# DIAMOND PRICE PREDICTION

By ~ Areeb Shaikh

# Data Overview

**Data Source :-**  **https://www.kaggle.com/datasets/shivam2503/diamonds**

Kaggle (Diamond Price Dataset with 50,000 records).

Features:-

1. Numerical: Carat, Dimensions (x, y, z).

2. Categorical: Cut (Fair, Good, Very Good, Premium, Ideal), Color (J to D), Clarity (I1 to IF).

3. **Target Variable:-** Price in USD.

# Project Overview

▶ **Problem Statement**

In the diamond industry, accurate pricing is a critical challenge due to the variability of diamond features such as carat, cut, color, clarity, and dimensions. The pricing process often involves manual evaluation, which can be subjective and inconsistent. This creates the need for a robust and data-driven solution to estimate diamond prices reliably and transparently.

▶ **Objective**

The objective of this project is to predict the price of diamonds based on key features such as carat, cut, color, clarity, and dimensions. By leveraging machine learning techniques, the project aims to build an accurate predictive model that can estimate diamond prices, providing insights for the jewelry industry and potential buyers
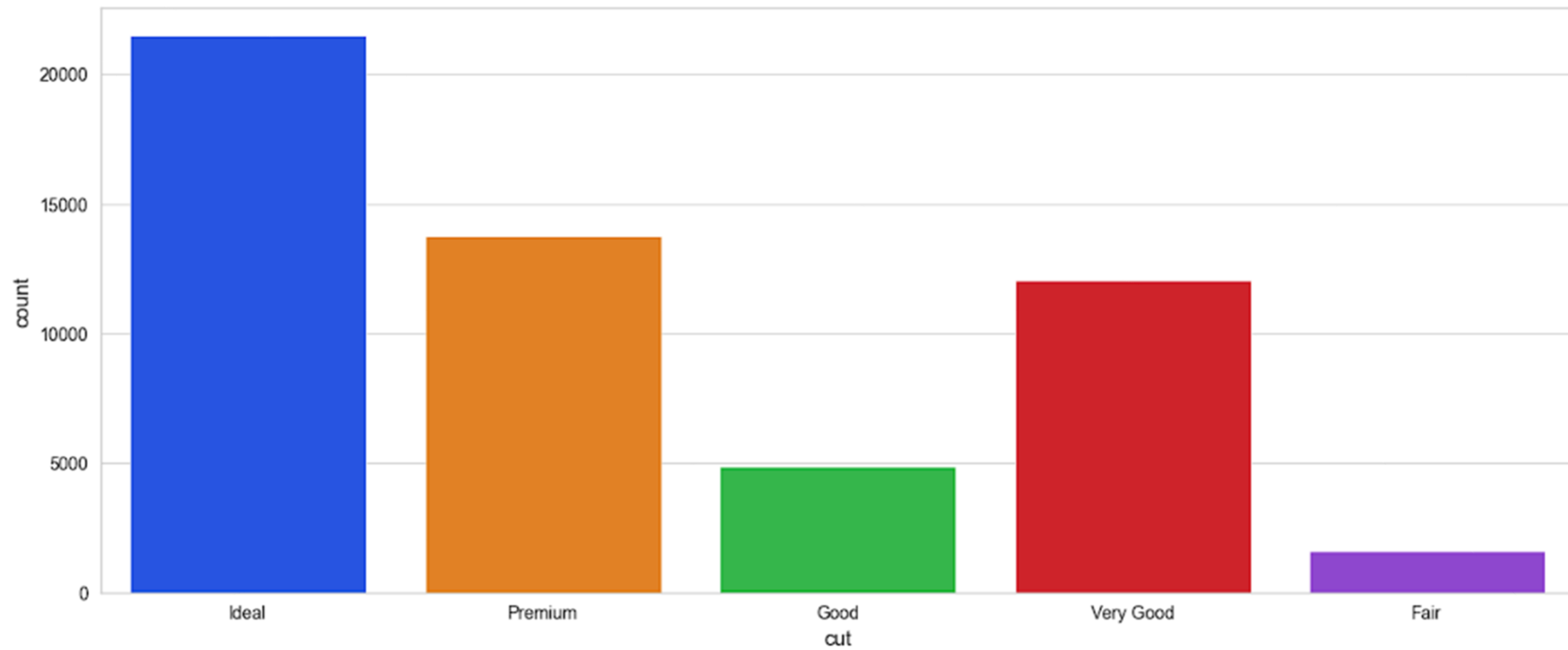
# Methodology

▶ **Data Preparation :-**

1) Perform exploratory data analysis (EDA) to identify patterns and trends.

➢ Univariate , Bivariate, Multivariate

2) Handle missing & Duplicate data (if any) and ensure consistent formatting.

3) Retain outliers as they reflect valid variability in diamond characteristics.

4) Feature Engineering

➢ Created New Attribute (Size) By Combining Existing Attribute (x, y, z)

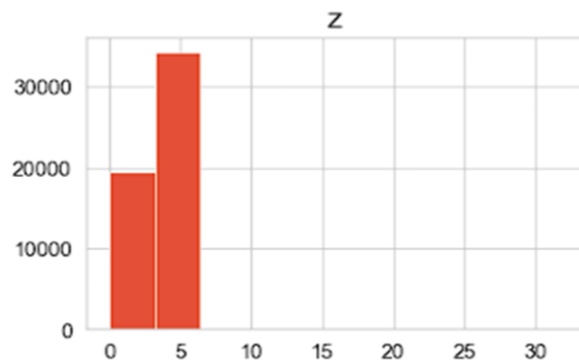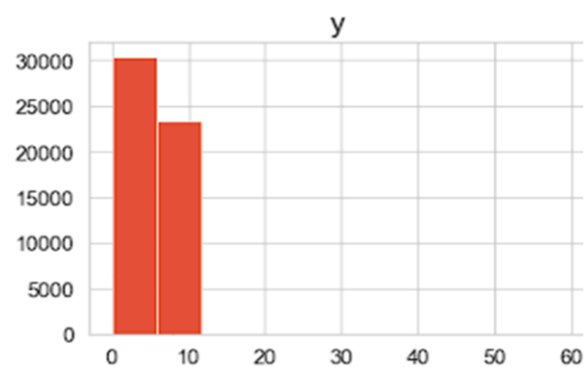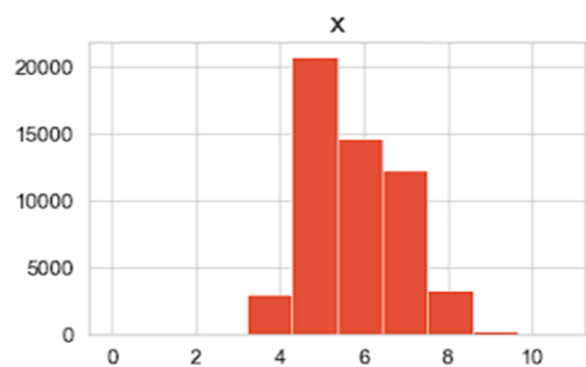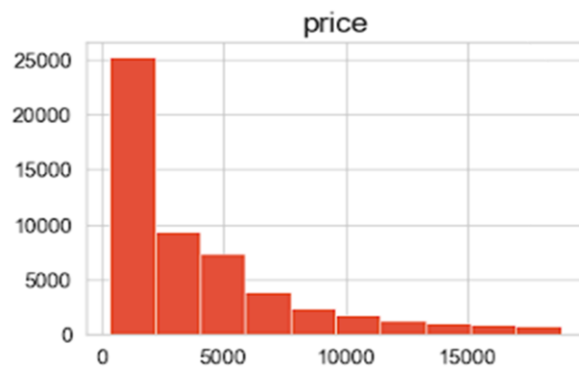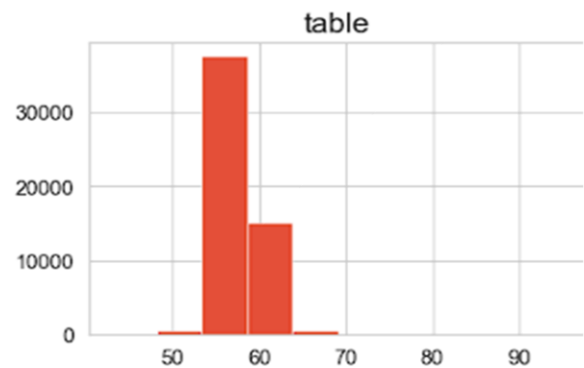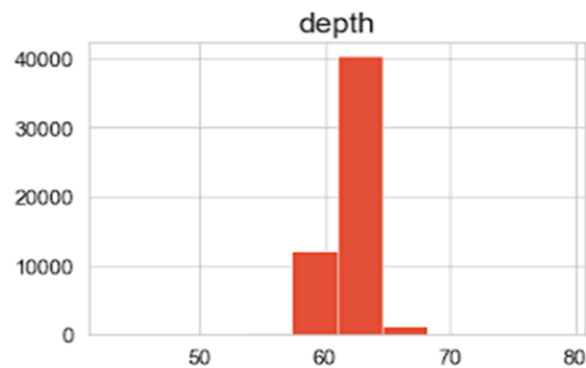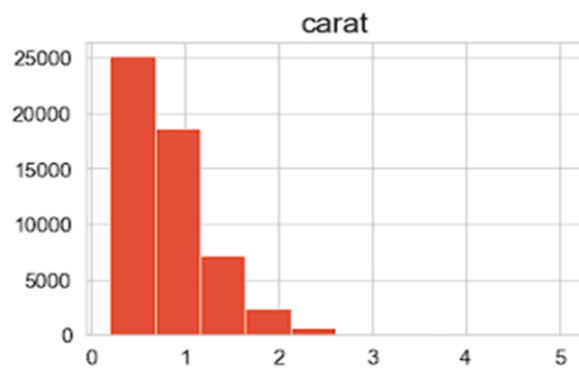5) Label Encoding To Categorical Features

# Insights & Visuals

▶ **Used visualizations**

1) Bar Plot (Categorical Distribution)

2) Histogram (Numerical Distribution)

3) Scatter Plot (Carat vs Price)

4) Scatter Plot (Size vs Price)

5) Box Plot

6) Scatter + Regression plot

7) Heatmap (Correlation Matrix)
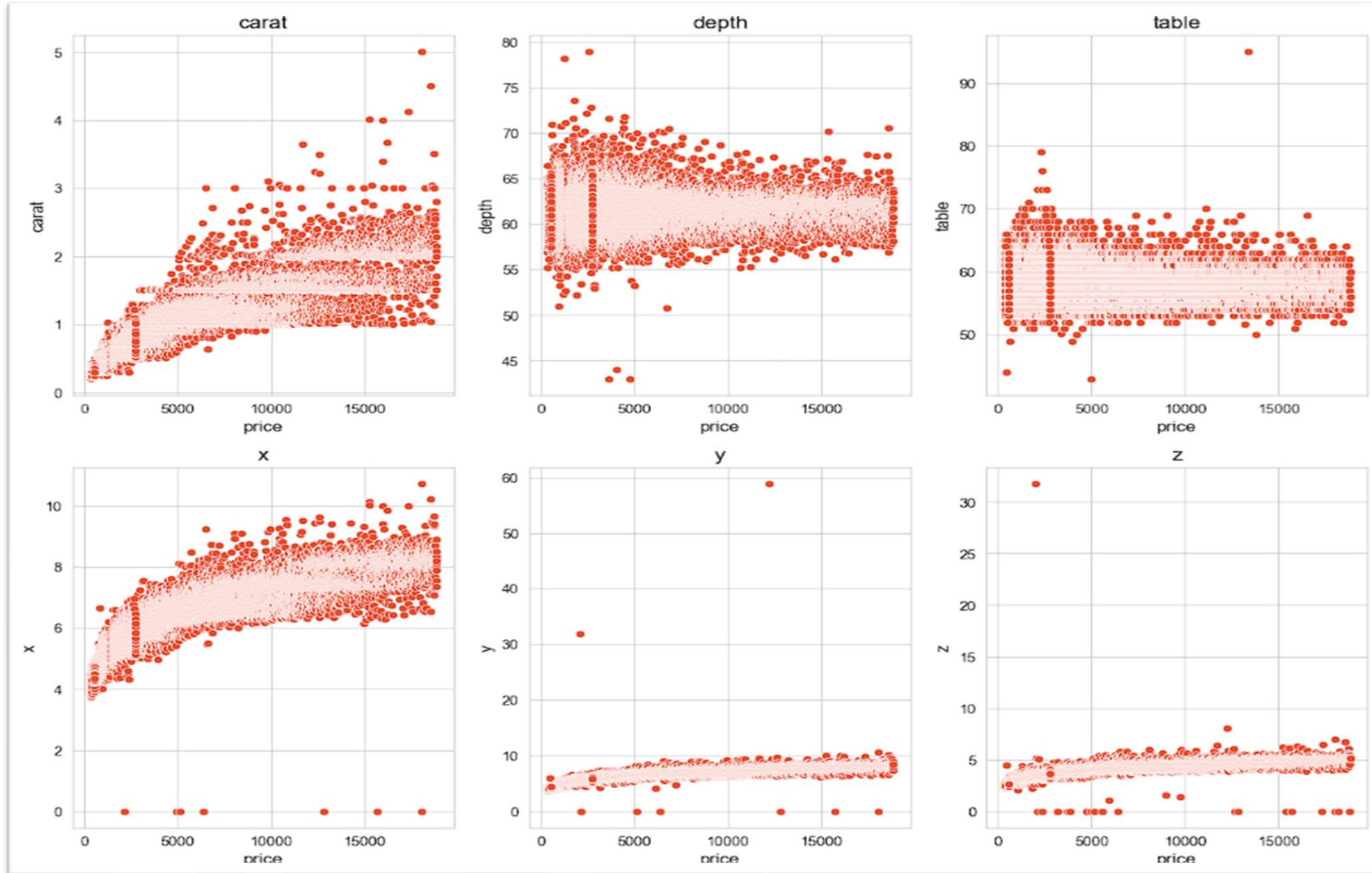
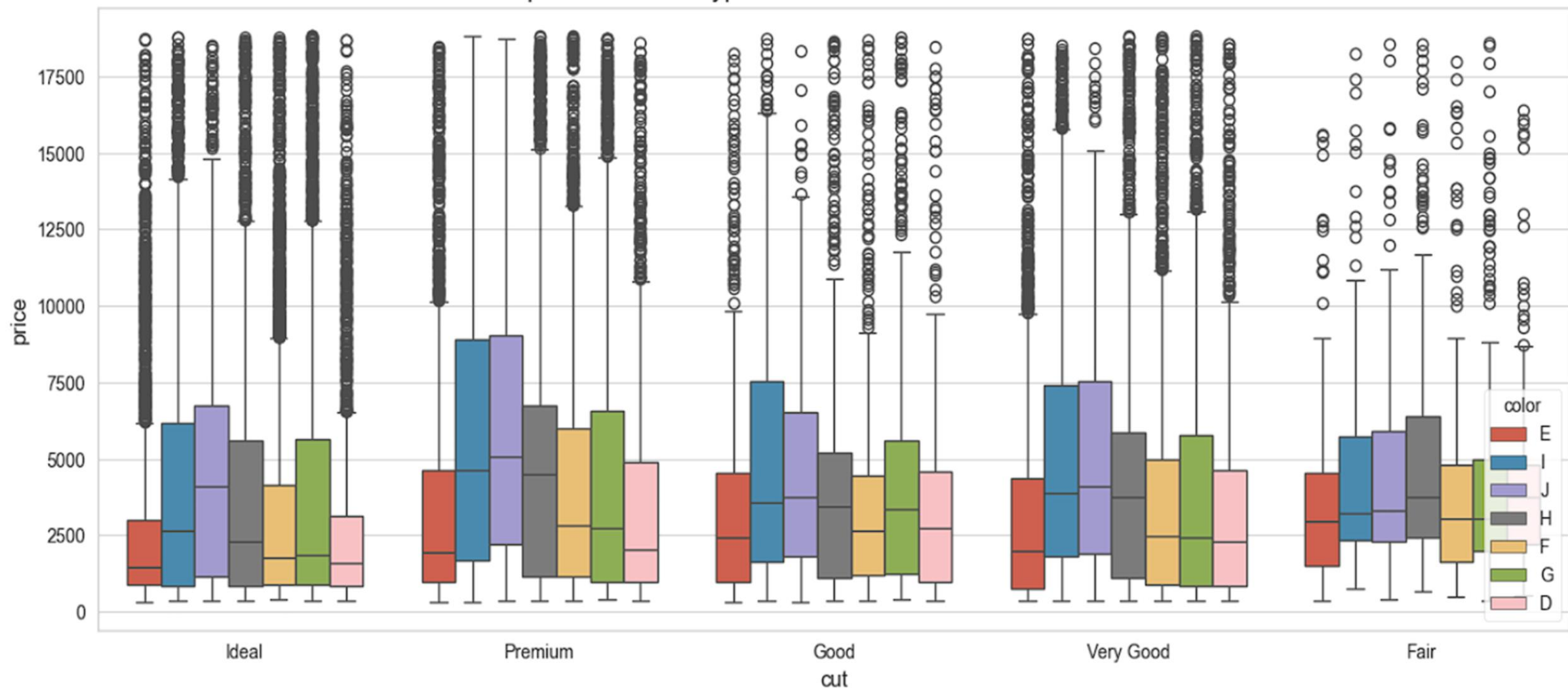8) Feature Importance Plot

Distribution Of Numerical Features

# Relationship of all the feature with Target variable

prices of all the types of diamonds based on their color
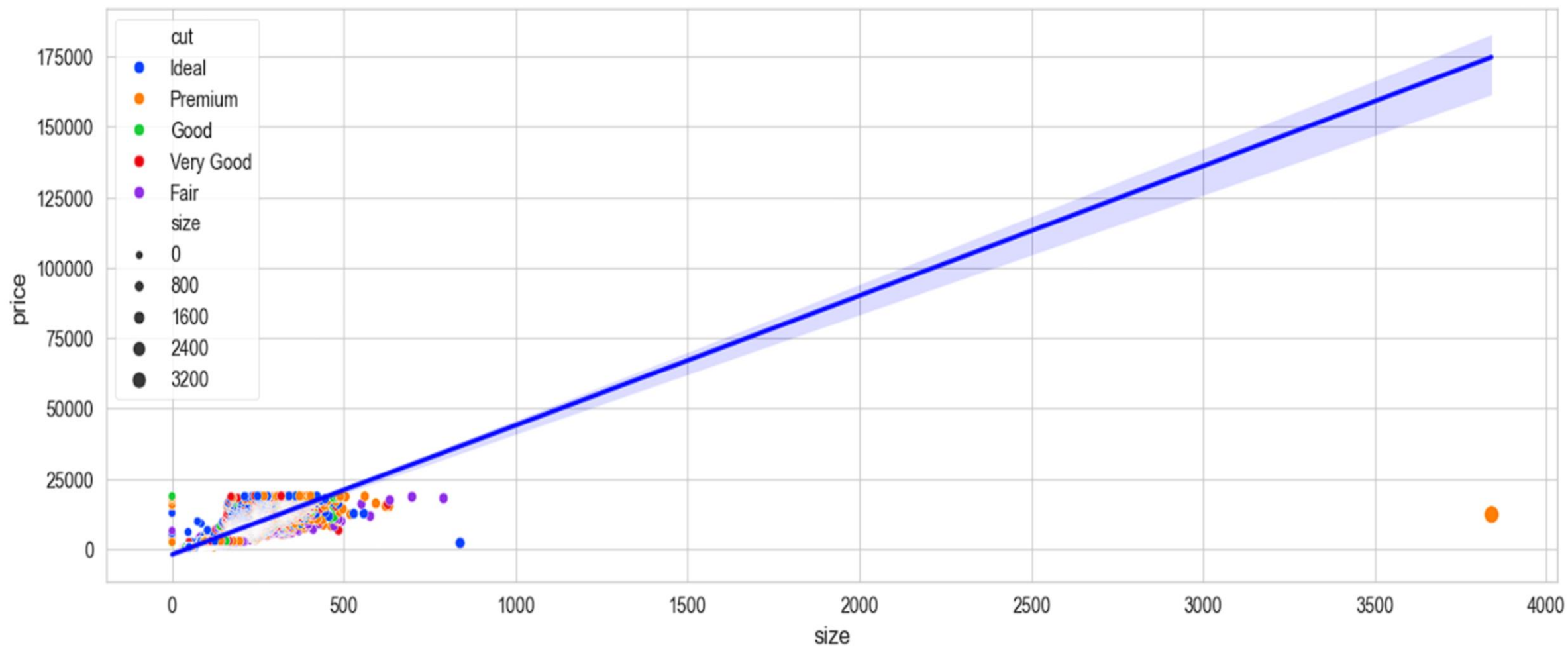
prices of all the types of diamonds based on their clarity

Scatter Plot of Carat vs Price with Depth Size and Cut Color

Scatter Plot of Size vs Price with Size and Cut Color

# Heatmap Of Correlation Matrix

Feature Importance

# Model Development

- **Tree Based Models** :-
1. Decision Tree
2. Random Forest
3. Ada Boost
4. XGB
- **Scaling-Based Models**
1. Linear Regression
2. KNN
3. SVM

# Challenges

1. **Outliers:-** Managing outliers to ensure they do not overly bias the model while retaining their valid variability.

2. **Categorical Features:-** Encoding qualitative attributes like cut, color, and clarity effectively for both tree-based and scaling-based models.

3. **Model Tuning:-** Balancing computation time during hyperparameter tuning for a large dataset.

# Model Performance Comparison

| | Model | Train R2 | Test R2 | Train MSE | Test MSE |
|---|---|---|---|---|---|
| 0 | Decision Tree | 0.886651 | 0.877203 | 1.814021e+06 | 1.842908e+06 |
| 1 | Random Forest | 0.897212 | 0.879838 | 1.645005e+06 | 1.803351e+06 |
| 2 | XGBoost | 0.887510 | 0.879653 | 1.800265e+06 | 1.806139e+06 |
| 3 | AdaBoost | 0.869742 | 0.866318 | 2.084621e+06 | 2.006263e+06 |
| 4 | Linear Regression | 0.852629 | 0.851238 | 2.358506e+06 | 2.232572e+06 |
| 5 | SVR | 0.855425 | 0.851829 | 2.313748e+06 | 2.223712e+06 |
| 6 | KNN | 0.898196 | 0.870080 | 1.629253e+06 | 1.949793e+06 |

# Actual VS Predicted Diamond Price



Actual vs. Predicted Diamond Prices

# Conclusion & Practical Application

## ▶ Conclusion

➢ **Impact of Models**:- Accurate diamond price predictions help jewelers, buyers, and insurers determine fair prices.

➢ **Best Model**: XG boost, due to its high R2 and low MSE on both training and test sets.

➢ **Good Alternatives**: Random Forest , Decision Tree and KNN also show strong performance.

➢ **Baseline Models**: Linear Regression and SVR provide good baselines but are outperformed by ensemble methods.

## ▶ Practical Application

1. **E-Commerce Platforms :-** Automating pricing recommendations for diamonds listed on online marketplaces.

2. **Inventory Management :-** Helping jewelers value their inventory based on consistent pricing.

3. **Fraud Detection :-** Identifying discrepancies in diamond pricing for authentication and fraud prevention.

4. **Customer Decision Support :-** Assisting customers in understanding and justifying diamond prices during purchases.

# Deployment Details

▶ **Description:** The diamond price prediction model has been deployed as a web application using the Streamlit framework, enabling users to interact with the model seamlessly.

▶ **Platform**: Streamlit cloud platform.

▶ **Deployment Process :**

1. Model trained and saved using Pickle.

2. Interactive web app built with Streamlit for a user-friendly interface.

3. Hosted on Streamlit's free cloud service.

▶ **Accessibility:**

1. Fully functional and accessible from any device via the link.

2. **Link :** https://diamond-price-predictions.streamlit.app/