1.1. : Read the data and do exploratory data analysis. Describe the data briefly. (Check the null values, Data types, shape, EDA). Perform Univariate and Bivariate Analysis.

Below provides five number summaries for all continuous variables such as carat, depth, table, x, y, z and price while providing statistics on unique value count, top value with its frequencies for each categorical variable such as cut, color and clarity.

	carat	cut	color	clarity	depth	table	х	У	z	price
count	26967	26967	26967	26967	26270	26967	26967	26967	26967	26967
unique	NaN	5	7	8	NaN	NaN	NaN	NaN	NaN	NaN
top	NaN	ldeal	G	SI1	NaN	NaN	NaN	NaN	NaN	NaN
freq	NaN	10816	5661	6571	NaN	NaN	NaN	NaN	NaN	NaN
mean	0.798375	NaN	NaN	NaN	61.74515	57.45608	5.729854	5.733569	3.538057	3939.518
std	0.477745	NaN	NaN	NaN	1.41286	2.232068	1.128516	1.166058	0.720624	4024.865
min	0.2	NaN	NaN	NaN	50.8	49	0	0	0	326
25%	0.4	NaN	NaN	NaN	61	56	4.71	4.71	2.9	945
50%	0.7	NaN	NaN	NaN	61.8	57	5.69	5.71	3.52	2375
75%	1.05	NaN	NaN	NaN	62.5	59	6.55	6.54	4.04	5360
max	4.5	NaN	NaN	NaN	73.6	79	10.23	58.9	31.8	18818

Data types of the independent variables

Independ ent variables	Data type
carat	float64
cut	object
color	object
clarity	object
depth	float64
table	float64
X	float64
У	float64
Z	float64

The predicted variable is numeric and continuous as expected.

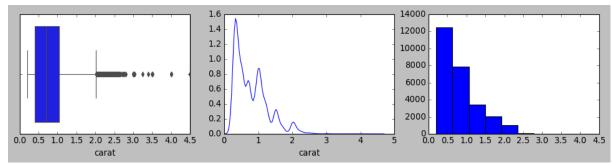
Shape:

26967 rows and 10 columns

Univariate analysis:

1. Univariate analysis for carat

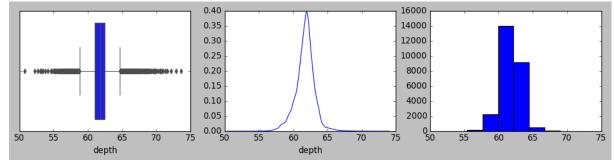
Mean is 0.798375, Median is 0.700000, Mode(s) are 0.3000 Column carat has outliers



Column carat is not normally distributed

2. Univariate analysis for depth

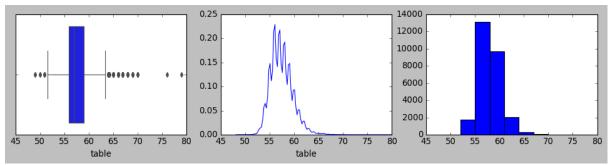
Mean is 61.745147, Median is nan, Mode(s) are 62.0000 Column depth has outliers



Column depth is normally distributed

3. Univariate analysis for table

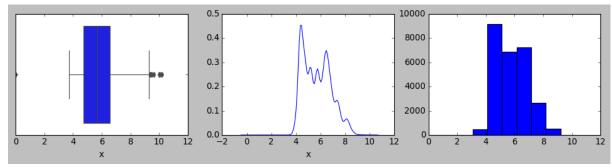
Mean is 57.456080, Median is 57.000000, Mode(s) are 56.0000 Column table has outliers



Column table is not normally distributed

4. Univariate analysis for x

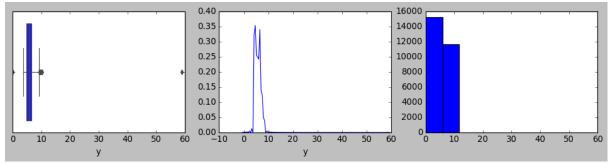
Mean is 5.729854, Median is 5.690000, Mode(s) are 4.3800 Column x has outliers



Column x is not normally distributed

5. Univariate analysis for y

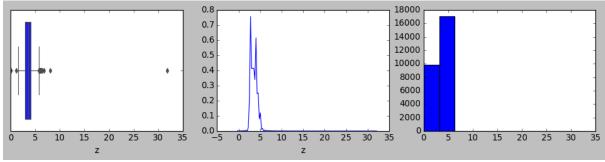
Mean is 5.733569, Median is 5.710000, Mode(s) are 4.3500 Column y has outliers



Column y is not normally distributed

6. Univariate analysis for z

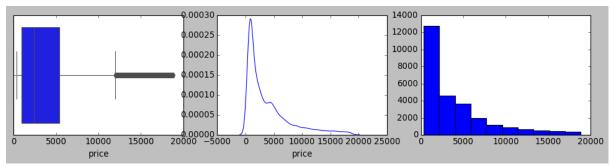
Mean is 3.538057, Median is 3.520000, Mode(s) are 2.6900 Column z has outliers



Column z is not normally distributed

7. Univariate analysis for price

Mean is 3939.518115, Median is 2375.000000, Mode(s) are 544.0000 Column price has outliers



Column price is not normally distributed

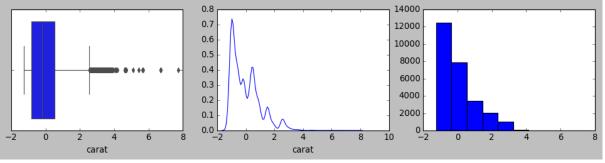
With the fact that independent variables have different units, scaling becomes necessary to remove the units variables are associated with so that the linear equation can be formed on the independent variables post standardization of them.

Z scoring based scaling of data would change the coefficient ,neutralize/remove the intercept while the accuracy score remains the same before and after. MSE would get scaled too.

Accordingly, below is the univariate analysis for the dataset that is standardized on z scores. To be noted that the below observation on the data is observed after encoding the categorical variables such as cut, color and clarity since in linear equation it demands every variable to be numeric to apply mathematical calculations on it.

1. Univariate analysis for carat

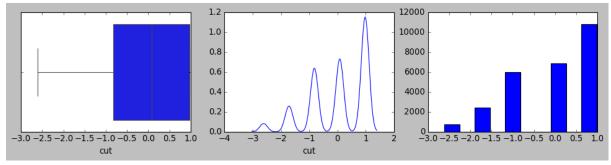
Mean is -0.000000, Median is -0.205920, Mode(s) are -1.0432 Column carat has outliers



Column carat is not normally distributed

2. Univariate analysis for cut

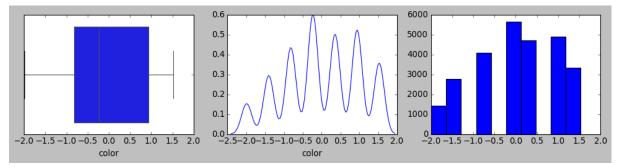
Mean is 0.000000, Median is 0.081246, Mode(s) are 0.9796 Column cut does not have outliers



Column cut is not normally distributed

3. Univariate analysis for color

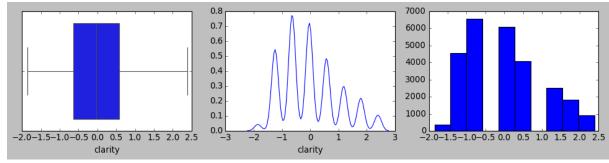
Mean is 0.000000, Median is -0.230890, Mode(s) are -0.2309 Column color does not have outliers



Column color is not normally distributed

4. Univariate analysis for clarity

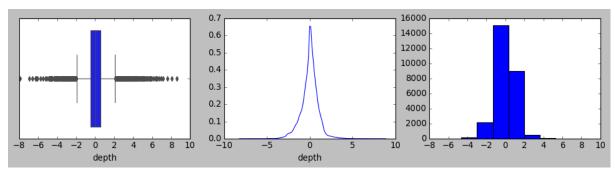
Mean is 0.000000, Median is -0.032241, Mode(s) are -0.6394 Column clarity does not have outliers



Column clarity is not normally distributed

5. Univariate analysis for depth

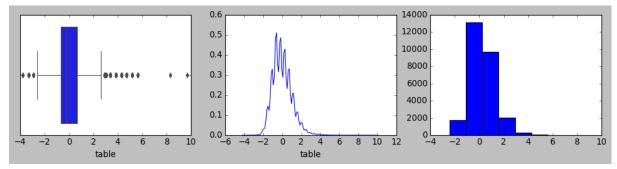
Mean is 0.000000, Median is 0.036839, Mode(s) are 0.0368 Column depth has outliers



Column depth is not normally distributed

6. Univariate analysis for table

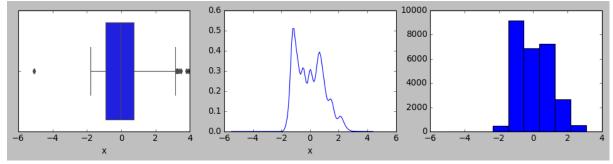
Mean is -0.000000, Median is -0.204334, Mode(s) are -0.6524 Column table has outliers



Column table is not normally distributed

7. Univariate analysis for x

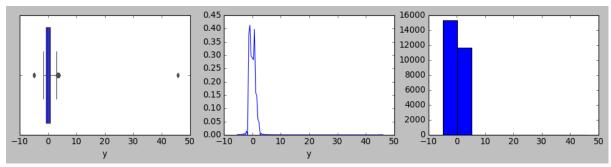
Mean is 0.000000, Median is -0.035316, Mode(s) are -1.1962 Column x has outliers



Column x is not normally distributed

8. Univariate analysis for y

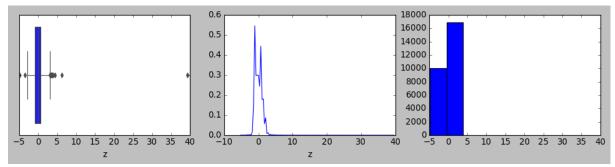
Mean is -0.000000, Median is -0.020213, Mode(s) are -1.1866 Column y has outliers



Column y is not normally distributed

9. Univariate analysis for z

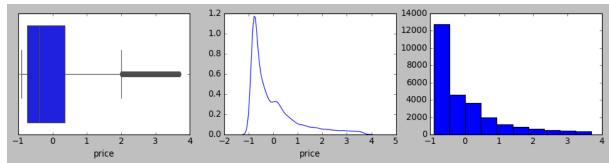
Mean is -0.000000, Median is -0.025058, Mode(s) are -1.1769 Column z has outliers



Column z is not normally distributed

10. Univariate analysis for price

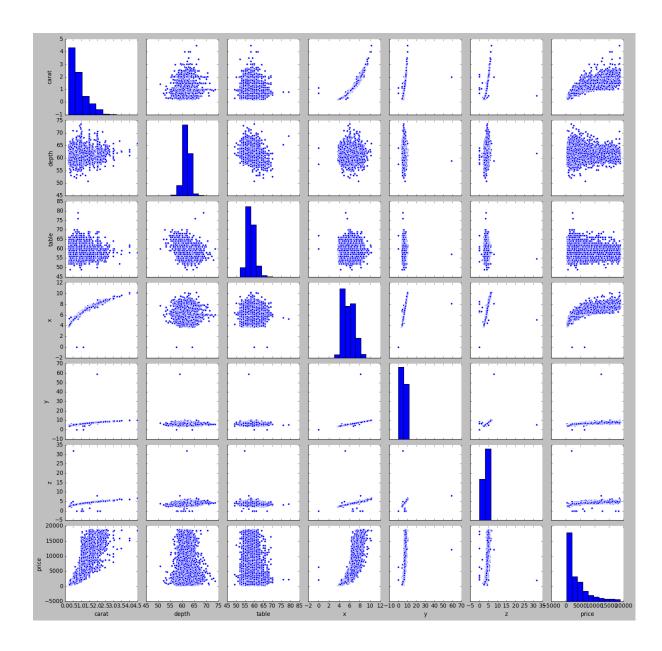
Mean is -0.000000, Median is -0.388720, Mode(s) are -0.8437 Column price has outliers



Column price is not normally distributed

Bi variate analysis:

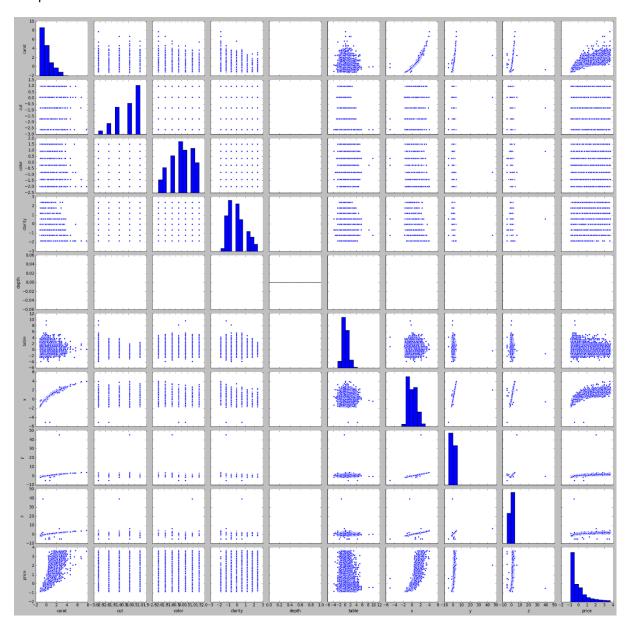
Pair plot for unscaled dataset:



Heatmap for unscaled dataset:

carat	1	0.035	0.18	0.98	0.94	0.94	0.92	- 0.90
depth	0.035	1	-0.29	-0.018	-0.024	0.097	-0.0022	- 0.75
table	0.18	-0.29	1	0.2	0.18	0.15	0.13	- 0.60
× -	0.98	-0.018	0.2	1	0.96	0.96	0.89	0.45
>-	0.94	-0.024	0.18	0.96	1	0.93	0.86	- 0.15
2 -	0.94	0.097	0.15	0.96	0.93	1	0.85	- 0.00
price	0.92	-0.0022	0.13	0.89	0.86	0.85	1 -	0.15
	carat	depth	table	×	ý	Z	price	

Pair plot for scaled dataset:



Heatmap for scaled dataset:

carat -	1	-0.14	-0.29	-0.36	'	0.18	0.98	0.94	0.94	0.92
cut	-0.14	1	0.027	0.18		-0.44	-0.13	-0.13	-0.15	-0.06
color	-0.29	0.027	1	-0.021		-0.024	-0.27	-0.26	-0.27	-0.17
clarity	-0.36	0.18	-0.021	1		-0.16	-0.38	-0.36	-0.36	-0.15
depth -										
table -	0.18	-0.44	-0.024	-0.16		1	0.2	0.18	0.15	0.13
x -	0.98	-0.13	-0.27	-0.38		0.2	1	0.96	0.96	0.89
у -	0.94	-0.13	-0.26	-0.36		0.18	0.96	1	0.93	0.86
Z -	0.94	-0.15	-0.27	-0.36		0.15	0.96	0.93	1	0.85
price -	0.92	-0.06	-0.17	-0.15		0.13	0.89	0.86	0.85	1
	carat	cut	color	clarity	depth	table	x	у	z	price

Observations:

It is noticed that price (target variable) depicts great correlation among carat, X, Y and Z independent variables across all combinations of them.

We also have to further observe if these variables are potential contributors for multi collinearity which is not good for explaining the coefficients. However, from a pure prediction standpoint these multi collinearity does not matter if we don't have to explain coefficients.

1.2: Impute null values if present, also check for the values which are equal to zero. Do they have any meaning or do we need to change them or drop them? Do you think scaling is necessary in this case?

There are null values in depth in 697 rows. There are various ways to handle missing values. Drop the rows, replace missing values with median values etc. It is a sizable data for us to retain as they have other independent variables contributing to the linear equation algorithm. Hence depth has been imputed with the median value.

Independent variables	Count of null values
carat	0
cut	0
color	0
clarity	0
depth	697
table	0

Х	0
У	0
Z	0

Below table suggests variables that have zeros as its value in the given dataset. Accordingly X, Y and Z are the variables that contain zeros across one or more rows.

Indepen dent variables	Zero values present (Yes/No)?
carat	No
cut	No
color	No
clarity	No
depth	No
table	No
X	Yes
У	Yes
Z	Yes

Below are the 9 rows that have 0's either in X, Y or Z. In the linear equation these rows would not contribute to prediction and hence no impact by retaining them. However dropping the rows could have adverse effects as other variables could potentially contribute with better coefficients in predicting price (target variable) which shall not be missed.

Index	carat	cut	color	clarity	depth	table	х	у	Z	price
5822	0.71	Good	F	SI2	64.1	60	0	0	0	2130
6035	2.02	Premium	Н	VS2	62.7	53	8.02	7.95	0	18207
6216	0.71	Good	F	SI2	64.1	60	0	0	0	2130
10828	2.2	Premium	I	SI1	61.2	59	8.42	8.37	0	17265
12499	2.18	Premium	Н	SI2	59.4	61	8.49	8.45	0	12631
12690	1.1	Premium	G	SI2	63	59	6.5	6.47	0	3696
17507	1.14	Fair	G	VS1	57.5	67	0	0	0	6381
18195	1.01	Premium	Н	l1	58.1	59	6.66	6.6	0	3167
23759	1.12	Premium	G	l1	60.4	59	6.71	6.67	0	2383

1.3: Encode the data (having string values) for Modelling. Data Split: Split the data into test and train (70:30). Apply Linear regression. Performance Metrics: Check the performance of Predictions on Train and Test sets using Rsquare, RMSE.

Please check below for the unique values across categorical variables and the respective count of number of rows:

```
CUT: 5 unique values
Fair
             781
            2441
Good
Very Good 6030
           6899
Premium
        10816
Ideal
Name: cut, dtype: int64
COLOR: 7 unique values
    1443
    2771
Ι
    3344
D
    4102
Η
F
    4729
Ε
    4917
    5661
G
Name: color, dtype: int64
CLARITY: 8 unique values
       365
I1
ΙF
       894
VVS1
       1839
VVS2
       2531
VS1
       4093
SI2
      4575
VS2
      6099
SI1
       6571
Name: clarity, dtype: int64
```

Since the linear algorithms require predictor variables to be numeric we shall assess if we need to convert them into codes or perform one hot encoding. Since in this case cut, color and clarity have values that indicate/depict the order we shall convert them into codes with highest value given to the variables as per below approach.

Based on the metadata clarified in the question, we wanted to assign the highest value depicting premium quality , best colour and best clarity for cut, color and clarity variables.

```
• Cut : Fair: 1, Good: 2, Very Good: 3, Premium: 4, Ideal: 5
```

• color: D:7, E:6, F:5, G:4, H:3, I:2, J:1

clarity: FL:11, IF:10, VVS1:9, VVS2:8, VS1:7, VS2:6, SI1:5, SI2:4, I1:3, I2:2, I3:1

Also, given that outliers can impact the model output while arriving at the linear equation and its coefficients, they shall be treated. In this case carat is treated for outliers beyond maximum whisker while depth, table, X, Y, Z and price are treated for outliers beyond both sides of the whisker. Below is the five number summary post the dataset has been treated for the outlier.

	carat	cut	color	clarity	depth	table	х	у	z	price
count	26967	2.70E+04	2.70E+04	2.70E+04	0	26967	26967	26967	26967	26967
mean	-0.01001	4.80E-16	1.49E-16	6.18E-17	NaN	-0.009131	0.000044	-0.00146	-0.00111	-0.050091
std	0.967962	1.00E+00	1.00E+00	1.00E+00	NaN	0.966443	0.998695	0.959443	0.967622	0.862377
min	-1.25252	-2.61E+00	-1.99E+00	-1.85E+00	NaN	-2.668463	-3.34946	-2.66846	-3.25843	-0.897815
25%	-0.83388	-8.17E-01	-8.17E-01	-6.39E-01	NaN	-0.652358	-0.90373	-0.87782	-0.88544	-0.744018
50%	-0.20592	8.12E-02	-2.31E-01	-3.22E-02	NaN	-0.204334	-0.03532	-0.02021	-0.02506	-0.38872
75%	0.526701	9.80E-01	9.41E-01	5.75E-01	NaN	0.691712	0.726761	0.691601	0.696552	0.352933
max	2.567575	9.80E-01	1.53E+00	2.40E+00	NaN	2.707817	3.172495	3.045731	3.069541	1.998361

Subsequently the split data into training and testing data at the ratio of 70:30 has the following five number summaries.

Train data:

	carat	cut	color	clarity	depth	table	х	у	z
count	18876	18876	18876	18876	18876	18876	18876	18876	18876
mean	-0.011	-0.00803	-0.001492	0.000503	0.004303	-0.007321	-0.00092	-0.00197	-0.00192
std	0.967735	0.999851	0.999856	0.995418	0.875865	0.964069	0.998108	0.959029	0.967608
min	-1.25252	-2.61367	-1.98943	-1.85372	-1.969594	-2.668463	-3.34946	-2.66846	-3.25843
25%	-0.83388	-0.81706	-0.81707	-0.6394	-0.463659	-0.652358	-0.91259	-0.87782	-0.88544
50%	-0.20592	0.081246	-0.23089	-0.03224	0.038319	-0.204334	-0.03532	-0.02879	-0.02506
75%	0.526701	0.97955	0.94147	0.574919	0.540298	0.691712	0.726761	0.691601	0.696552
max	2.567575	0.97955	1.52765	2.3964	2.046232	2.707817	3.172495	3.045731	3.069541

Test data:

i cot data.													
	carat	cut	color	clarity	depth	table	х	у	z				
count	8091	8091	8091	8091	8091	8091	8091	8091	8091				
mean	-0.00771	0.018739	0.003481	-0.00117	-0.000625	-0.013353	0.00228	-0.00026	0.000797				
std	0.968546	1.000221	1.00045	1.010731	0.87004	0.972004	1.000121	0.960467	0.967711				
min	-1.25252	-2.61367	-1.98943	-1.85372	-1.969594	-2.668463	-3.34946	-2.66846	-3.25843				
25%	-0.83388	-0.81706	-0.81707	-0.6394	-0.463659	-0.652358	-0.90373	-0.86924	-0.88544				
50%	-0.20592	0.081246	-0.23089	-0.03224	0.038319	-0.204334	-0.02645	-0.02021	-0.02506				
75%	0.526701	0.97955	0.94147	0.574919	0.540298	0.691712	0.7179	0.691601	0.696552				
max	2.567575	0.97955	1.52765	2.3964	2.046232	2.707817	3.163634	2.972834	3.069541				

Based on the regression model that has been built based on the above training and testing data set below are the statistical summary using the statsmodel library.

OLS Regression Results											
Dep. Variabl	.e:		ice R-squa			0.932 0.932					
Model:			_	Adj. R-squared:							
Method:		Least Squa		F-statistic:							
Date:	Su	n, 25 Oct 2	020 Prob (F-statistic)):	0.00					
Time:		11:00	:22 Log-Li	kelihood:		1293.0					
No. Observat		18	876 AIC:			-2566.					
Df Residuals	:	18	866 BIC:			-2488.					
Df Model:			9								
Covariance T	ype:	nonrob	ust								
	coef	std err	t	P> t	[0.025	0.975]					
Tatanant	0.0370	0.002	-22,466	0.000	-0.040	-0.034					
Intercept carat	-0.0370 1.0464	0.002	109.827	0.000	1.028	1.065					
cut color	0.0311	0.002	15.385	0.000	0.027	0.035					
	0.1161	0.002	66.597	0,000	0.113	0.119					
clarity	0.1772	0.002	96.577	0.000	0.174	0.181					
depth	-0.0034	0.003	-1.120	0.263	-0.009	0.003					
table	-0.0086	0.002	-3.946	0.000	-0.013	-0.004					
X	-0.4008	0.036	-11.137	0.000	-0.471	-0.330					
У	0.3530	0.037	9.426	0.000	0.280	0.426					
Z	-0.0427	0.016	-2.728	0.006	-0.073	-0.012					
Omnibus:	=======	2632.	707 Dunkin	-Watson:	.======	1.984					
	۸.										
Prob(Omnibus Skew:) •			-Bera (JB):		9984.946 0.00					
Skew: Kurtosis:			669 Prob(J 302 Cond.	,							
KUPTOSIS:		6.	SWZ Cond.	NO.		62.0					
========			========	========							

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Based on the null hypothesis that states that the independent variables have no co relationship with target variables from the universe, please note the coefficients of the predictor variables from the sample should be having p value lesser than the significance value of 5%. Based on the stats model based statistical observation predictor variables such as carat, cut, color, clarity, table, x, y and z shall be considered as a function of predicted variable (price) as their p value is much lesser than 5%. However for the predictor variable "depth" the p value indicates that the null hypothesis needs to be rejected and hence this will not contribute to the efficiency of the model. Hence "depth" shall be descoped from the linear equation.

Subsequently VIF has been evaluated to check multi collinearity among predictor variables. Any multi collinearity among predictor variables would mean that despite prediction itself need not be impacted with any low accuracy, the interpretation of the coefficients would be wrong without treating multi collinearity. Any VIF score between 1 to 5 is an acceptable range to retain the predictor variables in linear equation. Please find below the VIF for all the predictor variables considered so far without "depth" variable.

VIF for color :1.1199180000646427 VIF for clarity :1.239446323097618 VIF for cut :1.5097266302884316 VIF for table :1.6171628502013768 VIF for depth :2.8080406229119284 VIF for carat :31.360238644953316 VIF for z :104.8133511789732 VIF for y :366.3698715201059 VIF for x :379.79254326971653 As we can notice that variables such as x, y and z depicts much higher VIF than 5, further treatment on the linear equation to remove them one by one in the order of variable with highest VIF value taken out first, we could fine removing x, y and z results in much better VIF values removing all multi collinearity as below. Also the train and test data has been fixed to remove depth, x, y and z variables on which the model is getting executed with required predictor variables.

OLS Regression Results

Dep. Variabl	e:		p	rice	R-sq	uared:		0.931
Model:				OLS	Adj.	R-squared:		0.931
Method:		Least	: Squ	iares	F-st	atistic:		5.075e+04
Date:		Sun, 25	0ct	2020	Prob	(F-statistic):		0.00
Time:			11:0	8:27	Log-	Likelihood:		1174.3
No. Observat	ions:		1	.8876	AIC:			-2337.
Df Residuals	:		1	.8870	BIC:			-2290.
Df Model:				5				
Covariance T	ype:	r	nonro	bust				
	coet	f std	err		t	P> t	[0.025	0.975]
						0.000		
carat	0.9438					0.000		
cut	0.031					0.000		
table	-0.0089					0.000		
	0.115					0.000		0.119
clarity	0.182	3 0.	.002	100	.529	0.000	0.179	0.186
- "								
Omnibus:						in-Watson:		1.985
Prob(Omnibus):					que-Bera (JB):		7997.882
Skew:				.645				0.00
Kurtosis:			- 5	.917	Cond	. No.		1.86
=========	=======							

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Accordingly please find the revised VIF value for these remaining variables that has been considered in the revised linear equation.

VIF for color :1.1182180001886444
VIF for clarity :1.1987705393921508
VIF for cut :1.4942796531076148
VIF for table :1.580701029953695
VIF for depth :1.3231554005711066
VIF for carat :1.3031075405289636

We could notice the R squared value as well as adjusted R squared value has not changed much while we have the optimized set of predictor variables to make the linear equation to derive their coefficients.

Also the adjusted R square value (coefficient of determinant) is almost the same as R squared value indicating there is no statistical fluke indicating that there is no inconsistency on distribution between sample and the universe.

Please find below the coefficient of the predictor variables:

carat	0.9438
cut	0.0311
color	0.1155
clarity	0.1823

```
table -0.0089
```

Also, the intercept is -0.0384

Please find below the accuracy metrics/score:

- R Squared: 0.93 (Same as adjusted R Squared)
- Root mean square error: 0.2288135274586359

1.4: Inference: Basis on these predictions, what are the business insights and recommendations.

```
(-0.04) * Intercept + (0.94) * carat + (0.03) * cut + (-0.01) * table + (0.12) * color + (0.18) * clarity
```

Based on the coefficients making the above equation it can be understood that carat, x and y has the highest weight in that order towards predicting price with x having negative coefficient. Color and clarity also have reasonable coefficients that contribute to the predicting power for price.

So when carat increases by 1 unit price increases by 0.94 units while when cut increases by 1 unit price decreases by 0.01 unit and so on.

2.1: Data Ingestion: Read the dataset. Do the descriptive statistics and do null value condition check, write an inference on it. Perform Univariate and Bivariate Analysis. Do exploratory data analysis.

Please find below the data types for all the columns in the given dataset.

Column	Туре
Holliday Package	objec
	t
Salary	int64
age	int64
educ	int64
no_young_childre	
n	int64
no_older_childre	
n	int64
foreign	objec
10161811	t

Below provides five number summaries for all continuous variables such as Salary, age, educ, no_young_children and no_older_children while providing statistics on unique value count, top value with its frequencies for each categorical variable such as foreign and Holliday_Package.

	Holliday_Pack age	Salary	age	educ	no_young_child ren	no_older_child ren	foreig n
count	872	872	872	872	872	872	872
uniqu e	2	NaN	NaN	NaN	NaN	NaN	2
top	no	NaN	NaN	NaN	NaN	NaN	no
freq	471	NaN	NaN	NaN	NaN	NaN	656
mean	NaN	47729.17 2	39.955 28	9.307 34	0.311927	0.982798	NaN
std	NaN	23418.66 85	10.551 68	3.036 26	0.61287	1.086786	NaN
min	NaN	1322	20	1	0	0	NaN
25%	NaN	35324	32	8	0	0	NaN
50%	NaN	41903.5	39	9	0	1	NaN
75%	NaN	53469.5	48	12	0	2	NaN
max	NaN	236961	62	21	3	6	NaN

Null value check for all the columns in the dataset

Independent variables	Count of null value s
Holliday_Package	0
Salary	0
age	0
educ	0
no_young_childre n	0
no_older_childre n	0
foreign	0

There are no duplicates in the dataset.

Below is the normalized distribution of the predicted variable (Holliday_Package).

No: 54.01% Yes: 45.99%

The values across Yes and No are almost equally distributed.

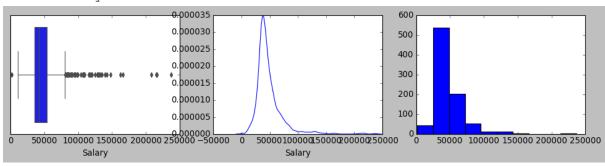
Univariate analysis:

Below analysis did not count in "Holidday_Package" and "foreign" columns as they are categorical and non numeric.

1. Univariate analysis for Salary

Mean is 47729.172018, Median is 41903.500000, Mode(s) are 32197.0000

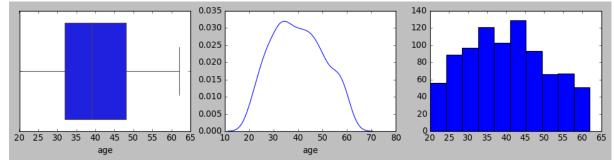
Column Salary has outliers



Column Salary is not normally distributed

2. Univariate analysis for age

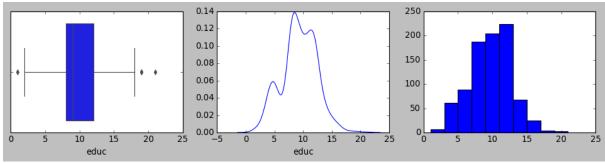
Mean is 39.955275, Median is 39.000000, Mode(s) are 44.0000 Column age does not have outliers



Column age is not normally distributed

3. Univariate analysis for educ

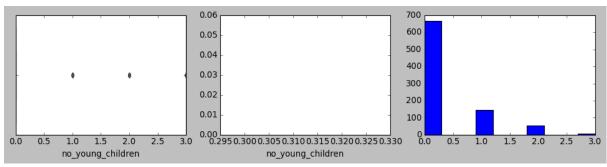
Mean is 9.307339, Median is 9.000000, Mode(s) are 8.0000 Column educ has outliers



Column educ is not normally distributed

4. Univariate analysis for no young children

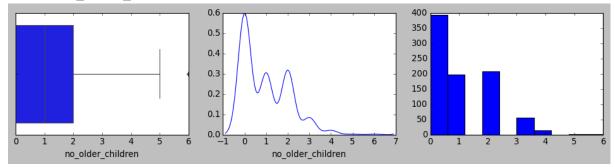
Mean is 0.311927, Median is 0.000000, Mode(s) are 0.0000 Column no_young_children has outliers



Column no young children is not normally distributed

5. Univariate analysis for no_older_children

Mean is 0.982798, Median is 1.000000, Mode(s) are 0.0000 Column no older children has outliers

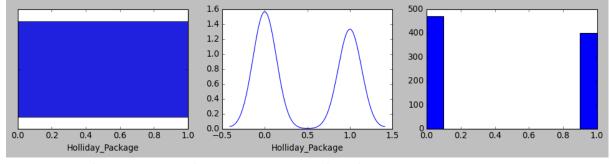


Column no_older_children is not normally distributed

Below analysis is done on the dataset after converting "Holidday_Package" into binary code (yes=1 and no=0) and "foreign" into dummies (foreign yes and foreign no).

Univariate analysis for Holliday_Package

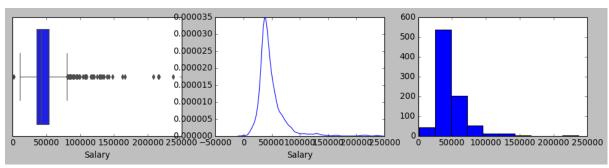
Mean is 0.459862, Median is 0.000000, Mode(s) are 0.0000 Column Holliday Package does not have outliers



Column Holliday Package is not normally distributed

2. Univariate analysis for Salary

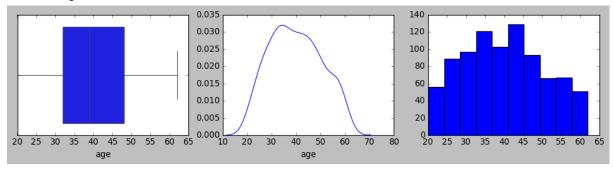
Mean is 47729.172018, Median is 41903.500000, Mode(s) are 32197.0000 Column Salary has outliers



Column Salary is not normally distributed

3. Univariate analysis for age

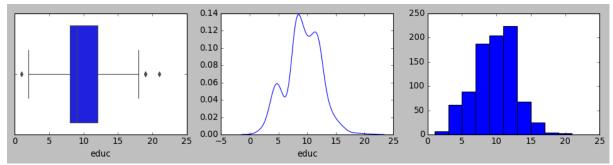
Mean is 39.955275, Median is 39.000000, Mode(s) are 44.0000 Column age does not have outliers



Column age is not normally distributed

4. Univariate analysis for educ

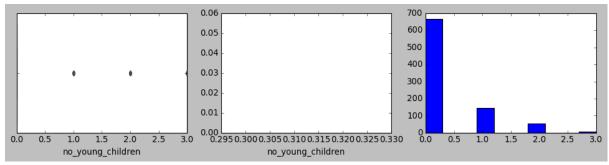
Mean is 9.307339, Median is 9.000000, Mode(s) are 8.0000 Column educ has outliers



Column educ is not normally distributed

5. Univariate analysis for $no_young_children$

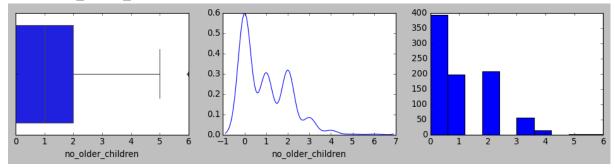
Mean is 0.311927, Median is 0.000000, Mode(s) are 0.0000 Column no_young_children has outliers



Column no young children is not normally distributed

6. Univariate analysis for no_older_children

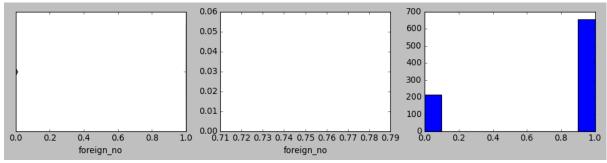
Mean is 0.982798, Median is 1.000000, Mode(s) are 0.0000 Column no older children has outliers



Column no_older_children is not normally distributed

7. Univariate analysis for foreign no

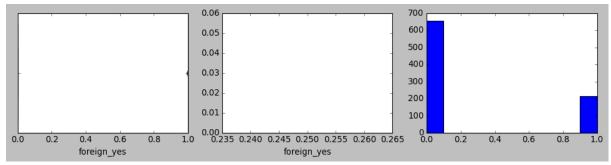
Mean is 0.752294, Median is 1.000000, Mode(s) are 1.0000



Column foreign_no is not normally distributed

8. Univariate analysis for foreign yes

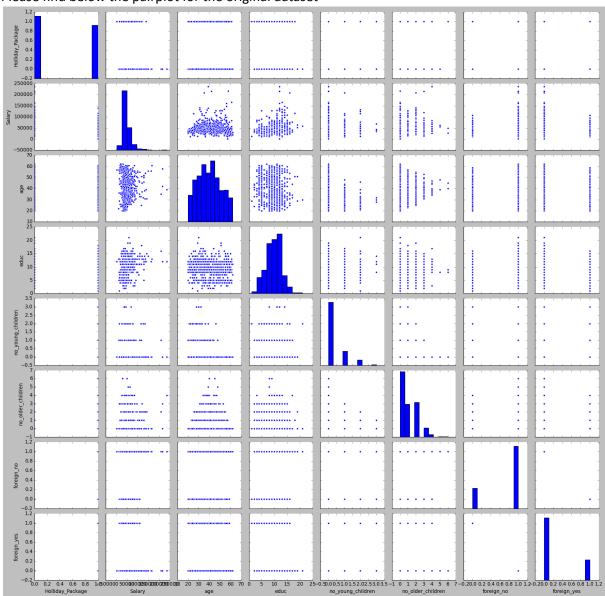
Mean is 0.247706, Median is 0.000000, Mode(s) are 0.0000 Column foreign yes has outliers



Column foreign_yes is not normally distributed

Bivariate analysis:

Please find below the pairplot for the original dataset



Please find below the heatmap for the original dataset



The above bi variate analysis clearly shows that there isn't a strong co relationship between many variables except for some traces of positive correlation observed between foreign(yes and no) and educ and some negative correlation between age and no_young_children.

Also the target variable (Holliday_Package) does not show a great deal of co relationship with other predictor variables while the predictor variables such as foreign(no and yes) followed by salary and no young children show minor negative co relationship with the predicted variable (target).

2.2: Do not scale the data. Encode the data (having string values) for Modelling. Data Split: Split the data into train and test (70:30). Apply Logistic Regression and LDA (linear discriminant analysis).

As salary alone showed a good number of outliers beyond maximum whisker it has been treated to match the maximum whisker accordingly along with minimal outliers beyond minimum whisker also being imputed.

Please find below the unique values across categorical variables and the respective count of number of rows. Also they are converted into numeric columns by

HOLLIDAY PACKAGE: 2 unique values

Yes: 401 rows No: 471 rows

FOREIGN: 2 unique values

Yes: 216 rows No: 656 rows

Also the categorical variable Holliday_Package has been type casted with binary values of 1 and 0 for 'yes' and 'no' respectively and hence making them align with reflecting on if the customer opted for vacation package(as 1) or not (as 0) respectively.

Also, the 'foreign' variable which does not have numerical significance has been converted to flag based (one hot encoding) variables with numeric binaries as their values across two

new variables namely foreign_yes and foreign_no depicting whether the customer is foreigner or not.

Please check below for the 5 number summary for the training data (70% of the sample) after the categorical object types have been converted to numeric as mentioned before.

	Salary	age	educ	no_young _children	no_old er_chil dren	foreign_no	foreign_yes
count	610	610	610	610	610	610	610
mean	46007.03	39.4557	9.372131	0.337705	0.9852	0.744262	0.255738
std	16024.82	10.4373	3.057341	0.64382	1.0542	0.436633	0.436633
min	8105.75	20	1	0	0	0	0
25%	35370.25	31	8	0	0	0	0
50%	42249.5	38	9	0	1	1	0
75%	53836.25	47	12	0	2	1	1
max	80687.75	62	21	3	5	1	1

Please check below for the 5 number summary for the testing data after the categorical object types have been converted to numeric as mentioned before.

	Salary	age	educ	no_young _children	er chii	foreign_no	foreign_yes
count	262	262	262	262	262	262	262
mean	44680.08	41.1183	9.156489	0.251908	0.9771	0.770992	0.229008
std	14903.91	10.7435	2.986934	0.530214	1.1611	0.420998	0.420998
min	14119	20	2	0	0	0	0
25%	35079.25	33	8	0	0	1	0
50%	40908.5	42	9	0	1	1	0
75%	52435.25	49	11	0	2	1	0
max	80687.75	62	19	3	6	1	1

Based on the regression model that has been built based on the above training and testing data set below are the statistical summary using the statsmodel library.

OLS Regression Results

Dep. Variable: Model:	Holliday		R-squared: Adj. R-squar	and.		.163 .154
Method:	Least		F-statistic			.9.54
Date:			Prob (F-stat		_	
Time:	500, 21	22:40:24				6.65
No. Observations:		610	AIC:			87.3
Df Residuals:		603	BIC:			18.2
Df Model:		6				
Covariance Type:	r	nonrobust				
	coef	std err	t	P> t	[0.025	0.975]
T-1						
					0.626	
•			-3.198		-6.6e-06	
•					-0.017	
educ	0.0121		1.667		-0.002	
no_older_children	-0.0096	0.019	-0.494	0.621	-0.048	0.029
no_young_children	-0.2714	0.036	-7.473	0.000	-0.343	-0.200
foreign_yes	0.5310	0.044	12.155	0.000	0.445	0.617
foreign_no	0.2607	0.053	4.953	0.000	0.157	0.364
0			Donald a United			
Omnibus:	-		Durbin-Watso			.015
Prob(Omnibus):			Jarque-Bera	(JB):	51	
Skew:			Prob(JB):		5.44	
Kurtosis:		1.610	Cond. No.		2.57	e+20

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The smallest eigenvalue is 2.19e-29. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

Based on the null hypothesis that states that the independent variables in the sample have no co relationship with target variable from the universe, please note the coefficients of predictor variables such as 'no_older_children' and 'educ' from the above statistical summary with p_value above significance value of 5% states that null will fly. This means those two variables cannot be part of predictor variables. However for the rest of the predictor variables such as Salary, age, no_young_children, foreign_yes and foreign_no the p_values is 0% and hence null hypothesis is rejected, which means their samples does not reflect the universe and potentially they are co related to the target variable "Holliday_Package".

Please find below the revised statistical summary post to the removal of "no_older_children" and "educ" as predictor variables. We are also proceeding for logistical regression and linear discriminant analysis model building by removing these predictor variables from the training and testing datasets.

OLS Regression Results

						====
Dep. Variable:	Holliday		R-squared:		e	.158
Model:		OLS	Adj. R-squar	red:	e	.153
Method:	Least	Squares	F-statistic:		2	8.45
Date:	Sat, 24 (Oct 2020	Prob (F-stat	istic):	1.12	e-21
Time:		22:35:12	Log-Likeliho	ood:	-38	8.28
No. Observations:		610	AIC:		7	86.6
Df Residuals:		605	BIC:		8	08.6
Df Model:		4				
Covariance Type:	ne	onrobust				
	coef	std err	t	P> t	[0.025	0.975]
Intercept	0.8415	0.069	12.247	0.000	0.707	0.976
· ·			-2.945			
•			-6.015			
no young children						
			13.526			
foreign no		0.042			0.220	
=======================================					=========	====
Omnibus:	1	5088.996	Durbin-Watso	n:	2	.027
Prob(Omnibus):		0.000	Jarque-Bera	(JB):	52	.756
Skew:			Prob(JB):		3.50	e-12
Kurtosis:		1.601				e+20
		========				====

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The smallest eigenvalue is 1.25e-28. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

Based on the revised predictor variables yielding above coefficients from statistical summary we could reject null hypothesis completely for the predictor variables Salary, age, no_young_children, foreign_yes and foreign_no predictor variables to proceed with resulting linear equation. However the coefficients of this linear equation is subject to the risk of multi collinearity and hence we check on the VIF for each of the predictor variables to check if they stay within the range of 1 to 5. Please find below the VIF accordingly:

VIF for Salary : 1.061323780479254 VIF for age : 1.3767750106591812 VIF for no_young_children : 1.370513314562445 VIF for foreign_yes : 7.9420154825877445 VIF for foreign_no : 25.500575817690464

Based on further analysis from the above VIF we could realize that foreign_no and Salary had to be removed in that sequence to bring the VIF within acceptable range as below:

Accordingly the revised summary from statsmodel library for OLS regression is as below:

OLS Regression Results

						==
Dep. Variable:	Holliday	Package	R-squared:		0.1	46
Model:		OLS	Adj. R-square	ed:	0.1	42
Method:	Least	Squares	F-statistic:		34.	60
Date:	Sun, 25 0	ct 2020	Prob (F-stati	stic):	1.19e-	20
Time:	1	2:48:18	Log-Likelihoo	od:	-392.	62
No. Observations:		610	AIC:		793	.2
Df Residuals:		606	BIC:		810	.9
Df Model:		3				
Covariance Type:	no	nrobust				
			t		_	0.975]
Intercept			10.697			1.162
age	-0.0127	0.002	-6.104	0.000	-0.017	-0.009
no_young_children	-0.2577	0.034	-7.615	0.000	-0.324	-0.191
foreign_yes	0.2685	0.043	6.246	0.000	0.184	0.353
						==
Omnibus:	8	034.241	Durbin-Watson	1:	2.0	54
Prob(Omnibus):			Jarque-Bera ((JB):	54.9	03
Skew:			Prob(JB):		1.20e-	
Kurtosis:		1.571	Cond. No.		20	6.
						==

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Based on the above summary it can be observed the coefficient of determinant (R squared) is slightly lesser than the original summary built before removing the multi collinearity. Higher the R squared value better the model is.

However, we shall proceed with running logistic regression and linear determinant model towards evaluating the model scores as per the built training and test data set and compare the efficiency of these two models.

2.3: Performance Metrics: Check the performance of Predictions on Train and Test sets using Accuracy, Confusion Matrix, Plot ROC curve and get ROC_AUC score for each model Final Model: Compare Both the models and write inference which model is best/optimized.

Please check below for model accuracy, confusion matrix, ROC curve and AUC score for training and testing data across logistic regression and linear discriminant analysis models.

Model accuracy:

Model accuracy	Training	Testing
Logistic Regression	0.66393443	0.65267176
LDA	0.65737705	0.65267176

Confusion matrix

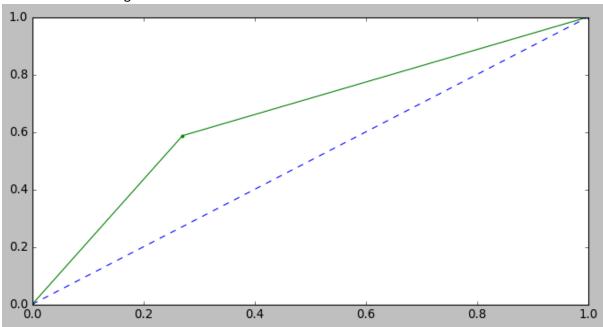
Logist	ic regression					LDA	
Training dataset	Predicted Negative	Predicted Positive	Training dataset		Training datas	Predicted Negative	Predicte
Actual Negative	240	89	Actual Negative	9	Actual Negation	245	5
Actual Negative	116	165	Actual Negative	5	Actual Negativ	125	5
	Predicted	Predicted				Predicted	
Testing dataset	Negative	Positive	Testing dataset		Testing datase	Negative	Predict
Actual Negative	107	35	Actual Negative	5	Actual Negativ	110)
Actual Negative	56	64	Actual Negative	1	Actual Negativ	59	

Classification report:

Logistic regression					LDA					
classifica	tion repo	ort for t	raining d	ata set:	classifica	classification report for training data se				
	precision	recall	f1-score	support		precision	recall	f1-score	sup	
0	0.67	0.73	0.7	329	0	0.66	0.74	0.7		
1	0.65	0.59	0.62	281	1	0.65	0.56	0.6		
accuracy			0.66	610	accuracy			0.66		
macro avg	0.66	0.66	0.66	610	macro avg	0.66	0.65	0.65		
weighted avg	0.66	0.66	0.66	610	weighted avg	0.66	0.66	0.65		
classific	ation rep	ort for t	esting da	ata set:	classific	ation report	for test	ing data	set	
	precision	recall	f1-score	support		precision	recall	f1-score	supp	
0	0.66	0.75	0.7	142	C	0.65	0.77	0.71		
1	0.65	0.53	0.58	120	1	0.66	0.51	0.57		
accuracy			0.65	262	accuracy			0.65		
macro avg	0.65	0.64	0.64	262	macro avg	0.65	0.64	0.64		
weighted avg	0.65	0.65	0.65	262	weighted avg	0.65	0.65	0.65		

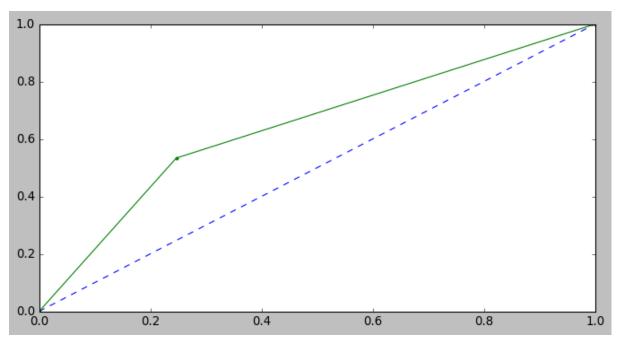
ROC/AUC from Logistic regression model:

ROC curve for training data



AUC for the Training Data: 0.658

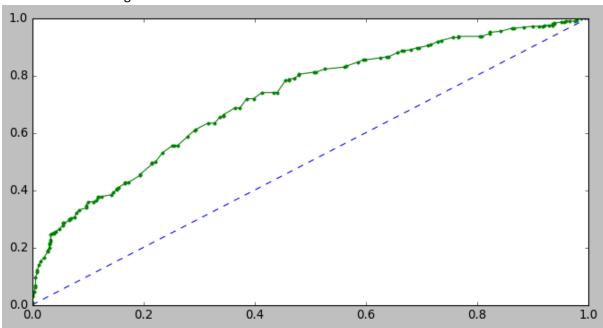
ROC curve for testing data



AUC for the Test Data: 0.643

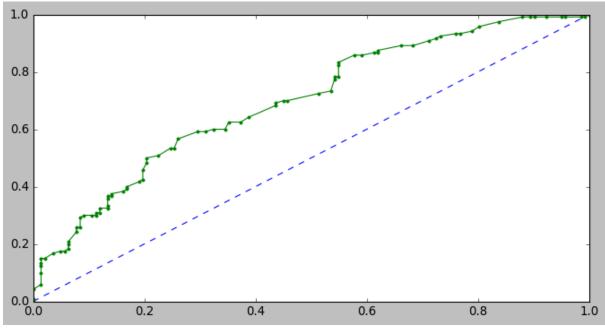
ROC/AUC from LDA:

ROC curve for training data



AUC for the Training Data: 0.720

ROC curve for testing data



AUC for the Test Data: 0.698

Accordingly, the below linear equation has been extracted upon which logistic regression model has applied sigmoid function to transform the raw numbers returned by the equation to probability to form the S curve wherein the probability ranges between 0 and 1.

```
(0.98) * Intercept + (-0.01) * age + (-0.26) * no_young_children + (0.27) * foreign yes
```

When we compare the performance of the model across accuracy scores and the classification report, we can notice that there isn't much difference on the performance across Logistic regression and LDA from the tuned models. There seems to be a very minor efficiency difference between the accuracy scores from the fact that logistic regression testing data outcome performed closer to the training outcome when compared to LDA.

2.4:Inference: Basis on these predictions, what are the insights and recommendations

Based on the analysis done on the data, it has been seen that Holliday_Package does not have a noticeable co relationship with any of the predictor variables. However, the distribution of Holliday_Package between customers who chose the vacation versus not is almost equal and hence the model has enough data to classify them well.

From an interpretation and insight perspective, there were few multicollinearity observed in the presence of Salary and non foreigners while education and number of older children data from the sample did not reflect the real data from the universe. This lead to limiting the predictive power of Holliday_Package to three aspects of the customer information on if it's a foreign customer and based on number of young children they had and the age they belong to. Out of them, foreign customers has the highest weightage in being able to sign up for the holiday package while number of young children had the next highest weightage clarifying that customer with lesser young children we able to choose the vacation package. There is also a slight tendency that customers signing up increases when the age decreases meaning mostly younger customers out of the lot has higher tendency to sign up for the vacation. This indicates potential customers could be parents of teenage students who are young adults while the parents themselves could be young enough within the age group of earning potentials coming from abroad.

Hence the recommendation would be to focus on offering family package deals with discounts during school or college holidays. This could enable upselling with that customer segment. Alternatively, for customers with younger children new vacation plans can be rolled out that are children friendly where they can spend more quality time in great resorts located near children attractions with reliable child care services that can generate interest level for that customer segment and hence provide increased market share opportunities.