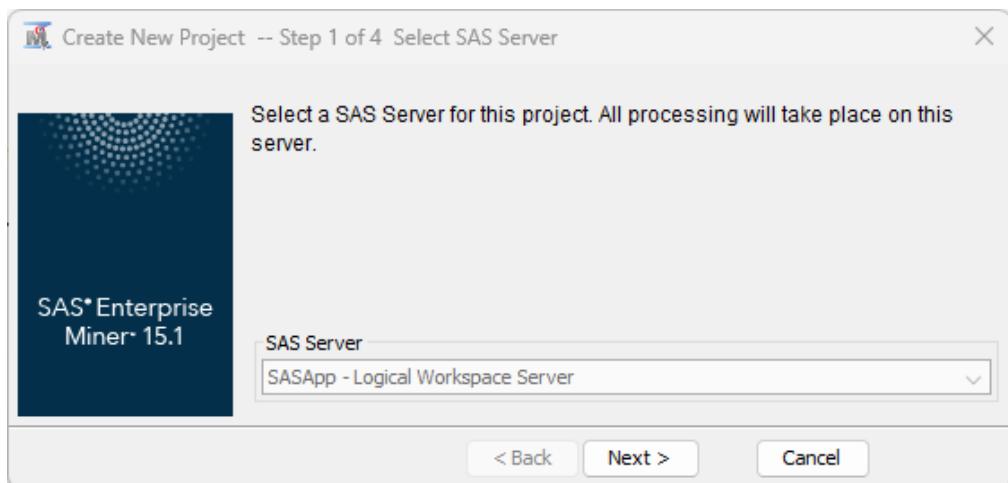
### 7.1 SAMPLE

#### 7.1.1 Create a SAS Enterprise Miner Project

- Select **File** → **New** → **Project** from the main menu or select **New Project** directly

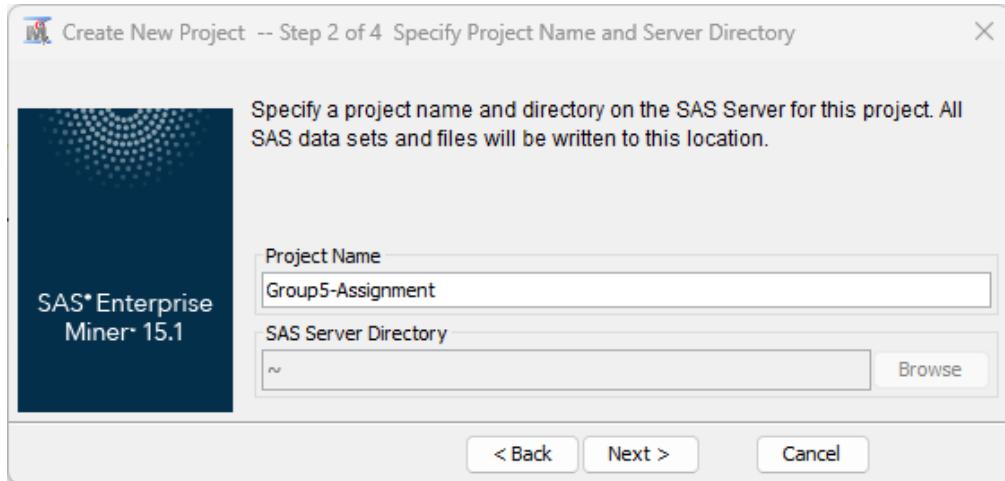


- Step 1 of 4 - Select SAS Server. We will use the default server. Select **Next >**

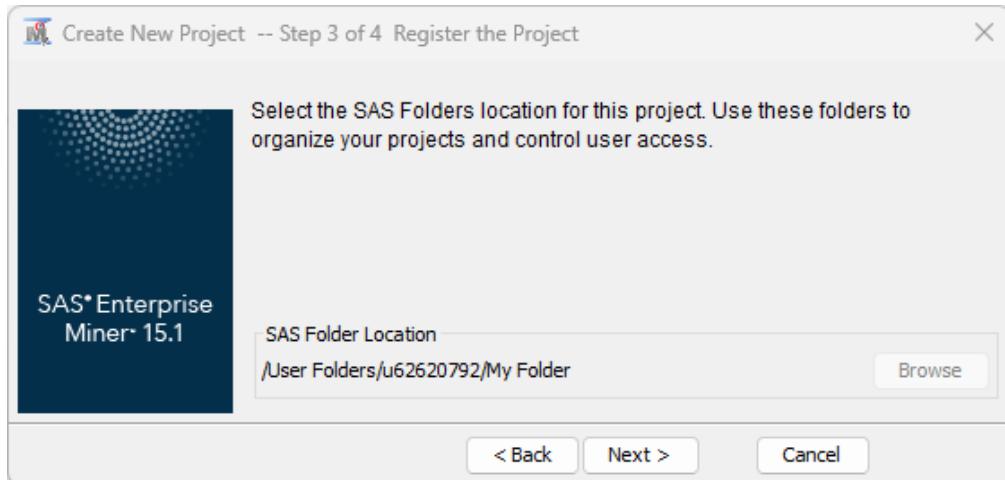


- Step 2 of 4 - Specify Project Name and Server Directory.

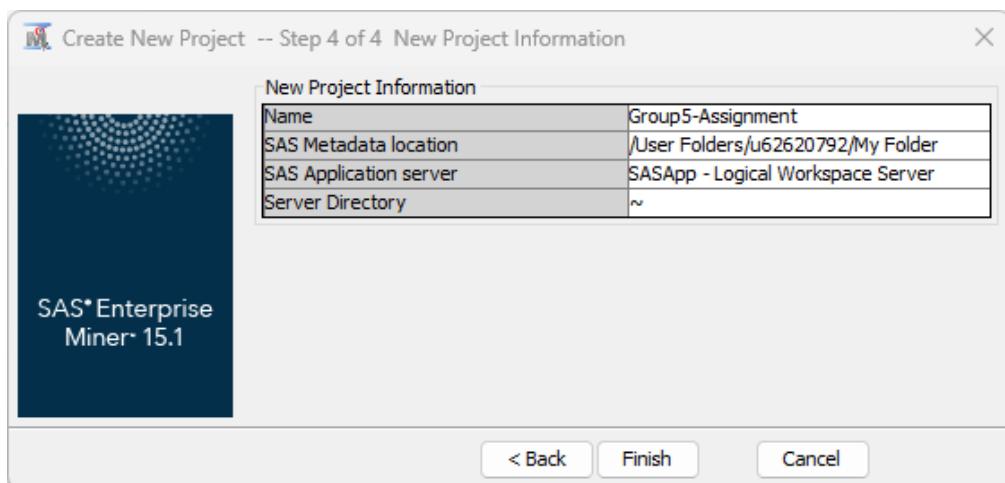
Type Project Name: **Group5-Assignment** and Select **Next >**



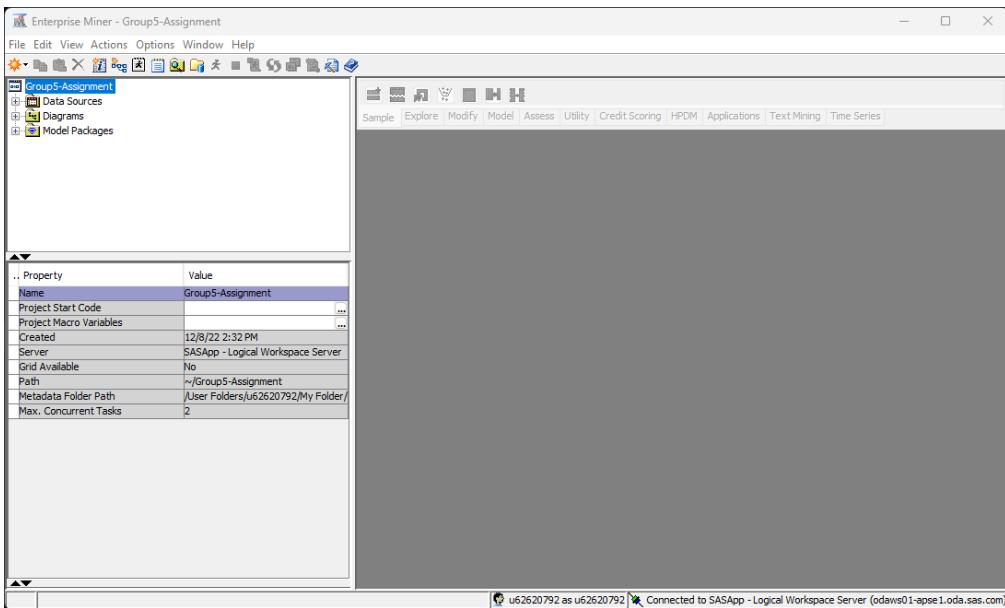
- Step 3 of 4 - Register the Project. Select **Next >**



- Step 4 of 4 - New Project Information. Select **Finish**

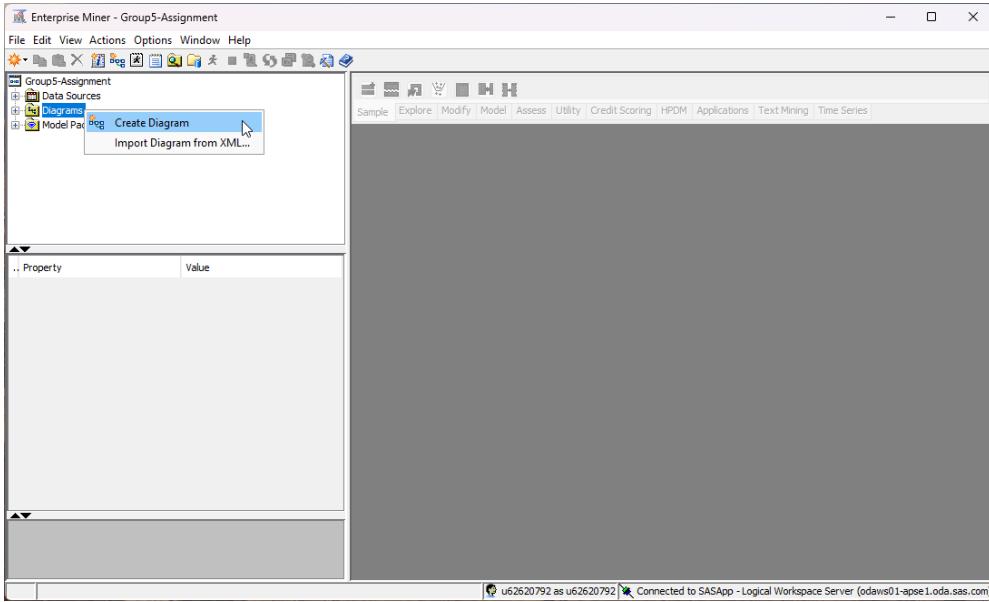


- The SAS Enterprise Miner Project is created successfully.

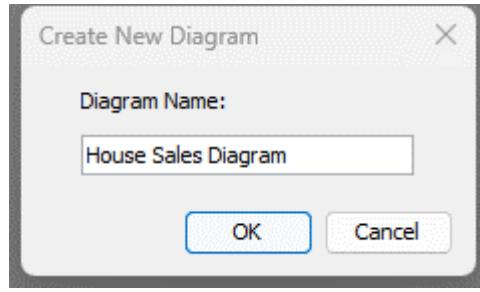


### 7.1.2 Create a SAS Enterprise Miner Diagram

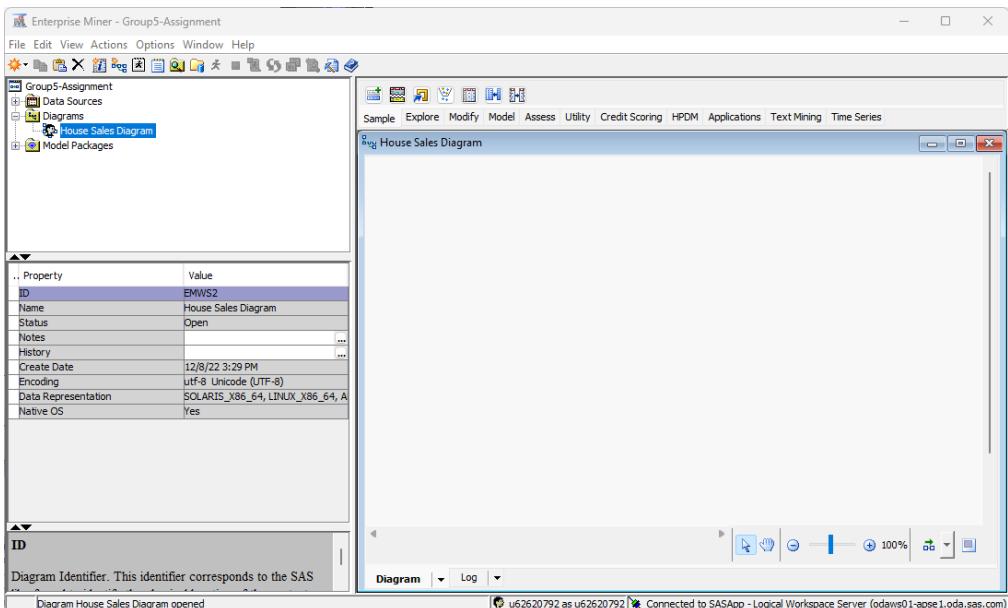
- Select **File → New → Diagram...** from the Main Menu or  
Select Right-click **Diagrams** directly, Select **Create Diagram**



- Type Diagram Name: **House Sales Diagram** and Select **OK**

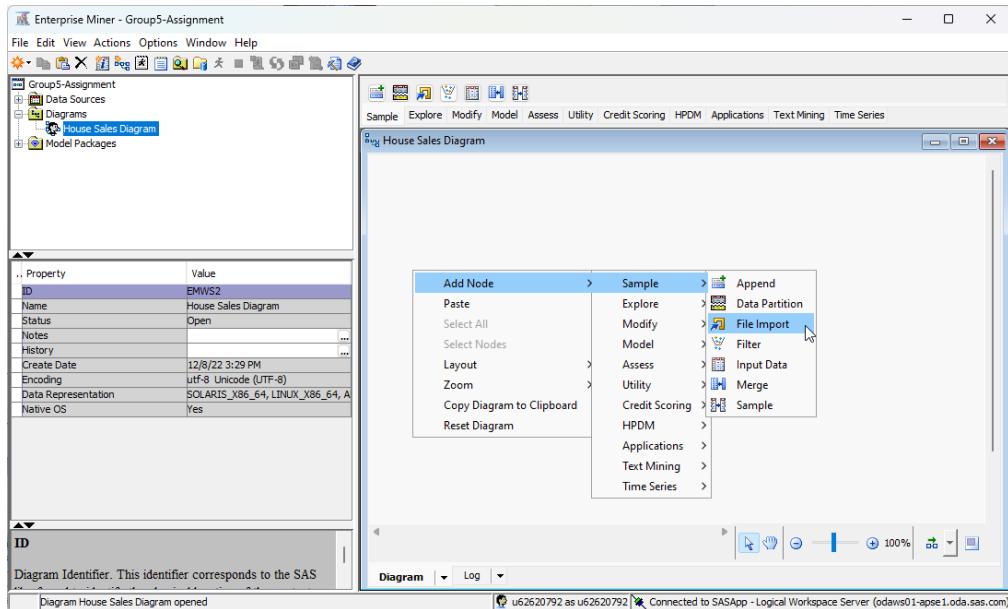


- **House Sales Diagram** is created successfully



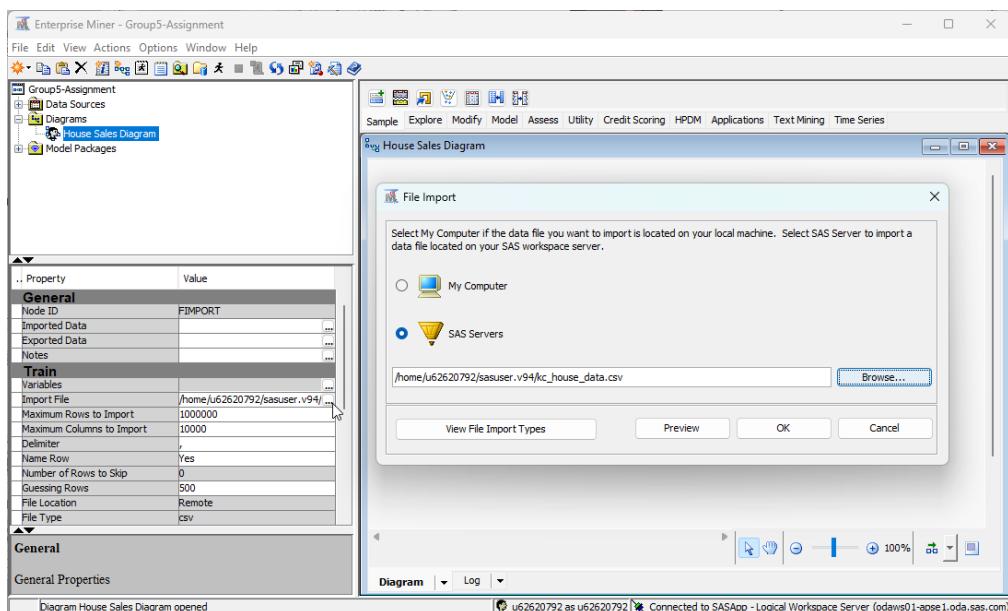
- Create a Sample Node to import the dataset into the SAS

**Right click on House Sales Diagram Window and select Add Note → Sample → File Import**

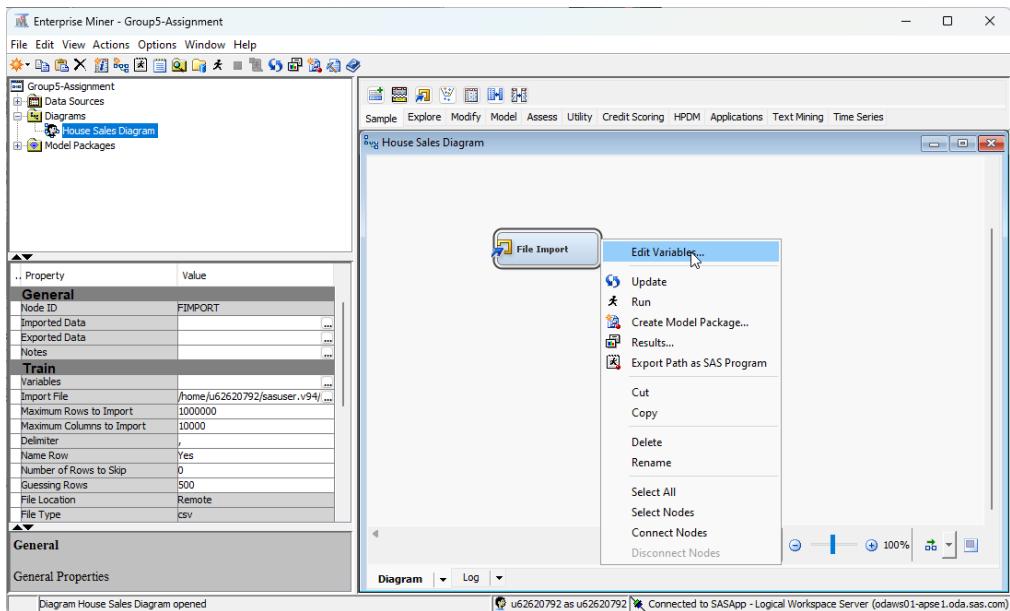


- Select Import File under Train on the left sidebar.

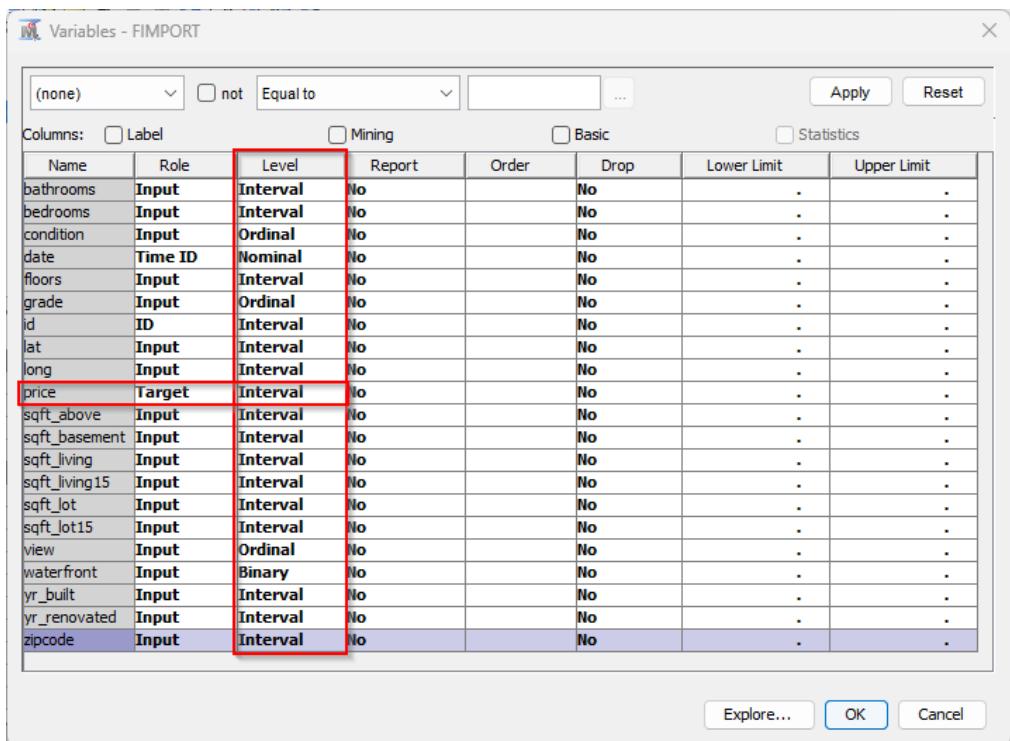
Browse the dataset directory, and select **OK**



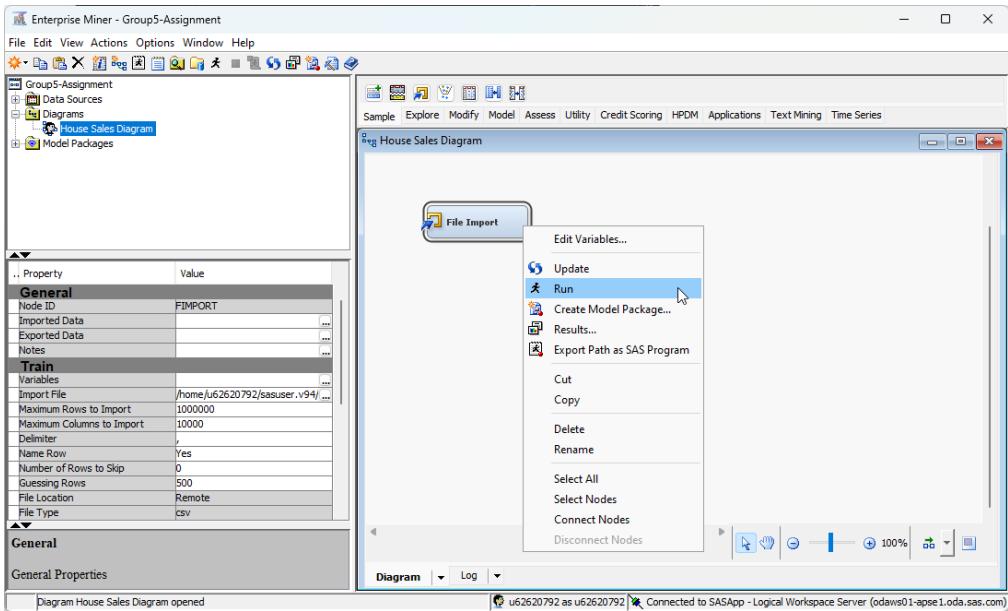
- Right-click the **File Import** node and select Edit Variables...



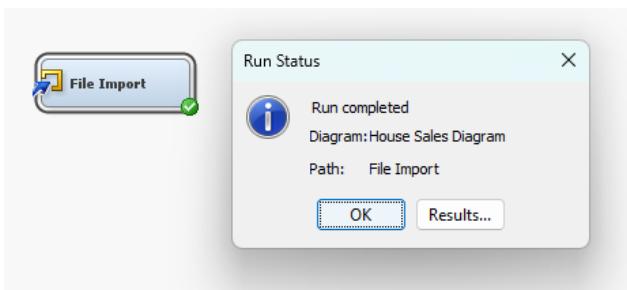
- Set *price* as the role of Target and update levels of the attributes because SAS Enterprise Miner will default to set all the numeric variables as interval and character variables as Nominal variables. After updating all attributes, select OK



- Right-click the **File Import** node and select **Run**



- Once the run has completed, a green arrow is displayed on the **File Import** node, indicating that the file was imported successfully. Select **Results...**



- **Results Window** shown. Once verified, close the window.

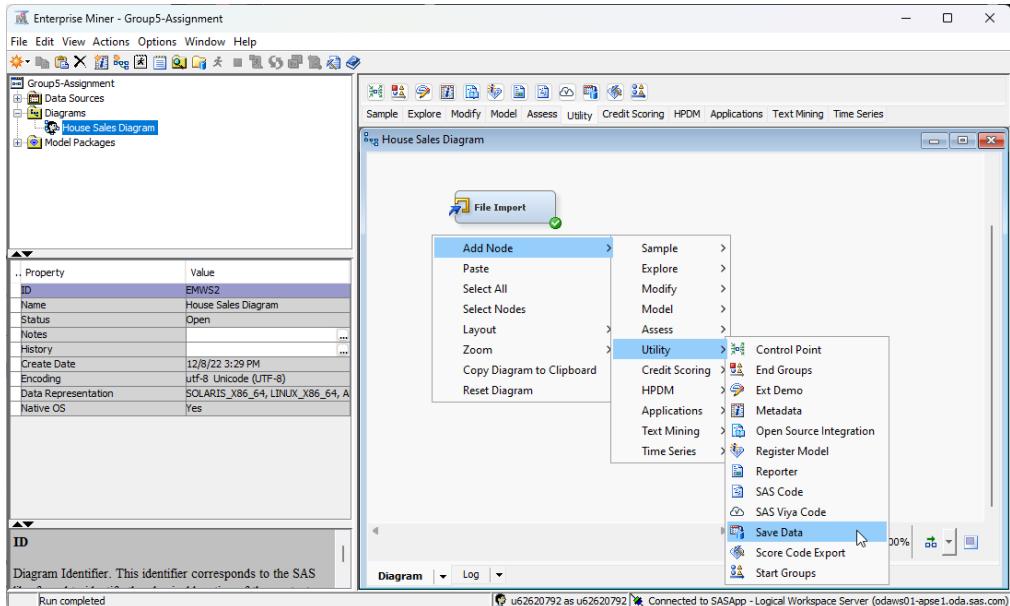
File Edit View Window

```

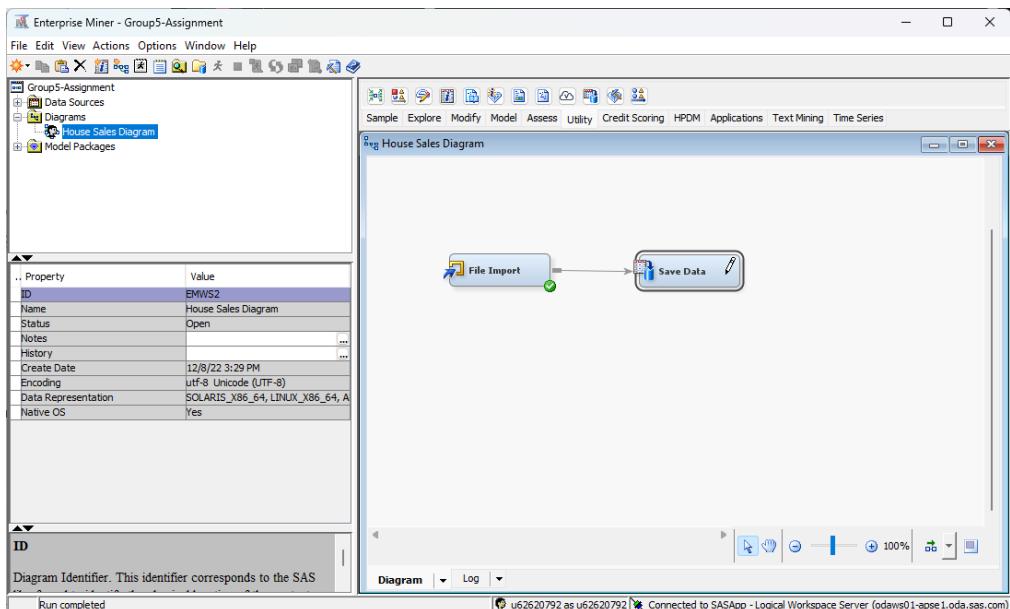
59
60      Alphabetic List of Variables and Attributes
61
62      #  Variable       Type    Len   Format   Informat
63
64      5  bathrooms     Num     8    BEST12.  BEST32.
65      4  bedrooms      Num     8    BEST12.  BEST32.
66      11 condition     Num     8    BEST12.  BEST32.
67      2  date          Char    15   $15.    $15.
68      8  floors         Num     8    BEST12.  BEST32.
69      12 grade          Num     8    BEST12.  BEST32.
70      1  id             Num     8    BEST12.  BEST32.
71      18 lat            Num     8    BEST12.  BEST32.
72      19 long           Num     8    BEST12.  BEST32.
73      3  price           Num     8    BEST12.  BEST32.
74      13 sqft_above     Num     8    BEST12.  BEST32.
75      14 sqft_basement Num     8    BEST12.  BEST32.
76      6  sqft_living    Num     8    BEST12.  BEST32.
77      20 sqft_living15 Num     8    BEST12.  BEST32.
78      7  sqft_lot        Num     8    BEST12.  BEST32.
79      21 sqft_lot15     Num     8    BEST12.  BEST32.
80      10 view            Num     8    BEST12.  BEST32.
81      9  waterfront      Num     8    BEST12.  BEST32.
82      15 yr_builtin     Num     8    BEST12.  BEST32.
83      16 yr_renovated   Num     8    BEST12.  BEST32.
84      17 zipcode         Num     8    BEST12.  BEST32.
85
86
87  *-----*
88  * Score Output
89  *-----*
90
91
92  *-----*
93  * Report Output
94  *-----*
95
96
97
98
99  Exported Attributes for TRAIN Port
100
101      Measurement   Frequency
102      Role          Level      Count
103
104      ID            INTERVAL   1
105      INPUT          BINARY     1
106      INPUT          INTERVAL   14
107      INPUT          ORDINAL   3
108      TARGET          INTERVAL   1
109      TIMEID         NOMINAL   1
110

```

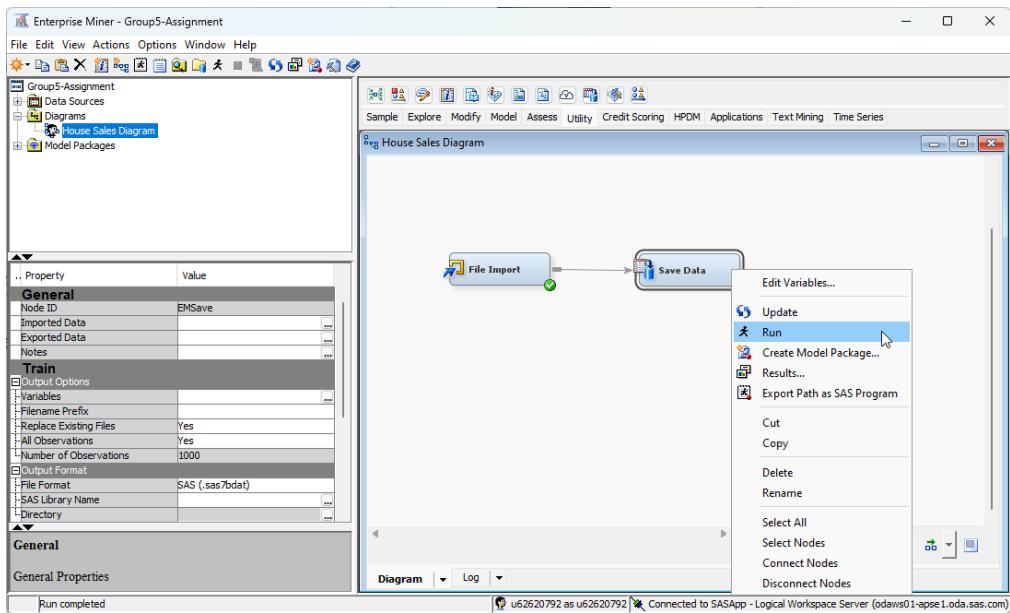
- Right click on House Sales Diagram Window and select Add Note → Utility → Save Data



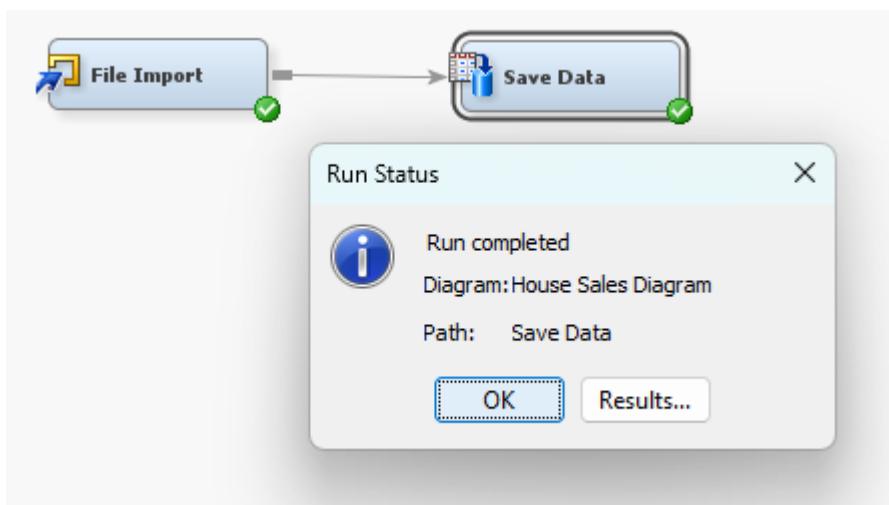
- Connect the Two Nodes



- Right-click the **Save Data** node and select **Run**

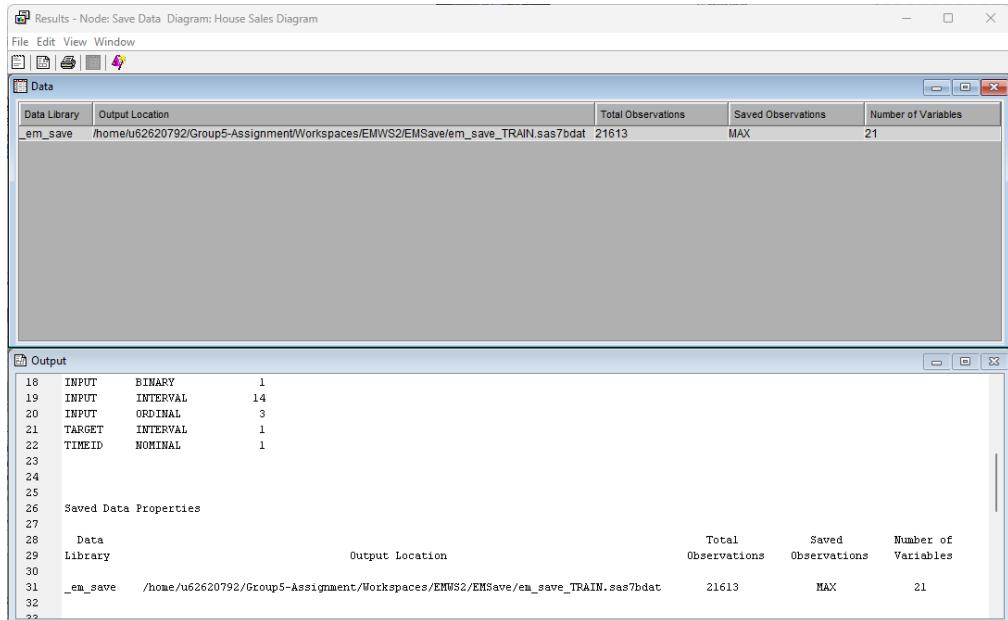


- Once the run has completed, a green arrow is displayed on the **Save Data** node, indicating that the file was imported successfully. Select **Results...**



- **Results Window** shown. The **Output Location** of the .sas7bdat file is identified.

Once verified, close the window.



The screenshot shows the SAS Results window with two main panes: Data and Output.

**Data Pane:**

Data Library	Output Location	Total Observations	Saved Observations	Number of Variables
_em_save	/home/u62620792/Group5-Assignment/Workspaces/EMWS2/EMSave/_em_save_TRAIN.sas7bdat	21613	MAX	21

**Output Pane:**

```

18 INPUT    BINARY      1
19 INPUT    INTERVAL    14
20 INPUT    ORDINAL     3
21 TARGET   INTERVAL    1
22 TIMEID   NOMINAL    1
23
24
25
26 Saved Data Properties
27
28 Data
29 Library          Output Location           Total       Saved      Number of
30                           Observations  Observations  Variables
31 _em_save        /home/u62620792/Group5-Assignment/Workspaces/EMWS2/EMSave/_em_save_TRAIN.sas7bdat  21613      MAX         21
32
33

```

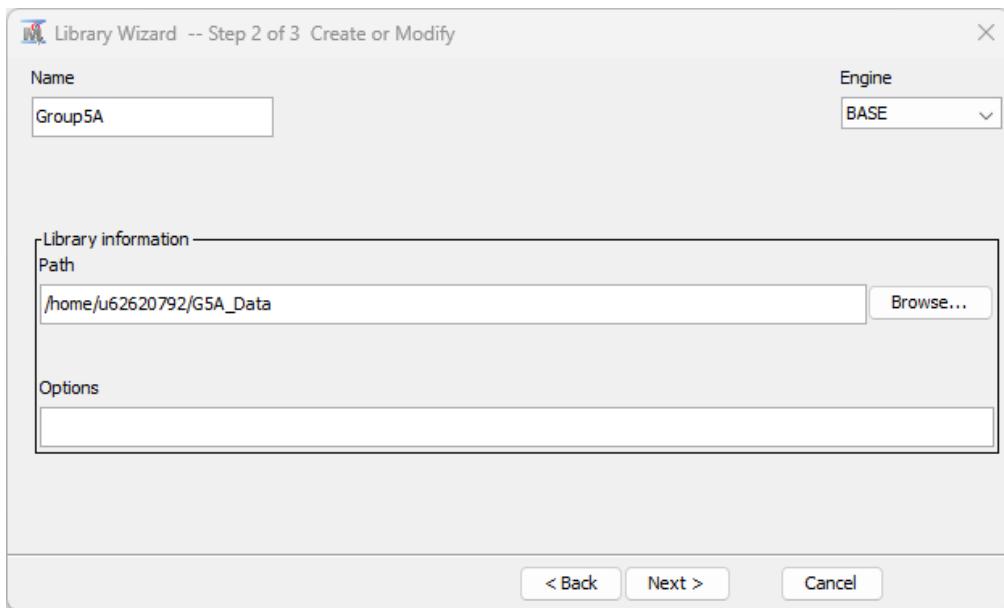
### 7.1.3 Create a SAS Enterprise Miner Library

- Select **File** → **New** → **Library** from the Main Menu.

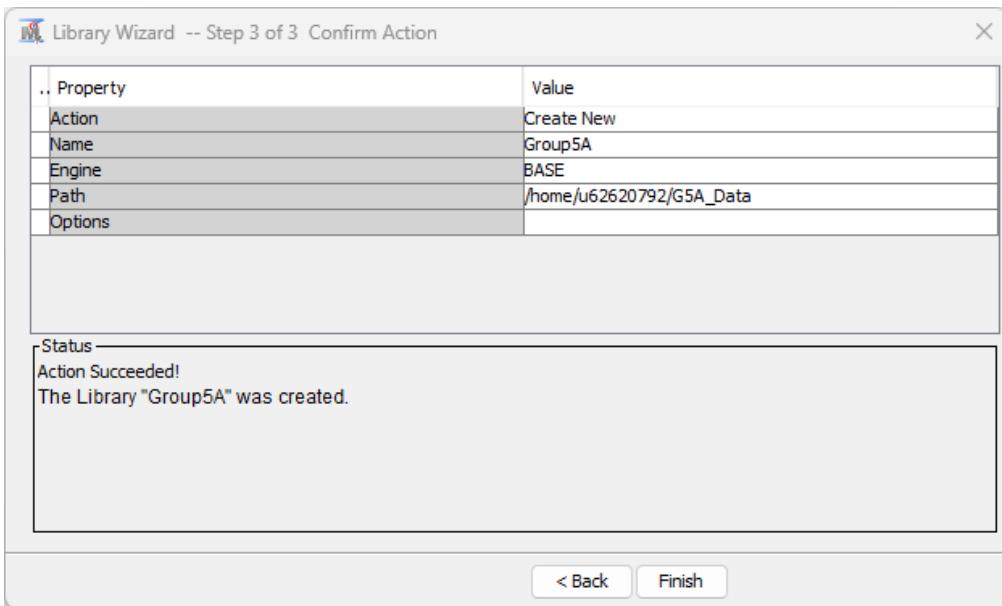
Select the **Create New Library** and **Next >**



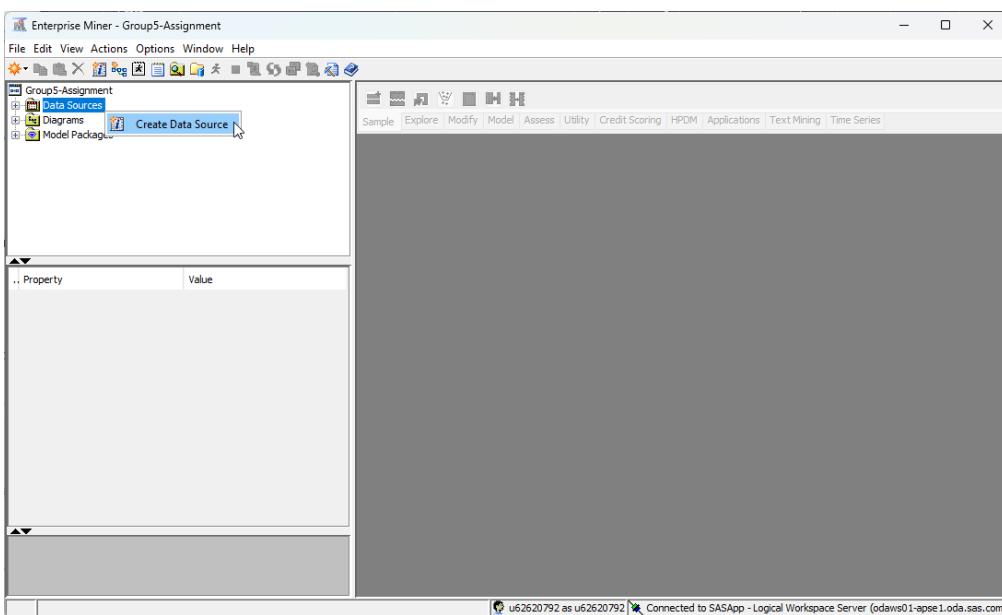
- Type the Library Name: **Group5A**. Select **Browse...** and input Library Path. Then select **Next >**



- The Library *Group5A* was created. Select **Finish**.



- Right-click the **Data Sources** and select **Create Data Source**.



- Step 1 of 8 – Metadata Source. Select **SAS Table** as source and then, select **Next >**



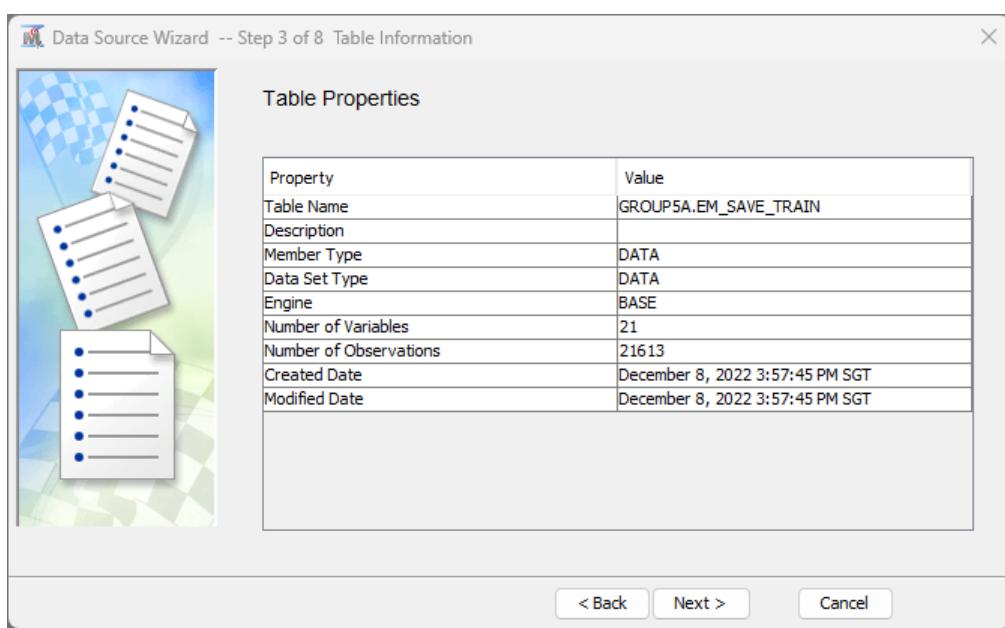
- Step 2 of 8 – Select a SAS Table. **Browse** the table and select target dataset. Then, Select **OK**

Name	Type
Em_save_train	Table

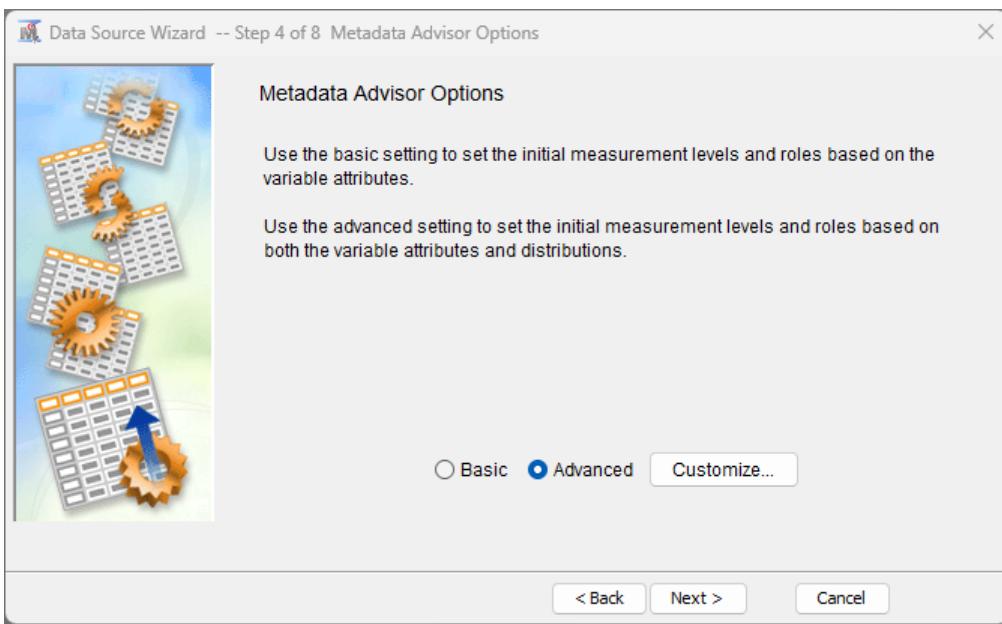
- Select Next >



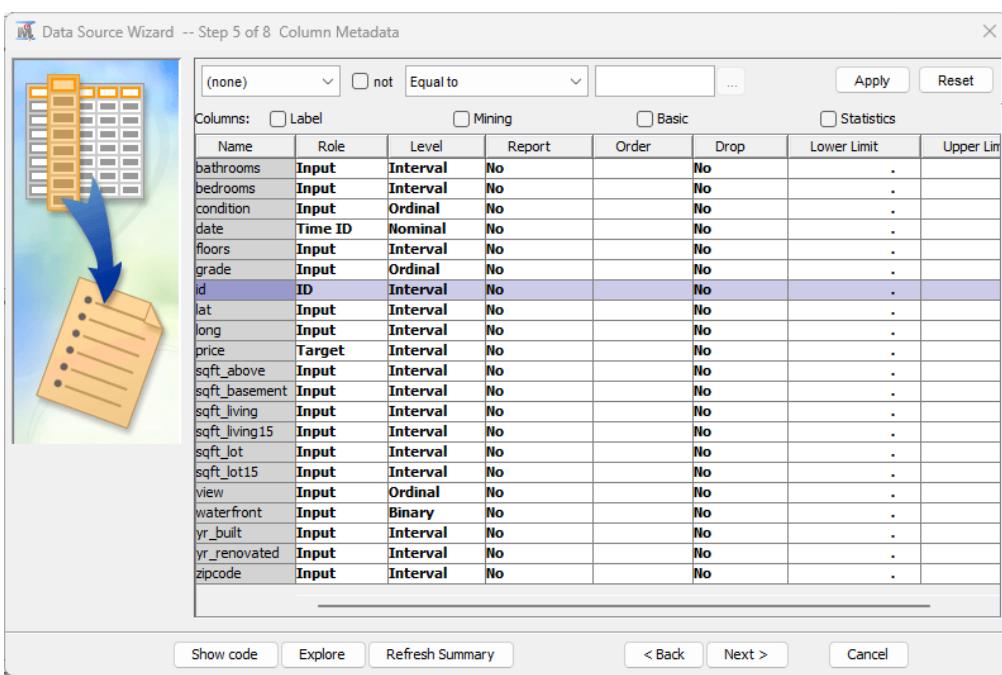
- Step 3 of 8 – Table Information. Select Next >



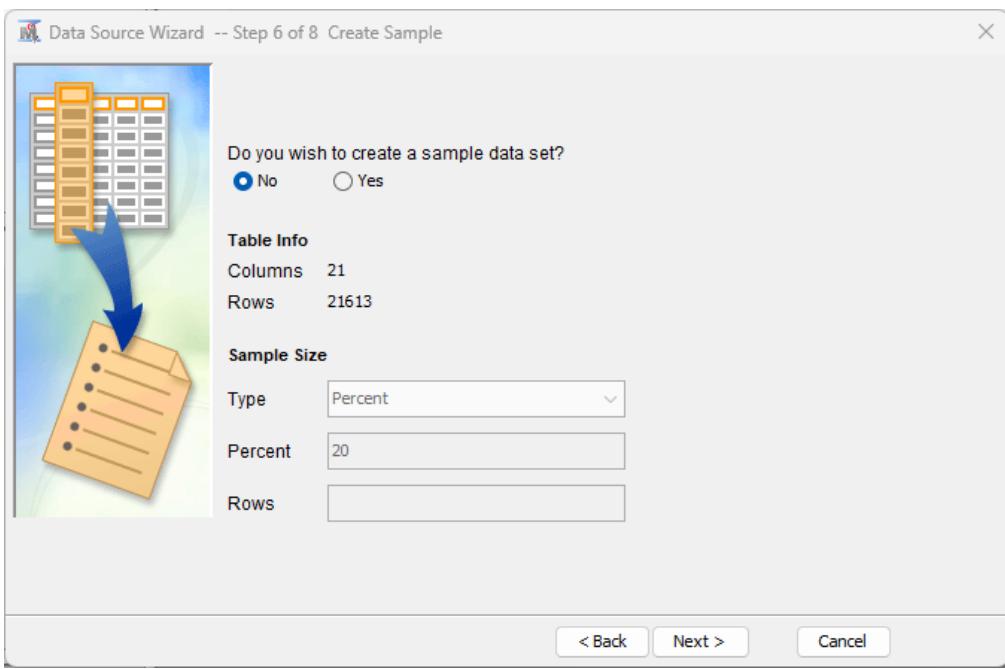
- Step 4 of 8 – Metadata Advisor. Select **Advanced** and **Next >**



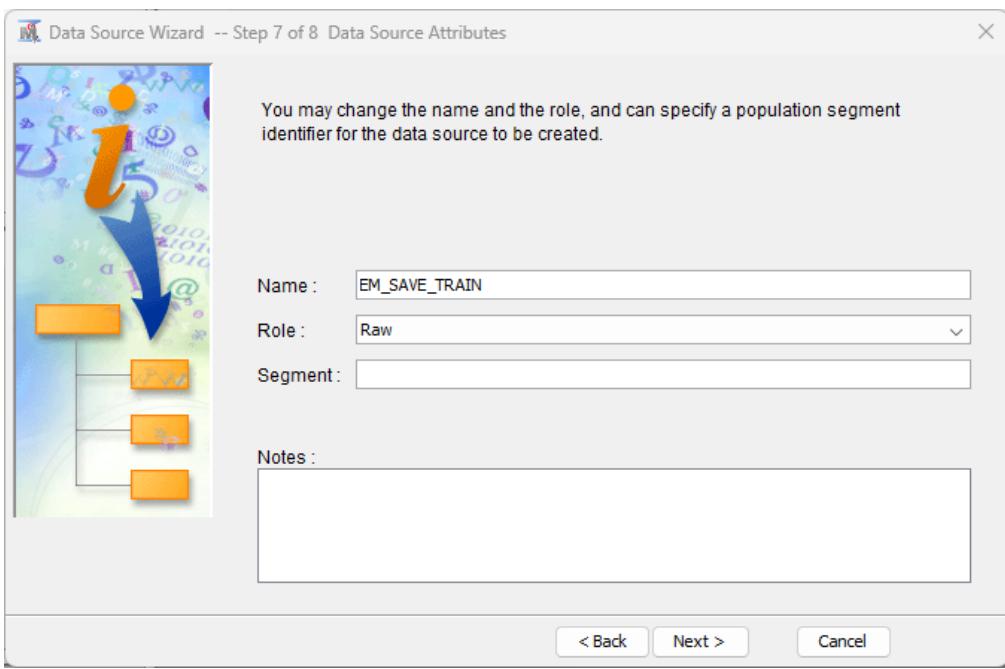
- Step 5 of 8 – Column Metadata. Update attributes on Roles and Levels. Select **Next >**



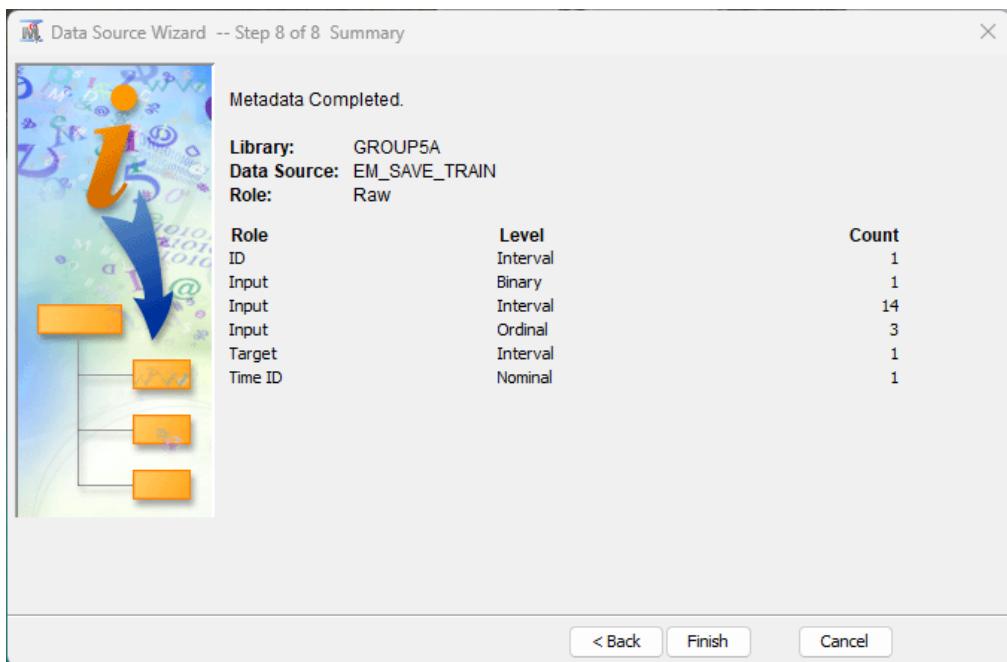
- Step 6 of 8 – Create Sample. Select Next >



- Step 7 of 8 – Data Source Attributes. Select Next >



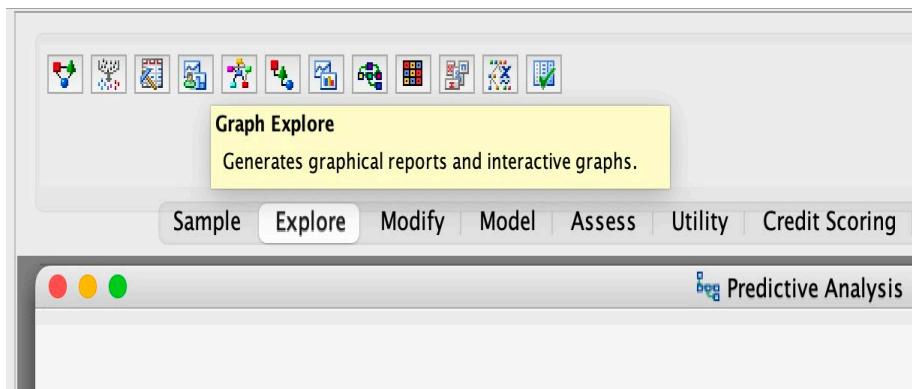
- Step 8 of 8 – Summary. Select Finish



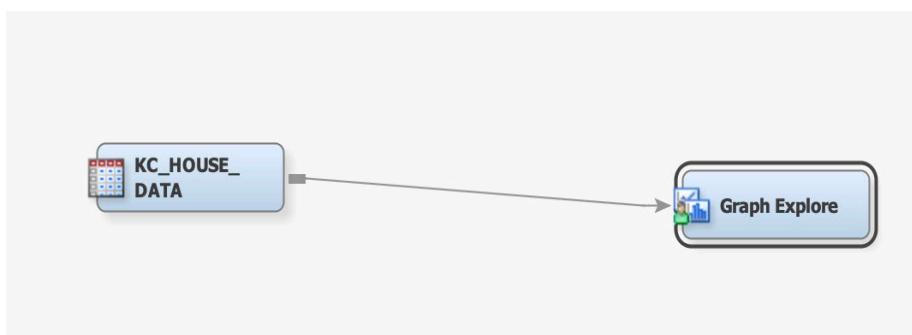
## 7.2 EXPLORE

### 7.2.1 Create Histogram, Pie Chart and Boxplot

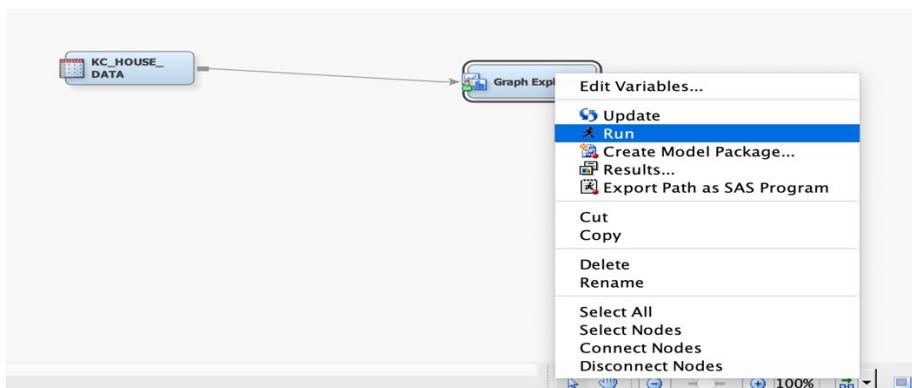
- Click Explore. Drag Graph Explore to the diagram



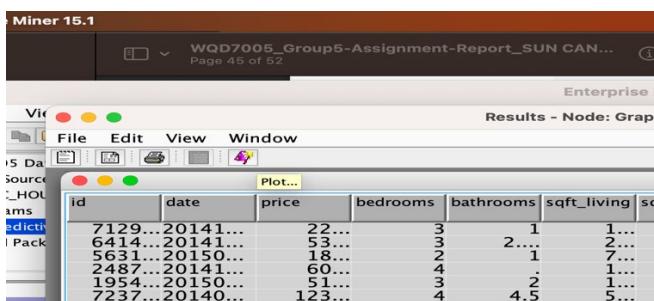
- Connect the Graph Explore Node to data source



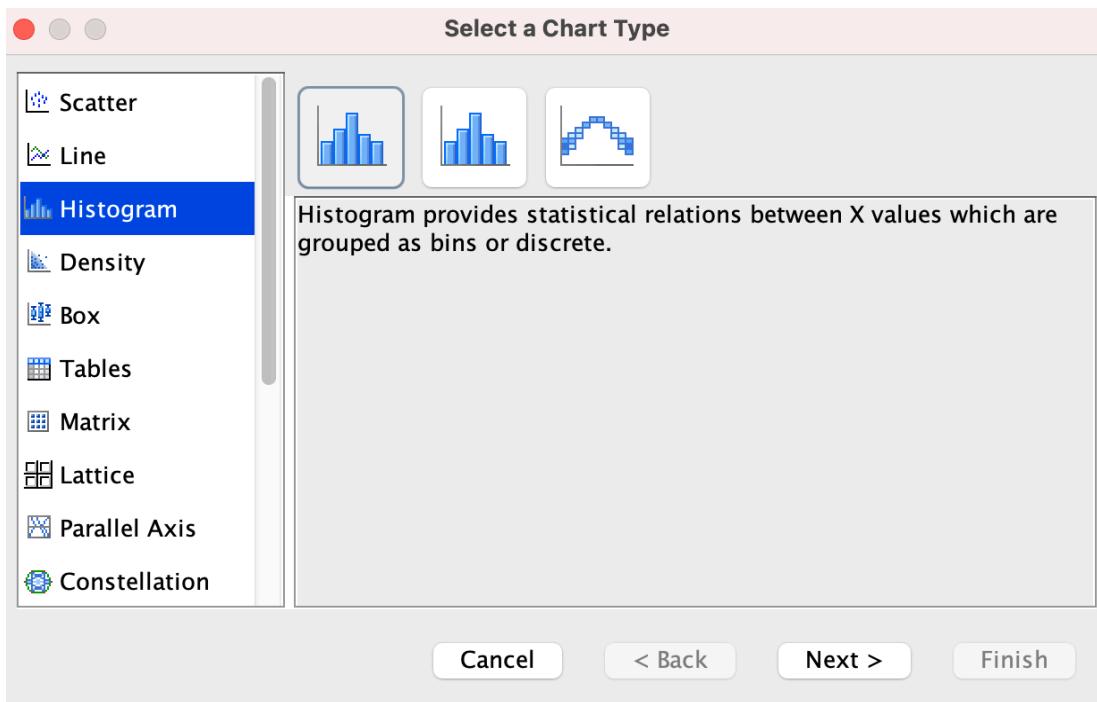
- Right click on Graph Explore Node, select Run and then Results



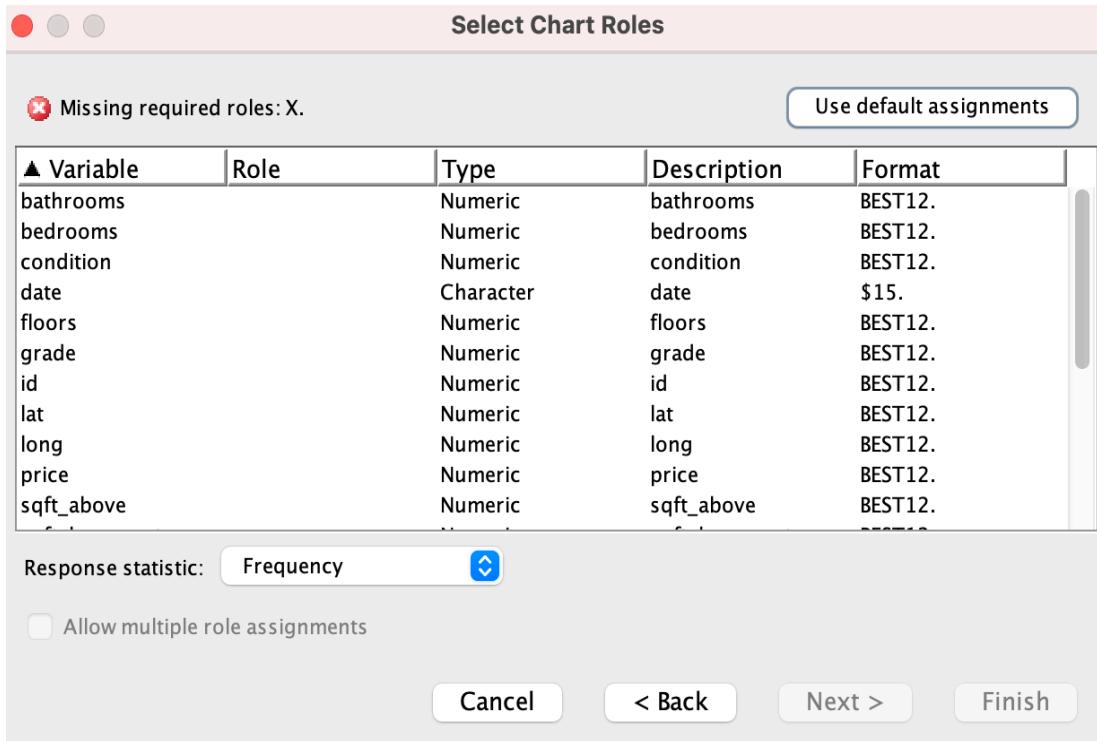
- Select Plot



- Select Histogram, Boxplot or Pie Chart

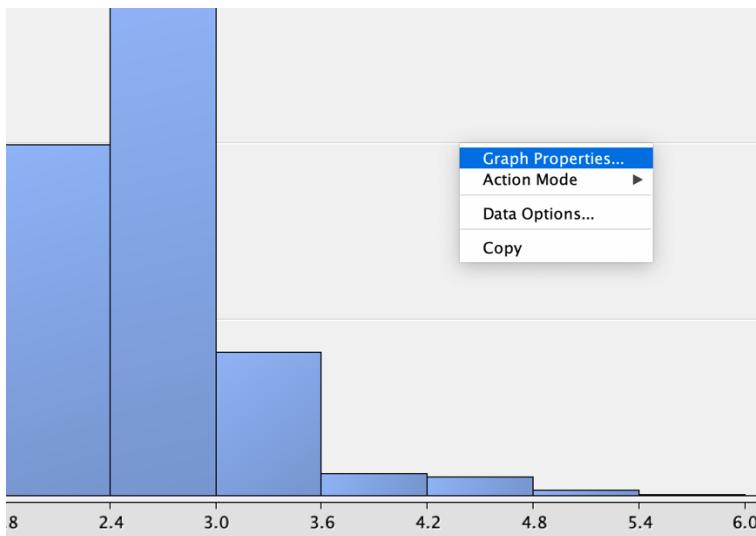


- Specify the variable and click Finish

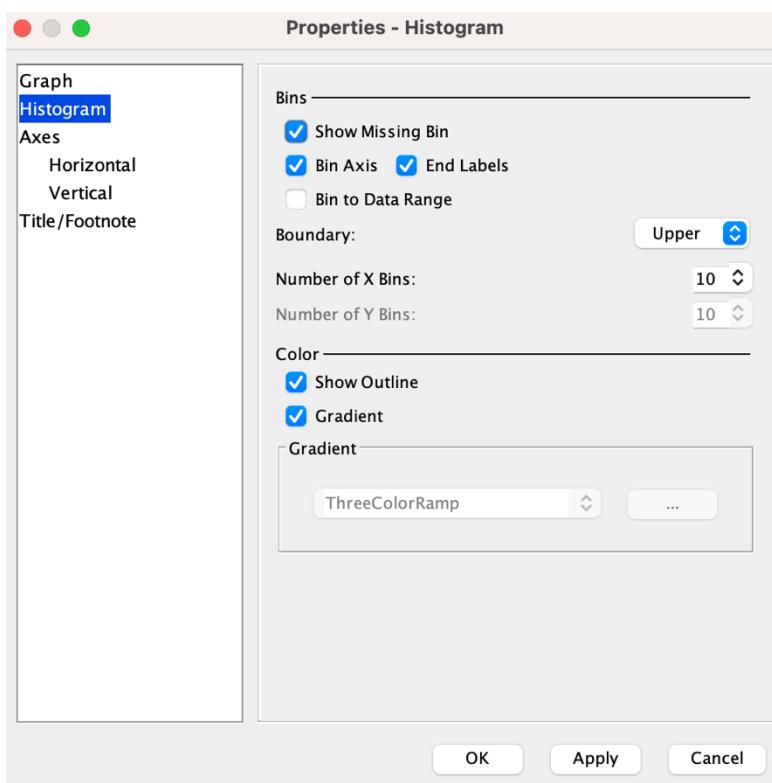


## 7.2.2 Add Missing Bin and Changing Graph Properties - Number of Bin in Histogram

- Right click on histogram and click graph properties

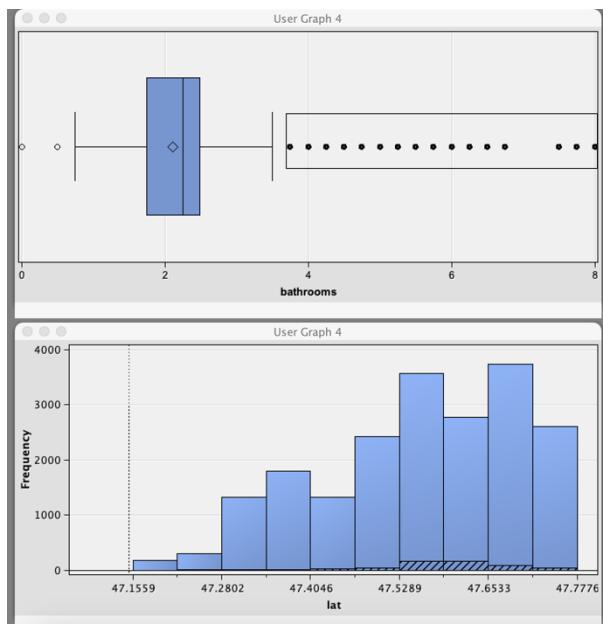


- Click “show missing bin” to display missing bin. Change the number of bin in histogram by changing the value of “Number of X Bins”.



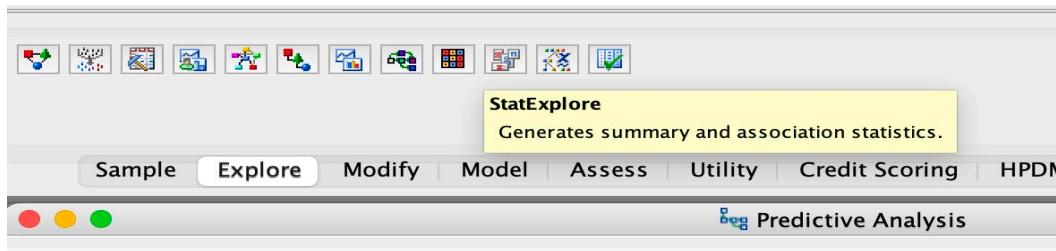
### 7.2.3 Display Variable Association

- Drag across the points on the diagram

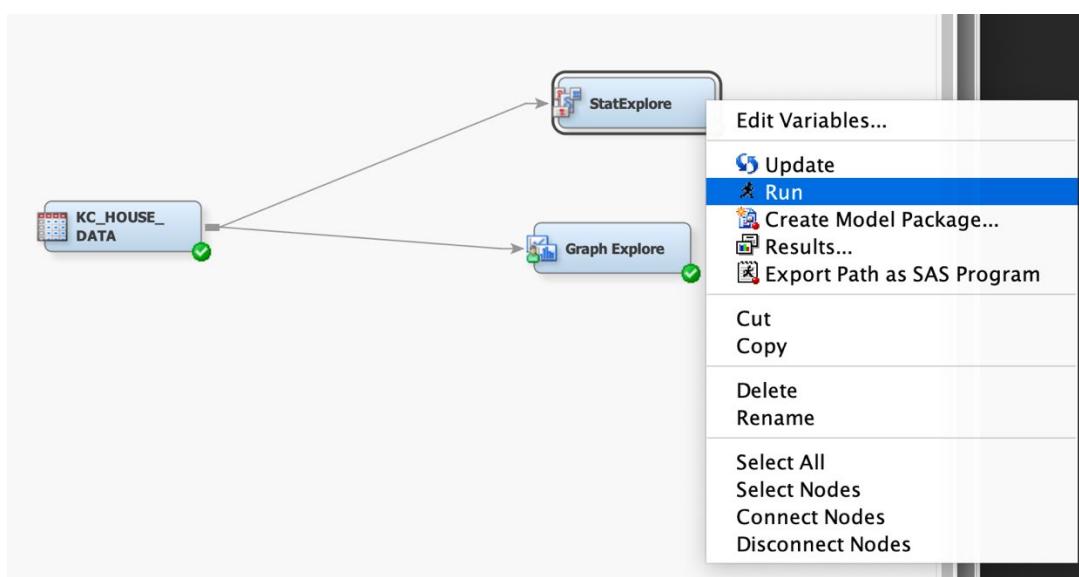


### 7.2.4 Create Summary Statistics

- Drag StatExplorer to the diagram

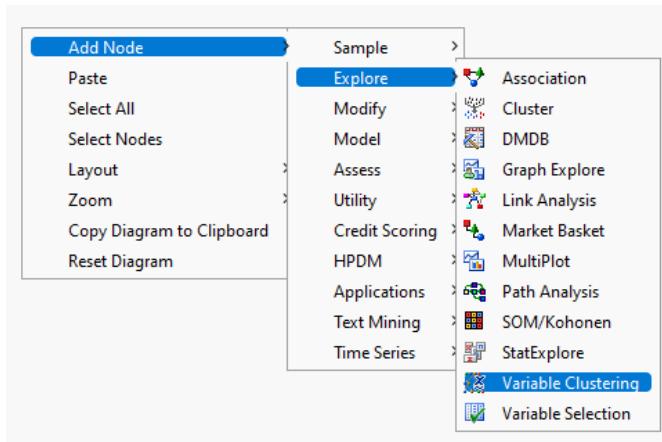


Connect the StatExplorer to the data source. Right click Run and Result

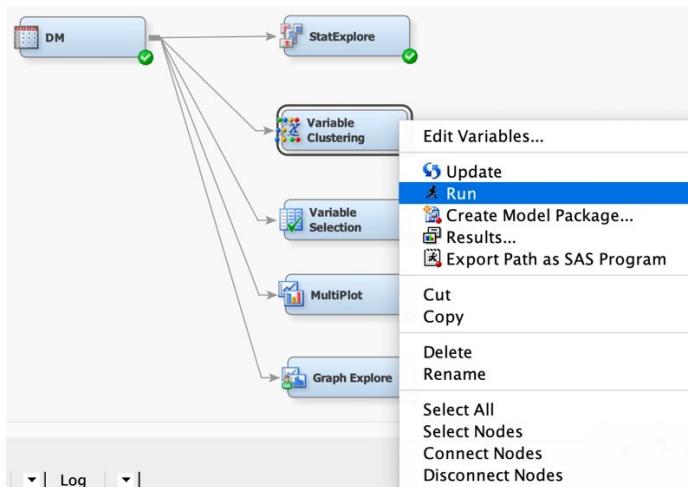


### 7.2.5 Perform Variable Clustering

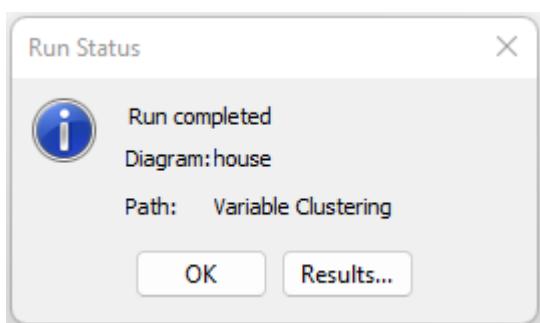
- Select Add Node → Explore → Variable Clustering.



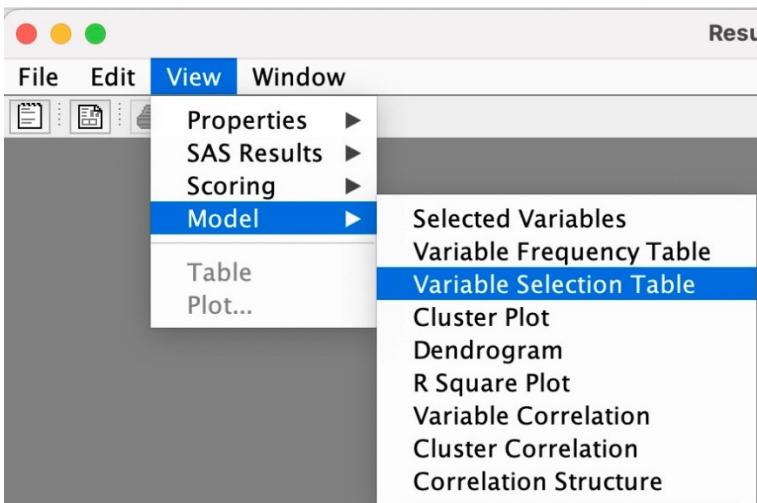
- Connect the data source and variable clustering and click Run.



- Once the run completed, click Results...



- Select **View**, **Model** and click on **Variable Selection Table**.

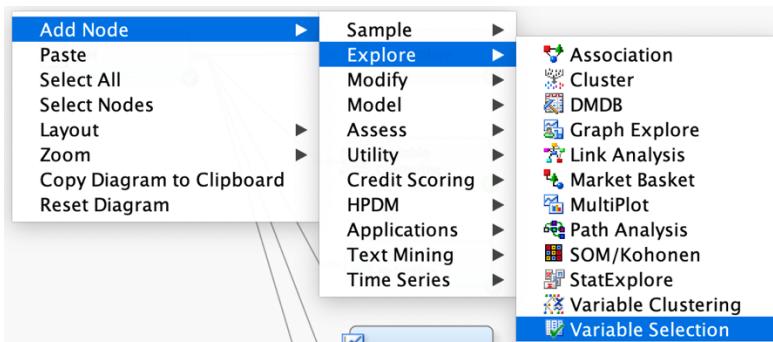


The screenshot shows the 'Variable Selection Table' results window. The title bar says 'Results - Node: Variable Clustering Diagram: DM'. The table has columns: Cluster, Variable, Label, R-Square With Own Cluster Component, Next Closest Cluster, R-Square with Next Cluster Component, Type, 1-R2 Ratio, and Variable Selected. The data rows show various variables from 'CLUS1' to 'CLUS4' assigned to clusters 'CLUS1', 'CLUS2', 'CLUS3', and 'CLUS4' respectively, along with their respective statistics and selection status.

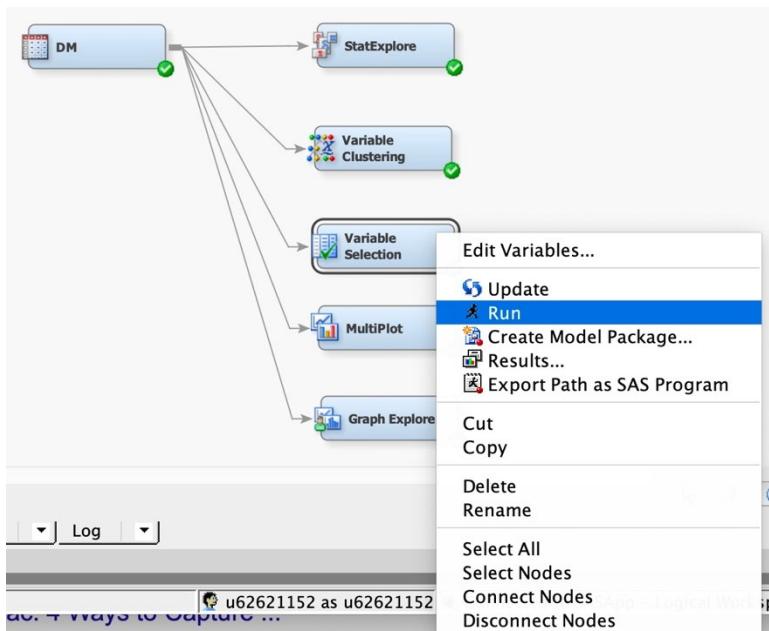
Cluster	Variable	Label	R-Square With Own Cluster Component	Next Closest Cluster	R-Square with Next Cluster Component	Type	1-R2 Ratio	Variable Selected
CLUS1	CLUS1	Cluster 1	0.0671...	1CLUS5	0.0671...	Cluster...	0.0671...	YES
CLUS1	SQFT_L...	sqft_liv...	0.8686...	CLUS4	0.1892...	Variable	0.1619...	NO
CLUS1	SQFT_A...	sqft_ab...	0.8446...	CLUS5	0.0877...	Variable	0.1703...	NO
CLUS1	BATHR...	bathro...	0.7268...	CLUS4	0.0804...	Variable	0.2970...	NO
CLUS1	SQFT_L...	sqft_liv...	0.6487...	CLUS5	0.0836...	Variable	0.3833...	NO
CLUS1	BEDRO...	bedroo...	0.4204...	CLUS4	0.0918...	Variable	0.6381...	NO
CLUS1	FLOORS	floors	0.3124...	CLUS4	0.0603...	Variable	0.7317...	NO
CLUS2	CLUS2	Cluster 2	0.0519...	1CLUS5	0.0519...	Cluster...	0.0519...	YES
CLUS2	SQFT_L...	sqft_lot...	0.8592...	CLUS5	0.0403...	Variable	0.1466...	NO
CLUS2	SQFT_L...	sqft_lot...	0.8592...	CLUS5	0.0492...	Variable	0.1480...	NO
CLUS3	CLUS3	Cluster 3	0.0026...	1CLUS1	0.0026...	Cluster...	0.0026...	YES
CLUS3	YR_BUILTyr_built		0.0026...	1CLUS1	0.0026...	Variable	0.0026...	NO
CLUS4	CLUS4	Cluster 4	0.0460...	1CLUS1	0.0460...	Cluster	0.0460...	NOFS

## 7.2.6 Perform Variable Selection

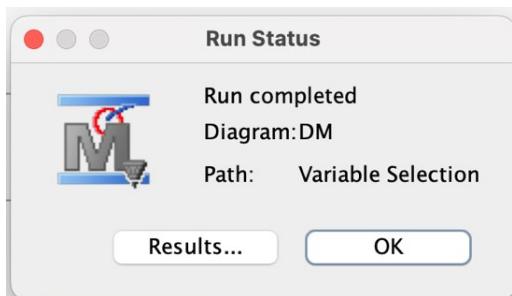
- Select **Add Node** → **Explore** → **Variable Selection**.



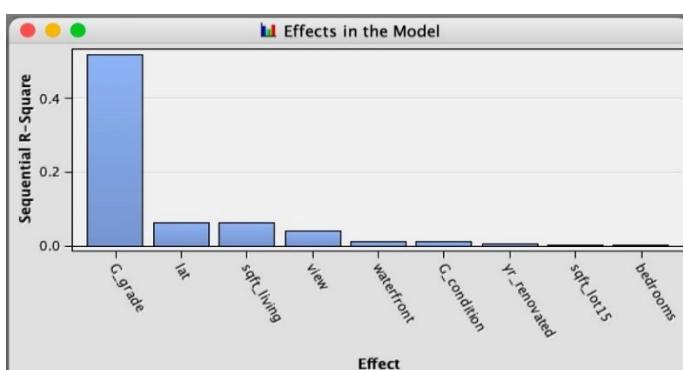
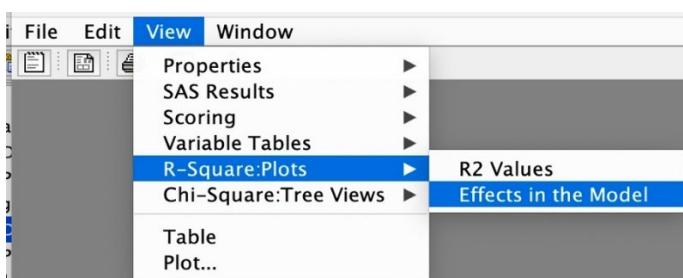
- Connect the data source and variable clustering and click **Run**.



- Once the run completed, click **Results...**

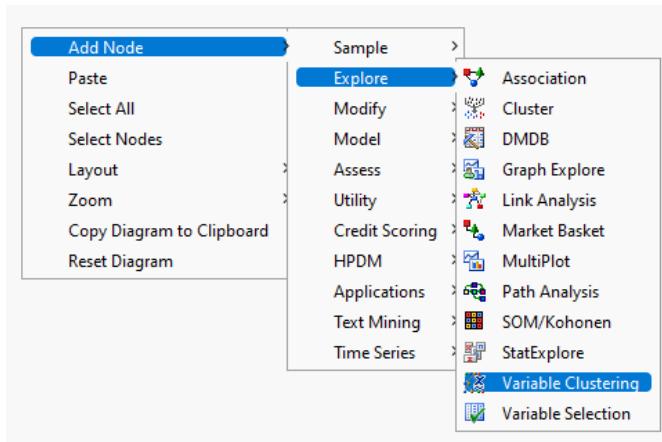


- Select **View**, **R-Square: Plots** and click on **Effects in the Model**.

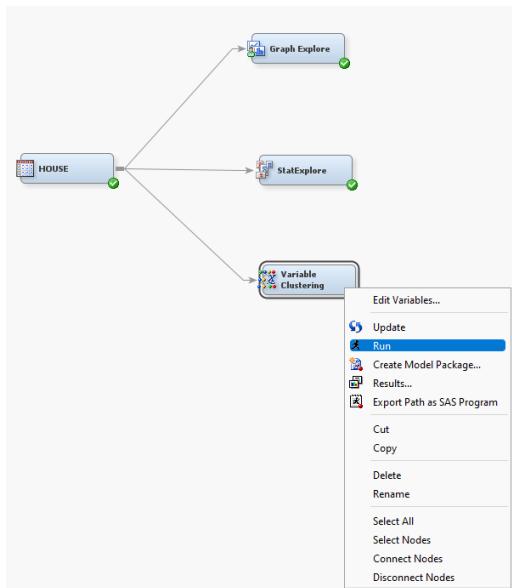


### 7.2.7 Perform Variable Correlation

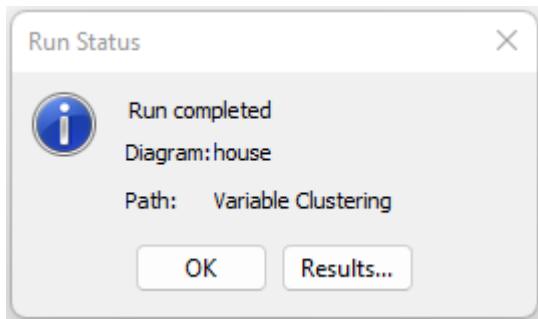
- Select Add Node → Explore → Variable Clustering.



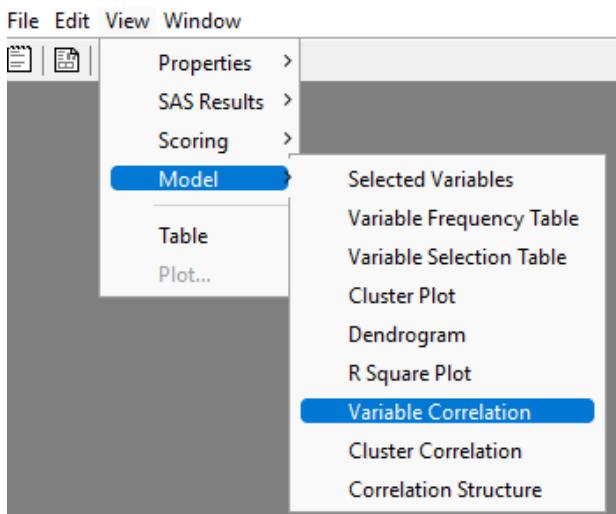
- Connect the data source and variable clustering and click Run.



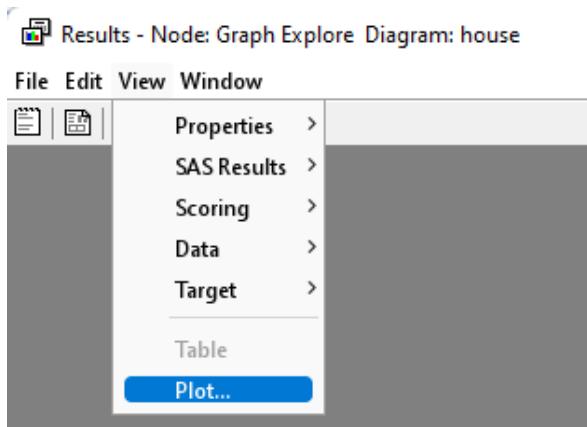
- Once the run completed, click Results...



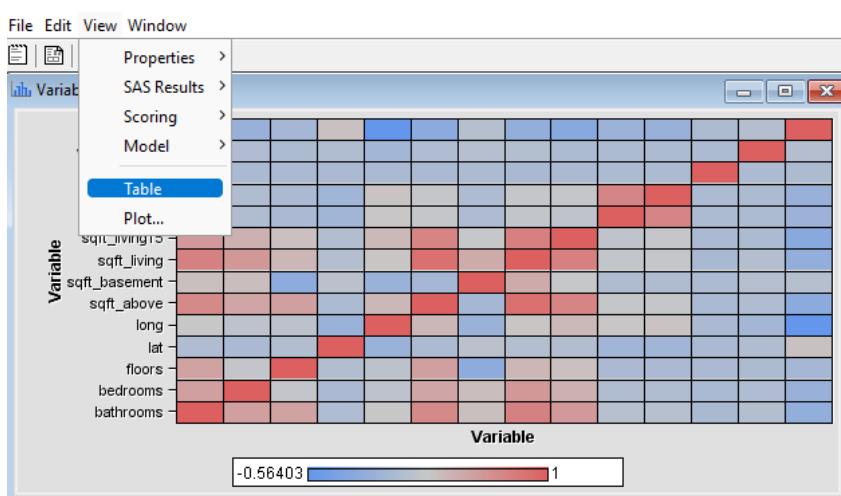
- Select **View**, **Model** and click on **Variable Correlation**.



- Select **View** and click on **Plot...**

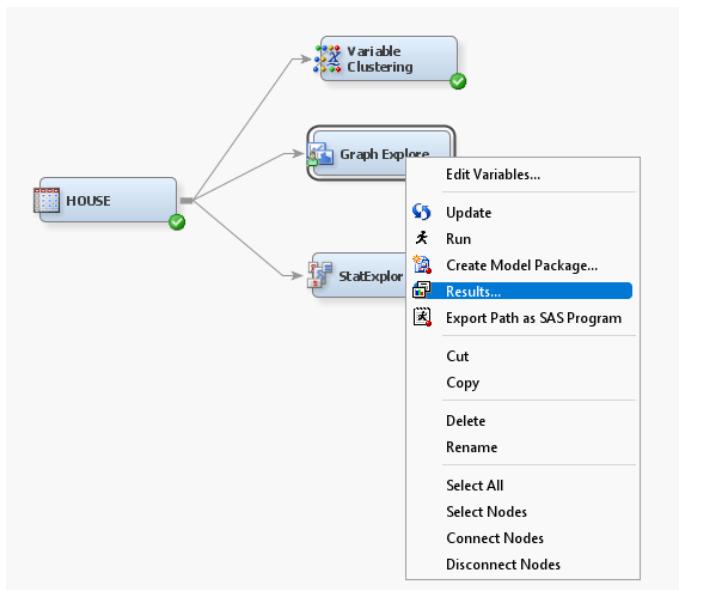


- Then, click on the correlation matrix plot before selecting **View** and click on **Table** to obtain correlation table.

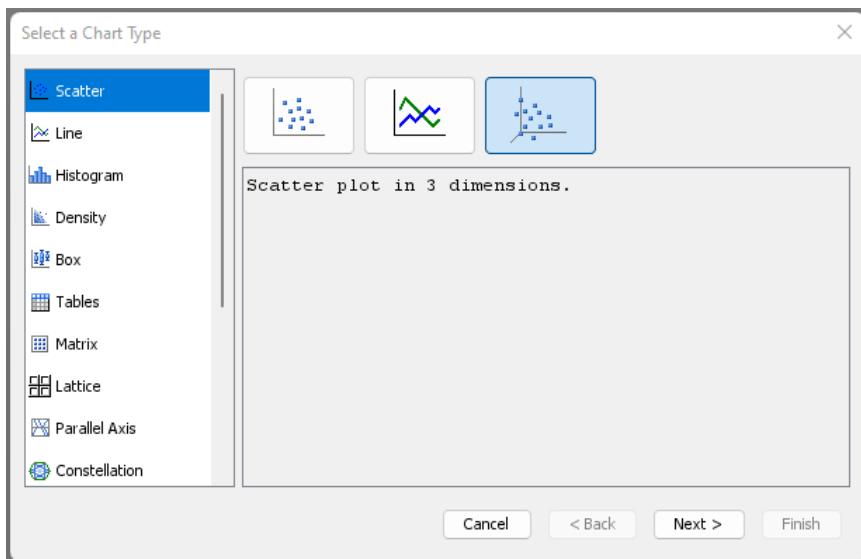


### 7.2.8 Interesting Visualizations

- Right-click the **Graph Explore** node and select **Results**.



- Select **3D Scatter** and click **Next**.



- Select at least 3 variables as X, Y and Z. Then, click **Finish**.

Select Chart Roles

▲ Variable	Role	Type	Description	Format
date		Character	date	\$15.
floors		Numeric	floors	BEST12.
grade		Numeric	grade	BEST12.
id		Numeric	id	BEST12.
lat		Numeric	lat	BEST12.
long		Numeric	long	BEST12.
price	X	Numeric	price	BEST12.
sqft_above	Z	Numeric	sqft_above	BEST12.
sqft_basement		Numeric	sqft_basement	BEST12.
sqft_living	Y	Numeric	sqft_living	BEST12.
sqft_living15		Numeric	sqft_living15	BEST12.
sqft_lot		Numeric	sqft_lot	BEST12.
sqft_lot15		Numeric	sqft_lot15	BEST12.
view		Numeric	view	BEST12.

Allow multiple role assignments

### 7.2.9 Correlation Table

<b>Variable</b>	<b>Variable</b>	<b>Correlation</b>
bathrooms	bathrooms	1.0000
bedrooms	bathrooms	0.5159
floors	bathrooms	0.5007
lat	bathrooms	0.0246
long	bathrooms	0.2231
sqft above	bathrooms	0.6854
sqft basement	bathrooms	0.2837
sqft living	bathrooms	0.7547
sqft living15	bathrooms	0.5687
sqft lot	bathrooms	0.0877
sqft lot15	bathrooms	0.0872
yr built	bathrooms	-0.0145
yr renovated	bathrooms	0.0507
zipcode	bathrooms	-0.2039
bathrooms	bedrooms	0.5159
bedrooms	bedrooms	1.0000
floors	bedrooms	0.1754
lat	bedrooms	-0.0090
long	bedrooms	0.1295
sqft above	bedrooms	0.4776
sqft basement	bedrooms	0.3030
sqft living	bedrooms	0.5766
sqft living15	bedrooms	0.3916
sqft lot	bedrooms	0.0317
sqft lot15	bedrooms	0.0292
yr built	bedrooms	-0.0131
yr renovated	bedrooms	0.0188
zipcode	bedrooms	-0.1528
bathrooms	floors	0.5007
bedrooms	floors	0.1754
floors	floors	1.0000
lat	floors	0.0497
long	floors	0.1253
sqft above	floors	0.5238
sqft basement	floors	-0.2457
sqft living	floors	0.3539
sqft living15	floors	0.2798
sqft lot	floors	-0.0052
sqft lot15	floors	-0.0113
yr built	floors	-0.0108

yr_renovated	floors	0.0063
zipcode	floors	-0.0591
bathrooms	lat	0.0246
bedrooms	lat	-0.0090
floors	lat	0.0497
lat	lat	1.0000
long	lat	-0.1355
sqft_above	lat	-0.0008
sqft_basement	lat	0.1105
sqft_living	lat	0.0526
sqft_living15	lat	0.0488
sqft_lot	lat	-0.0857
sqft_lot15	lat	-0.0864
yr_built	lat	-0.0106
yr_renovated	lat	0.0294
zipcode	lat	0.2670
bathrooms	long	0.2231
bedrooms	long	0.1295
floors	long	0.1253
lat	long	-0.1355
long	long	1.0000
sqft_above	long	0.3438
sqft_basement	long	-0.1447
sqft_living	long	0.2402
sqft_living15	long	0.3346
sqft_lot	long	0.2295
sqft_lot15	long	0.2544
yr_built	long	-0.0092
yr_renovated	long	-0.0684
zipcode	long	-0.5640
bathrooms	sqft_above	0.6854
bedrooms	sqft_above	0.4776
floors	sqft_above	0.5238
lat	sqft_above	-0.0008
long	sqft_above	0.3438
sqft_above	sqft_above	1.0000
sqft_basement	sqft_above	-0.0519
sqft_living	sqft_above	0.8766
sqft_living15	sqft_above	0.7319
sqft_lot	sqft_above	0.1835
sqft_lot15	sqft_above	0.1940
yr_built	sqft_above	-0.0134

yr_renovated	sqft_above	0.0232
zipcode	sqft_above	-0.2612
bathrooms	sqft_basement	0.2837
bedrooms	sqft_basement	0.3030
floors	sqft_basement	-0.2457
lat	sqft_basement	0.1105
long	sqft_basement	-0.1447
sqft_above	sqft_basement	-0.0519
sqft_basement	sqft_basement	1.0000
sqft_living	sqft_basement	0.4351
sqft_living15	sqft_basement	0.2004
sqft_lot	sqft_basement	0.0153
sqft_lot15	sqft_basement	0.0173
yr_built	sqft_basement	-0.0043
yr_renovated	sqft_basement	0.0713
zipcode	sqft_basement	0.0748
bathrooms	sqft_living	0.7547
bedrooms	sqft_living	0.5766
floors	sqft_living	0.3539
lat	sqft_living	0.0526
long	sqft_living	0.2402
sqft_above	sqft_living	0.8766
sqft_basement	sqft_living	0.4351
sqft_living	sqft_living	1.0000
sqft_living15	sqft_living	0.7564
sqft_lot	sqft_living	0.1728
sqft_lot15	sqft_living	0.1833
yr_built	sqft_living	-0.0141
yr_renovated	sqft_living	0.0553
zipcode	sqft_living	-0.1995
bathrooms	sqft_living15	0.5687
bedrooms	sqft_living15	0.3916
floors	sqft_living15	0.2798
lat	sqft_living15	0.0488
long	sqft_living15	0.3346
sqft_above	sqft_living15	0.7319
sqft_basement	sqft_living15	0.2004
sqft_living	sqft_living15	0.7564
sqft_living15	sqft_living15	1.0000
sqft_lot	sqft_living15	0.1446
sqft_lot15	sqft_living15	0.1832
yr_built	sqft_living15	-0.0119

yr_renovated	sqft_living15	-0.0027
zipcode	sqft_living15	-0.2790
bathrooms	sqft_lot	0.0877
bedrooms	sqft_lot	0.0317
floors	sqft_lot	-0.0052
lat	sqft_lot	-0.0857
long	sqft_lot	0.2295
sqft_above	sqft_lot	0.1835
sqft_basement	sqft_lot	0.0153
sqft_living	sqft_lot	0.1728
sqft_living15	sqft_lot	0.1446
sqft_lot	sqft_lot	1.0000
sqft_lot15	sqft_lot	0.7186
yr_built	sqft_lot	-0.0012
yr_renovated	sqft_lot	0.0076
zipcode	sqft_lot	-0.1296
bathrooms	sqft_lot15	0.0872
bedrooms	sqft_lot15	0.0292
floors	sqft_lot15	-0.0113
lat	sqft_lot15	-0.0864
long	sqft_lot15	0.2544
sqft_above	sqft_lot15	0.1940
sqft_basement	sqft_lot15	0.0173
sqft_living	sqft_lot15	0.1833
sqft_living15	sqft_lot15	0.1832
sqft_lot	sqft_lot15	0.7186
sqft_lot15	sqft_lot15	1.0000
yr_built	sqft_lot15	-0.0017
yr_renovated	sqft_lot15	0.0078
zipcode	sqft_lot15	-0.1472
bathrooms	yr_built	-0.0145
bedrooms	yr_built	-0.0131
floors	yr_built	-0.0108
lat	yr_built	-0.0106
long	yr_built	-0.0092
sqft_above	yr_built	-0.0134
sqft_basement	yr_built	-0.0043
sqft_living	yr_built	-0.0141
sqft_living15	yr_built	-0.0119
sqft_lot	yr_built	-0.0012
sqft_lot15	yr_built	-0.0017
yr_built	yr_built	1.0000

yr_renovated	yr_built	-0.0023
zipcode	yr_built	0.0092
bathrooms	yr_renovated	0.0507
bedrooms	yr_renovated	0.0188
floors	yr_renovated	0.0063
lat	yr_renovated	0.0294
long	yr_renovated	-0.0684
sqft_above	yr_renovated	0.0232
sqft_basement	yr_renovated	0.0713
sqft_living	yr_renovated	0.0553
sqft_living15	yr_renovated	-0.0027
sqft_lot	yr_renovated	0.0076
sqft_lot15	yr_renovated	0.0078
yr_built	yr_renovated	-0.0023
yr_renovated	yr_renovated	1.0000
zipcode	yr_renovated	0.0644
bathrooms	zipcode	-0.2039
bedrooms	zipcode	-0.1528
floors	zipcode	-0.0591
lat	zipcode	0.2670
long	zipcode	-0.5640
sqft_above	zipcode	-0.2612
sqft_basement	zipcode	0.0748
sqft_living	zipcode	-0.1995
sqft_living15	zipcode	-0.2790
sqft_lot	zipcode	-0.1296
sqft_lot15	zipcode	-0.1472
yr_built	zipcode	0.0092
yr_renovated	zipcode	0.0644
zipcode	zipcode	1.0000

## 7.3 MODIFY

### 7.3.1 Modify Inconsistent Data

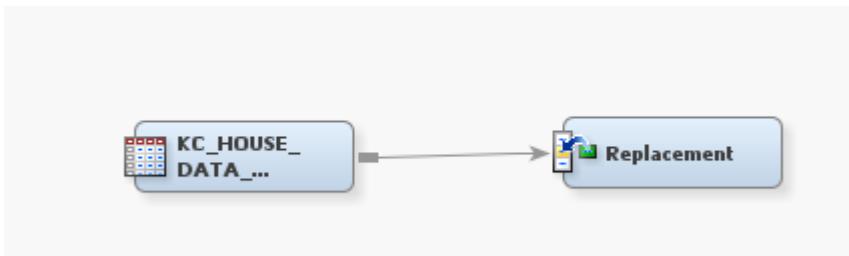
- In Talend, go to yr\_built variable

- Apply Extract parts of the text function on yr\_built and set the End Index to 4, yr\_built substring will be generated, with all the data are consistent

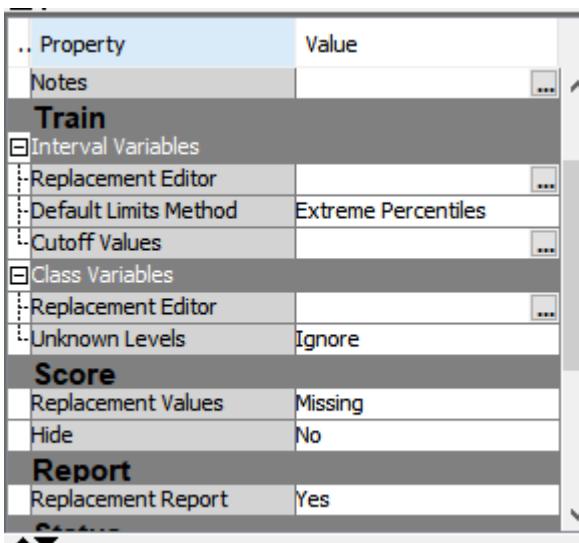
- Remove yr\_built column and rename yr\_built\_substring column to yr\_built column

### 7.3.2 Modify Noisy Data and Incomplete Data

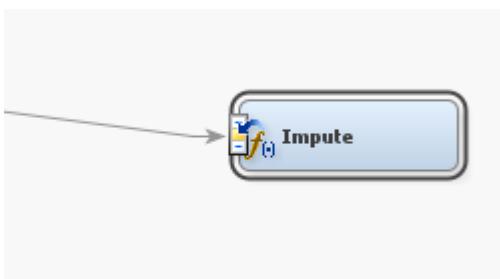
- Add Replacement Node



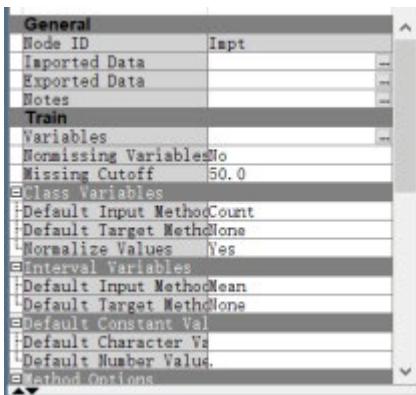
- Select Interval Variable >> Replacement Editor >> Variables that have noisy data >> Yes
- Select Default Limits Method >> Extreme Percentiles
- Select Score >> Replacement Values >> Missing



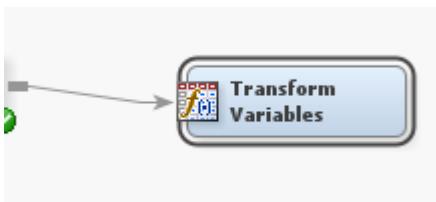
- Add Impute Node



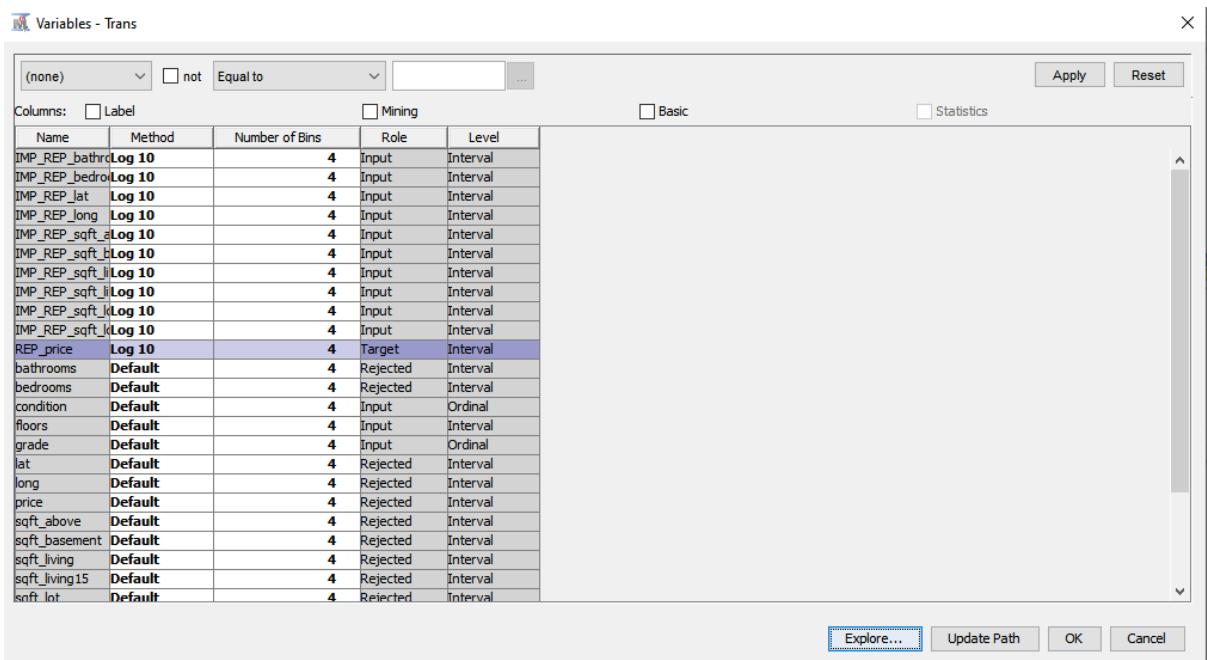
- Select Interval Variables >> Default Input Method >> Mean



- Add Transform Node



- Select Transform variables >> Train >> Variables >> Method >> Log 10



### 7.3.3 Data Partition

- Add Data Partition Node to Split Training Data and Validation Data



- Under Data Set Allocations, set 50% for Training and 50% for Validation

Data Set Allocations	
Training	50.0
Validation	50.0
Test	0.0

## 7.4 MODEL

### 7.4.1 Decision Tree Model

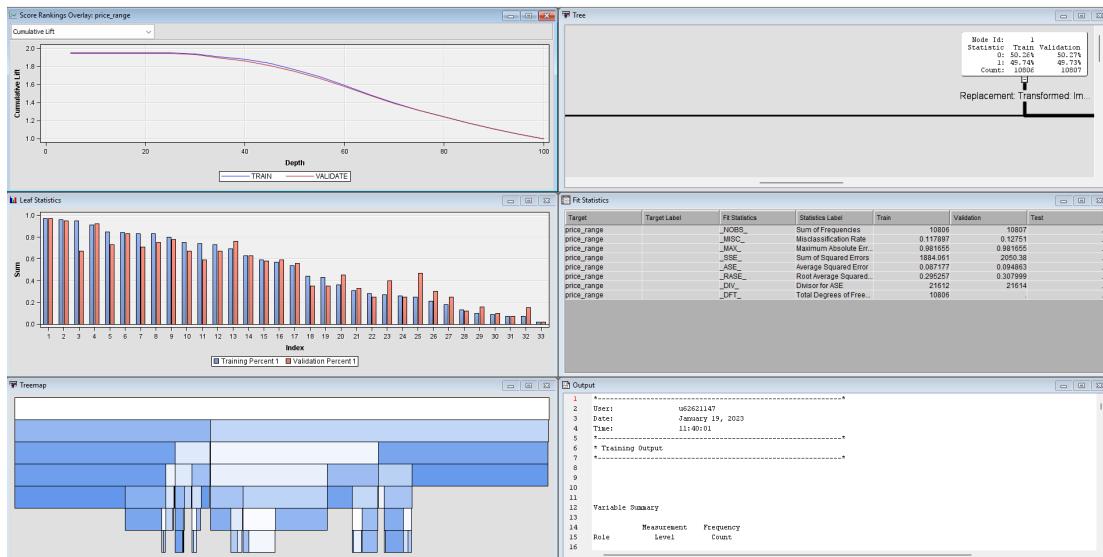
- Add and connect the decision tree node to the diagram.



- The parameter is set as below to build a decision tree model.

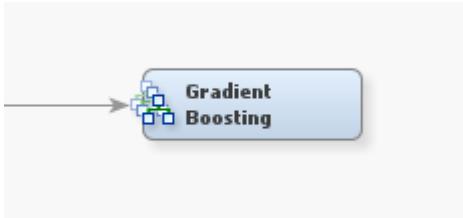
General	
Node ID	Tree
Imported Data	[...]
Exported Data	[...]
Notes	[...]
Train	
Variables	[...]
Interactive	[...]
Import Tree Model	No
Tree Model Data Set	[...]
Use Frozen Tree	No
Use Multiple Targets	No
Splitting Rule	
Interval Target Criterion	Variance
Nominal Target Criterion	Entropy
Ordinal Target Criterion	Entropy
Significance Level	0.2
Missing Values	Use in search
Use Input Once	No
Maximum Branch	2
Maximum Depth	6
Minimum Categorical Size	30

- Run the decision tree node with the above parameters and the outputs are as below:



## 7.4.2 Gradient Boosting Model

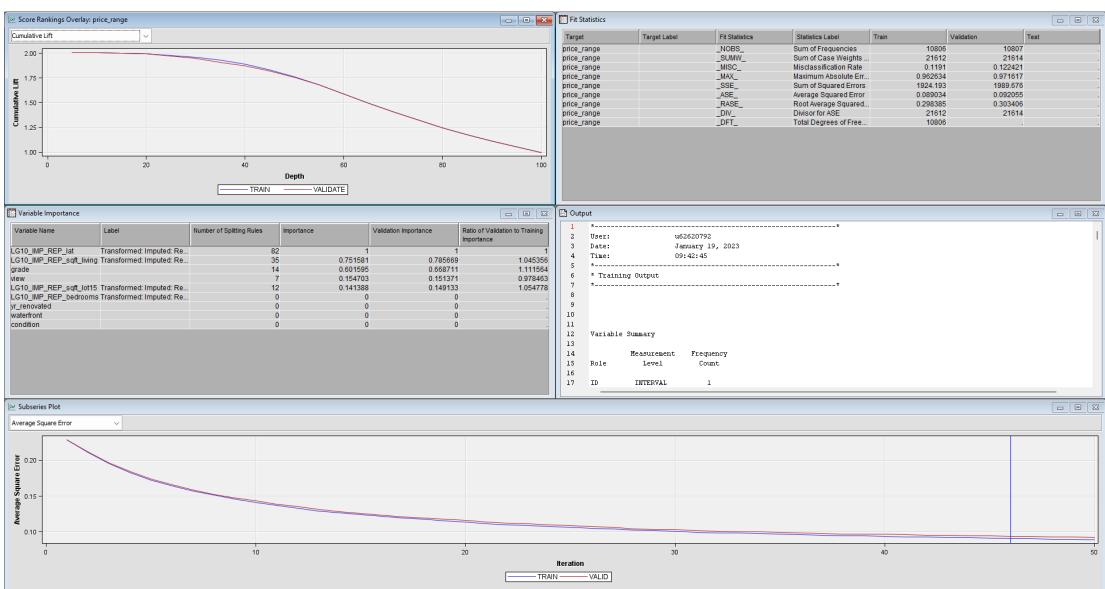
- Add Gradient Boosting Node into the diagram to create a Gradient Boosting Model.



- Set preferred parameters in the Gradient Boosting node:

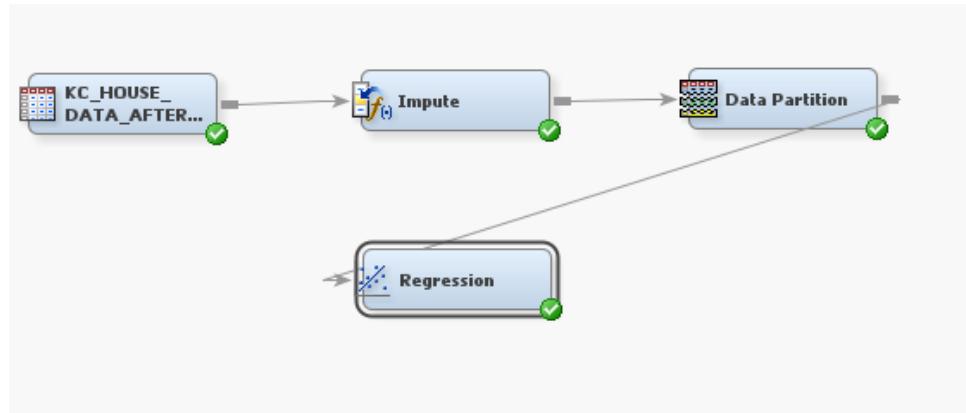
General	
Node ID	Boost
Imported Data	[...]
Exported Data	[...]
Notes	[...]
Train	
Variables	[...]
Series Options	
N Iterations	50
Seed	12345
Shrinkage	0.1
Train Proportion	60
Splitting Rule	
Huber M-Regression	No
Maximum Branch	2
Maximum Depth	2
Minimum Categorical Size	5
Reuse Variable	1
Categorical Bins	30
Interval Bins	100
Missing Values	Use in search
Performance	Disk
Node	
Leaf Fraction	0.001
Number of Surrogate Rules	0
Split Size	.

- After setting the parameters, run Gradient Boosting Node and view the result:



### 7.4.3 Logistics Regression Model

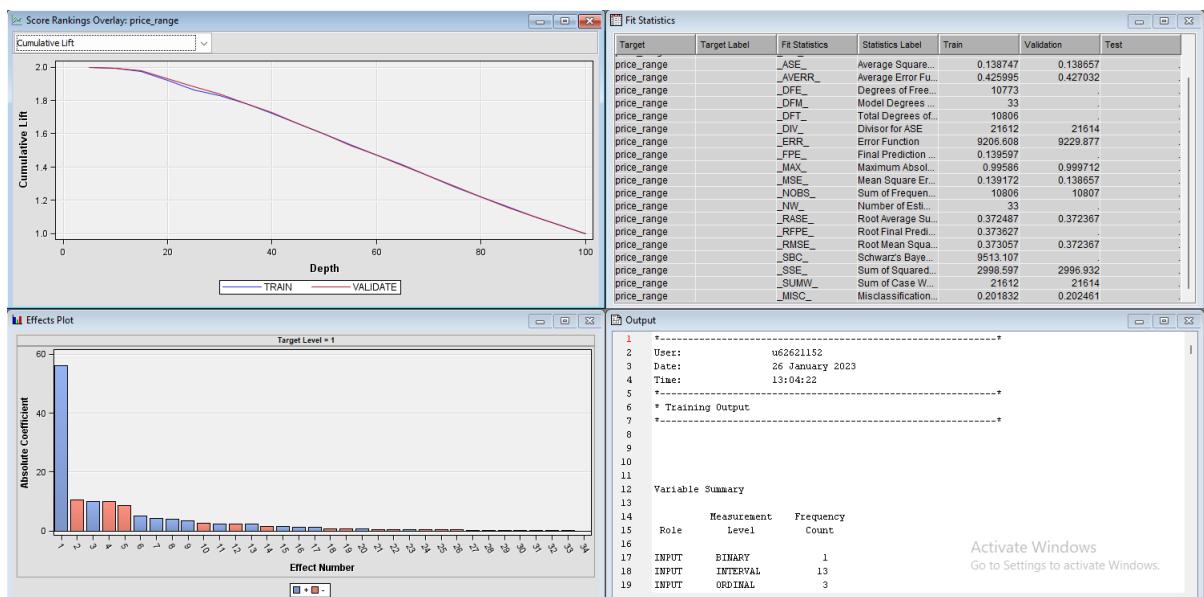
- Drag and drop the logistic regression model into the diagram and connect the previous nodes.



- Select the logistic regression as the regression type in the node and set the specified parameters:

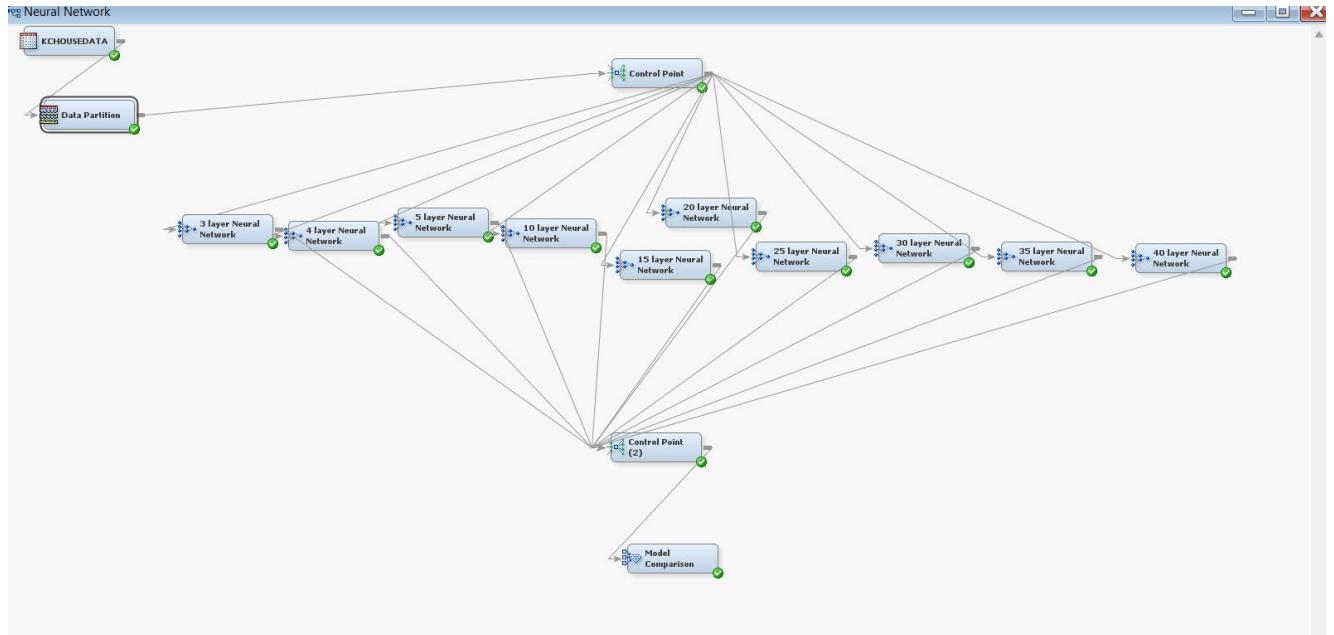
... Property	Value
Variables	
Equation	
Main Effects	Yes
Two-Factor Interactions	No
Polynomial Terms	No
Polynomial Degree	2
User Terms	No
Term Editor	
Class Targets	
Regression Type	Logistic Regression
Link Function	Logit
Model Options	
Suppress Intercept	No
Input Coding	Deviation
Model Selection	
Selection Model	Stepwise
Selection Criterion	Default
Use Selection Defaults	Yes
Selection Options	
Optimization Options	

- After setting the parameters, run the logistic regression and view the results as below:



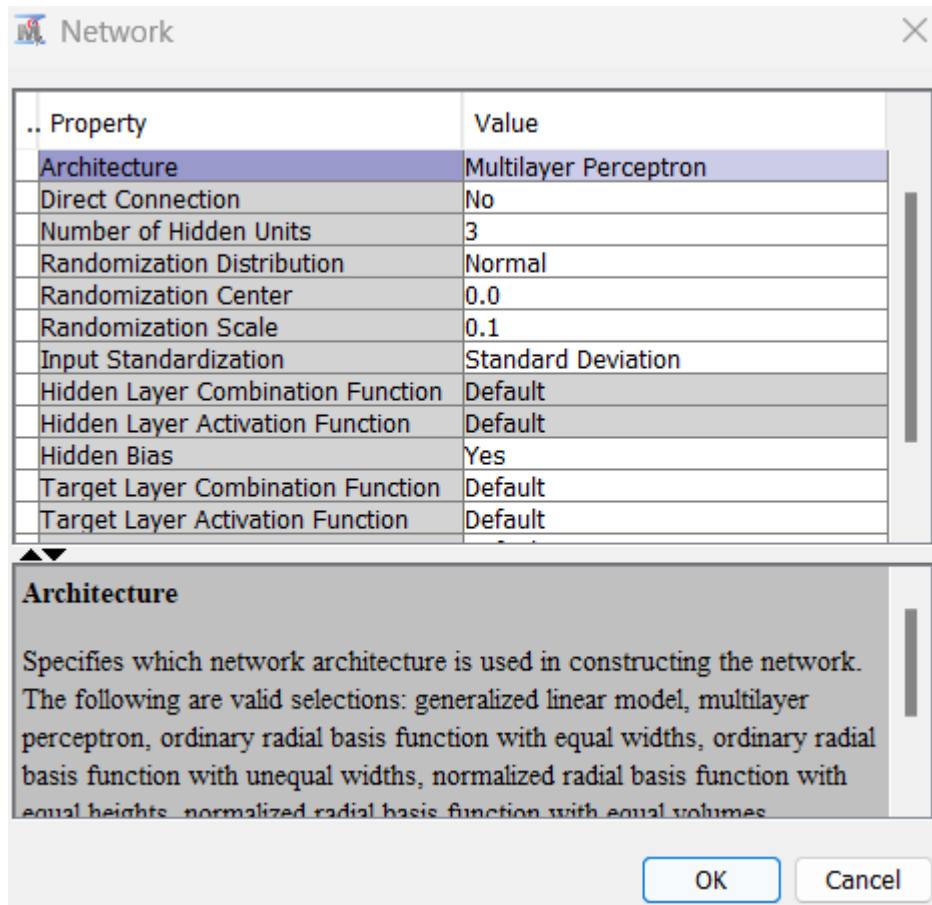
#### 7.4.4 Neural Network Model

- Create neural network model by dragging and dropping the different neural network nodes with various hidden layers into the diagram and connect the previous nodes.

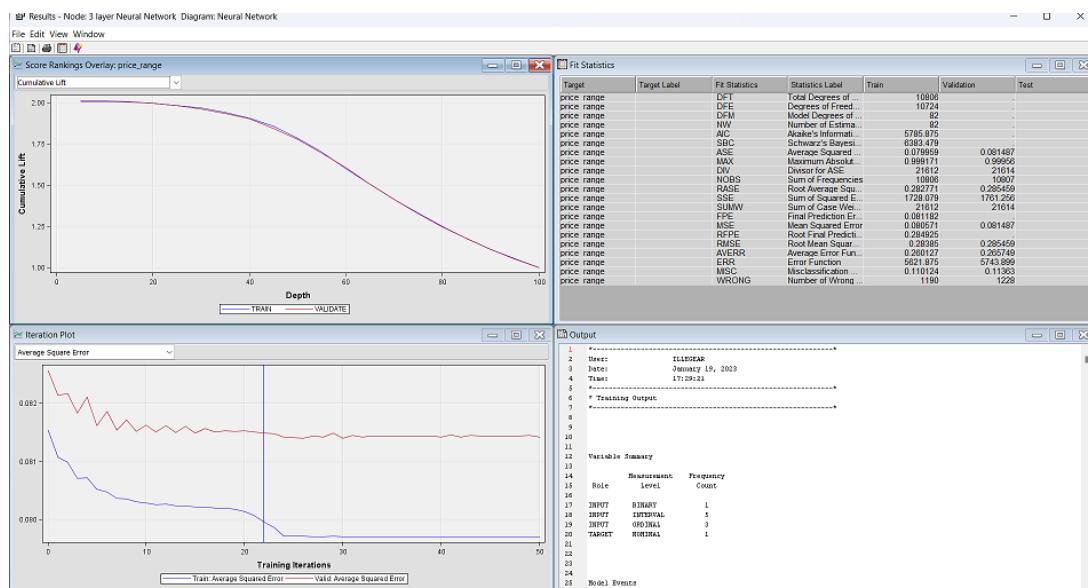


- Set preferred parameters in each neural network node

.. Property	Value
<b>General</b>	
Node ID	Neural
Imported Data	...
Exported Data	...
Notes	...
<b>Train</b>	
Variables	...
Continue Training	No
Network	...
Optimization	...
Initialization Seed	12345
Model Selection Criterion	Profit/Loss
Suppress Output	No
<b>Score</b>	
Hidden Units	No
Residuals	Yes
Standardization	No
<b>Status</b>	
Create Time	1/19/23 10:51 AM
Run ID	42640b96-3ed1-4ddd-b4
Last Error	
Last Status	Complete
Last Run Time	1/19/23 5:28 PM
Run Duration	0 Hr. 0 Min. 28.12 Sec.
Grid Host	
User-Added Node	No



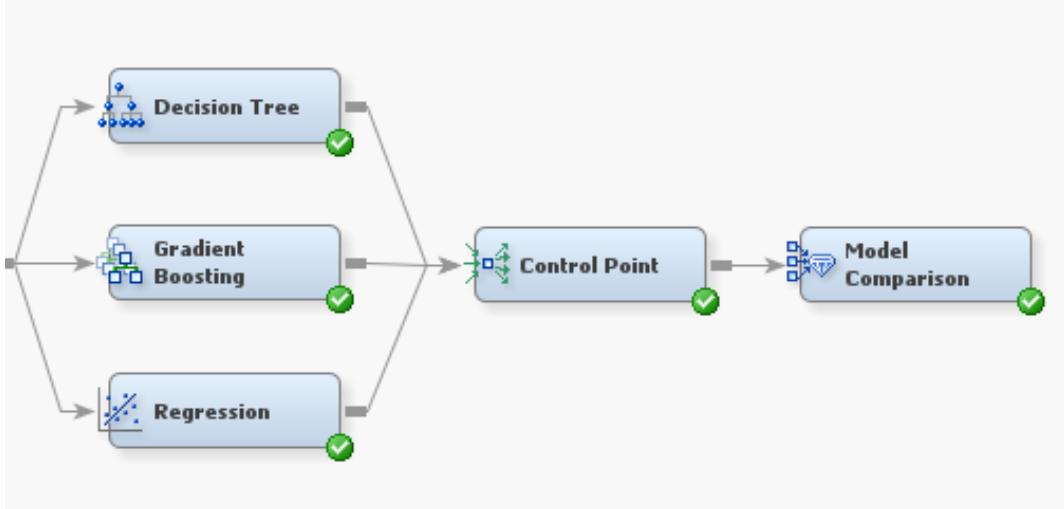
- After setting the parameters, run the neural network node and view the results as follows:



## 7.5 ASSESS

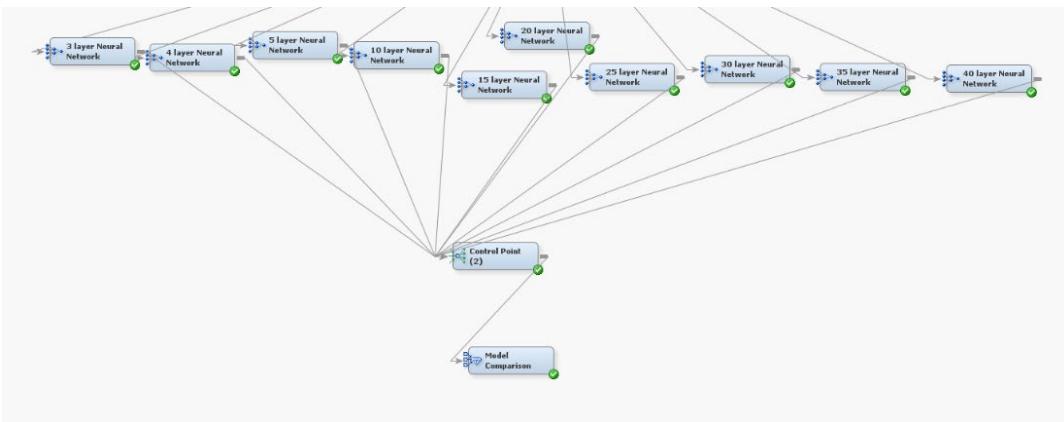
### 7.5.1 Compare between Decision Tree, Logistic Regression, and Gradient Boosting

Link Decision Tree, Logistic Regression and Gradient Nodes to a Control Point Node from Utility Tab and subsequently link it a Model Comparison Node from the Assess Tab.



### 7.5.2 Compare Neural Network

Link all the neural network nodes to the control point node and link the control point node to model comparison node.



### 7.5.3 Compare the Best Model

With two Model Comparison Nodes available, link both Model Comparison Nodes to a new Control Point Node and subsequently link it to a new Model Comparison Node.

