

How many properties are listed for sale, and in which cities?

```
In [5]: 1 properties_for_sale = df[df['purpose'] == 'For Sale']
        2
        3 total_properties_for_sale = len(properties_for_sale)
        4 print("\nProperties listed for sale:")
        5 print(f"Total properties for sale: {total_properties_for_sale}")
        6
        7 properties_by_city = properties_for_sale.groupby('city').size().reset_index()
        8 print("\nProperties for sale in each city:")
        9 print(properties_by_city)
       10
```

Properties listed for sale:

Total properties for sale: 70947

Properties for sale in each city:

	city	properties_count
0	Faisalabad	1611
1	Islamabad	8794
2	Karachi	27210
3	Lahore	26221
4	Rawalpindi	7111

Location Analysis:

Which locations have the highest and lowest average property prices?

```
In [6]: 1 average_price_by_location = df.groupby('location')['price'].mean().sort_values()
        2 highest_avg_price_location = average_price_by_location.idxmax()
        3 lowest_avg_price_location = average_price_by_location.idxmin()
        4
        5 print(f"Highest average price location: {highest_avg_price_location} (Average Price: Rs{average_price_by_location[highest_avg_price_location]:.2f})")
        6 print(f"Lowest average price location: {lowest_avg_price_location} (Average Price: Rs{average_price_by_location[lowest_avg_price_location]:.2f})")
        7
        8
        9
```

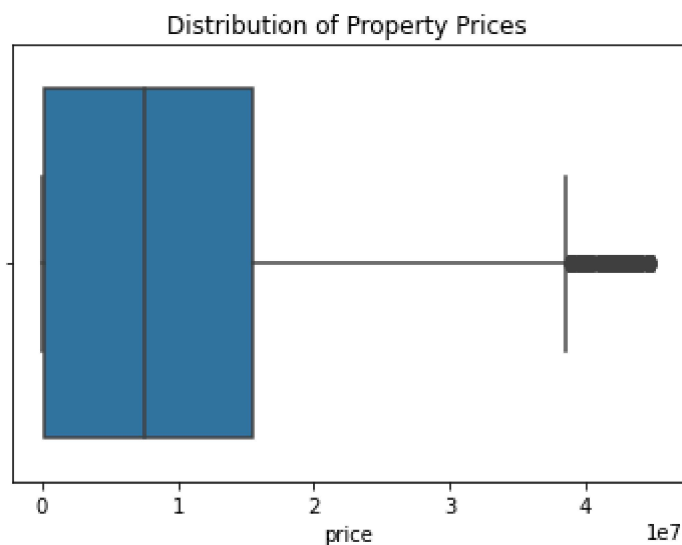
Highest average price location: TECH Society (Average Price: Rs42500000.00)

Lowest average price location: Beaumont Road (Average Price: Rs16000.00)

Price Analysis

Are there outliers or high-value properties in the dataset?

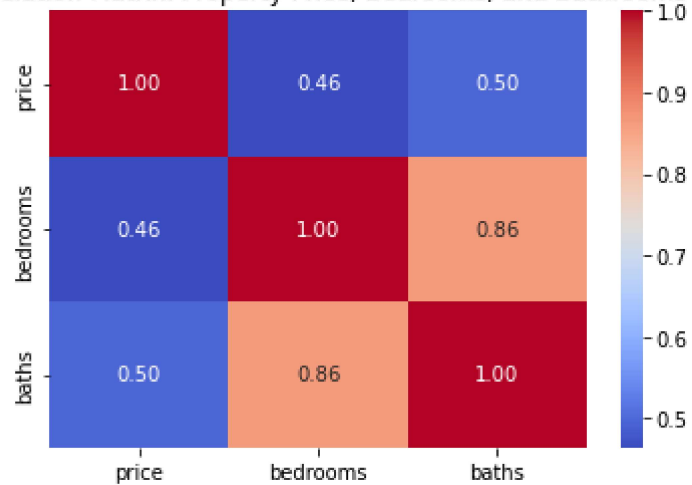
```
In [10]: 1 # outliers
2 sns.boxplot(x='price', data=df)
3 plt.title("Distribution of Property Prices")
4 plt.show()
5
6 correlation_matrix = df[['price', 'bedrooms', 'baths']].corr()
7
8 print("\nCorrelation Matrix:")
9 print(correlation_matrix)
10
11
12 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
13 plt.title("Correlation Matrix: Property Price, Bedrooms, and Bathrooms")
14 plt.show()
15
```



Correlation Matrix:

	price	bedrooms	baths
price	1.000000	0.464393	0.496222
bedrooms	0.464393	1.000000	0.863885
baths	0.496222	0.863885	1.000000

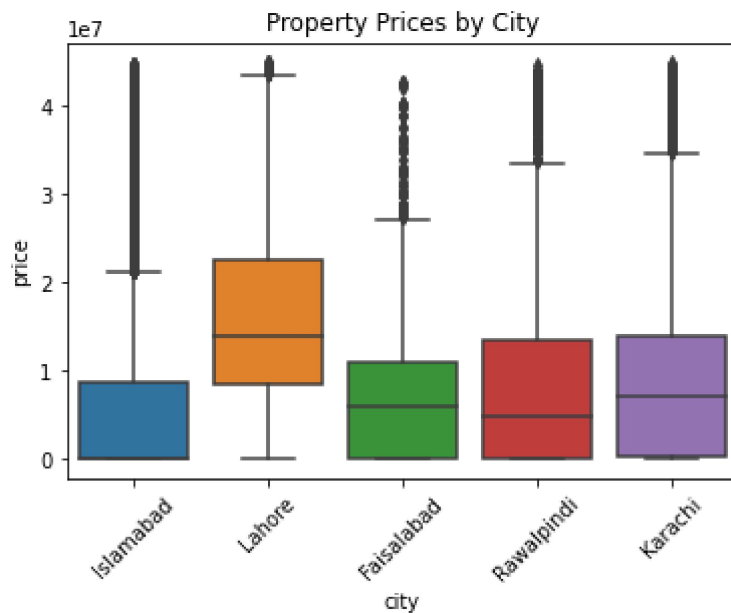
Correlation Matrix: Property Price, Bedrooms, and Bathrooms



City Comparison

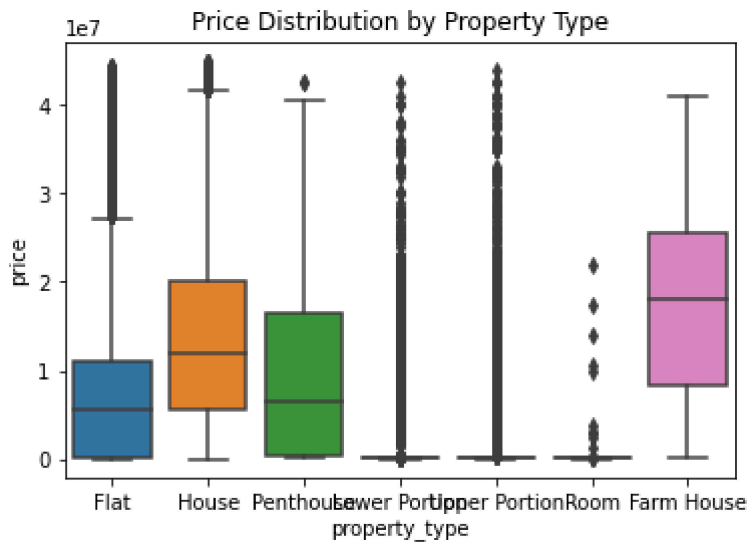
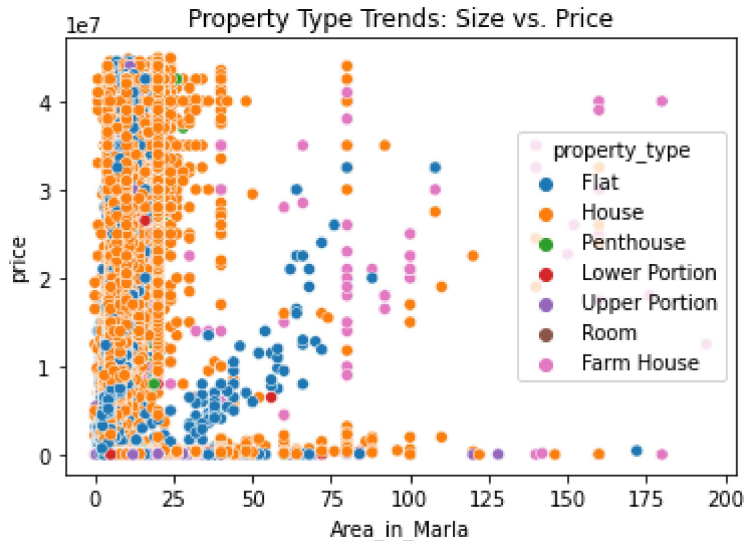
Question: How do property prices differ between cities?

```
In [11]: 1 sns.boxplot(x='city', y='price', data=df)
2 plt.title("Property Prices by City")
3 plt.xticks(rotation=45)
4 plt.show()
```



Are there trends or patterns specific to flats, houses, or other property types?

```
In [16]: 1 sns.scatterplot(x='Area_in_Marla', y='price', hue='property_type', data=df)
2          plt.title("Property Type Trends: Size vs. Price")
3          plt.show()
4
5          sns.boxplot(x='property_type', y='price', data=df)
6          plt.title("Price Distribution by Property Type")
7          plt.show()
8
```



Predictive Modeling:

```
In [21]: 1 # One-hot encoding
2 df = pd.get_dummies(df, columns=['location', 'property_type', 'city', 'pur
3
4
5 X = df.drop(columns=['price'])
6 y = df['price']
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
9
10
11 linear_reg_model = LinearRegression()
12 linear_reg_model.fit(X_train, y_train)
13
14
15 # Predictions
16 linear_reg_predictions = linear_reg_model.predict(X_test)
17
18
19 import pandas as pd
20 from sklearn.model_selection import train_test_split
21 from sklearn.linear_model import LinearRegression
22 from sklearn.metrics import r2_score, mean_squared_error
23
24 df = pd.read_csv('house_prices.csv')
25
26 # Select features and target variable
27 X = df[['location', 'bedrooms', 'Area_in_Marla']]
28 y = df['price']
29
30
31 X = pd.get_dummies(X, columns=['location'], drop_first=True)
32
33 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
34
35 model = LinearRegression()
36 model.fit(X_train, y_train)
37
38 predictions = model.predict(X_test)
39
40 # Evaluation
41 r2 = r2_score(y_test, predictions)
42 rmse = mean_squared_error(y_test, predictions, squared=False) # RMSE
43 print(f"R-squared (R2): {r2:.4f}")
44 print(f"Root Mean Squared Error (RMSE): {rmse:.2f}")
45
```

R-squared (R2): -922385505381862784.0000

Root Mean Squared Error (RMSE): 10275059879753562.00

```
In [22]: 1 # Random Forest Regression model
2 random_forest_model = RandomForestRegressor()
3 random_forest_model.fit(X_train, y_train)
4
```

Out[22]: RandomForestRegressor()

```
In [23]: 1 random_forest_predictions = random_forest_model.predict(X_test)
```

```
In [24]: 1 # Evaluate performance
2 def evaluate_model(predictions, model_name):
3     r2 = r2_score(y_test, predictions)
4     rmse = sqrt(mean_squared_error(y_test, predictions))
5     print(f"\nPerformance of {model_name} Model:")
6     print(f"R-squared (R2): {r2:.4f}")
7     print(f"Root Mean Squared Error (RMSE): {rmse:.2f}")
8
```

```
In [25]: 1 evaluate_model(random_forest_predictions, random_forest_model)
```

Performance of RandomForestRegressor() Model:
R-squared (R2): 0.5474
Root Mean Squared Error (RMSE): 7197282.44