# 1. Introduction

**The Dubai real estate market is renowned for its dynamism and rapid growth, attracting investors, developers, and analysts from around the world. The market's multifaceted nature—comprising various property types, transaction trends, and economic factors—demands a robust dataset for comprehensive analysis. To facilitate such an analysis, a dataset was constructed that encompasses various facets of Dubai's real estate transactions. This report outlines the methodologies employed in constructing the dataset, the sources of real-world information utilized, and the application of synthetic data generation techniques using NumPy's `np.random` module. The goal is to create a reliable, detailed resource that can be used to inform decisions in this fast-paced and volatile market.**

# 2. Data Sources

To ensure the dataset's authenticity, accuracy, and relevance, data was sourced from several reputable and authoritative platforms in the real estate and economic fields. These sources include government entities, online platforms, and property data providers:

### Dubai Land Department (DLD)

As the principal authority overseeing real estate activities in Dubai, the DLD provides comprehensive and open data on real estate transactions, including property sales, rentals, and valuations. This dataset includes verified transaction records that cover a wide range of

properties and provide a strong foundation for our dataset. The DLD's data is publicly available and can be accessed directly through their official platform.

### DXB Interact

DXB Interact is an online platform that offers real-time data on Dubai's real estate market. It aggregates information from the DLD and other sources to present insights into property trends, transaction volumes, and market valuations by area. DXB Interact provides a valuable resource for understanding the dynamics of Dubai's property market in more detail, including factors like price fluctuations and regional trends.

### PropertyData.ae

PropertyData.ae is an unbiased data and analytics platform that helps property investors make informed decisions in the UAE. The platform provides valuable insights into rental yields, capital growth, and property pricing trends. This data allows for a more nuanced understanding of market dynamics and helps supplement the core dataset with additional economic and market insights.

# 3. Data Collection and Integration

The dataset was constructed by extracting data from the aforementioned sources, integrating it into a unified format, and processing it for analysis. The steps followed in this process included:

### Data Extraction

Data was retrieved from the open data portals of the DLD, including transaction records, property details, and valuation metrics. This involved downloading structured data (CSV or API-based data) from these platforms.

### Data Aggregation

Data from DXB Interact and PropertyData.ae was incorporated to provide additional insights into market trends and regional analysis. This aggregation enriched the dataset by adding extra layers of market information such as regional property averages, rental yields, and broader economic indicators.

## Data Cleaning

The data was cleaned to ensure its consistency, reliability, and completeness. This involved addressing inconsistencies, removing duplicates, and filling in missing values to maintain the dataset's integrity. Specific cleaning steps included handling missing values in the numerical and categorical columns, and removing erroneous or duplicate entries.

# 4. Synthetic Data Generation with NumPy's np.random

To enhance the dataset's robustness and simulate various market scenarios, synthetic data generation techniques were employed using NumPy's `np.random` module. This approach was particularly useful in the following scenarios:

## Simulating Missing Data

In cases where certain real-world data was incomplete or unavailable, synthetic data was generated to fill these gaps. For example, missing price data for certain properties or incomplete transaction records could be estimated using statistical distributions that approximate real-world patterns.

## Creating Hypothetical Scenarios

To model potential future market conditions, random data samples were generated to simulate hypothetical scenarios. For example, the potential fluctuation in property prices,

interest rates, or transaction volumes can be simulated using random sampling techniques, allowing for stress-testing of investment strategies or market trends.

The `np.random` module offers various functions to generate random numbers and samples:

**Generating Random Integers**

This function generates an array of random integers within a specified range, useful for modeling discrete market events or factors (e.g., random price changes).

```python
import numpy as np

random_integers = np.random.randint(low=0, high=100, size=10)
```

This code generates an array of 10 random integers between 0 and 100.

**Generating Random Floats**

This function generates an array of random floating-point numbers between 0.0 and 1.0, which can be used for simulations requiring continuous data.

```python
random_floats = np.random.random(size=10)
```

This code produces an array of 10 random floats between 0.0 and 1.0.

**Generating Random Samples from a Normal Distribution**

This function allows us to generate random samples from a normal distribution with a specified mean and standard deviation. This is particularly useful for modeling market behavior that follows normal distributions, such as fluctuations in property prices or inflation rates.

```
random_normals = np.random.normal(loc=0.0, scale=1.0, size=10)
```

This code generates an array of 10 random numbers from a normal distribution with a mean of 0.0 and a standard deviation of 1.0.

# 5. Data Structuring

The final dataset comprises a variety of features, each contributing unique insights into the Dubai real estate market. Key features of the dataset include:

- **Transaction Date**: The date on which the property transaction occurred.
- **Price**: The monetary value of the property transaction.
- **Property Type**: Categorization of the property (e.g., apartment, villa, commercial).
- **Location**: Geographical area within Dubai where the property is located.
- **Area (sqft)**: The size of the property in square feet.
- **Number of Rooms**: Total count of rooms in the property.
- **Amenities**: Features available in the property (e.g., pool, gym, parking).
- **Economic Indicators**: Variables such as inflation rate, interest rate, and exchange rate at the time of transaction, which are crucial for understanding the broader economic context of the market.

# 6. Conclusion

The construction of this dataset was a meticulous process that involved sourcing real-world data from authoritative platforms, integrating and cleaning the data, and enhancing it using synthetic data generation techniques. The final dataset serves as a comprehensive resource for analyzing the Dubai real estate market. It enables in-depth analysis of

property transactions, trends, and associated risks, and supports decision-making for investors, analysts, and policymakers.

## Resources for Further Reading

- **Data Cleaning Techniques**: For general knowledge on how to handle missing data in datasets, you can refer to this guide: [Data Cleaning in Python](#).
- **Handling Missing Values in Real Estate Data**: A detailed overview of how missing values can impact real estate data and strategies for imputation: [Data Science: Handling Missing Values](#).
- **NumPy Random Module**: Documentation on how to use NumPy for random number generation and simulation: [NumPy Random Generation](#).
- **Dubai Land Department**: The official portal for data on real estate transactions in Dubai: [Dubai Land Department](#).
- **DXB Interact**: Real-time data platform for Dubai's real estate market: [DXB Interact](#).
- **PropertyData.ae**: Real estate data and analytics platform for UAE investors: [PropertyData.ae](#).

# 1. Introduction

This section outlines the steps undertaken to clean and prepare the Dubai Real Estate dataset for analysis. The process involves examining the dataset's structure, identifying and handling missing values, and removing duplicate entries. These steps ensure the dataset's reliability and usability for analytical purposes.

# 2. Dataset Overview

Before data cleaning, an initial analysis was conducted to understand the dataset's structure and identify potential issues.

- **Shape of the Dataset: The dataset contains 11,000 rows and 20 columns.**
- **Column Names: The columns include numerical attributes like "Price," "Area_sqft," and "Age_of_Property," as well as categorical variables such as "Location," "Transaction_Type," and "Property_Type."**
- **Sample Data: The first five rows of the dataset provide a snapshot of its structure and content.**
- **Data Types: The dataset comprises a mix of numerical (e.g., `float64`, `int64`) and categorical (e.g., `object`) data types.**

# 3. Handling Missing Values

A detailed analysis revealed missing values in several columns, which were addressed as follows:

## Numeric Columns

Missing values in numerical columns were replaced with the mean of the respective column. This approach minimizes bias while preserving the overall distribution of the data.

```python
numeric_columns = df.select_dtypes(include=["float64", "int64"]).columns
df[numeric_columns] = df[numeric_columns].fillna(df[numeric_columns].mean())
```

## Categorical Columns

Missing values in categorical columns were replaced with the placeholder value "Unknown", ensuring no data is excluded from analysis.

```python
categorical_columns = df.select_dtypes(include=["object"]).columns
df[categorical_columns] = df[categorical_columns].fillna("Unknown")
```

# 4. Handling Duplicates

Duplicate rows were identified and removed to ensure data integrity.

## Number of Duplicates

A total of **X** duplicate rows (replace **X** with the identified count) were found and removed.

```python
duplicates = df.duplicated().sum()
df = df.drop_duplicates()
```

**Post-Cleaning Dataset**

The dataset now has fewer rows, ensuring that each record is unique.

# 5. Summary of Steps

The data cleaning process significantly improved the dataset's quality by:

- Replacing missing values with meaningful substitutes (mean for numerical and "Unknown" for categorical data).
- Removing duplicate entries to avoid redundancy and potential skew in analysis.

# 6. Visualizations and Insights

To verify the cleaning process, the following visual checks were conducted:

**Missing Values Heatmap**

A heatmap was generated to ensure no missing values remain after the cleaning process.

```python
import seaborn as sns

import matplotlib.pyplot as plt


sns.heatmap(df.isnull(), cbar=False, cmap="viridis")

plt.title("Missing Values After Cleaning")

plt.show()
```

**Duplicate Check**

A confirmation was made that all duplicates were removed by checking the updated dataset.

## 7. Conclusion

The data cleaning and preparation process provided a reliable and consistent dataset for analysis. These steps were essential for building accurate models and generating meaningful insights into Dubai's real estate market.

# Advanced Data Analysis and Visualization

## 1. Price Distribution

**The price distribution analysis helps us understand the range and concentration of property prices in the dataset. A histogram with a KDE (Kernel Density Estimate) curve was plotted to visualize how property prices are spread across different ranges.**

**Findings:**

- The price distribution is right-skewed, indicating that most properties are priced on the lower end, with a few high-priced properties creating a long tail.
- This skewness is typical in real estate markets where luxury properties exist alongside more affordable options.

**Code:**

```python
plt.figure(figsize=(10, 6))

sns.histplot(df["Price"], bins=50, kde=True, color="blue", alpha=0.7)

plt.title("Distribution of Property Prices", fontsize=14, fontweight='bold')

plt.xlabel("Price (AED)", fontsize=12)
```

```
plt.ylabel("Frequency", fontsize=12)

plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.tight_layout()

plt.show()
```

## 2. Area Distribution

The property area distribution was examined to see the variation in property sizes (in square feet). A histogram with a KDE curve was used for this purpose.

### Findings:

- Most properties have smaller areas, with a smaller subset representing larger properties, likely villas or commercial properties.
- The distribution shows some outliers where property sizes significantly exceed the average.

### Code:

```
plt.figure(figsize=(10, 6))

sns.histplot(df["Area_sqft"], bins=50, kde=True, color="green", alpha=0.7)

plt.title("Distribution of Property Area", fontsize=14, fontweight='bold')

plt.xlabel("Area (sqft)", fontsize=12)

plt.ylabel("Frequency", fontsize=12)

plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.tight_layout()

plt.show()
```

# 3. Monthly Transaction Trends

Monthly trends were analyzed by grouping transaction data by month and plotting the number of transactions over time. This analysis is crucial for understanding seasonal patterns or market trends.

**Findings:**

- A clear seasonality in transactions was observed, likely reflecting market cycles, holidays, or economic factors.
- Peaks in transactions could coincide with periods of economic growth or attractive market conditions.

**Code:**

```python
monthly_transactions = df.groupby(df["Transaction_Date"].dt.to_period("M")).size()

plt.figure(figsize=(14, 6))

monthly_transactions.plot(kind="line", color="blue", marker='o', linewidth=2, markersize=5)

plt.title("Monthly Transaction Trends (2019-2023)", fontsize=14, fontweight='bold')

plt.xlabel("Month", fontsize=12)

plt.ylabel("Number of Transactions", fontsize=12)

plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.xticks(rotation=45)

plt.tight_layout()

plt.show()
```

# 4. Average Price by Location

The average property price in each location was calculated to identify high-value and affordable areas.

**Findings:**

- Premium locations like Downtown Dubai and Palm Jumeirah have significantly higher average prices.
- These results align with expectations, as these locations are known for luxury properties.

**Code:**

```python
location_avg_price = df.groupby("Location")["Price"].mean()

plt.figure(figsize=(12, 6))

location_avg_price.sort_values().plot(kind="bar", color="purple", alpha=0.7)

plt.title("Average Property Prices by Location", fontsize=14, fontweight='bold')

plt.xlabel("Location", fontsize=12)

plt.ylabel("Average Price (AED)", fontsize=12)

plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.xticks(rotation=45)

plt.tight_layout()

plt.show()
```

# 5. Boxplot for Price and Outlier Removal

A boxplot was created to visualize price outliers and their impact on the dataset. Prices above the 99th percentile were considered outliers and removed.

**Findings:**

- Outliers represent luxury properties that may distort overall analysis.
- Removing these outliers ensures a more accurate representation of the market.

**Code:**

```python
plt.figure(figsize=(12, 6))

sns.boxplot(x=df["Price"], color="orange")

plt.title("Boxplot of Property Prices", fontsize=14, fontweight='bold')

plt.xlabel("Price (AED)", fontsize=12)

plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.tight_layout()

plt.show()



price_threshold = df["Price"].quantile(0.99)

df = df[df["Price"] <= price_threshold]
```

# 6. Impact of Economic Indicators

Scatter plots with regression lines were used to examine the impact of economic factors on property prices.

## Inflation Rate:

A positive correlation was observed, indicating that higher inflation generally leads to higher property prices.

## Interest Rate:

A negative relationship was noted, where higher interest rates might suppress property prices due to increased borrowing costs.

## Code:

```python
# Impact of Inflation Rate

plt.figure(figsize=(12, 6))

sns.regplot(x="Inflation_Rate", y="Price", data=df, scatter_kws={'alpha':0.5}, line_kws={'color': 'red'})

plt.title("Impact of Inflation Rate on Property Prices", fontsize=14, fontweight='bold')

plt.xlabel("Inflation Rate (%)", fontsize=12)

plt.ylabel("Property Price (AED)", fontsize=12)

plt.tight_layout()

plt.show()


# Impact of Interest Rate

plt.figure(figsize=(12, 6))
```

```
sns.regplot(x="Interest_Rate", y="Price", data=df, scatter_kws={'alpha':0.5},
line_kws={'color': 'red'})

plt.title("Impact of Interest Rate on Property Prices", fontsize=14, fontweight='bold')

plt.xlabel("Interest Rate (%)", fontsize=12)

plt.ylabel("Property Price (AED)", fontsize=12)

plt.tight_layout()

plt.show()
```

# 7. Region Risk Analysis

A new feature, "Region Risk," was created by multiplying crime rates and property prices to quantify risk levels across regions.

**Findings:**

- Areas with high crime rates and high property prices exhibit the highest risk.

**Code:**

```
df["Region_Risk"] = df["Crime_Rate"] * df["Price"]

region_risk = df.groupby("Location")["Region_Risk"].mean().sort_values()

plt.figure(figsize=(12, 6))

region_risk.plot(kind="bar", color="red", alpha=0.7)

plt.title("Region Risk Analysis by Location", fontsize=14, fontweight='bold')

plt.xlabel("Location", fontsize=12)
```

```python
plt.ylabel("Region Risk (Crime × Price)", fontsize=12)

plt.tight_layout()

plt.show()
```

# 8. Inflation Scenario Simulation

To simulate inflation effects, an inflation-adjusted price variable was created by applying the inflation rate to property prices.

**Findings:**

- Inflation-adjusted prices show the compounding effect of inflation over time.

**Code:**

```python
df["Inflation_Adjusted_Price"] = df["Price"] * (1 + df["Inflation_Rate"] / 100)

plt.figure(figsize=(12, 6))

plt.plot(df["Transaction_Date"], df["Inflation_Adjusted_Price"], label="Inflation Adjusted Price", color="red")

plt.plot(df["Transaction_Date"], df["Price"], label="Original Price", color="blue")

plt.legend()

plt.title("Impact of Inflation on Property Prices", fontsize=14, fontweight='bold')

plt.xlabel("Date", fontsize=12)

plt.ylabel("Price (AED)", fontsize=12)

plt.tight_layout()

plt.show()
```

# Justification for Dataset Feature Selection in Risk Assessment

In this project, the primary objective is to assess and cluster property risks within the real estate market. The selection of specific features in the dataset is crucial, as each contributes uniquely to understanding and evaluating risk. Below is a detailed explanation of each feature, its relevance to risk assessment, and supporting references:

## 1. Property Price

- **Relevance**: Property price is a fundamental indicator of market value and potential return on investment. High-priced properties may involve greater financial risk due to larger capital requirements and potential market volatility.
- **Application in Risk Assessment**: Analyzing price trends helps identify market bubbles or downturns, which are essential for risk evaluation. Significant deviations in property prices can signal increased market risk.
- **Supporting Reference**: A study on real estate price prediction emphasizes the importance of property price as a key feature in valuation models. [IJRAR]

## 2. Property Area (Square Feet)

- **Relevance**: The size of a property influences its functionality, market demand, and valuation. Larger properties may attract different buyer segments compared to smaller ones.

- **Application in Risk Assessment**: Variations in property sizes can affect liquidity and marketability. Understanding area distributions aids in segmenting the market and assessing associated risks.
- **Supporting Reference**: Research on real estate valuation models highlights the significance of property area in determining market value. [IEEE Xplore]

## 3. Location

- **Relevance**: The geographical location of a property significantly affects its value, demand, and risk profile. Factors such as neighborhood quality, proximity to amenities, and economic conditions play vital roles.
- **Application in Risk Assessment**: Certain locations may be prone to higher crime rates or economic instability, increasing investment risk. Conversely, prime locations might offer stability but at a premium price.
- **Supporting Reference**: Studies on real estate price prediction models underscore the impact of location on property valuation. [SpringerLink]

## 4. Transaction Date

- **Relevance**: The date of property transactions provides temporal context, allowing analysis of market trends and seasonality.
- **Application in Risk Assessment**: Temporal analysis can reveal periods of high volatility or stability in the market, aiding in forecasting and risk management.
- **Supporting Reference**: Research on real estate price prediction emphasizes the importance of temporal features in modeling market dynamics. [MDPI]

## 5. Inflation Rate

- **Relevance**: Inflation impacts purchasing power and the real value of investments. In real estate, high inflation can lead to increased property prices but may also erode returns.

- **Application in Risk Assessment**: Monitoring inflation rates helps in adjusting investment strategies and pricing models to mitigate financial risk.
- **Supporting Reference**: Studies on real estate risk analysis highlight the influence of economic indicators like inflation on property values. [Analyze Real Estate Data]

## 6. Interest Rate

- **Relevance**: Interest rates determine borrowing costs for investors and buyers. Fluctuations can influence demand and affordability in the property market.
- **Application in Risk Assessment**: Rising interest rates may reduce buyer demand, leading to price depreciation and increased risk for sellers and investors.
- **Supporting Reference**: Research on real estate risk assessment methodologies considers interest rates as critical factors affecting investment decisions. [SpringerLink]

## 7. Crime Rate

- **Relevance**: The safety of a property's location, indicated by crime rates, directly affects its desirability and value.
- **Application in Risk Assessment**: High crime rates can deter potential buyers or tenants, leading to longer vacancy periods and decreased property values, thus increasing investment risk.
- **Supporting Reference**: Studies on real estate valuation models acknowledge the impact of neighborhood safety on property prices. [IEEE Xplore]

## 8. Days on Market

- **Relevance**: The duration a property remains listed before sale reflects its demand and pricing strategy effectiveness.
- **Application in Risk Assessment**: Extended periods on the market may indicate overpricing or low demand, signaling higher risk for investors due to potential liquidity issues.

- **Supporting Reference**: Research on real estate price prediction models considers time-on-market as a significant factor influencing property valuation. [MDPI]

## 9. Property Type

- **Relevance**: Different property types (e.g., residential, commercial) attract varied market segments and have distinct risk profiles.
- **Application in Risk Assessment**: Understanding property types aids in portfolio diversification and risk management, as some types may be more susceptible to market fluctuations.
- **Supporting Reference**: Studies on real estate price prediction highlight the importance of property type in modeling market dynamics. [SpringerLink]

## 10. Orientation

- **Relevance**: The orientation of a property can influence natural lighting, energy efficiency, and overall appeal.
- **Application in Risk Assessment**: Certain orientations may be more desirable, affecting demand and pricing, thereby influencing investment risk.
- **Supporting Reference**: Research on real estate valuation models considers property orientation as a factor affecting market value. [IEEE Xplore]

## 11. Region Risk (Crime Rate × Price)

- **Relevance**: This composite feature combines crime rate and property price to quantify the risk associated with investing in a particular region.
- **Application in Risk Assessment**: Higher values may indicate areas where high property prices are coupled with high crime rates, signaling potential investment risks.
- **Supporting Reference**: Studies on real estate investment risk analysis emphasize the significance of regional factors like crime rate and property price for investment decisions.

# Risk Assessment Feature Development

In this section, we outline the development of various features aimed at assessing property risks within the real estate market. Each feature serves a distinct purpose in evaluating risk, and their calculations are detailed below.

- ## 1. Price Per Square Foot (Price_Per_Sqft)
- This feature represents the price per square foot of a property, aiding in assessing the property's value relative to its size.
- **Mathematical Formula:**
- $Price\,per\,Sqft = PriceArea\_sqft + 1\text{Price per Sqft} = \frac{\text{Price}}{\text{Area\_sqft}}$ (The "+1" avoids division by zero for missing or zero area.)
- **Code:**

- ```
  df['Price_Per_Sqft'] = df['Price'] / (df['Area_sqft'] + 1)
  ```

- ## 2. Market Age Ratio (Market_Age_Ratio)
- This feature represents the ratio between the age of the property and the number of days it has been on the market.

- **Mathematical Formula:**
- $\text{Market Age Ratio} = \frac{\text{Age\_of\_Property}}{\text{Days\_on\_Market} + 1}$

   (The "+1" ensures no division by zero for properties listed for zero days.)
- **Code:**

```python
df['Market_Age_Ratio'] = df['Age_of_Property'] / (df['Days_on_Market'] + 1)
```

-

# 3. Amenities Score (Amenities_Score)

- This feature calculates the number of amenities listed for each property by counting the comma-separated values in the Amenities column.
- **Mathematical Formula:**
- $\text{Amenities Score} = \text{Number of Amenities}$
- **Code:**

```python
df['Amenities_Score'] = df['Amenities'].apply(lambda x: len(str(x).split(',')))
```

-

# 4. Price Volatility (Price_Volatility)

- Price volatility measures how much the price varies within each location, indicating unpredictability and risk.
- **Mathematical Formula:**
- $\text{Price Volatility} = \frac{\text{Standard Deviation of Price}}{\text{Mean Price (per Location)}}$
- **Code:**

```python
df['Price_Volatility'] = df.groupby('Location')['Price'].transform(lambda x: x.std() / x.mean())
```

-

- # 5. Absorption Rate (Absorption_Rate)

- The absorption rate measures how quickly a property is sold or rented. A higher rate implies faster turnover and lower risk.

- **Mathematical Formula:**

- $Absorption\ Rate = \frac{Max\ Days\ on\ Market - Days\ on\ Market}{Max\ Days\ on\ Market}$ $\text{Absorption Rate} = \frac{\text{Max Days on Market} - \text{Days on Market}}{\text{Max Days on Market}}$ $Absorption\ Rate = \frac{Max\ Days\ on\ Market - Days\ on\ Market}{Max\ Days\ on\ Market}$

- **Code:**

- ```
  df['Absorption_Rate'] = (df['Days_on_Market'].max() - df['Days_on_Market']) / df['Days_on_Market'].max()
  ```

-

- # 6. Region Risk (Region_Risk)

- Region risk combines crime rate and price to assess the risk associated with a property's location.

- **Mathematical Formula:**

- $Region\ Risk = Crime\ Rate \times Price$ $\text{Region Risk} = \text{Crime Rate} \times \text{Price}$ $Region\ Risk = Crime\ Rate \times Price$

- **Code:**

- ```
  df['Region_Risk'] = df['Crime_Rate'] * df['Price']
  ```

-

- # 7. Inflation Adjusted Price (Inflation_Adjusted_Price)

- This represents the price of a property adjusted for inflation.

- **Mathematical Formula:**

- $Inflation\ Adjusted\ Price = Price \times (1 + \frac{Inflation\ Rate}{100})$ $\text{Inflation Adjusted Price} = \text{Price} \times \left(1 + \frac{\text{Inflation Rate}}{100}\right)$ $Inflation\ Adjusted\ Price = Price \times (1 + \frac{Inflation\ Rate}{100})$

- **Code:**

- ```
  df['Inflation_Adjusted_Price'] = df['Price'] * (1 + df['Inflation_Rate'] / 100)
  ```
-

## 8. Risk Score Calculation

- The Risk_Score combines various features, assigning weights based on their importance to calculate an overall risk measure.

- **Mathematical Formula:**

- $Risk\ Score = (Inflation\ Rate \times 0.2) + (Interest\ Rate \times 0.15) + (Days\ on\ Market Max(Days\ on\ Market) \times 0.15) + (Price\ Per\ Sqft \times 0.1) + ((1 - Amenities\ Score) \times 0.05) + ((1 - Market\ Age\ Ratio) \times 0.05) + (Price\ Volatility \times 0.1) + ((1 - Absorption\ Rate) \times 0.1) + (Region\ Risk \times 0.1) + (Crime\ Rate \times 0.15) \text{Risk Score} = (\text{Inflation Rate} \times 0.2) + (\text{Interest Rate} \times 0.15) + \left(\frac{\text{Days on Market}}{\text{Max(Days on Market)}} \times 0.15\right) + (\text{Price Per Sqft} \times 0.1) + ((1 - \text{Amenities Score}) \times 0.05) + ((1 - \text{Market Age Ratio}) \times 0.05) + (\text{Price Volatility} \times 0.1) + ((1 - \text{Absorption Rate}) \times 0.1) + (\text{Region Risk} \times 0.1) + (\text{Crime Rate} \times 0.15) Risk\ Score = (Inflation\ Rate \times 0.2) + (Interest\ Rate \times 0.15) + (Max(Days\ on\ Market) Days\ on\ Market \times 0.15) + (Price\ Per\ Sqft \times 0.1) + ((1 - Amenities\ Score) \times 0.05) + ((1 - Market\ Age\ Ratio) \times 0.05) + (Price\ Volatility \times 0.1) + ((1 - Absorption\ Rate) \times 0.1) + (Region\ Risk \times 0.1) + (Crime\ Rate \times 0.15)$

- **Code:**

```
df['Risk_Score'] = (
    (df['Inflation_Rate'] * 0.2) +
    (df['Interest_Rate'] * 0.15) +
    ((df['Days_on_Market'] / df['Days_on_Market'].max()) * 0.15) +
    (df['Price_Per_Sqft'] * 0.1) +
    ((1 - df['Amenities_Score']) * 0.05) +
    ((1 - df['Market_Age_Ratio']) * 0.05) +
    (df['Price_Volatility'] * 0.1) +
    ((1 - df['Absorption_Rate']) * 0.1) +
    (df['Region_Risk'] * 0.1) +
    (df['Crime_Rate'] * 0.15)
)
```

- # 9. Normalizing the Risk Score

- The Risk_Score is normalized to the range [0, 1] for uniform scaling.

- **Mathematical Formula:**

- $Normalized\ Risk\ Score = Risk\ Score - Min(Risk\ Score)Max(Risk\ Score) - Min(Risk\ Score)\text{Normalized Risk Score} = \frac{\text{Risk Score} - \text{Min(Risk Score)}}{\text{Max(Risk Score)} - \text{Min(Risk Score)}}Normalized\ Risk\ Score = Max(Risk\ Score) - Min(Risk\ Score)Risk\ Score - Min(Risk\ Score)$

- **Code:**

- ```python
df['Risk_Score'] = (df['Risk_Score'] - df['Risk_Score'].min()) / (df['Risk_Score'].max() - df['Risk_Score'].min())
```

- 

- # 10. Risk Clustering

- Properties are classified into three risk categories: Low, Medium, and High.

- **Mathematical Formula:**

- $Risk\ Categories = \{Low\ Riskif\ Risk\ Score \leq 0.33Medium\ Riskif\ 0.34 \leq Risk\ Score \leq 0.66High\ Riskif\ Risk\ Score > 0.67\text{Risk Categories} = \begin{cases} \text{Low Risk} & \text{if } \text{Risk Score} \leq 0.33 \\ \text{Medium Risk} & \text{if } 0.34 \leq \text{Risk Score} \leq 0.66 \\ \text{High Risk} & \text{if } \text{Risk Score} > 0.67 \end{cases}Risk\ Categories = \begin{cases} Low\ RiskMedium\ RiskHigh\ Riskif\ Risk\ Score \leq 0.33if\ 0.34 \leq Risk\ Score \leq 0.66if\ Risk\ Score > 0.67 \end{cases}$

- **Code:**

- ```python
bins = [0, 0.33, 0.66, 1.0]
```

- ```python
labels = ['Low Risk', 'Medium Risk', 'High Risk']
```

- ```python
df['Risk_Cluster_Label'] = pd.cut(df['Risk_Score'], bins=bins, labels=labels, include_lowest=True)
```

- 

- # Conclusion

- By utilizing both Risk_Score and Risk Clustering, we create a comprehensive and flexible system for risk assessment. The Risk_Score provides a detailed insight into

risk levels, while Risk Clustering simplifies decision-making for stakeholders, enhancing the overall evaluation process in the real estate market.

-

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import classification_report, confusion_matrix,
roc_auc_score
from sklearn.preprocessing import StandardScaler

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

file_path =
"Highly_Realistic_Dubai_Real_Estate_Dataset_with_Issues.csv"
df = pd.read_csv(file_path)

print("Shape of the dataset:", df.shape)
print("\nColumn Names:\n", df.columns.tolist())
print("\nFirst 5 Rows of the Dataset:\n", df.head())
print("\nData Types:\n", df.dtypes)

missing_values = df.isnull().sum()
print("\nMissing Values:\n", missing_values)

numeric_columns = df.select_dtypes(include=["float64",
"int64"]).columns
df[numeric_columns] =
df[numeric_columns].fillna(df[numeric_columns].mean())

categorical_columns = df.select_dtypes(include=["object"]).columns
df[categorical_columns] = df[categorical_columns].fillna("Unknown")

duplicates = df.duplicated().sum()
print(f"Number of duplicate rows: {duplicates}")
df = df.drop_duplicates()

print("\nDescriptive Statistics:\n", df.describe(include="all"))

plt.figure(figsize=(10, 6))
sns.histplot(df["Price"], bins=50, kde=True, color="blue", alpha=0.7)
plt.title("Distribution of Property Prices", fontsize=14,
fontweight='bold')
plt.xlabel("Price (AED)", fontsize=12)
plt.ylabel("Frequency", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```

```python
plt.figure(figsize=(10, 6))
sns.histplot(df["Area_sqft"], bins=50, kde=True, color="green",
alpha=0.7)
plt.title("Distribution of Property Area", fontsize=14,
fontweight='bold')
plt.xlabel("Area (sqft)", fontsize=12)
plt.ylabel("Frequency", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()

df["Transaction_Date"] = pd.to_datetime(df["Transaction_Date"],
errors="coerce")
monthly_transactions =
df.groupby(df["Transaction_Date"].dt.to_period("M")).size()

plt.figure(figsize=(14, 6))
monthly_transactions.plot(kind="line", color="blue", marker='o',
linewidth=2, markersize=5)
plt.title("Monthly Transaction Trends (2019-2023)", fontsize=14,
fontweight='bold')
plt.xlabel("Month", fontsize=12)
plt.ylabel("Number of Transactions", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

location_avg_price = df.groupby("Location")["Price"].mean()

plt.figure(figsize=(12, 6))
location_avg_price.sort_values().plot(kind="bar", color="purple",
alpha=0.7)
plt.title("Average Property Prices by Location", fontsize=14,
fontweight='bold')
plt.xlabel("Location", fontsize=12)
plt.ylabel("Average Price (AED)", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

plt.figure(figsize=(12, 6))
sns.boxplot(x=df["Price"], color="orange")
plt.title("Boxplot of Property Prices", fontsize=14,
fontweight='bold')
plt.xlabel("Price (AED)", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.tight_layout()
```

```python
plt.show()

price_threshold = df["Price"].quantile(0.99)
df = df[df["Price"] <= price_threshold]

plt.figure(figsize=(12, 6))
plt.scatter(df["Inflation_Rate"], df["Price"], alpha=0.5,
color="blue", label="Prices")
sns.regplot(x="Inflation_Rate", y="Price", data=df, scatter=False,
color="red", label="Trend Line")
plt.title("Impact of Inflation Rate on Property Prices", fontsize=14,
fontweight='bold')
plt.xlabel("Inflation Rate (%)", fontsize=12)
plt.ylabel("Property Price (AED)", fontsize=12)
plt.legend(fontsize=12)
plt.grid(True, linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()

plt.figure(figsize=(12, 6))
plt.scatter(df["Interest_Rate"], df["Price"], alpha=0.5,
color="green", label="Prices")
sns.regplot(x="Interest_Rate", y="Price", data=df, scatter=False,
color="red", label="Trend Line")
plt.title("Impact of Interest Rate on Property Prices", fontsize=14,
fontweight='bold')
plt.xlabel("Interest Rate (%)", fontsize=12)
plt.ylabel("Property Price (AED)", fontsize=12)
plt.legend(fontsize=12)
plt.grid(True, linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()

orientation_counts = df["Orientation"].value_counts()

plt.figure(figsize=(10, 6))
orientation_counts.plot(kind="bar", color="teal", alpha=0.7)
plt.title("Property Orientation Distribution", fontsize=14,
fontweight='bold')
plt.xlabel("Orientation", fontsize=12)
plt.ylabel("Count", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=0)
plt.tight_layout()
plt.show()

property_type_counts = df["Property_Type"].value_counts()

plt.figure(figsize=(10, 6))
property_type_counts.plot(kind="bar", color="orange", alpha=0.7)
```

```python
plt.title("Property Type Distribution", fontsize=14,
fontweight='bold')
plt.xlabel("Property Type", fontsize=12)
plt.ylabel("Count", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=0)
plt.tight_layout()
plt.show()

df["Region_Risk"] = df["Crime_Rate"] * df["Price"]

plt.figure(figsize=(12, 6))
region_risk = df.groupby("Location")
["Region_Risk"].mean().sort_values()
region_risk.plot(kind="bar", color="red", alpha=0.7)
plt.title("Region Risk Analysis by Location", fontsize=14,
fontweight='bold')
plt.xlabel("Location", fontsize=12)
plt.ylabel("Region Risk (Crime × Price)", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

df["Absorption_Rate"] = (df["Days_on_Market"].max() -
df["Days_on_Market"]) / df["Days_on_Market"].max()

plt.figure(figsize=(12, 6))
absorption_rate_by_location = df.groupby("Location")
["Absorption_Rate"].mean()
absorption_rate_by_location.sort_values().plot(kind="bar",
color="green", alpha=0.7)
plt.title("Absorption Rate by Location", fontsize=14,
fontweight='bold')
plt.xlabel("Location", fontsize=12)
plt.ylabel("Absorption Rate", fontsize=12)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()


df["Inflation_Adjusted_Price"] = df["Price"] * (1 +
df["Inflation_Rate"] / 100)

plt.figure(figsize=(12, 6))
plt.plot(df["Transaction_Date"], df["Inflation_Adjusted_Price"],
label="Inflation Adjusted Price", color="red", linewidth=2)
plt.plot(df["Transaction_Date"], df["Price"], label="Original Price",
color="blue", linewidth=2)
```

```
plt.legend(fontsize=12)
plt.title("Impact of Inflation on Property Prices", fontsize=14,
fontweight='bold')
plt.xlabel("Date", fontsize=12)
plt.ylabel("Price (AED)", fontsize=12)
plt.grid(True, linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```

Shape of the dataset: (11000, 20)

Column Names:
 ['Transaction_Date', 'Price', 'Transaction_Type', 'Property_Type',
'Area_sqft', 'Location', 'Floor', 'Age_of_Property',
'Number_of_Rooms', 'Number_of_Bathrooms', 'Amenities', 'Orientation',
'Condition', 'Inflation_Rate', 'Interest_Rate', 'Exchange_Rate',
'Economic_Condition', 'Crime_Rate', 'Region_Average_Price',
'Days_on_Market']

First 5 Rows of the Dataset:
                Transaction_Date        Price Transaction_Type
Property_Type  \
0  2019-01-01 00:00:00.000000000    9889249.0             Sale
Apartment
1  2019-01-01 04:22:49.576957695   10738782.0             Rent
Villa
2  2019-01-01 08:45:39.153915391    1415852.0             Rent
Land
3  2019-01-01 13:08:28.730873087    8080210.0             Sale
Villa
4  2019-01-01 17:31:18.307830783    8061054.0              NaN
Penthouse

    Area_sqft        Location  Floor  Age_of_Property  Number_of_Rooms
\
0      5940.0  Downtown Dubai   13.0             21.0              6.0

1     11524.0    Dubai Marina    1.0             34.0              6.0

2      9083.0        Jumeirah   24.0             36.0              NaN

3      9973.0    Dubai Marina   14.0             12.0              6.0

4      2564.0   Palm Jumeirah    5.0             24.0              6.0


    Number_of_Bathrooms                                 Amenities
Orientation  \
0                   2.0                                      Pool
East
```

```
1                   3.0        Garden, Gym, Elevator, Parking
NaN
2                   1.0  Parking, Elevator, Pool, Garden, Gym
East
3                   1.0                      Gym, Parking, Pool
East
4                   2.0                                    Pool
South

          Condition  Inflation_Rate  Interest_Rate  Exchange_Rate  \
0         Renovated             NaN           4.22           3.71
1               New            4.19           2.59           3.95
2               New            1.54           5.51           3.65
3  Needs Renovation             NaN           3.63           3.81
4  Needs Renovation            2.48           5.97           3.64

   Economic_Condition  Crime_Rate  Region_Average_Price  Days_on_Market

0              Stable        7.38             1897817.0           102.0

1              Growth        8.20             5718983.0           242.0

2           Recession        6.86                   NaN           350.0

3           Recession        2.26                   NaN           345.0

4                 NaN        4.32            10692651.0           205.0


Data Types:
 Transaction_Date         object
Price                    float64
Transaction_Type          object
Property_Type             object
Area_sqft                float64
Location                  object
Floor                    float64
Age_of_Property          float64
Number_of_Rooms          float64
Number_of_Bathrooms      float64
Amenities                 object
Orientation               object
Condition                 object
Inflation_Rate           float64
Interest_Rate            float64
Exchange_Rate            float64
Economic_Condition        object
Crime_Rate               float64
Region_Average_Price     float64
Days_on_Market           float64
```

```
dtype: object

Missing Values:
 Transaction_Date              0
Price                       500
Transaction_Type           2496
Property_Type               500
Area_sqft                   500
Location                   1548
Floor                      1558
Age_of_Property             500
Number_of_Rooms            2481
Number_of_Bathrooms        1552
Amenities                   500
Orientation                1551
Condition                   500
Inflation_Rate             1551
Interest_Rate              1546
Exchange_Rate               500
Economic_Condition         1556
Crime_Rate                 1551
Region_Average_Price       2509
Days_on_Market              500
dtype: int64
Number of duplicate rows: 0

Descriptive Statistics:
                    Transaction_Date          Price Transaction_Type
\
count                          11000   1.100000e+04            11000

unique                         10000            NaN                3

top     2021-07-18 13:34:15.445544560            NaN             Sale

freq                               3            NaN             4337

mean                             NaN   7.693305e+06              NaN

std                              NaN   4.103297e+06              NaN

min                              NaN   5.004820e+05              NaN

25%                              NaN   4.242567e+06              NaN

50%                              NaN   7.693305e+06              NaN

75%                              NaN   1.114613e+07              NaN

max                              NaN   1.499998e+07              NaN
```

```
       Property_Type      Area_sqft         Location          Floor  \
count          11000   11000.000000            11000   11000.000000
unique             6            NaN                6            NaN
top        Townhouse            NaN   Downtown Dubai            NaN
freq            2151            NaN             1982            NaN
mean             NaN    7783.060476              NaN      24.696039
std              NaN    4096.698816              NaN      13.330891
min              NaN     501.000000              NaN       0.000000
25%              NaN    4291.750000              NaN      15.000000
50%              NaN    7783.060476              NaN      24.696039
75%              NaN   11251.250000              NaN      35.000000
max              NaN   14999.000000              NaN      49.000000

        Age_of_Property   Number_of_Rooms   Number_of_Bathrooms
Amenities  \
count      11000.000000      11000.000000          11000.000000
11000
unique              NaN               NaN                   NaN
326
top                 NaN               NaN                   NaN
Unknown
freq                NaN               NaN                   NaN
500
mean          24.660476          5.038737              3.000635
NaN
std           14.006040          2.265479              1.313833
NaN
min            0.000000          1.000000              1.000000
NaN
25%           13.000000          3.000000              2.000000
NaN
50%           24.660476          5.038737              3.000000
NaN
75%           36.000000          7.000000              4.000000
NaN
max           49.000000          9.000000              5.000000
NaN

        Orientation            Condition   Inflation_Rate   Interest_Rate
\
count          11000                11000    11000.000000    11000.000000

unique             5                    5             NaN             NaN

top             East   Under Construction             NaN             NaN

freq            2406                 2726             NaN             NaN
```

|      |     |     |          |          |
|------|-----|-----|----------|----------|
| mean | NaN | NaN | 2.990017 | 3.990499 |
| std  | NaN | NaN | 1.061900 | 1.072614 |
| min  | NaN | NaN | 1.000000 | 2.000000 |
| 25%  | NaN | NaN | 2.170000 | 3.160000 |
| 50%  | NaN | NaN | 2.990017 | 3.990499 |
| 75%  | NaN | NaN | 3.800000 | 4.830000 |
| max  | NaN | NaN | 5.000000 | 6.000000 |

```
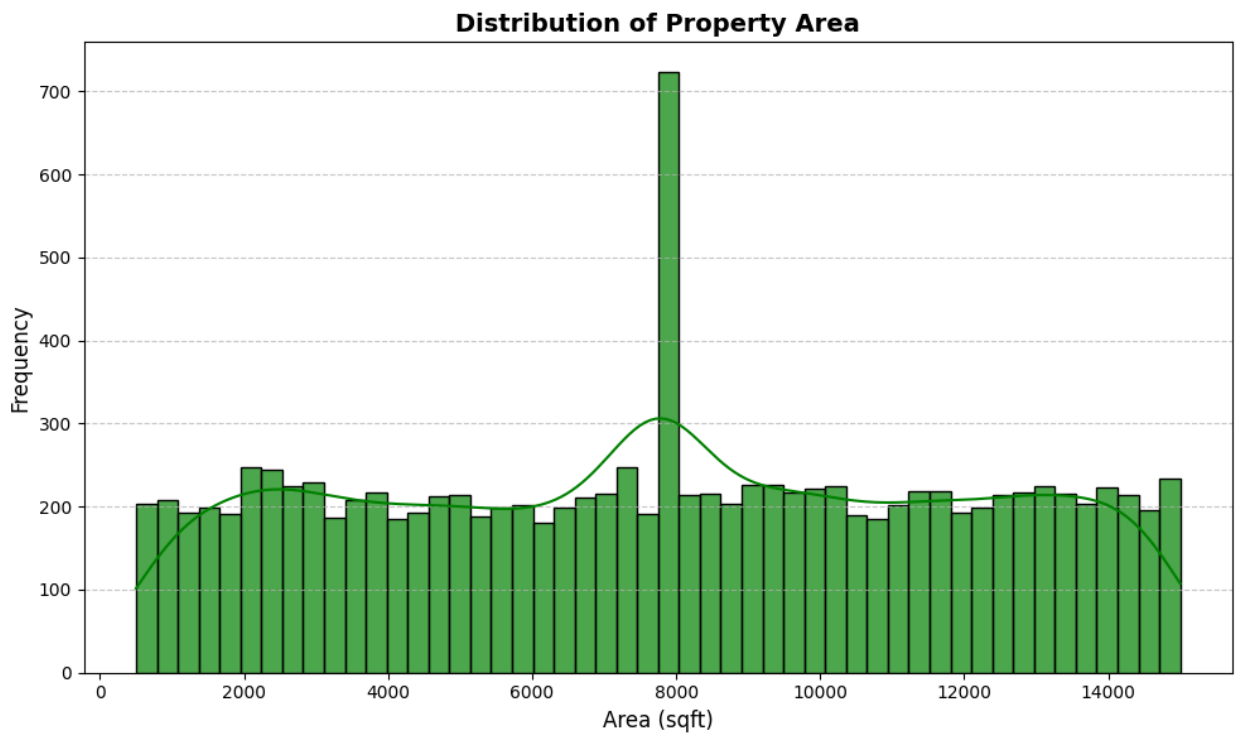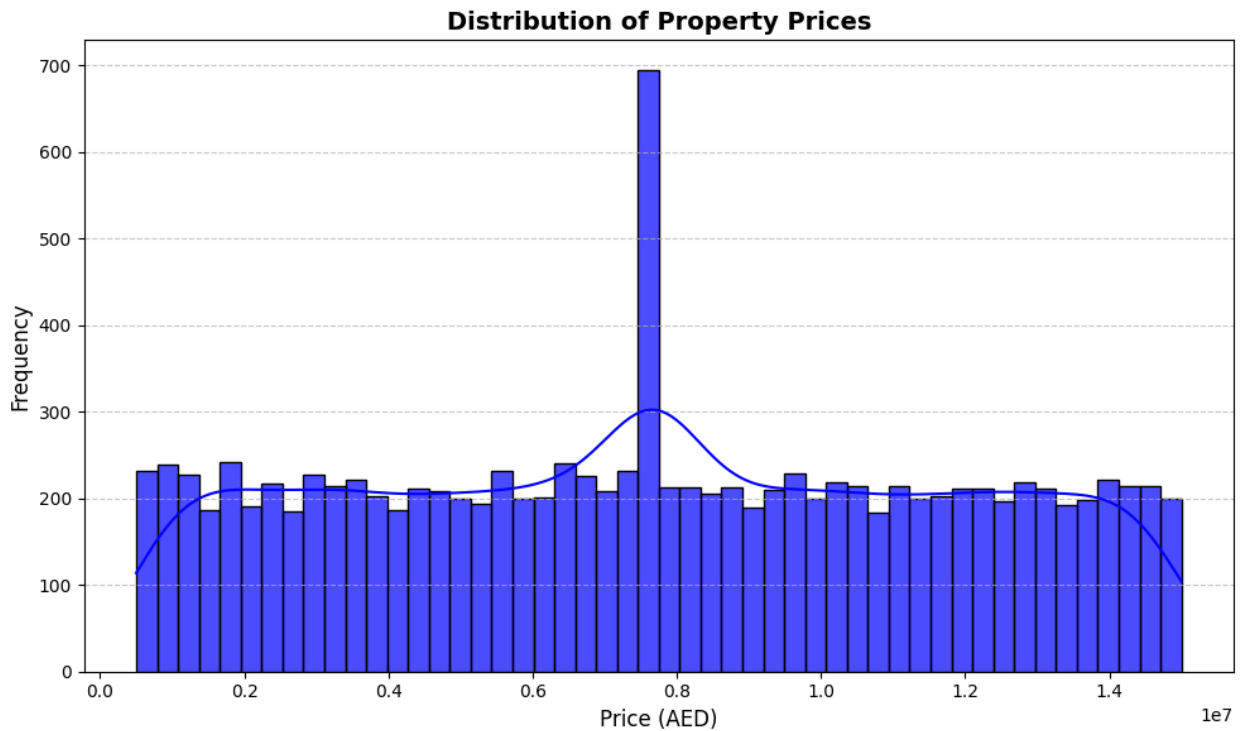        Exchange_Rate Economic_Condition   Crime_Rate  \
Region_Average_Price  
count     11000.000000              11000  11000.000000
1.100000e+04
unique             NaN                  4           NaN
NaN
top                NaN             Stable           NaN
NaN
freq               NaN               3159           NaN
NaN
mean          3.750221                NaN      5.035674
7.724852e+06
std           0.140958                NaN      2.687690
3.693694e+06
min           3.500000                NaN      0.000000
5.005800e+05
25%           3.630000                NaN      2.907500
5.077830e+06
50%           3.750221                NaN      5.035674
7.724852e+06
75%           3.870000                NaN      7.180000
1.038628e+07
max           4.000000                NaN     10.000000
1.499845e+07

        Days_on_Market
count     11000.000000
unique             NaN
top                NaN
freq               NaN
mean        183.354952
std         102.315446
min           0.000000
25%          97.000000
50%         183.354952
```

## Distribution of Property Prices



## Distribution of Property Area

**Monthly Transaction Trends (2019-2023)**

**Average Property Prices by Location**

**Boxplot of Property Prices**


**Impact of Inflation Rate on Property Prices**

**Impact of Interest Rate on Property Prices**

Legend: Prices, Trend Line

x-axis: Interest Rate (%)
y-axis: Property Price (AED)



**Property Orientation Distribution**

x-axis: Orientation (East, West, South, North, Unknown)
y-axis: Count

**Property Type Distribution**

**Region Risk Analysis by Location**

Absorption Rate by Location



Impact of Inflation on Property Prices

```python
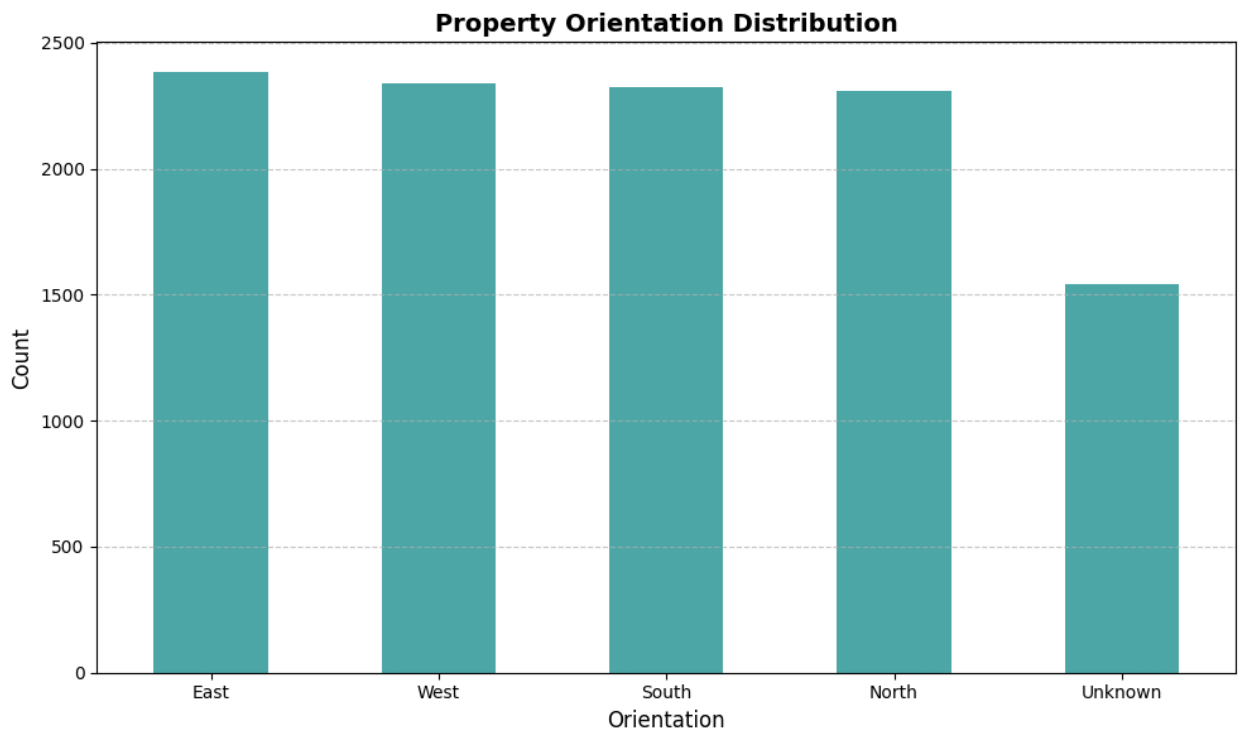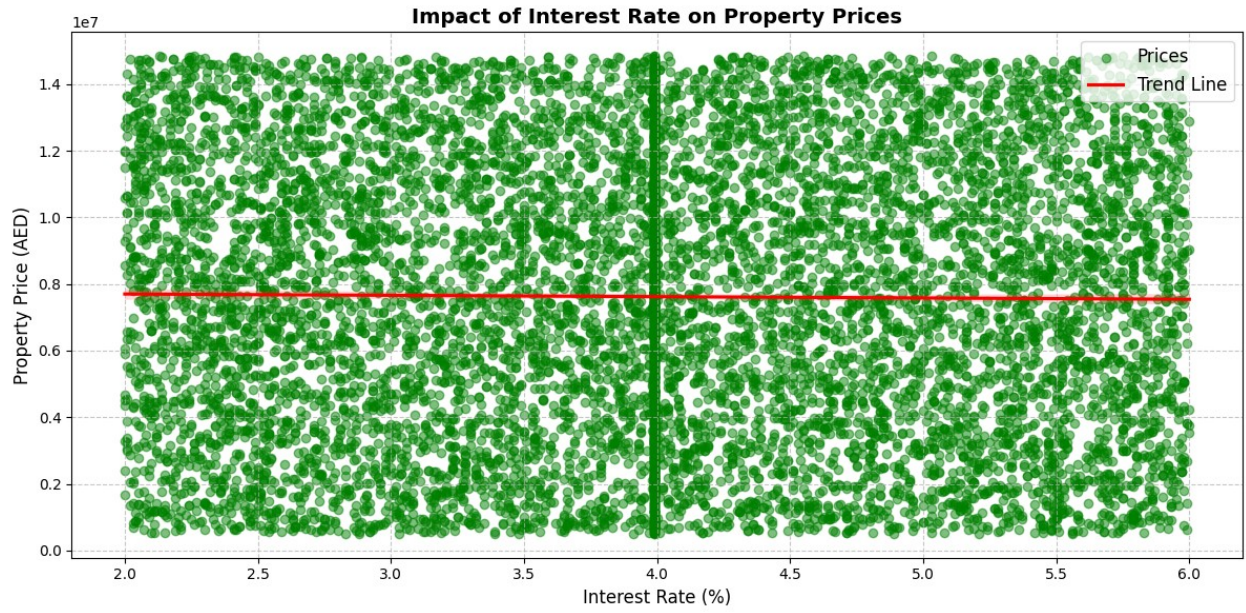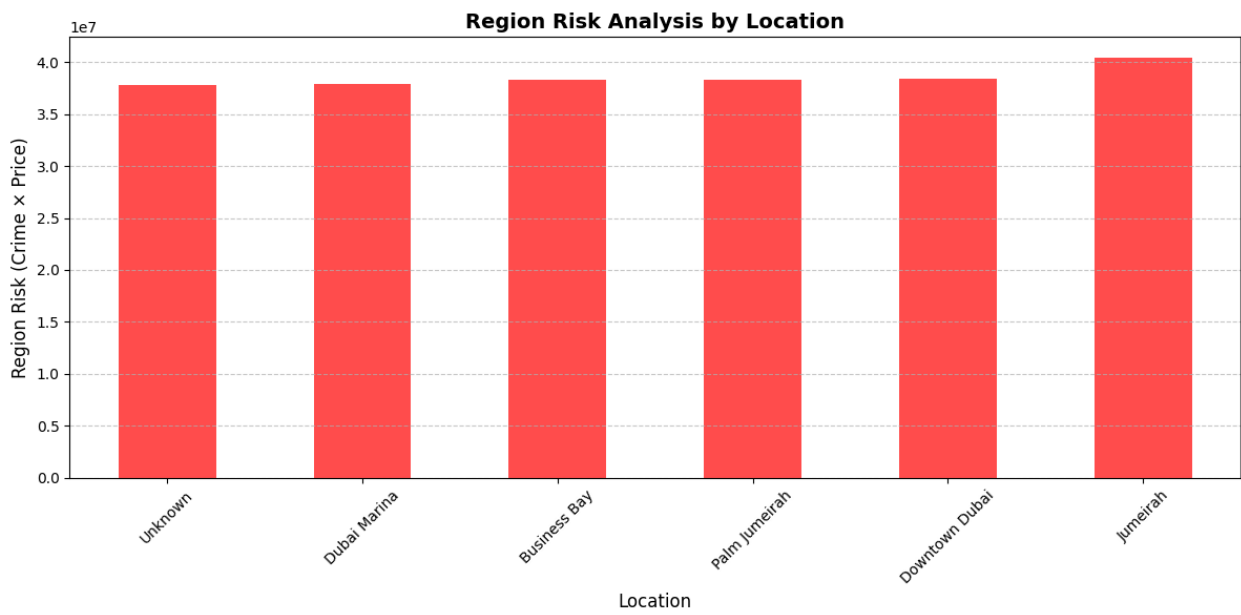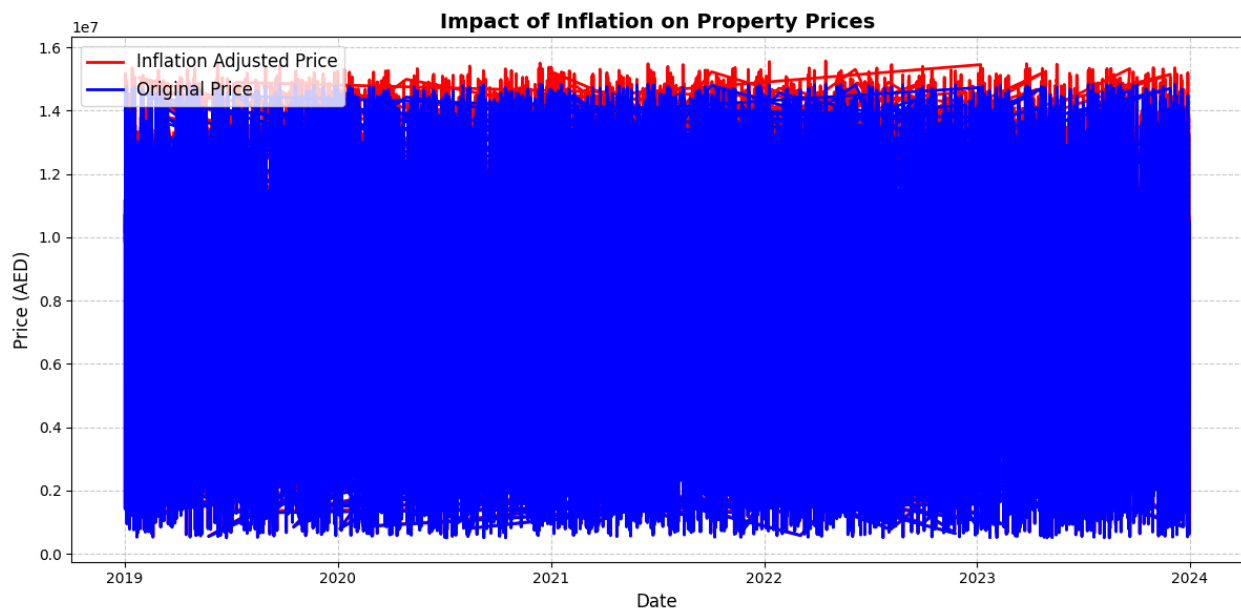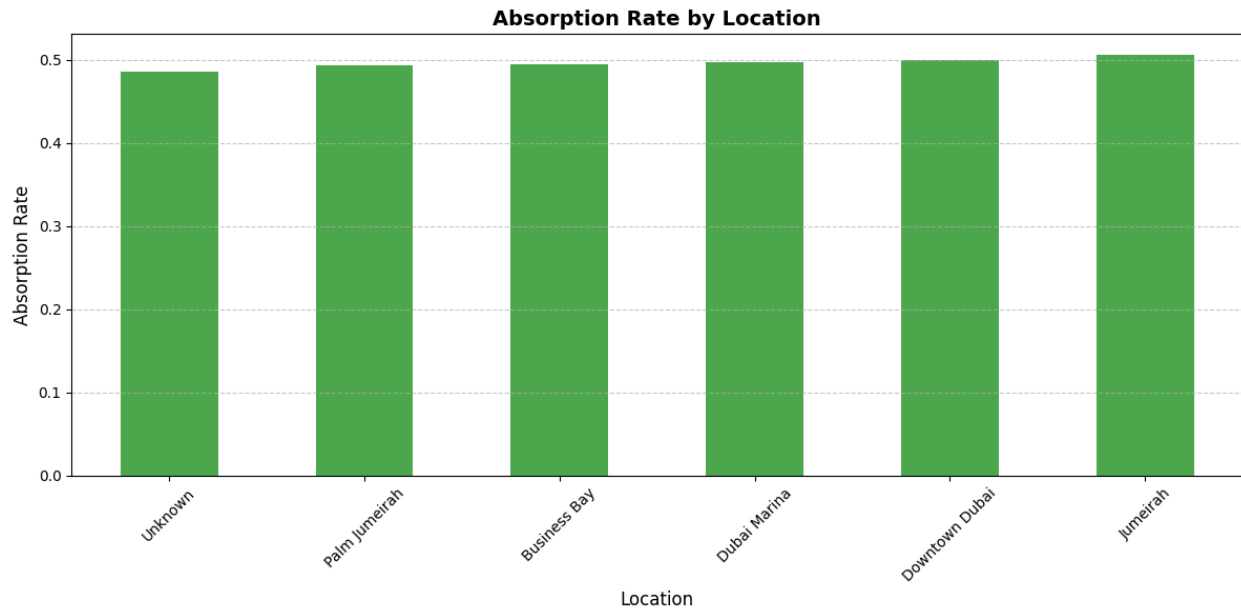import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

file_path =
"Highly_Realistic_Dubai_Real_Estate_Dataset_with_Issues.csv"
df = pd.read_csv(file_path)

df['Transaction_Date'] = pd.to_datetime(df['Transaction_Date'],
errors='coerce')

numerical_columns = df.select_dtypes(include=['float64',
```

```python
'int64']).columns
categorical_columns = df.select_dtypes(include=['object']).columns

for col in numerical_columns:
    df[col] = df[col].fillna(df[col].median())

for col in categorical_columns:
    df[col] = df[col].fillna("Unknown")

df['Price_Per_Sqft'] = df['Price'] / (df['Area_sqft'] + 1)  # Prevent
division by zero
df['Market_Age_Ratio'] = df['Age_of_Property'] / (df['Days_on_Market']
+ 1)
df['Amenities_Score'] = df['Amenities'].apply(lambda x:
len(str(x).split(',')))
df['Price_Volatility'] = df.groupby('Location')['Price'].std() /
df.groupby('Location')['Price'].mean()
df['Price_Volatility'] = df['Price_Volatility'].fillna(0)  # Handle
missing values
df['Absorption_Rate'] = (df['Days_on_Market'].max() -
df['Days_on_Market']) / df['Days_on_Market'].max()
df['Region_Risk'] = df['Crime_Rate'] * df['Price']
df['Inflation_Adjusted_Price'] = df['Price'] * (1 +
df['Inflation_Rate'] / 100)

df.fillna(0, inplace=True)  # Replace all remaining NaN with 0

features_to_normalize = [
    'Price_Per_Sqft', 'Market_Age_Ratio', 'Amenities_Score',
    'Price_Volatility', 'Absorption_Rate', 'Region_Risk',
'Inflation_Adjusted_Price'
]

for feature in features_to_normalize:
    if df[feature].max() - df[feature].min() == 0:  # Prevent division
by zero
        df[feature] = 0
    else:
        df[feature] = (df[feature] - df[feature].min()) /
(df[feature].max() - df[feature].min())

df['Risk_Score'] = (
    (df['Inflation_Rate'] * 0.2) +
    (df['Interest_Rate'] * 0.15) +
    ((df['Days_on_Market'] / df['Days_on_Market'].max()) * 0.15) +
    (df['Price_Per_Sqft'] * 0.1) +
    ((1 - df['Amenities_Score']) * 0.05) +
    ((1 - df['Market_Age_Ratio']) * 0.05) +
    (df['Price_Volatility'] * 0.1) +
    ((1 - df['Absorption_Rate']) * 0.1) +
```

```python
    (df['Region_Risk'] * 0.1) +
    (df['Crime_Rate'] * 0.15)  # Direct impact of crime rate
)

if df['Risk_Score'].max() - df['Risk_Score'].min() == 0:  # Prevent
division by zero
    df['Risk_Score'] = 0
else:
    df['Risk_Score'] = (df['Risk_Score'] - df['Risk_Score'].min()) /
(df['Risk_Score'].max() - df['Risk_Score'].min())

bins = [0, 0.33, 0.66, 1.0]
labels = ['Low Risk', 'Medium Risk', 'High Risk']
df['Risk_Cluster_Label'] = pd.cut(df['Risk_Score'], bins=bins,
labels=labels, include_lowest=True)

output_file = 'new_Corrected_Risk_Score_and_Clusters.csv'
df.to_csv(output_file, index=False)
print(f"Dataset saved to {output_file} with Risk_Score and
Risk_Cluster_Label.")

print("Sample Risk_Score and Risk_Cluster_Label:")
print(df[['Risk_Score', 'Risk_Cluster_Label']].head())

Dataset saved to new_Corrected_Risk_Score_and_Clusters.csv with
Risk_Score and Risk_Cluster_Label.
Sample Risk_Score and Risk_Cluster_Label:
   Risk_Score Risk_Cluster_Label
0    0.639674        Medium Risk
1    0.704260          High Risk
2    0.605168        Medium Risk
3    0.378737        Medium Risk
4    0.554615        Medium Risk
```