

What is Diamonds Prices Dataset?

This document explores a dataset containing prices and attributes for approximately 54,000 round-cut diamonds. There are 53,940 diamonds in the dataset with 10 features (carat, cut, color, clarity, depth, table, price, x, y, and z). Most variables are numeric in nature, but the variables cut, color, and clarity are ordered factor variables with the following levels.#

About the currency for the price column: it is Price (\$)

And About the columns x,y, and z they are diamond measurements as ((x: length in mm, y: width in mm, z: depth in mm))

```
In [ ]: # There are 28 columns
#carat:weight of the diamond
#cut:quality of the cut (Fair, Good, Very Good, Premium, Ideal)
#color:Diamond colour, from J (worst) to D (best)
#clarity:how clear the diamond is (I1 (worst), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (best))
#x:length in mm
#y:width in mm
#z:Depth in mm
#price:Total depth percentage = z / mean(x, y) = 2 * z / (x + y) (43--79)
#table:Width of top of diamond relative to widest point (43--96)
#price:Price in US dollars (326--16,823)
```

importing the libraries

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [4]: df = pd.read_csv("D:/Datasets/Diamonds Prices2022.csv")
df.head()
```

```
Out[4]: Unnamed: 0  carat    cut  color  clarity  depth  table  price      x      y      z
0          0.23    Ideal    E   SI2    61.5   55.0   326    3.95    3.88    2.43
1          2.01  Premium    E   SI1    59.8   61.0   326    3.99    3.84    2.31
2          3.23    Good    E   VS1    56.9   65.0   327    4.09    4.07    2.31
3          4.29  Premium    I   VS2    62.4   58.0   334    4.20    4.23    2.61
4          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
5          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
```

Data Preprocessing

```
In [5]: df.shape
Out[5]: (53943, 11)
```

```
In [6]: #checking for null values
df.isnull().sum()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53943 entries, 0 to 53942
Data columns (total 11 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Unnamed: 0  53943 non-null   int64
 1   carat       53943 non-null   float64
 2   cut         53943 non-null   object
 3   color       53943 non-null   object
 4   clarity     53943 non-null   object
 5   depth       53943 non-null   float64
 6   table       53943 non-null   float64
 7   price       53943 non-null   int64
 8   x           53943 non-null   float64
 9   y           53943 non-null   float64
10  z           53943 non-null   float64
dtypes: float64(4), int64(2), object(5)
memory usage: 4.5+ MB
```

```
In [7]: #checking descriptive Values
df.describe()
```

```
Out[7]: Unnamed: 0    carat    cut  color  clarity  depth  table  price      x      y      z
count  53943.000000  53943.000000  53943.000000  53943.000000  53943.000000  53943.000000  53943.000000  53943.000000  53943.000000  53943.000000
mean    26970.200000    0.709395    61.740322    57.467261    61.922694    57.731658    57362.66    3.731658    3.736266    3.538720
std    15572.147122    0.472998    1.426266    2.235458    3.969338447    1.121730    11421.03    1.142103    1.142103    0.705678
min      1.000000    0.200000    43.000000    43.000000    326.000000    0.000000    0.000000    0.000000    0.000000    0.000000
25%   13486.500000    0.400000    61.000000    56.000000    56.000000    47.000000    47300.00    4.730000    4.730000    2.910000
50%   26972.000000    0.700000    61.800000    57.000000    57.000000    5.700000    57100.00    5.710000    5.710000    3.530000
75%   40487.500000    1.040000    62.500000    58.000000    58.000000    52.400000    65400.00    6.540000    6.540000    4.040000
max   53943.000000    5.010000    69.000000    96.000000    18823.000000    10.740000    58.900000    58.900000    31.800000
```

```
In [9]: #values count for categorical values
df["cut"].value_counts()
```

```
Out[9]: Ideal      21951
Premium    13793
Very Good  12083
Good       4986
Fair       3619
Name: cut, dtype: int64
```

```
In [10]: df["color"].value_counts()
```

```
Out[10]: G      11392
E       9709
F       9543
H       8394
D       8775
I       8422
J       2889
Name: color, dtype: int64
```

```
In [11]: df["clarity"].value_counts()
```

```
Out[11]: SI1      13867
VS2     12759
SI2     9194
VS1     8171
VVS2     5666
VVS1     3855
IF       1790
I1       741
Name: clarity, dtype: int64
```

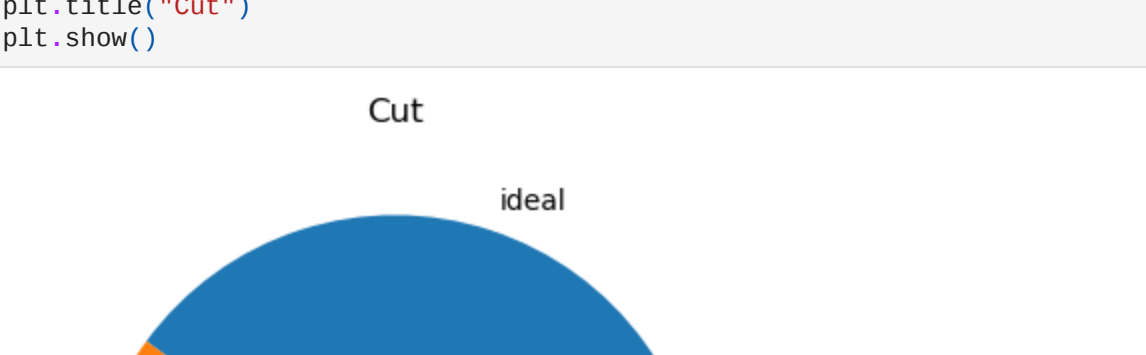
```
In [12]: df.head(10)
```

```
Out[12]: Unnamed: 0  carat    cut  color  clarity  depth  table  price      x      y      z
0          0.23    Ideal    E   SI2    61.5   55.0   326    3.95    3.88    2.43
1          2.01  Premium    E   SI1    59.8   61.0   326    3.99    3.84    2.31
2          3.23    Good    E   VS1    56.9   65.0   327    4.09    4.07    2.31
3          4.29  Premium    I   VS2    62.4   58.0   334    4.20    4.23    2.63
4          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
5          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
6          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
7          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
8          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
9          5.01    Good    J   SI2    63.3   58.0   335    4.34    4.35    2.75
```

Exploratory Data Analysis

```
In [13]: sns.histplot(df["price"],bins=20)
```

```
Out[13]: <AxesSubplot: xlabel='price', ylabel='Count'>
```



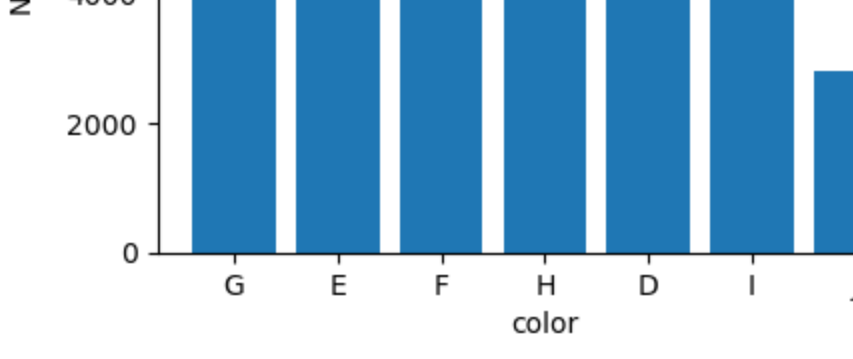
```
In [14]: sns.histplot(df["carat"],bins=20)
```

```
Out[14]: <AxesSubplot: xlabel='carat', ylabel='Count'>
```

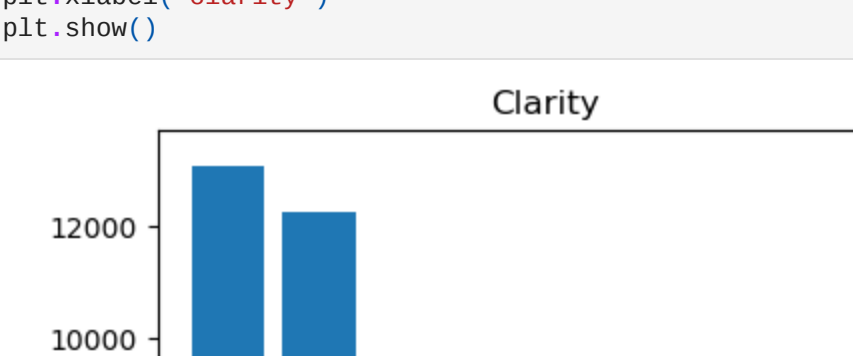


Most of the diamonds are less than 1 carat in weight

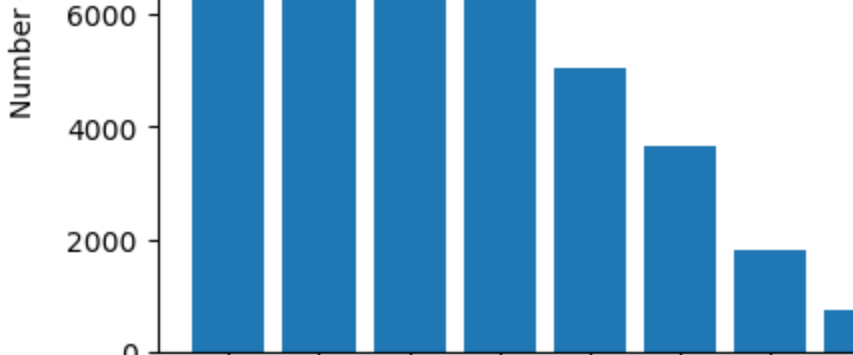
```
In [16]: plt.figure(figsize=(5,5))
plt.plot(df["cut"].value_counts(),label=["Ideal","Fair","Good","Very Good","Premium"])
plt.title("Cut")
plt.show()
```



```
In [17]: plt.figure(figsize=(5,5))
plt.bar(df["color"].value_counts().index,df["color"].value_counts())
plt.title("Number of diamonds")
plt.xlabel("color")
plt.show()
```

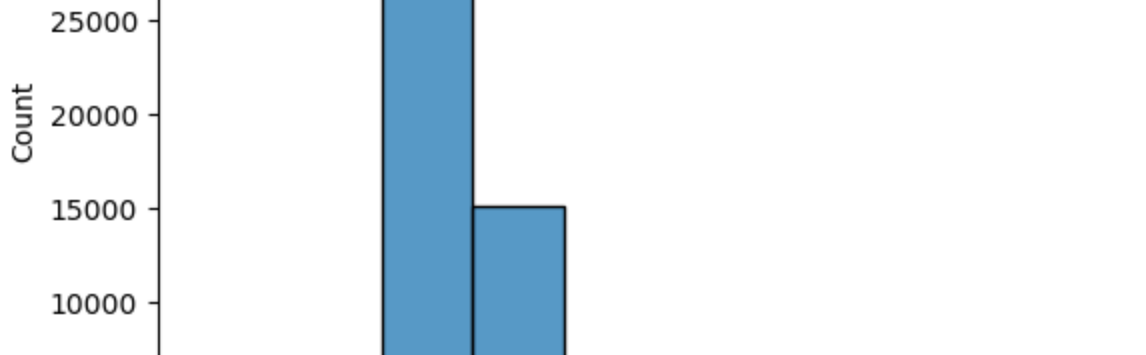


```
In [18]: plt.figure(figsize=(5,5))
plt.bar(df["clarity"].value_counts().index,df["clarity"].value_counts())
plt.title("Clarity")
plt.xlabel("Number of Diamonds")
plt.xlabel("clarity")
plt.show()
```



```
In [19]: sns.histplot(df["table"],bins=10)
```

```
Out[19]: <AxesSubplot: xlabel='table', ylabel='Count'>
```



Comparing Diamond's Features With Price

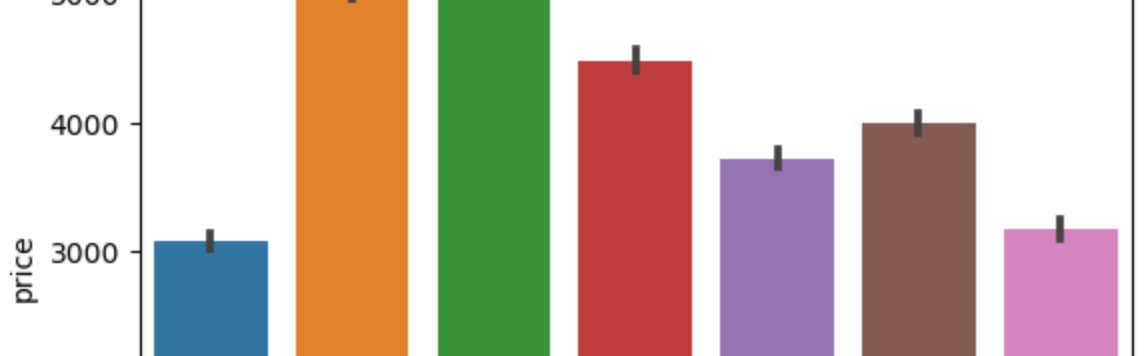
```
In [20]: sns.barplot(x="cut",y="price",data=df)
```

```
Out[20]: <AxesSubplot: xlabel='cut', ylabel='price'>
```



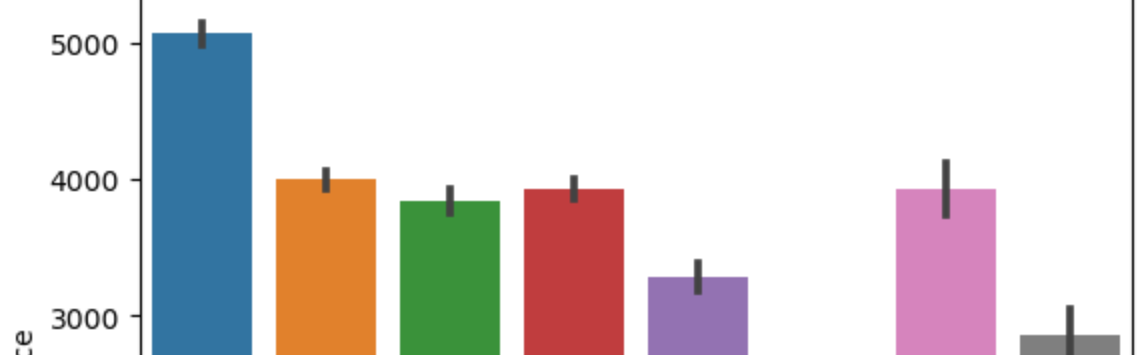
```
In [21]: sns.barplot(x="color",y="price",data=df)
```

```
Out[21]: <AxesSubplot: xlabel='color', ylabel='price'>
```



```
In [22]: sns.barplot(x="clarity",y="price",data=df)
```

```
Out[22]: <AxesSubplot: xlabel='clarity', ylabel='price'>
```



j color and I1 color are the worst features of the diamond,however when the data is plotted onthe bar graph,it seen the price of diamonds with J and I1 clarity is higher.

Than the price of diamonds of D color and IF clarity,which is opposite to each other what i expected

Data Preprocessing2

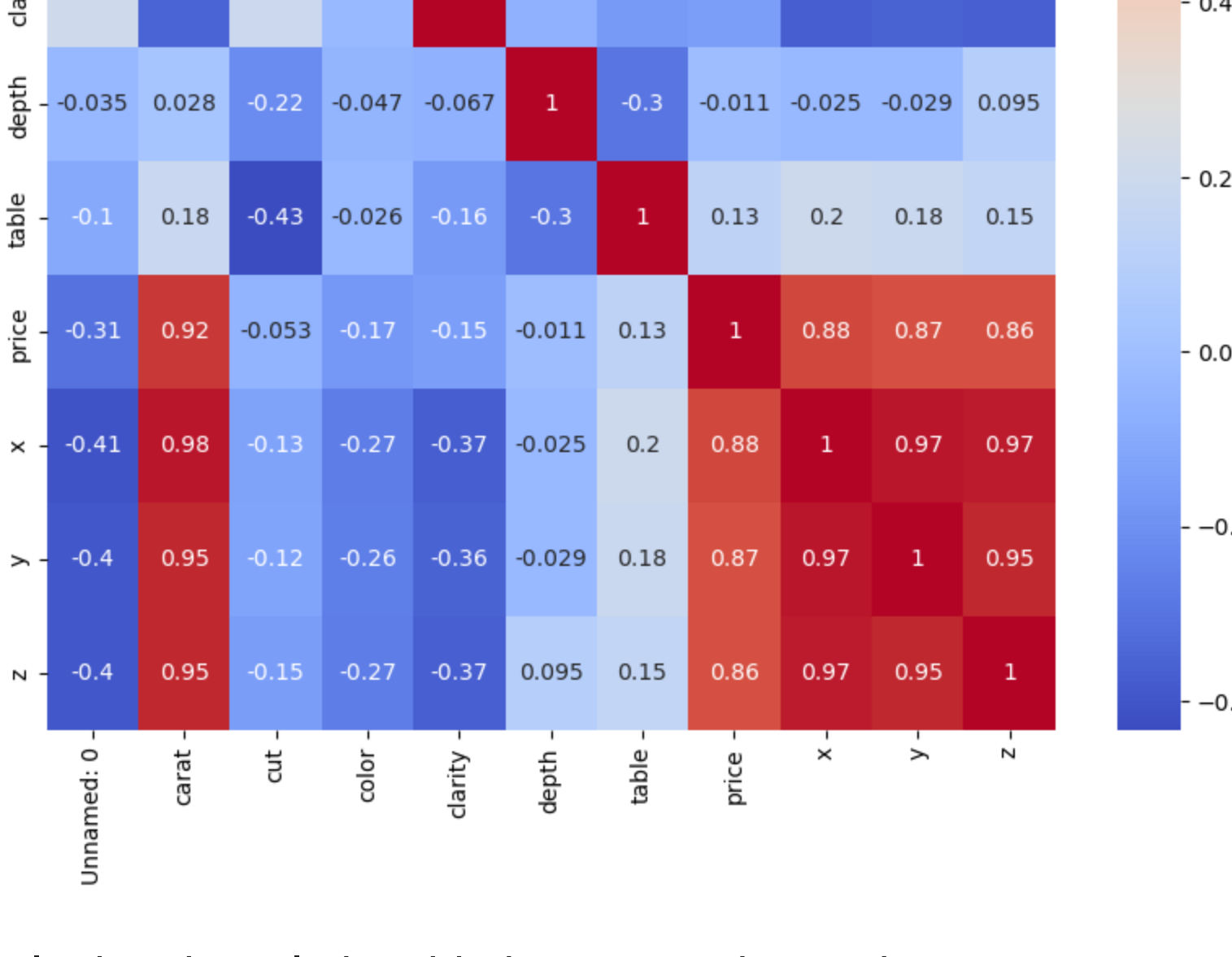
```
In [24]: df["cut"] = df["cut"].map({"Ideal":1,"Premium":2,"Very Good":3,"Good":4,"Fair":5})
df["color"] = df["color"].map({"D":1,"E":2,"F":3,"G":4,"H":5,"I":6,"J":7})
df["clarity"] = df["clarity"].map({"SI1":1,"SI2":2,"VS1":3,"VS2":4,"VVS1":5,"VVS2":6,"IF":7,"I1":8})
```

Coorelation

```
In [25]: #coorelation matrix
df.corr()
```

```
Out[25]: Unnamed: 0    carat    cut  color  clarity  depth  table  price      x      y      z
Unnamed: 0    1.000000  0.377670  0.096462  0.095149  0.206529  0.034896  0.101764  0.405405  0.398809  0.399135
carat          0.377670  1.000000  -0.134964  -0.291439  -0.325033  0.020234  0.181602  0.921581  0.975063  0.951721  0.953387
cut            0.096462  -0.134964  1.000000  0.020906  0.189171  -0.210037  -0.433387  -0.053487  -0.125664  -0.121461  -0.149320
color          0.095149  -0.291439  0.020906  1.000000  0.025646  -0.047316  -0.026457  -0.175200  -0.270281  -0.263579  -0.268226
clarity        0.206529  -0.325033  0.189171  -0.025646  1.000000  0.067355  -0.160328  -0.146791  -0.371996  -0.368417  -0.366946
depth          -0.034896  0.020234  -0.210037  -0.047316  0.067355  1.000000  -0.295788  -0.010630  -0.025289  -0.028340  0.046927
table          -0.101764  0.181602  -0.433387  -0.160328  -0.295788  1.000000  0.127118  0.193333  0.183750  0.150915
price          0.405675  0.921591  -0.053487  -0.175200  -0.146791  0.010630  1.000000  0.884433  0.885419  0.861249
x              0.398809  0.951721  -0.121461  -0.263579  -0.368417  -0.371996  0.884433  1.000000  0.970771  0.970771
y              0.399135  0.951721  -0.121461  -0.263579  -0.368417  -0.371996  0.885419  0.970771  1.000000  0.962005
z              0.399135  0.951721  -0.149320  -0.268226  -0.366946  0.046927  0.861249  0.970771  0.962005  1.000000
```

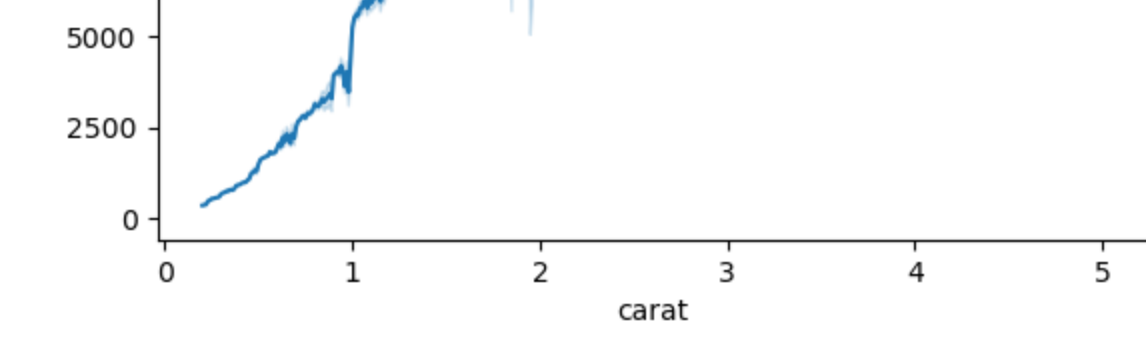
```
In [27]: #plotting the coorelation heatmap
plt.figure(figsize=(10,10))
sns.heatmap(df.corr(),annot=True,cmap="coolwarm")
plt.title("Coorelation Heatmap")
plt.show()
```



Plotting the relationship between price and carat

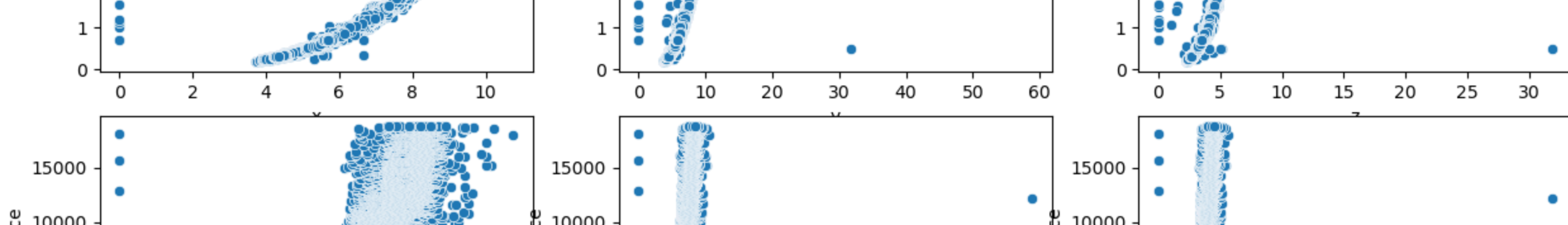
```
In [29]: sns.lnplot(x="carat",y="price",data=df)
```

```
Out[29]: <AxesSubplot: xlabel='carat', ylabel='price'>
```



From the lneploit it is quite to clear the price of diamonds is increasing with the increase of the carat of the diamonds.

```
In [31]: fig,ax=plt.subplots(2,3,figsize=(15,5))
ax=sns.scatterplot(x="y",y="carat",data=df,ax=ax[0,0])
sns.scatterplot(x="y",y="carat",data=df,ax=ax[0,1])
sns.scatterplot(x="y",y="carat",data=df,ax=ax[0,2])
sns.scatterplot(x="y",y="price",data=df,ax=ax[1,0])
sns.scatterplot(x="y",y="price",data=df,ax=ax[1,1])
sns.scatterplot(x="y",y="price",data=df,ax=ax[1,2])
plt.show()
```



```
In [31]: # Majority of the diamonds have x values between 4 and 6 , y values between 4 and 3, and z values between 2 and 6, Diamonds with other Dimensions are very rare.
```

Train And testing

```
In [32]: from sklearn.model_selection import train_test_split
x_train,x_train,y_test,y_train = train_test_split(df.drop('price',axis=1),df['price'],test_size=0.2,random_state=42)
```

Model Building

Decision Tree Regressor

```
In [33]: from sklearn.tree import DecisionTreeRegressor
dt = DecisionTreeRegressor()
```

```
Out[33]: DecisionTreeRegressor()
```

```
In [34]: #training the model
dt.fit(x_train,y_train)
```

```
Out[34]: <sklearn.metrics._mean_squared_error>
```

```
In [35]: #predicting the test set
dt_pred = dt.predict(x_test)
```

Random Forest Regressor

```
In [36]: from sklearn.ensemble import RandomForestRegressor
rf = RandomForestRegressor()
```

```
Out[36]: RandomForestRegressor()
```

```
In [37]: #training the model
rf.fit(x_train,y_train)
```

```
Out[37]: 0.9999750994313683
```

```
In [38]: #predicting the test set
rf_pred = rf.predict(x_test)
```

Model Evaluation

```
In [39]: from sklearn.metrics import mean_squared_error,mean_absolute_error
```

Decision Tree Regressor

```
In [40]: #distribution plot for actual and predicted values
ax = sns.distplot(y_test,hist=False,color='r',label='Actual Value')
```

```
sns.distplot(dt_pred,hist=False,color='b',label='Fitted Values',ax=ax)
```

```
sns.distplot(rf_pred,hist=False,color='b',label='Fitted Values',ax=ax)
```

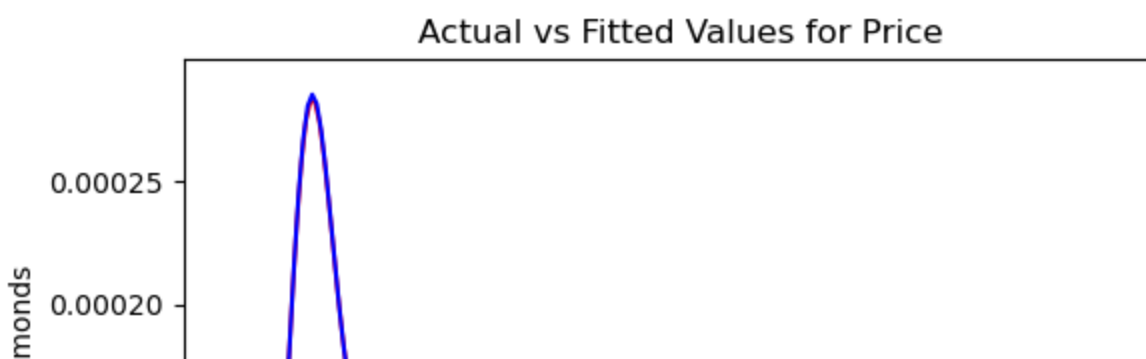
```
plt.xlabel('Actual vs Fitted values for Price')
```

```
plt.ylabel('Proportion of Diamonds')
```

```
plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).



```
In [41]: print("Decision Tree Regressor RMSE: ",np.sqrt(mean_squared_error(y_test,dt_pred)))
print("Random Forest Regressor Accuracy: ",r2_score(x_test,y_test))
print("Decision Tree Regressor MAE: ",mean_absolute_error(y_test,dt_pred))
```

```
Decision Tree Regressor RMSE: 14.6944827889647
Random Forest Regressor Accuracy: 0.9999999999999999
Decision Tree Regressor MAE: 10.18446398569931
```

Random Forest Regressor

```
In [42]: #distribution plot for actual and predicted values
ax = sns.distplot(y_test,hist=False,color='r',label='Actual Value')
```

```
sns.distplot(rf_pred,hist=False,color='b',label='Fitted Values',ax=ax)
```

```
sns.distplot(rf_pred,hist=False,color='b',label='Fitted Values',ax=ax)
```

```
plt.xlabel('Actual vs Fitted values for Price')
```

```
plt.ylabel('Proportion of Diamonds')
```

```
plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).


```
In [43]: print("Random Forest Regressor RMSE: ",np.sqrt(mean_squared_error(y_test,rf_pred)))
print("Random Forest Regressor Accuracy: ",r2_score(x_test,y_test))
print("Random Forest Regressor MAE: ",mean_absolute_error(y_test,rf_pred))
```

```
Random Forest Regressor RMSE: 84.844681294595
Random Forest Regressor Accuracy: 0.9999999999999999
Random Forest Regressor MAE: 8.24149498523235
```

Conclusion

Both the models having almost same accuracy. However, the Random Forest Regressor model is slightly better than the Decision Tree Regressor model.

There is something interesting about the data. The price of the diamonds with J color and I1 clarity is higher than the price of the diamonds with D color and IF clarity which couldn't be explained by the models. This could be because of the other factors that affect the price of the diamond.

```
In [ ]:
```