

FODS PROJECT PHASE - II

# SPARKLE VALUATE

## DIAMOND PRICE PREDICTION

---

Senthamizh Vani A P - 211701049

Meenakshi R - 211701031

# CONTENTS

---

1. ABSTRACT
2. UNDERSTANDING THE DATASETS
3. DATA PREPROCESSING
4. EXPLORATORY DATA ANALYSIS
5. VISUALISATION
6. MODEL BUILDING
7. MODEL EVALUATION



# ABSTRACT

---

The "Diamond Price Prediction" project aims to develop a data-driven model to accurately estimate the price of diamonds based on key features such as carat, cut, color, clarity, and other relevant attributes. Using advanced machine learning algorithms, the model analyzes patterns and trends in historical diamond pricing data to generate precise price predictions. This system can assist gemologists, jewelers, and buyers in making informed decisions, reducing the risk of overpricing or underpricing. By leveraging feature engineering and robust validation techniques, this project aspires to enhance transparency and efficiency in the diamond marketplace, bridging the gap between subjective evaluation and objective assessment.

# UNDERSTANDING THE DATASET

---

## Diamond Price Prediction dataset

- The aim of this analysis is to predict the price of diamonds based on their characteristics.
- The dataset used for this analysis is the Diamond dataset from kaggle.
- The dataset contains 50000 observation and 10 variables.



# UNDERSTANDING THE DATASET

---

## Dataset Attributes

- **Carat:** The weight of the diamond, one of the most important factors influencing the price. Larger diamonds are generally more valuable.
- **Cut:** The quality of the diamond's cut, affecting its symmetry and brilliance. Cut categories often include Excellent, Good, Fair, etc.
- **Color:** A measure of how colorless a diamond is, usually graded on a scale from D (colorless) to Z (light yellow or brown).
- **Clarity:** The presence of internal or external flaws (inclusions or blemishes), with a scale ranging from Flawless (FL) to Included (I1, I2, I3).
- **Depth:** The total depth percentage of the diamond, calculated as the depth divided by the average diameter.
- **Price:** The target variable, usually in USD, representing the price of the diamond.

# TYPES OF DATA

- **Categorical data**

Attributes such as cut and color represent categorical data because they denote specific categories. These attributes can be used to group and compare attributes of the diamond.

```
#Loading the dataset
df = pd.read_csv('diamonds.csv')
df.head()
```

	carat	cut	color	clarity	depth	table	price	x	y	z
0	0.23	Ideal	E	SI2	61.5	55.0	326	3.95	3.98	2.43
1	0.21	Premium	E	SI1	59.8	61.0	326	3.89	3.84	2.31
2	0.23	Good	E	VS1	56.9	65.0	327	4.05	4.07	2.31
3	0.29	Premium	I	VS2	62.4	58.0	334	4.20	4.23	2.63
4	0.31	Good	J	SI2	63.3	58.0	335	4.34	4.35	2.75

# TYPES OF DATA

- **Numerical data**
- In diamond price prediction, numerical data plays a crucial role in determining the value of diamonds based on specific measurable attributes. These attributes are typically numeric and quantitative, which allow machine learning models to predict diamond prices based on patterns.

```
#values count of categorical variables
print(df.cut.value_counts(),'\n',df.color.value_counts(),'\n',df.clarity.value_counts())
```

```
cut
Ideal      19938
Premium    12806
Very Good  11204
Good       4557
Fair        1495
Name: count, dtype: int64
color
G      10452
E      9085
F      8864
H      7711
D      6224
I      5058
J      2606
Name: count, dtype: int64
clarity
SI1     12115
VS2     11404
SI2     8519
VS1     7579
VVS2    4694
VVS1    3369
IF      1632
I1      688
Name: count, dtype: int64
```

# RELATIONSHIP BETWEEN THE FEATURES

---

1. **Carat Weight:** Heavier diamonds generally have higher prices, but the relationship is not linear.
2. Cut Quality: Higher quality cuts (Fair, Good, Very Good, Ideal, Premium) tend to increase the diamond's price.
3. **Color Grade:** Diamonds with less color (closer to D) are more valuable
4. Correlation Matrix: A matrix showing the correlation coefficients between different features (e.g., carat, cut, color) to identify strong relationships.
- .
5. **Feature vs. Target Analysis:** Analyzing how each feature (e.g., carat weight, cut quality) relates to the target variable (price) helps in understanding the impact of each feature on the price.
6. Feature vs. Price Analysis: Analyzing how each feature (e.g., carat weight, cut quality) relates to the target variable (price) helps in model building.
7. **Predictive Modeling:** Using machine learning models like Linear Regression, Random Forest, and Gradient Boosting to predict diamond prices with high accuracy.

# DATASET DESCRIPTION

```
#checking descriptive statistics  
df.describe()
```

	carat	depth	table	price	x	y	z
count	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000
mean	0.799444	61.753006	57.457830	3944.805440	5.734403	5.737956	3.541056
std	0.475173	1.431088	2.232092	3997.938105	1.123077	1.145579	0.707065
min	0.200000	43.000000	43.000000	326.000000	0.000000	0.000000	0.000000
25%	0.400000	61.000000	56.000000	951.000000	4.710000	4.720000	2.910000
50%	0.700000	61.800000	57.000000	2410.000000	5.700000	5.710000	3.530000
75%	1.040000	62.500000	59.000000	5351.000000	6.540000	6.540000	4.040000
max	5.010000	79.000000	95.000000	18823.000000	10.740000	58.900000	31.800000

# DATA PREPROCESSING

Steps involved in Data Preprocessing:

- Data cleaning
- Identifying and removing outliers
- Encoding categorical variables

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 10 columns):
 #   Column   Non-Null Count  Dtype  
---  --  
 0   carat    50000 non-null   float64
 1   cut      50000 non-null   object 
 2   color    50000 non-null   object 
 3   clarity  50000 non-null   object 
 4   depth    50000 non-null   float64
 5   table    50000 non-null   float64
 6   price    50000 non-null   int64  
 7   x        50000 non-null   float64
 8   y        50000 non-null   float64
 9   z        50000 non-null   float64
dtypes: float64(6), int64(1), object(3)
memory usage: 3.8+ MB
```

```
df.shape
```

```
(50000, 10)
```

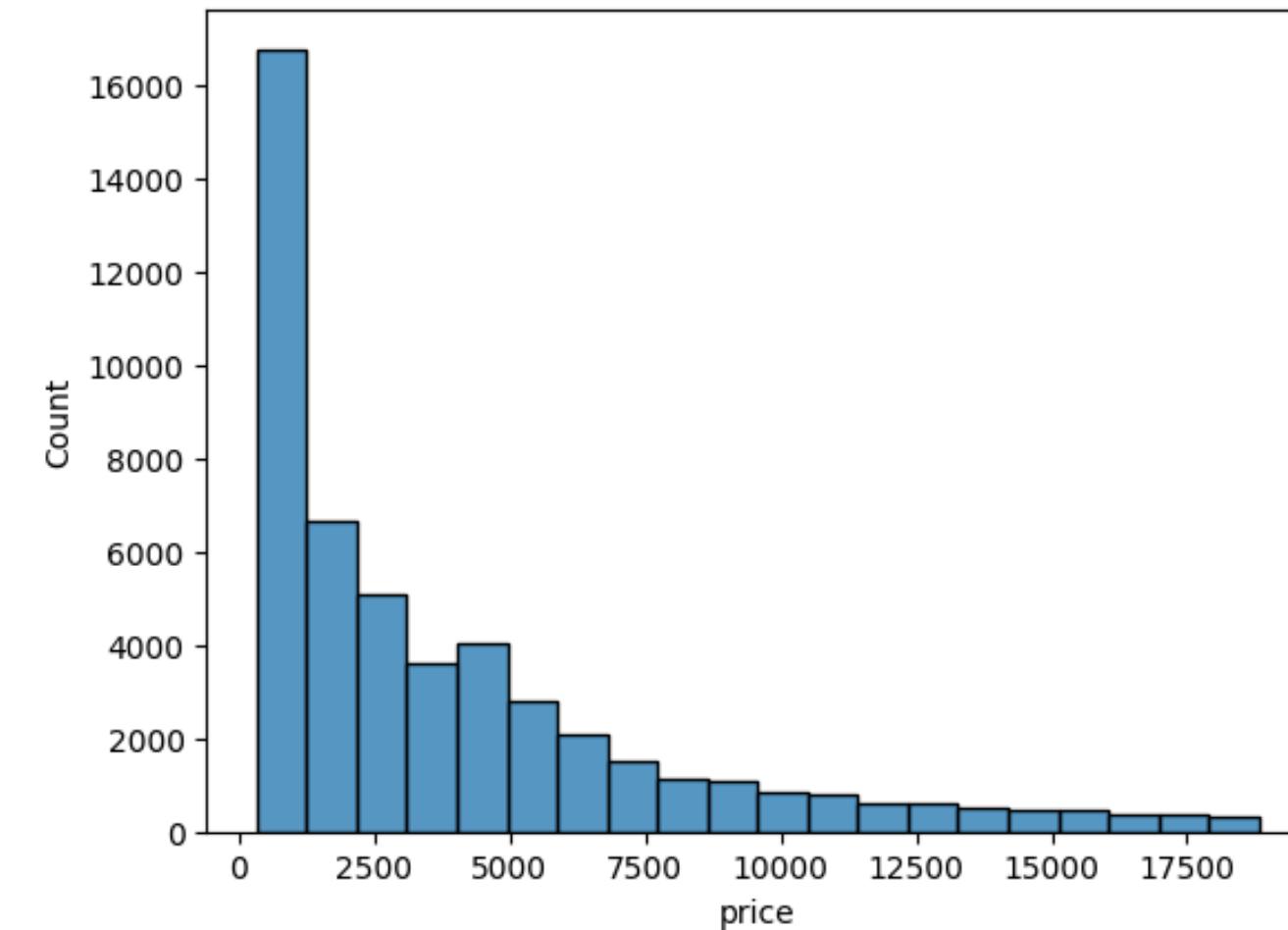
```
#checking for null values
df.info()
```

# EXPLORATORY DATA ANALYSIS

EDA helps you understand the structure, patterns, and potential issues in the data before building a model. For this project, the EDA process involves several stages, including data understanding, visualizations, and statistical analyses.

```
sns.histplot(df['price'], bins = 20)
```

```
<Axes: xlabel='price', ylabel='Count'>
```

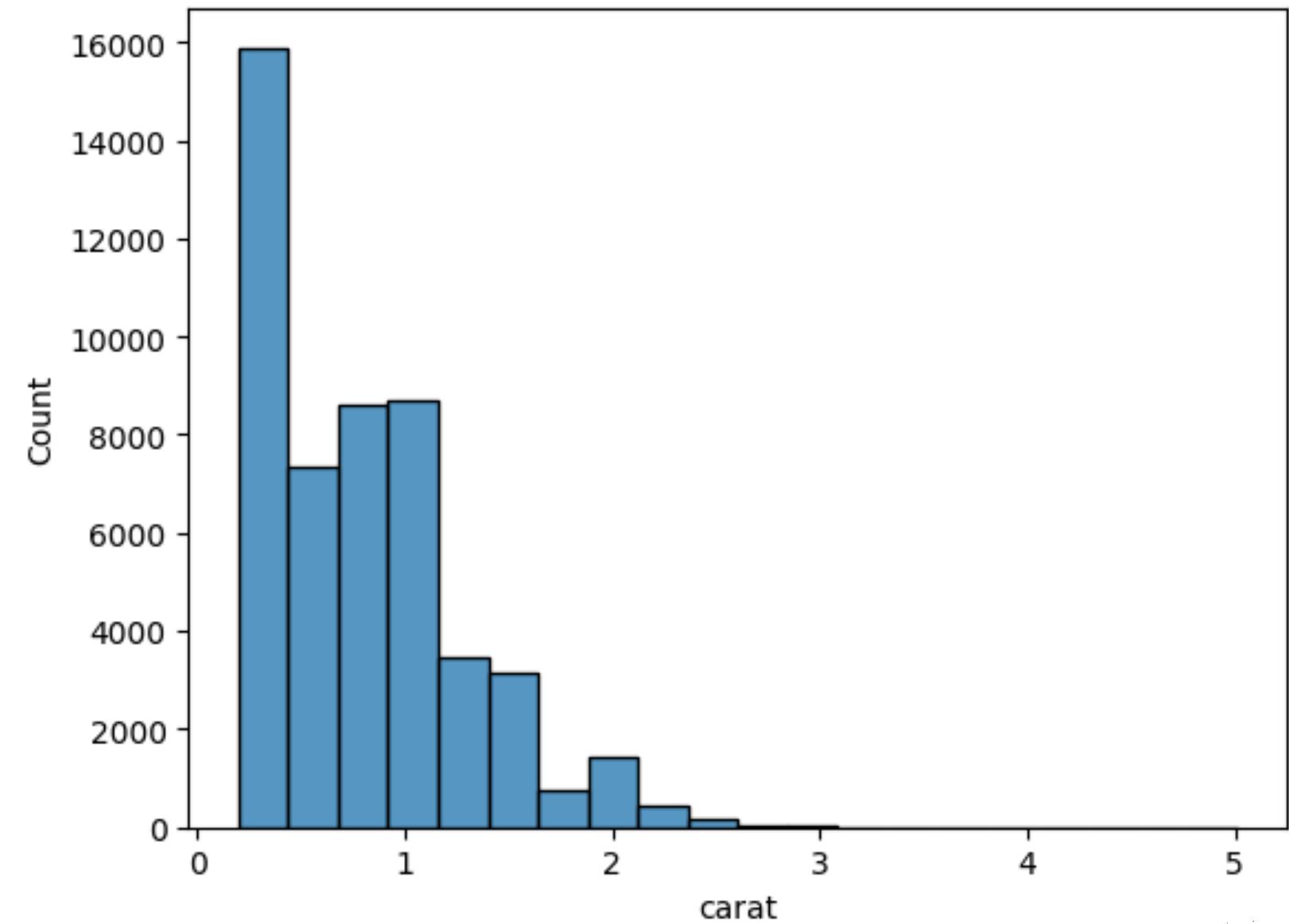


# EXPLORATORY DATA ANALYSIS

---

```
sns.histplot(df['carat'], bins=20)
```

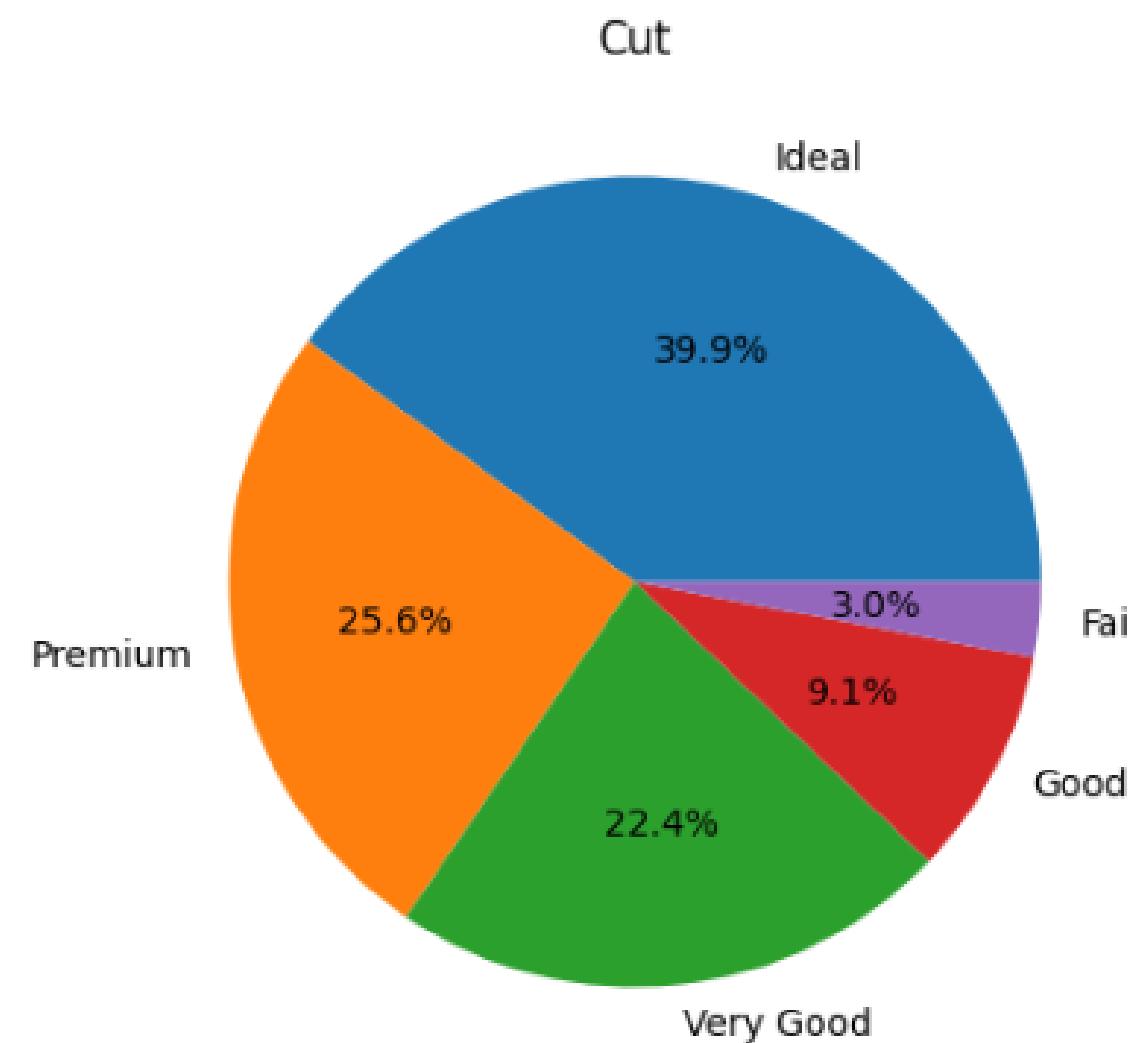
```
<Axes: xlabel='carat', ylabel='Count'>
```



# EXPLORATORY DATA ANALYSIS

## Composition of cut in diamond

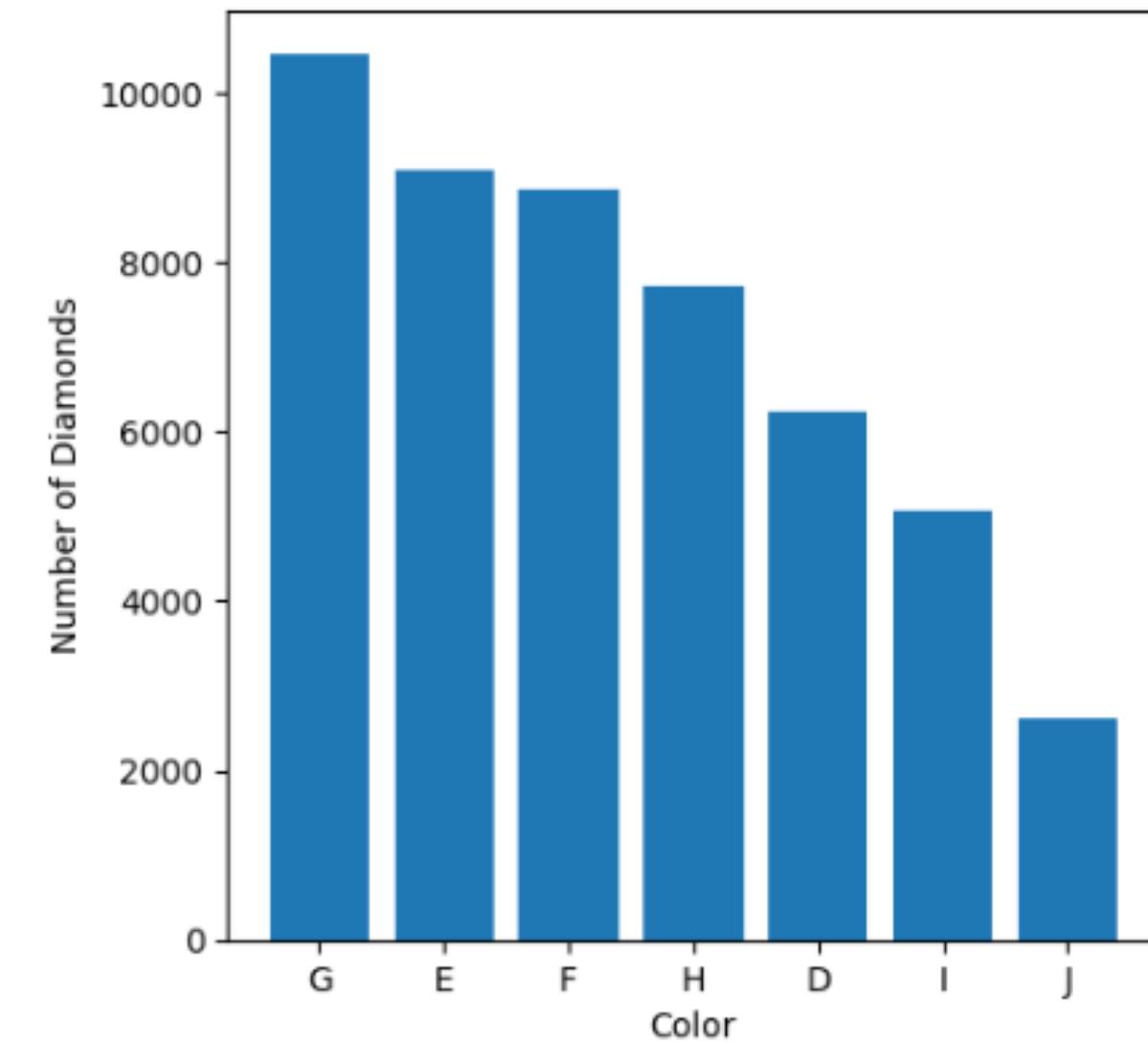
```
plt.figure(figsize=(5,5))
plt.pie(df['cut'].value_counts(),labels=['Ideal','Premium','Very Good','Good','Fair'],autopct='%1.1f%%')
plt.title('Cut')
plt.show()
```



# EXPLORATORY DATA ANALYSIS

relationship between number of diamonds vs color

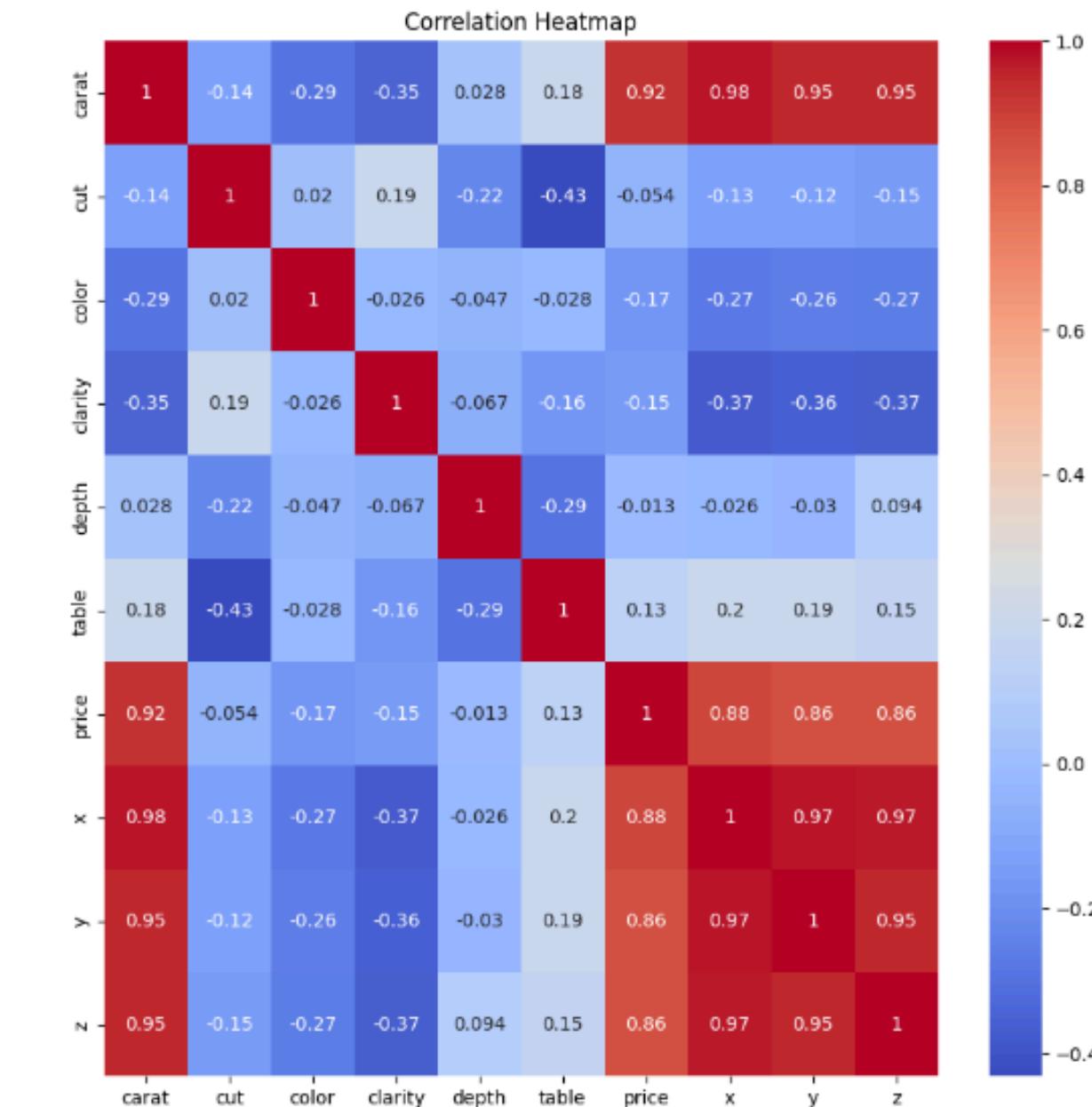
```
plt.figure(figsize=(5,5))
plt.bar(df['color'].value_counts().index,df['color'].value_counts())
plt.ylabel("Number of Diamonds")
plt.xlabel("Color")
plt.show()
```



# EXPLORATORY DATA ANALYSIS

## Correlation matrix

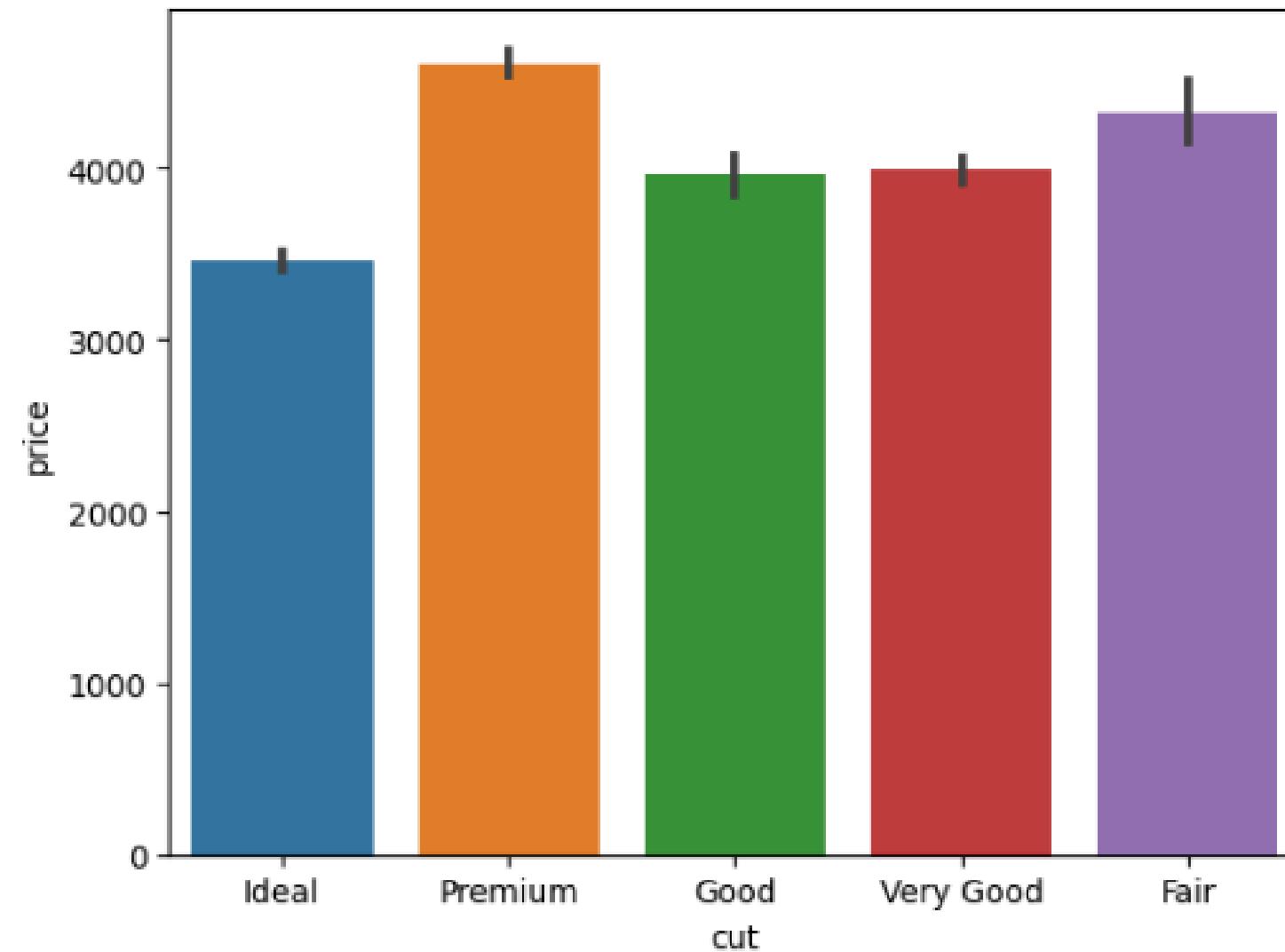
```
#Plotting the correlation heatmap
plt.figure(figsize=(10,10))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')
plt.show()
```



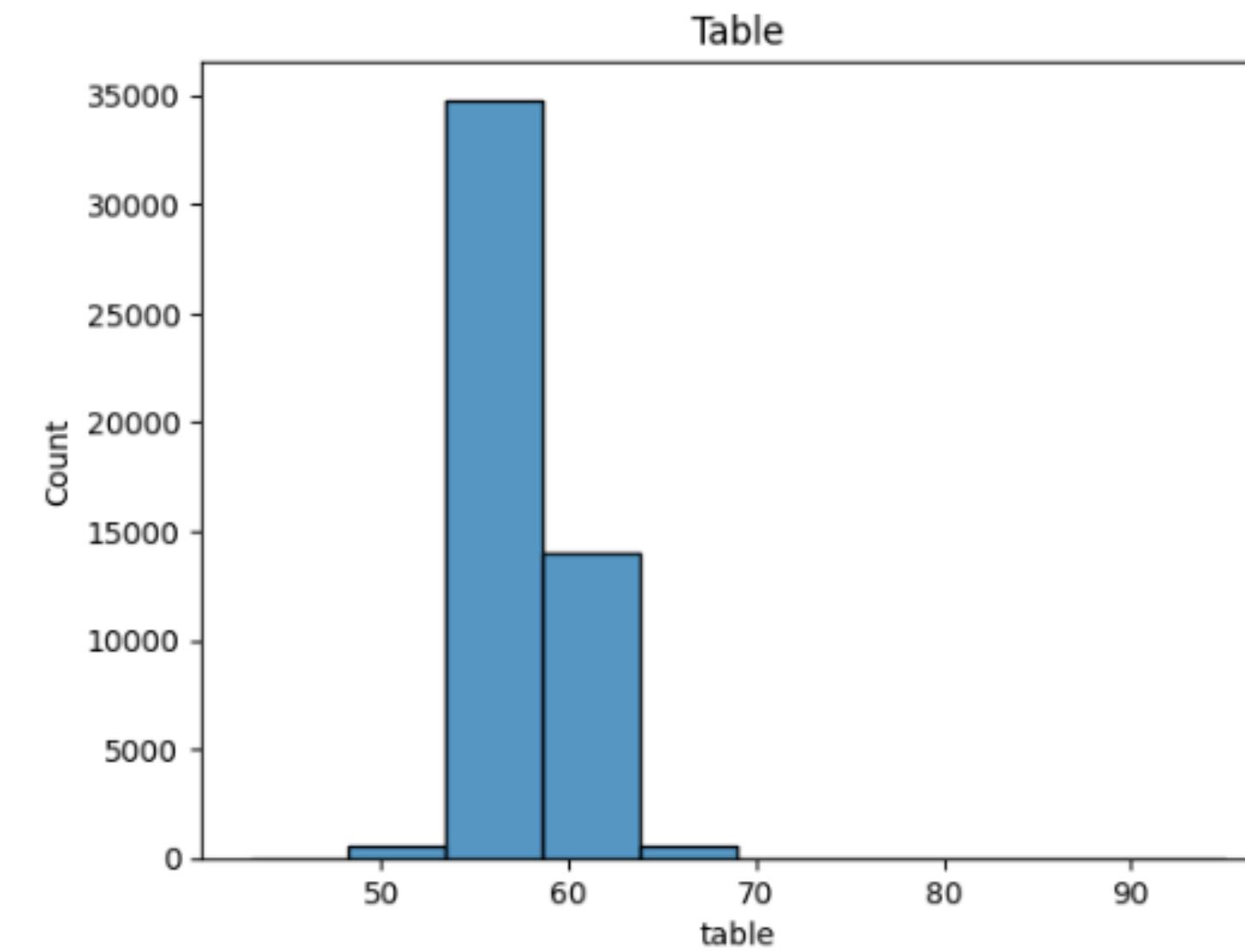
# EXPLORATORY DATA ANALYSIS

## Comparing Diamond's features with Price

```
- sns.barplot(x='cut',y='price',data=df)
- <Axes: xlabel='cut', ylabel='price'>
```

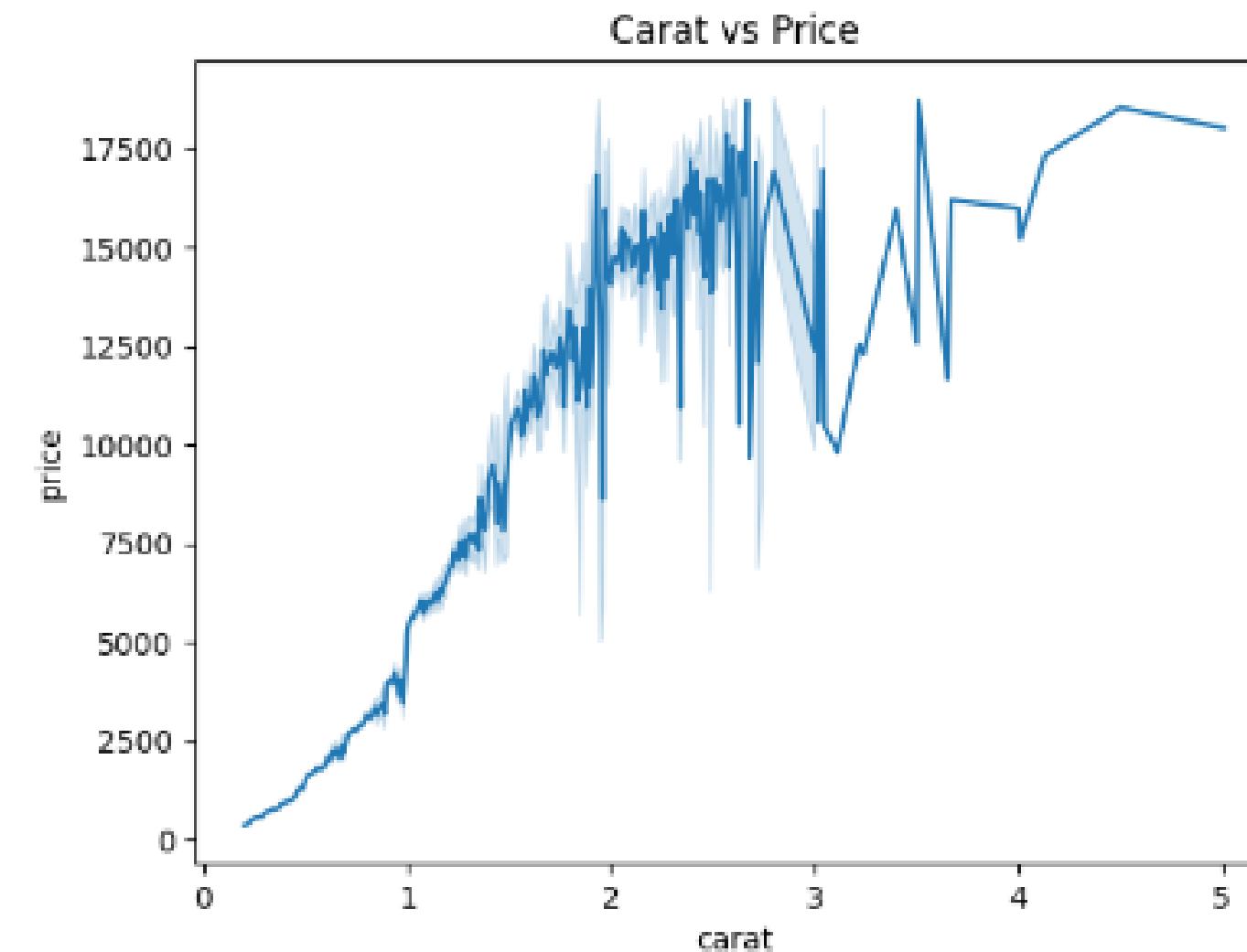


```
sns.histplot(df['table'],bins=10)
plt.title('Table')
plt.show()
```



# VISUALISATION

From the lineplot it is quite clear that the price of the diamond increases with the increase in the carat of the diamond. However, diamonds with less carat also have high price. This is because of the other factors that affect the price of the diamond.



Plotting the relationship between Price and Carat

```
sns.lineplot(x="carat",y="price",data=df)  
plt.title('Carat vs Price')  
plt.show()
```

# MODEL BUILDING

---

The model uses data science techniques to predict diamond prices based on features like carat, cut, color, and clarity. Steps include data preprocessing (encoding, scaling), exploratory analysis, and applying regression models (e.g., Linear Regression, Random Forests, XGBoost). Feature engineering and hyperparameter tuning ensure accurate, reliable predictions for market use

# MODEL BUILDING

## DECISION TREE REGRESSOR

```
: ▾ DecisionTreeRegressor  
DecisionTreeRegressor()
```

```
: #training the model  
dt.fit(x_train,y_train)  
#train accuracy  
dt.score(x_train,y_train)
```

```
: 0.999995617234543
```

```
: from sklearn.ensemble import RandomForestRegressor  
rf = RandomForestRegressor()  
rf
```

```
: ▾ RandomForestRegressor  
RandomForestRegressor()
```

```
: #training the model  
rf.fit(x_train,y_train)  
#train accuracy  
rf.score(x_train,y_train)
```

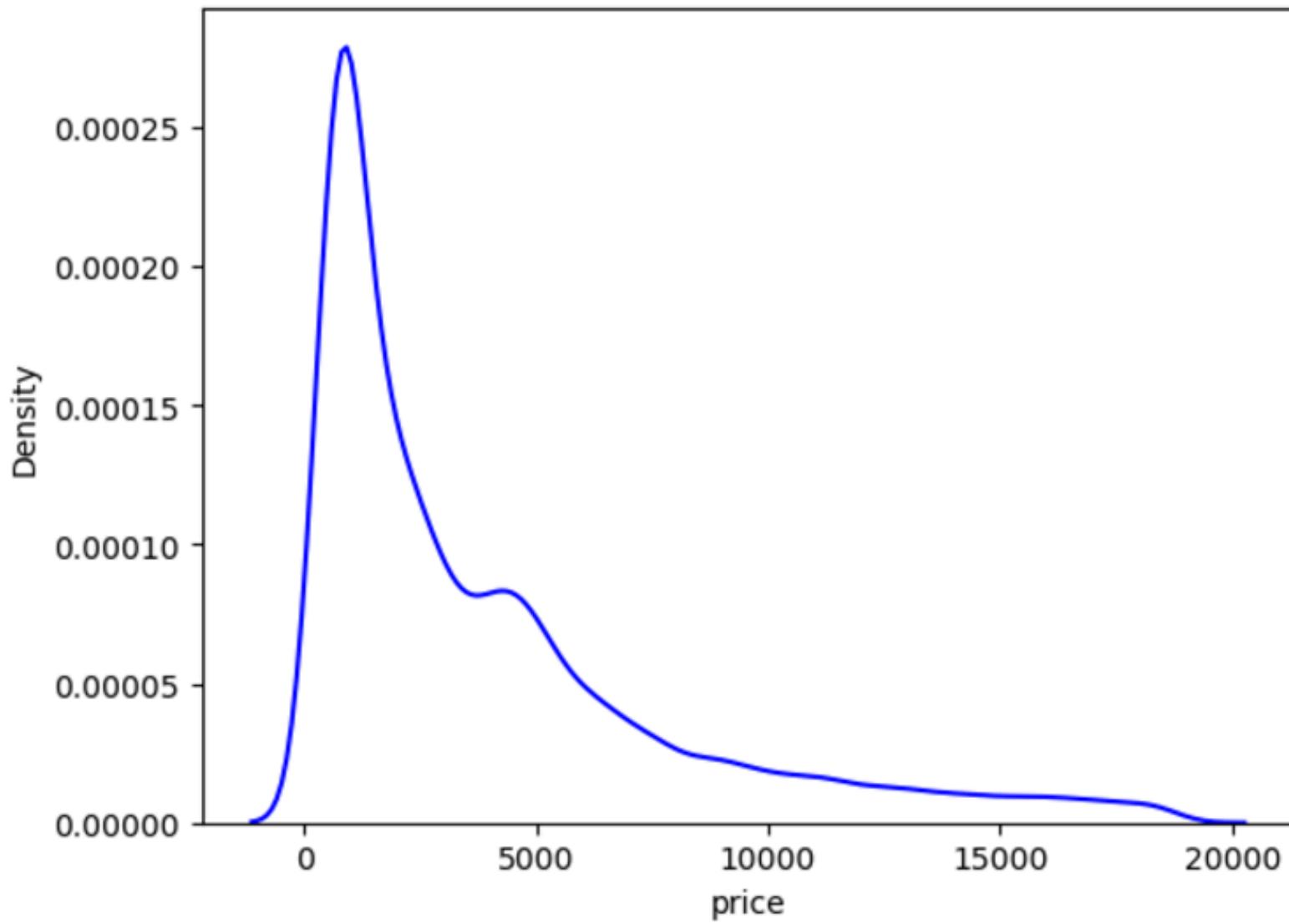
```
: 0.9970489911076822
```

The Decision Tree Regressor predicts diamond prices by creating a tree-based model that splits data on feature thresholds, capturing nonlinear relationships and providing interpretable, efficient predictions for complex datasets.

# MODEL EVALUATION

---

## Decision Tree Regressor



# CONCLUSION

---

Both the models have almost the 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.

**THANK YOU!**