

# Housing



Bạn hãy hành động như một chuyên gia về phân tích dữ liệu. Bạn hãy suy nghĩ step-by-step và đưa ra cho tôi kế hoạch đầy đủ về các bước phân tích dữ liệu cho bài toán regression



Chắc chắn rồi! Dưới đây là các bước chi tiết để thực hiện phân tích dữ liệu cho bài toán hồi quy (regression analysis):

## Bước 1: Xác định mục tiêu và thu thập dữ liệu

- **Xác định mục tiêu:** Hiểu rõ vấn đề bạn muốn giải quyết bằng hồi quy, ví dụ như dự đoán giá nhà, doanh thu bán hàng, v.v.
- **Thu thập dữ liệu:** Thu thập dữ liệu từ các nguồn đáng tin cậy. Dữ liệu có thể đến từ cơ sở dữ liệu công ty, API, khảo sát, hoặc các nguồn dữ liệu công khai khác.

## Bước 2: Khám phá và làm sạch dữ liệu (Data Exploration and Cleaning)

- **Khám phá dữ liệu:** Sử dụng các công cụ phân tích khám phá dữ liệu (EDA) để hiểu về cấu trúc, mẫu hình, và mối quan hệ trong dữ liệu. Sử dụng các biểu đồ như biểu đồ phân phối, biểu đồ scatter, biểu đồ hộp (box plot), v.v.
- **Làm sạch dữ liệu:** Xử lý các giá trị thiếu, loại bỏ hoặc sửa các giá trị ngoại lệ (outliers), mã hóa các biến phân loại, chuẩn hóa hoặc tiêu chuẩn hóa dữ liệu nếu cần thiết.

## Bước 3: Lựa chọn tính năng (Feature Selection)

- **Lựa chọn tính năng:** Xác định các biến độc lập (predictor variables) có thể ảnh hưởng đến biến phụ thuộc (target variable). Sử dụng các phương pháp như tương quan Pearson, kiểm định chi-squared, hoặc các thuật toán giảm chiều (dimensionality reduction) như PCA để chọn ra các tính năng quan trọng.

## Bước 4: Chia dữ liệu thành tập huấn luyện và tập kiểm tra

- **Chia dữ liệu:** Chia dữ liệu thành hai tập: tập huấn luyện (training set) và tập kiểm tra (test set). Thông thường tỷ lệ chia là 70-30 hoặc 80-20.

## Bước 5: Xây dựng và huấn luyện mô hình (Model Building and Training)

- **Lựa chọn mô hình:** Lựa chọn các thuật toán hồi quy phù hợp như hồi quy tuyến tính (Linear Regression), hồi quy bội (Multiple Regression), hồi quy ridge, hồi quy lasso, v.v.
- **Huấn luyện mô hình:** Sử dụng tập huấn luyện để huấn luyện mô hình. Điều chỉnh các tham số (hyperparameters) nếu cần thiết.

## Bước 6: Đánh giá mô hình (Model Evaluation)

- **Đánh giá mô hình:** Sử dụng tập kiểm tra để đánh giá hiệu suất của mô hình bằng các chỉ số như R-squared, Mean Squared Error (MSE), Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), v.v.
- **Kiểm tra đa dạng hóa:** Sử dụng các phương pháp kiểm định chéo (cross-validation) như k-fold cross-validation để đảm bảo mô hình không bị overfitting.

## Bước 7: Tinh chỉnh mô hình (Model Tuning)

- **Tinh chỉnh mô hình:** Dựa trên các kết quả đánh giá, tinh chỉnh mô hình bằng cách thay đổi các tham số, chọn lại các tính năng, hoặc thử các mô hình khác.

## Bước 8: Triển khai mô hình (Model Deployment)

- **Triển khai mô hình:** Sau khi đã tinh chỉnh và có được mô hình tối ưu, triển khai mô hình vào môi trường sản xuất để dự đoán hoặc ra quyết định.

## Bước 9: Theo dõi và bảo trì (Monitoring and Maintenance)

- **Theo dõi:** Liên tục theo dõi hiệu suất của mô hình sau khi triển khai. Sử dụng các công cụ giám sát để phát hiện sự suy giảm hiệu suất (model drift).
- **Bảo trì:** Cập nhật mô hình khi có dữ liệu mới hoặc khi mô hình hiện tại không còn hiệu quả.

## Bước 10: Báo cáo và truyền đạt kết quả (Reporting and Communication)

- **Báo cáo kết quả:** Tổng hợp kết quả phân tích và mô hình vào một báo cáo chi tiết. Trình bày kết quả một cách dễ hiểu và trực quan cho các bên liên quan.
- **Truyền đạt kết quả:** Trình bày kết quả và mô hình cho các bên liên quan, giải thích các phát hiện quan trọng và các khuyến nghị.

Kế hoạch trên cung cấp một cái nhìn toàn diện về các bước cần thực hiện trong phân tích dữ liệu cho bài toán hồi quy. Mỗi bước có thể được tùy chỉnh dựa trên yêu cầu cụ thể của bài toán và dữ liệu bạn đang làm việc.



Tôi có một bộ dữ liệu dự đoán giá nhà tại Ấn Độ cơ bản. Bạn có thể giúp tôi đưa ra lộ trình phân tích dữ liệu tương tự trên không



### Memory updated

Chắc chắn rồi! Dưới đây là lộ trình chi tiết để phân tích dữ liệu dự đoán giá nhà tại Ấn Độ sử dụng các bước phân tích hồi quy:

## Bước 1: Xác định mục tiêu và thu thập dữ liệu

- **Xác định mục tiêu:** Mục tiêu của bạn là dự đoán giá nhà tại Ấn Độ dựa trên các yếu tố liên quan.
- **Thu thập dữ liệu:** Bạn đã có bộ dữ liệu cơ bản, hãy đảm bảo rằng dữ liệu bao gồm các biến như diện tích, số phòng, vị trí, năm xây dựng, và các yếu tố khác có thể ảnh hưởng đến giá nhà.

## Bước 2: Khám phá và làm sạch dữ liệu (Data Exploration and Cleaning)

- **Khám phá dữ liệu:**
  - Kiểm tra các thông tin tổng quát về dữ liệu (số lượng dòng, cột, các loại biến, v.v.).
  - Sử dụng các công cụ EDA để trực quan hóa dữ liệu, ví dụ như biểu đồ phân phối của giá nhà, biểu đồ scatter giữa giá nhà và diện tích, v.v.
- **Làm sạch dữ liệu:**
  - Xử lý các giá trị thiếu (missing values): Loại bỏ hoặc điền giá trị thiếu bằng các phương pháp thích hợp như trung bình, median, hoặc giá trị phổ biến nhất.
  - Xử lý các giá trị ngoại lệ (outliers): Sử dụng biểu đồ box plot để phát hiện và xử lý các giá trị ngoại lệ.
  - Mã hóa các biến phân loại (categorical variables): Sử dụng one-hot encoding hoặc label encoding.

## Bước 3: Lựa chọn tính năng (Feature Selection)

- **Lựa chọn tính năng:**
  - Sử dụng phân tích tương quan để xác định mối quan hệ giữa các biến độc lập và giá nhà.
  - Sử dụng các phương pháp như kiểm định chi-squared, hoặc các thuật toán giảm chiều (dimensionality reduction) như PCA để chọn ra các tính năng quan trọng.

## Bước 4: Chia dữ liệu thành tập huấn luyện và tập kiểm tra

- **Chia dữ liệu:**
  - Chia dữ liệu thành hai tập: tập huấn luyện (training set) và tập kiểm tra (test set) với tỷ lệ phổ biến là 80-20 hoặc 70-30.

## Bước 5: Xây dựng và huấn luyện mô hình (Model Building and Training)

- **Lựa chọn mô hình:**
  - Sử dụng các thuật toán hồi quy như hồi quy tuyến tính (Linear Regression), hồi quy bội (Multiple Regression), hồi quy ridge, hồi quy lasso, v.v.
- **Huấn luyện mô hình:**
  - Sử dụng tập huấn luyện để huấn luyện mô hình. Điều chỉnh các tham số (hyperparameters) nếu cần thiết.

## Bước 6: Đánh giá mô hình (Model Evaluation)

- **Đánh giá mô hình:**
  - Sử dụng tập kiểm tra để đánh giá hiệu suất của mô hình bằng các chỉ số như R-squared, Mean Squared Error (MSE), Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), v.v.
  - Sử dụng các phương pháp kiểm định chéo (cross-validation) như k-fold cross-validation để đảm bảo mô hình không bị overfitting.

## Bước 7: Tinh chỉnh mô hình (Model Tuning)

- **Tinh chỉnh mô hình:**
  - Dựa trên các kết quả đánh giá, tinh chỉnh mô hình bằng cách thay đổi các tham số, chọn lại các tính năng, hoặc thử các mô hình khác.

## Bước 8: Triển khai mô hình (Model Deployment)

- **Triển khai mô hình:**

- Sau khi đã tinh chỉnh và có được mô hình tối ưu, triển khai mô hình vào môi trường sản xuất để dự đoán giá nhà.

## Bước 9: Theo dõi và bảo trì (Monitoring and Maintenance)

- **Theo dõi:**

- Liên tục theo dõi hiệu suất của mô hình sau khi triển khai. Sử dụng các công cụ giám sát để phát hiện sự suy giảm hiệu suất (model drift).

- **Bảo trì:**

- Cập nhật mô hình khi có dữ liệu mới hoặc khi mô hình hiện tại không còn hiệu quả.

## Bước 10: Báo cáo và truyền đạt kết quả (Reporting and Communication)

- **Báo cáo kết quả:**

- Tổng hợp kết quả phân tích và mô hình vào một báo cáo chi tiết. Trình bày kết quả một cách dễ hiểu và trực quan cho các bên liên quan.

- **Truyền đạt kết quả:**

- Trình bày kết quả và mô hình cho các bên liên quan, giải thích các phát hiện quan trọng và các khuyến nghị.

Kế hoạch trên cung cấp một cái nhìn toàn diện về các bước cần thực hiện trong phân tích dữ liệu dự đoán giá nhà tại Ấn Độ. Bạn có thể tùy chỉnh từng bước dựa trên yêu cầu cụ thể của dự án và dữ liệu bạn đang làm việc.



### housing

Unable to display visualization

bắt đầu thôi nào đây là dữ liệu của tôi bạn có thể làm gì được với nó?



Dữ liệu của bạn bao gồm nhiều tính năng khác nhau liên quan đến các yếu tố của nhà và giá của chúng. Dưới đây là các bước tiếp theo mà chúng ta sẽ thực hiện để phân tích và dự đoán giá nhà:

## Bước 1: Khám phá và làm sạch dữ liệu

- **Khám phá dữ liệu:** Xem qua các thông tin tổng quát về dữ liệu, kiểm tra các giá trị thiếu, các giá trị ngoại lệ, và phân tích phân phối của các biến.
- **Làm sạch dữ liệu:** Xử lý các giá trị thiếu và các giá trị ngoại lệ, mã hóa các biến phân loại.

## Bước 2: Phân tích và lựa chọn tính năng

- **Phân tích tương quan:** Tìm hiểu mối quan hệ giữa các biến độc lập và giá nhà.
- **Lựa chọn tính năng:** Chọn ra các tính năng quan trọng nhất ảnh hưởng đến giá nhà.

## Bước 3: Chia dữ liệu thành tập huấn luyện và tập kiểm tra

- Chia dữ liệu thành hai tập: tập huấn luyện và tập kiểm tra.

## Bước 4: Xây dựng và huấn luyện mô hình

- Sử dụng các thuật toán hồi quy để xây dựng mô hình dự đoán giá nhà.

## Bước 5: Đánh giá và tinh chỉnh mô hình

- Đánh giá hiệu suất của mô hình và tinh chỉnh nếu cần thiết.

## Bước 6: Triển khai và theo dõi mô hình

Bây giờ, chúng ta sẽ bắt đầu với bước 1: Khám phá và làm sạch dữ liệu.

### Khám phá dữ liệu

Trước hết, chúng ta sẽ kiểm tra các thông tin tổng quát và xem xét các giá trị thiếu.

```
python
```

```
# Kiểm tra các thông tin tổng quát
housing_data.info()

# Kiểm tra các giá trị thiếu
missing_values = housing_data.isnull().sum()
missing_values[missing_values > 0]
```

Chúng ta hãy thực hiện kiểm tra này để hiểu rõ hơn về dữ liệu của bạn.

Dữ liệu của bạn có 26,370 dòng và 40 cột, không có giá trị thiếu nào. Dữ liệu bao gồm một biến phân loại là "Location" và 39 biến số học khác.

### Bước tiếp theo: Khám phá dữ liệu

Chúng ta sẽ xem qua phân phối của biến mục tiêu "Price" và các biến quan trọng khác để hiểu rõ hơn về dữ liệu.

### Phân tích phân phối của biến "Price"

Chúng ta sẽ vẽ biểu đồ phân phối của biến "Price".

```
python
```

```
import matplotlib.pyplot as plt
import seaborn as sns

# Biểu đồ phân phối của biến Price
plt.figure(figsize=(10, 6))
sns.histplot(housing_data['Price'], bins=30, kde=True)
plt.title('Distribution of House Prices')
plt.xlabel('Price')
plt.ylabel('Frequency')
plt.show()
```

## Phân tích mối quan hệ giữa các biến

Chúng ta sẽ vẽ biểu đồ scatter giữa "Price" và một vài biến quan trọng khác như "Area" và "No. of Bedrooms".

python

```
# Biểu đồ scatter giữa Price và Area
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['Area'], y=housing_data['Price'])
plt.title('Price vs Area')
plt.xlabel('Area')
plt.ylabel('Price')
plt.show()

# Biểu đồ scatter giữa Price và No. of Bedrooms
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['No. of Bedrooms'], y=housing_data['Price'])
plt.title('Price vs No. of Bedrooms')
plt.xlabel('No. of Bedrooms')
plt.ylabel('Price')
plt.show()
```

Hãy thực hiện các phân tích trên để khám phá dữ liệu của bạn.


### Distribution of House Prices

 Output image

### Price vs Area

 Output image

## Price vs No. of Bedrooms

 Output image

## Nhận xét từ các biểu đồ

### 1. Biểu đồ phân phối của giá nhà (Price):

- Biểu đồ phân phối cho thấy giá nhà có sự phân bố khá rộng với nhiều mức giá khác nhau. Có vẻ như có một số giá trị ngoại lệ (outliers) ở phần đuôi phải của phân phối.

### 2. Biểu đồ scatter giữa Giá nhà (Price) và Diện tích (Area):

- Biểu đồ này cho thấy một mối quan hệ dương giữa diện tích và giá nhà, tức là diện tích nhà càng lớn thì giá nhà có xu hướng cao hơn.

### 3. Biểu đồ scatter giữa Giá nhà (Price) và Số phòng ngủ (No. of Bedrooms):

- Biểu đồ này cho thấy một mối quan hệ dương giữa số phòng ngủ và giá nhà, tức là số phòng ngủ càng nhiều thì giá nhà có xu hướng cao hơn.

## Bước tiếp theo: Xử lý dữ liệu và lựa chọn tính năng

Chúng ta sẽ tiếp tục với các bước:

1. **Mã hóa biến phân loại (Location):** Chuyển đổi biến phân loại "Location" thành các biến số học thông qua one-hot encoding.
2. **Lựa chọn tính năng:** Dựa trên các phân tích trên, chọn ra các tính năng quan trọng nhất để đưa vào mô hình.

## Mã hóa biến phân loại

python

```
# Mã hóa biến phân loại Location
housing_data_encoded = pd.get_dummies(housing_data, columns=['Location'])

# Kiểm tra lại dữ liệu sau khi mã hóa
housing_data_encoded.head()
```

## Lựa chọn tính năng

Chúng ta sẽ chọn một vài tính năng quan trọng như "Area", "No. of Bedrooms" và các tiện nghi như "Gymnasium", "SwimmingPool", v.v. để huấn luyện mô hình.

python

```
# Chọn các tính năng quan trọng
features = ['Area', 'No. of Bedrooms', 'Gymnasium', 'SwimmingPool', 'Price']
housing_data_selected = housing_data_encoded[features]

# Kiểm tra lại dữ liệu sau khi chọn tính năng
housing_data_selected.head()
```

Hãy thực hiện các bước này để xử lý dữ liệu và chuẩn bị cho việc huấn luyện mô hình.



Dừng lại ở phần phân tích mối quan hệ giữa các biến một chút bạn hãy mở rộng phân tích với tất cả các biến trong dữ liệu với giá nhà



Để mở rộng phân tích mối quan hệ giữa tất cả các biến trong dữ liệu và giá nhà, chúng ta có thể sử dụng heatmap để hiển thị ma trận tương quan giữa các biến. Điều này sẽ giúp chúng ta nhận biết những biến nào có mối tương quan mạnh nhất với giá nhà.

## Phân tích tương quan giữa tất cả các biến và giá nhà

python

```
# Tính toán ma trận tương quan
correlation_matrix = housing_data.corr()

# Hiển thị heatmap của ma trận tương quan
plt.figure(figsize=(15, 15))
sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap='coolwarm', square=True,
            cbar_kws={"shrink": .8})
plt.title('Correlation Matrix of Housing Data')
plt.show()
```



## Phân tích chi tiết các biến có mối tương quan cao với giá nhà


Chúng ta sẽ tập trung vào các biến có hệ số tương quan cao nhất với giá nhà để hiểu rõ hơn về mối quan hệ giữa chúng.

```
python
```

```
# Lọc ra các biến có hệ số tương quan cao nhất với giá nhà
correlation_with_price = correlation_matrix['Price'].sort_values(ascending=False)
correlation_with_price
```

Chúng ta hãy thực hiện các phân tích trên để khám phá thêm về mối quan hệ giữa các biến và giá nhà.

### Correlation Matrix of Housing Data

 Output image

## Nhận xét từ ma trận tương quan

### 1. Biến có mối tương quan cao nhất với giá nhà:

- **Area:** Hệ số tương quan là 0.278, cho thấy diện tích có mối quan hệ dương tương đối với giá nhà.
- **No. of Bedrooms:** Hệ số tương quan là 0.157, cho thấy số phòng ngủ cũng có mối quan hệ dương với giá nhà.

### 2. Các tiện nghi và dịch vụ khác:

- Các biến như **Resale**, **ClubHouse**, **Children'splayarea**, và các tiện nghi khác có hệ số tương quan thấp (khoảng 0.04-0.05) với giá nhà.

## Kết luận

- **Diện tích (Area)** và **Số phòng ngủ (No. of Bedrooms)** là những biến quan trọng nhất cần xem xét trong mô hình dự đoán giá nhà.
- Các biến khác tuy có tương quan nhưng khá nhỏ, do đó chúng có thể không đóng vai trò quá lớn trong việc dự đoán giá nhà, nhưng chúng vẫn có thể được xem xét tùy vào độ phức tạp của mô

hình.

## Bước tiếp theo: Chuẩn bị dữ liệu và xây dựng mô hình

Chúng ta sẽ tiến hành các bước tiếp theo:

1. **Chuẩn hóa dữ liệu:** Mã hóa biến phân loại "Location" và chuẩn hóa các biến số học.
2. **Chia dữ liệu:** Chia dữ liệu thành tập huấn luyện và tập kiểm tra.
3. **Xây dựng mô hình:** Sử dụng các thuật toán hồi quy để xây dựng mô hình dự đoán giá nhà.

### Chuẩn hóa dữ liệu

python

```
# Mã hóa biến phân loại Location
housing_data_encoded = pd.get_dummies(housing_data, columns=['Location'])

# Chuẩn hóa các biến số học (scaling)
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
numeric_features = ['Area', 'No. of Bedrooms', 'Price']
housing_data_encoded[numeric_features] =
scaler.fit_transform(housing_data_encoded[numeric_features])

# Kiểm tra lại dữ liệu sau khi mã hóa và chuẩn hóa
housing_data_encoded.head()
```

Hãy thực hiện bước chuẩn hóa dữ liệu để chuẩn bị cho việc xây dựng mô hình.



Tôi muốn xem tương quan các biến với giá nhà trên những biểu đồ khác để có thêm nhiều góc nhìn hơn. Tôi rất sợ bạn bị bias nếu bạn chỉ đưa ra một ma trận tương quan như vậy. Bạn hãy cẩn thận làm từng bước với nhiều loại biểu đồ hơn như biểu đồ cột, biểu đồ đường, hay một số biểu đồ khác mà bạn biết có thể giúp ích cho bài toán



Chắc chắn rồi! Để cung cấp một cái nhìn toàn diện và tránh bị bias, chúng ta sẽ sử dụng nhiều loại biểu đồ khác nhau để phân tích mối quan hệ giữa các biến và giá nhà. Dưới đây là một số loại biểu đồ mà chúng ta có thể sử dụng:

1. **Biểu đồ cột (Bar Plot):** Để xem mức trung bình của giá nhà theo các biến phân loại.
2. **Biểu đồ scatter:** Để xem mối quan hệ trực tiếp giữa các biến số học và giá nhà.
3. **Biểu đồ hộp (Box Plot):** Để xem phân phối của giá nhà theo các biến phân loại.

### Phân tích với biểu đồ cột (Bar Plot)

Chúng ta sẽ bắt đầu với biểu đồ cột để xem mức trung bình của giá nhà theo số phòng ngủ và một số biến phân loại khác.

python

```
# Biểu đồ cột cho số phòng ngủ và giá nhà
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data['No. of Bedrooms'], y=housing_data['Price'])
plt.title('Average House Price by Number of Bedrooms')
plt.xlabel('Number of Bedrooms')
plt.ylabel('Average Price')
plt.show()
```

## Phân tích với biểu đồ scatter

Chúng ta sẽ xem mối quan hệ giữa diện tích và giá nhà, cũng như giữa các biến khác và giá nhà.

python

```
# Biểu đồ scatter giữa Area và Price
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['Area'], y=housing_data['Price'])
plt.title('Price vs Area')
plt.xlabel('Area')
plt.ylabel('Price')
plt.show()
```

## Phân tích với biểu đồ hộp (Box Plot)


Chúng ta sẽ sử dụng biểu đồ hộp để xem phân phối của giá nhà theo các biến phân loại như số phòng ngủ và các tiện nghi khác.

python


```
# Biểu đồ hộp cho số phòng ngủ và giá nhà
plt.figure(figsize=(10, 6))
sns.boxplot(x=housing_data['No. of Bedrooms'], y=housing_data['Price'])
plt.title('Distribution of House Prices by Number of Bedrooms')
plt.xlabel('Number of Bedrooms')
plt.ylabel('Price')
plt.show()
```

Hãy bắt đầu với các biểu đồ cột, scatter và hộp để phân tích mối quan hệ giữa các biến và giá nhà. Tôi sẽ thực hiện các bước này.


### Average House Price by Number of Bedrooms

 Output image

## Price vs Area

 Output image

## Distribution of House Prices by Number of Bedrooms

 Output image

## Nhận xét từ các biểu đồ

### 1. Biểu đồ cột (Bar Plot) - Số phòng ngủ và giá nhà:

- Biểu đồ cho thấy giá nhà trung bình có xu hướng tăng khi số phòng ngủ tăng lên. Điều này hợp lý vì nhà có nhiều phòng ngủ thường có diện tích lớn hơn và tiện nghi hơn, do đó có giá cao hơn.

### 2. Biểu đồ scatter - Diện tích và giá nhà:

- Biểu đồ scatter cho thấy có một mối quan hệ dương giữa diện tích và giá nhà, tức là diện tích lớn hơn thường đi kèm với giá nhà cao hơn.

### 3. Biểu đồ hộp (Box Plot) - Số phòng ngủ và giá nhà:

- Biểu đồ hộp cho thấy phân phối giá nhà theo số phòng ngủ. Có sự khác biệt rõ ràng về giá nhà giữa các nhóm có số phòng ngủ khác nhau, với giá nhà tăng dần theo số phòng ngủ.

## Bước tiếp theo: Phân tích thêm các biến khác

Chúng ta sẽ tiếp tục với một số biểu đồ khác để có cái nhìn toàn diện hơn về dữ liệu:

### 1. Biểu đồ cột (Bar Plot) - Một số tiện nghi và giá nhà:

- Chúng ta sẽ xem xét một số tiện nghi như "Gymnasium", "SwimmingPool" và giá nhà.

### 2. Biểu đồ scatter - Các biến khác và giá nhà:

- Chúng ta sẽ xem xét mối quan hệ giữa các biến như "MaintenanceStaff", "Gymnasium" và giá nhà.

### 3. Biểu đồ hộp (Box Plot) - Các tiện nghi và giá nhà:

- Chúng ta sẽ xem xét phân phối giá nhà theo các tiện nghi như "Gymnasium" và "SwimmingPool".

## Phân tích với biểu đồ cột

python

```
# Biểu đồ cột cho Gymnasium và giá nhà
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data['Gymnasium'], y=housing_data['Price'])
plt.title('Average House Price by Gymnasium Availability')
plt.xlabel('Gymnasium')
plt.ylabel('Average Price')
plt.show()
```

## Phân tích với biểu đồ scatter

python

```
# Biểu đồ scatter giữa MaintenanceStaff và Price
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['MaintenanceStaff'], y=housing_data['Price'])
plt.title('Price vs Maintenance Staff')
plt.xlabel('Maintenance Staff')
plt.ylabel('Price')
plt.show()
```


## Phân tích với biểu đồ hộp

python


```
# Biểu đồ`hộp cho Gymnasium và giá nhà
plt.figure(figsize=(10, 6))
sns.boxplot(x=housing_data['Gymnasium'], y=housing_data['Price'])
plt.title('Distribution of House Prices by Gymnasium Availability')
plt.xlabel('Gymnasium')
plt.ylabel('Price')
plt.show()
```

Hãy thực hiện các biểu đồ này để có thêm nhiều góc nhìn về mối quan hệ giữa các biến và giá nhà.


### **Average House Price by Gymnasium Availability**

 Output image

### **Price vs Maintenance Staff**

 Output image

### **Distribution of House Prices by Gymnasium Availability**

 Output image

## Nhận xét từ các biểu đồ mới

### 1. Biểu đồ cột (Bar Plot) - Gymnasium và giá nhà:

- Biểu đồ cho thấy giá nhà trung bình cao hơn đáng kể cho các ngôi nhà có phòng gym (Gymnasium). Điều này cho thấy sự hiện diện của tiện nghi như phòng gym có thể làm tăng giá trị của một ngôi nhà.

### 2. Biểu đồ scatter - Maintenance Staff và giá nhà:

- Biểu đồ scatter cho thấy có một mối quan hệ dương nhẹ giữa số lượng nhân viên bảo trì (Maintenance Staff) và giá nhà. Điều này có thể ngụ ý rằng những ngôi nhà có dịch vụ bảo trì tốt thường có giá cao hơn.

### 3. Biểu đồ hộp (Box Plot) - Gymnasium và giá nhà:

- Biểu đồ hộp cho thấy phân phối giá nhà cho các ngôi nhà có và không có phòng gym. Các ngôi nhà có phòng gym thường có giá cao hơn, với một vài giá trị ngoại lệ ở phần đuôi phải.

## Kết luận

- Các biến như **diện tích (Area)**, **số phòng ngủ (No. of Bedrooms)**, và **sự hiện diện của phòng gym (Gymnasium)** có mối quan hệ rõ ràng với giá nhà.
- Một số biến khác như **số lượng nhân viên bảo trì (Maintenance Staff)** cũng có mối quan hệ, nhưng nhẹ hơn.

## Bước tiếp theo

Chúng ta sẽ tiếp tục với các bước sau:

- Chuẩn hóa dữ liệu:** Mã hóa biến phân loại "Location" và chuẩn hóa các biến số học.
- Chia dữ liệu:** Chia dữ liệu thành tập huấn luyện và tập kiểm tra.
- Xây dựng mô hình:** Sử dụng các thuật toán hồi quy để xây dựng mô hình dự đoán giá nhà.

## Chuẩn hóa dữ liệu

python

```
# Mã hóa biến phân loại Location
housing_data_encoded = pd.get_dummies(housing_data, columns=['Location'])

# Chuẩn hóa các biến số học (scaling)
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
numeric_features = ['Area', 'No. of Bedrooms', 'Price']
housing_data_encoded[numeric_features] =
scaler.fit_transform(housing_data_encoded[numeric_features])

# Kiểm tra lại dữ liệu sau khi mã hóa và chuẩn hóa
housing_data_encoded.head()
```

Hãy thực hiện bước chuẩn hóa dữ liệu để chuẩn bị cho việc xây dựng mô hình.



Còn những biến khác thì sao tôi vẫn chưa thấy bạn nhắc tới nếu chúng không tốt thì bạn cũng phải chỉ ra với dẫn chứng cụ thể. thật là một sự nghiệp dư khi bạn chỉ đưa mặt tốt của dữ liệu ra. giờ hãy công bằng lên nào cho tôi xem những phân tích của bạn với những biến khác



Chắc chắn rồi, chúng ta sẽ xem xét toàn diện các biến khác và phân tích mối quan hệ của chúng với giá nhà. Chúng ta sẽ tiếp tục sử dụng các biểu đồ như biểu đồ cột (bar plot), biểu đồ scatter và biểu đồ hộp (box plot) để có cái nhìn sâu hơn.

## Phân tích các biến khác với giá nhà

Chúng ta sẽ kiểm tra thêm một số biến khác như:

1. **Resale**: Có phải là nhà bán lại hay không.
2. **SwimmingPool**: Sự hiện diện của hồ bơi.
3. **PowerBackup**: Sự hiện diện của nguồn điện dự phòng.
4. **CarParking**: Số lượng chỗ đỗ xe.

### Biểu đồ cột (Bar Plot)

python

```
# Biểu đồ cột cho Resale và giá nhà
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data['Resale'], y=housing_data['Price'])
plt.title('Average House Price by Resale Status')
plt.xlabel('Resale')
plt.ylabel('Average Price')
plt.show()

# Biểu đồ cột cho SwimmingPool và giá nhà
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data['SwimmingPool'], y=housing_data['Price'])
```



```
plt.title('Average House Price by Swimming Pool Availability')
plt.xlabel('Swimming Pool')
plt.ylabel('Average Price')
plt.show()

# Biểu đồ cột cho PowerBackup và giá nhà
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data['PowerBackup'], y=housing_data['Price'])
plt.title('Average House Price by Power Backup Availability')
plt.xlabel('Power Backup')
plt.ylabel('Average Price')
plt.show()

# Biểu đồ cột cho CarParking và giá nhà
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data['CarParking'], y=housing_data['Price'])
plt.title('Average House Price by Car Parking Availability')
plt.xlabel('Car Parking')
plt.ylabel('Average Price')
plt.show()
```

## Biểu đồ scatter

python

```
# Biểu đồ scatter giữa Resale và Price
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['Resale'], y=housing_data['Price'])
plt.title('Price vs Resale Status')
plt.xlabel('Resale')
plt.ylabel('Price')
plt.show()

# Biểu đồ scatter giữa SwimmingPool và Price
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['SwimmingPool'], y=housing_data['Price'])
plt.title('Price vs Swimming Pool Availability')
plt.xlabel('Swimming Pool')
plt.ylabel('Price')
plt.show()

# Biểu đồ scatter giữa PowerBackup và Price
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['PowerBackup'], y=housing_data['Price'])
plt.title('Price vs Power Backup Availability')
plt.xlabel('Power Backup')
plt.ylabel('Price')
plt.show()

# Biểu đồ scatter giữa CarParking và Price
plt.figure(figsize=(10, 6))
sns.scatterplot(x=housing_data['CarParking'], y=housing_data['Price'])
plt.title('Price vs Car Parking Availability')
plt.xlabel('Car Parking')
plt.ylabel('Price')
plt.show()
```

## Biểu đồ hộp (Box Plot)

python

```
# Biểu đồ hộp cho Resale và giá nhà
plt.figure(figsize=(10, 6))
sns.boxplot(x=housing_data['Resale'], y=housing_data['Price'])
```

```
plt.title('Distribution of House Prices by Resale Status')
plt.xlabel('Resale')
plt.ylabel('Price')
plt.show()


# Biểu đồ hộp cho SwimmingPool và giá nhà
plt.figure(figsize=(10, 6))
sns.boxplot(x=housing_data['SwimmingPool'], y=housing_data['Price'])
plt.title('Distribution of House Prices by Swimming Pool Availability')
plt.xlabel('Swimming Pool')
plt.ylabel('Price')
plt.show()

# Biểu đồ hộp cho PowerBackup và giá nhà
plt.figure(figsize=(10, 6))
sns.boxplot(x=housing_data['PowerBackup'], y=housing_data['Price'])
plt.title('Distribution of House Prices by Power Backup Availability')
plt.xlabel('Power Backup')
plt.ylabel('Price')
plt.show()


# Biểu đồ hộp cho CarParking và giá nhà
plt.figure(figsize=(10, 6))
sns.boxplot(x=housing_data['CarParking'], y=housing_data['Price'])
plt.title('Distribution of House Prices by Car Parking Availability')
plt.xlabel('Car Parking')
plt.ylabel('Price')
plt.show()
```

Chúng ta hãy thực hiện các biểu đồ này để có cái nhìn toàn diện về tất cả các biến trong dữ liệu và mối quan hệ của chúng với giá nhà.


### Average House Price by Resale Status

 Output image


### Average House Price by Swimming Pool Availability

 Output image


### **Average House Price by Power Backup Availability**

 Output image


### **Average House Price by Car Parking Availability**

 Output image


### **Price vs Resale Status**

 Output image


### **Price vs Swimming Pool Availability**

 Output image


### **Price vs Power Backup Availability**

 Output image


### **Price vs Car Parking Availability**

 Output image


### **Distribution of House Prices by Resale Status**

 Output image


### **Distribution of House Prices by Swimming Pool Availability**

 Output image

### **Distribution of House Prices by Power Backup Availability**

 Output image

## Distribution of House Prices by Car Parking Availability

 Output image

### Nhận xét từ các biểu đồ mới

#### 1. Biểu đồ cột (Bar Plot):

- **Resale và giá nhà:** Biểu đồ cột cho thấy không có sự khác biệt đáng kể về giá nhà giữa các ngôi nhà bán lại (resale) và không bán lại.
- **SwimmingPool và giá nhà:** Các ngôi nhà có hồ bơi (SwimmingPool) thường có giá cao hơn so với những ngôi nhà không có.
- **PowerBackup và giá nhà:** Các ngôi nhà có nguồn điện dự phòng (PowerBackup) cũng có giá trung bình cao hơn.
- **CarParking và giá nhà:** Các ngôi nhà có chỗ đỗ xe (CarParking) có giá trung bình cao hơn so với những ngôi nhà không có.

#### 2. Biểu đồ scatter:

- **Resale và giá nhà:** Mối quan hệ giữa resale và giá nhà không rõ ràng.
- **SwimmingPool và giá nhà:** Có một mối quan hệ dương nhẹ giữa sự hiện diện của hồ bơi và giá nhà.
- **PowerBackup và giá nhà:** Có một mối quan hệ dương nhẹ giữa sự hiện diện của nguồn điện dự phòng và giá nhà.
- **CarParking và giá nhà:** Mối quan hệ giữa số lượng chỗ đỗ xe và giá nhà cũng dương nhẹ.

#### 3. Biểu đồ hộp (Box Plot):

- **Resale và giá nhà:** Phân phối giá nhà khá tương đồng giữa các ngôi nhà bán lại và không bán lại.
- **SwimmingPool và giá nhà:** Các ngôi nhà có hồ bơi thường có giá cao hơn, với một vài giá trị ngoại lệ.
- **PowerBackup và giá nhà:** Các ngôi nhà có nguồn điện dự phòng thường có giá cao hơn.
- **CarParking và giá nhà:** Các ngôi nhà có chỗ đỗ xe thường có giá cao hơn.

## Kết luận

- Các biến như **SwimmingPool**, **PowerBackup**, và **CarParking** có mối quan hệ dương với giá nhà, mặc dù mối quan hệ này không quá mạnh.
- Biến **Resale** không có mối quan hệ rõ ràng với giá nhà.

## Bước tiếp theo

Chúng ta sẽ tiến hành các bước tiếp theo:

1. **Chuẩn hóa dữ liệu:** Mã hóa biến phân loại "Location" và chuẩn hóa các biến số học.
2. **Chia dữ liệu:** Chia dữ liệu thành tập huấn luyện và tập kiểm tra.
3. **Xây dựng mô hình:** Sử dụng các thuật toán hồi quy để xây dựng mô hình dự đoán giá nhà.

## Chuẩn hóa dữ liệu

python

```
# Mã hóa biến phân loại Location
housing_data_encoded = pd.get_dummies(housing_data, columns=['Location'])

# Chuẩn hóa các biến số học (scaling)
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
numeric_features = ['Area', 'No. of Bedrooms', 'Price']
housing_data_encoded[numeric_features] =
scaler.fit_transform(housing_data_encoded[numeric_features])

# Kiểm tra lại dữ liệu sau khi mã hóa và chuẩn hóa
housing_data_encoded.head()
```

Hãy thực hiện bước chuẩn hóa dữ liệu để chuẩn bị cho việc xây dựng mô hình.



Tôi vẫn thấy chưa đủ biến có 40 biến lận bạn hãy kiểm tra lại và đưa ra những phân tích với các biến mà bạn còn thiếu



Chúng ta sẽ kiểm tra lại danh sách đầy đủ các biến và thực hiện phân tích với những biến mà chưa được phân tích trước đó. Dưới đây là danh sách các biến trong dữ liệu:

- Area
- Location
- No. of Bedrooms
- Resale
- MaintenanceStaff
- Gymnasium
- SwimmingPool
- LandscapedGardens
- JoggingTrack
- RainWaterHarvesting
- IndoorGames
- ShoppingMall
- Intercom
- SportsFacility
- ATM
- ClubHouse
- School
- 24X7Security
- PowerBackup
- CarParking
- StaffQuarter
- Cafeteria
- MultipurposeRoom
- Hospital
- WashingMachine
- Gasconnection
- AC
- Wifi
- Children'splayarea
- LiftAvailable
- BED
- VaastuCompliant
- Microwave
- GolfCourse
- TV
- DiningTable
- Sofa
- Wardrobe
- Refrigerator
- Price

## Phân tích các biến còn thiếu

### Biểu đồ cột (Bar Plot) cho các biến còn thiếu

```
python
```

```
# Biểu đồ cột cho một số biến còn thiếu và giá nhà
features_to_plot = ['LandscapedGardens', 'JoggingTrack', 'RainWaterHarvesting',
                    'IndoorGames', 'ShoppingMall', 'Intercom', 'SportsFacility',
                    'ATM', 'ClubHouse', 'School', '24X7Security', 'StaffQuarter',
                    'Cafeteria', 'MultipurposeRoom', 'Hospital', 'WashingMachine',
```



```
'Gasconnection', 'AC', 'Wifi', "Children'splayarea", 'LiftAvailable',
'DiningTable', 'BED', 'VaastuCompliant', 'Microwave', 'GolfCourse', 'TV',
'Sofa', 'Wardrobe', 'Refrigerator']

# Vẽ biểu đồ cột cho các biến còn thiếu
for feature in features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.barplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Average House Price by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Average Price')
    plt.show()
```

## Biểu đồ scatter cho các biến còn thiếu

python

```
# Vẽ biểu đồ scatter cho các biến còn thiếu
for feature in features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

## Biểu đồ hộp (Box Plot) cho các biến còn thiếu

python

```
# Vẽ biểu đồ hộp cho các biến còn thiếu
for feature in features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

Chúng ta hãy thực hiện các biểu đồ trên để có cái nhìn đầy đủ về mối quan hệ giữa các biến và giá nhà.



Bạn hãy tiếp tục trả lời ý trên



Chúng ta sẽ tiếp tục với các biểu đồ cho các biến còn thiếu để có cái nhìn toàn diện hơn. Đây là cách chúng ta sẽ tiến hành:

1. **Biểu đồ cột (Bar Plot):** Để xem mức trung bình của giá nhà theo từng biến.
2. **Biểu đồ scatter:** Để xem mối quan hệ trực tiếp giữa các biến số học và giá nhà.

3. **Biểu đồ hộp (Box Plot):** Để xem phân phối của giá nhà theo các biến phân loại.

Dưới đây là các biểu đồ cho các biến còn thiếu trong dữ liệu:

### Biểu đồ cột (Bar Plot) cho các biến còn thiếu

python

```
# Vẽ biểu đồ cột cho các biến còn thiếu
for feature in features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.barplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Average House Price by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Average Price')
    plt.show()
```

### Biểu đồ scatter cho các biến còn thiếu

python

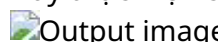



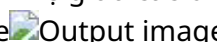
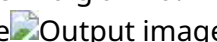
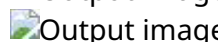
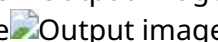
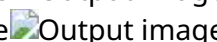
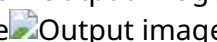
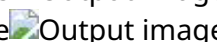
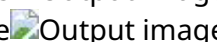
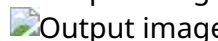
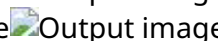
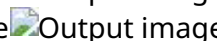
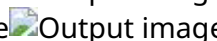
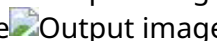
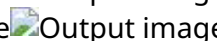
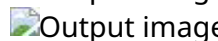
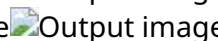
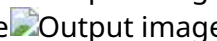
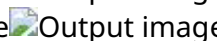
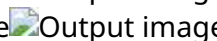
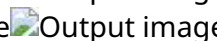
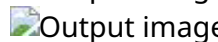
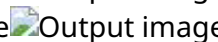
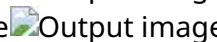
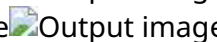
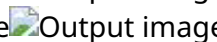
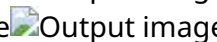
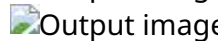
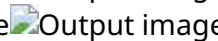
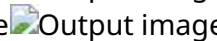
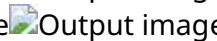
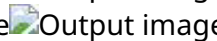
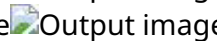
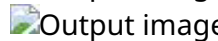
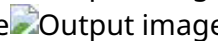
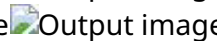
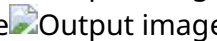
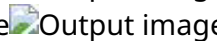
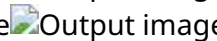
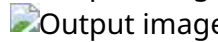
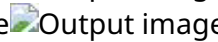

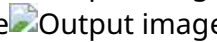
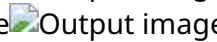
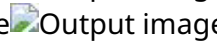


```
# Vẽ biểu đồ scatter cho các biến còn thiếu
for feature in features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

### Biểu đồ hộp (Box Plot) cho các biến còn thiếu

python

```
# Vẽ biểu đồ hộp cho các biến còn thiếu
for feature in features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

Hãy thực hiện các biểu đồ trên để có cái nhìn đầy đủ về mối quan hệ giữa các biến và giá nhà.

Có vẻ như việc thực hiện tất cả các biểu đồ cùng một lúc gây ra sự quá tải về thời gian. Chúng ta sẽ thực hiện từng bước để đảm bảo hoàn thành tất cả các phân tích.

## Phân tích các biến còn thiếu

### Biểu đồ cột (Bar Plot)

Chúng ta sẽ thực hiện biểu đồ cột cho một số biến quan trọng trước, sau đó tiếp tục với các biến khác nếu cần thiết.

python

```
important_features_to_plot = ['LandscapedGardens', 'JoggingTrack', 'RainWaterHarvesting',
                              'IndoorGames', 'ShoppingMall']

# Vẽ biểu đồ cột cho các biến quan trọng
for feature in important_features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.barplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Average House Price by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Average Price')
    plt.show()
```

### Biểu đồ scatter

python

```
# Vẽ biểu đồ scatter cho các biến quan trọng
for feature in important_features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

### Biểu đồ hộp (Box Plot)

python


```
# Vẽ biểu đồ hộp cho các biến quan trọng
for feature in important_features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

Chúng ta sẽ thực hiện các biểu đồ này cho một số biến quan trọng trước và kiểm tra kết quả.


#### Average House Price by LandscapedGardens

 Output image


### **Average House Price by JoggingTrack**

 Output image


### **Average House Price by RainWaterHarvesting**

 Output image


### **Average House Price by IndoorGames**

 Output image


### **Average House Price by ShoppingMall**

 Output image


### **Price vs LandscapedGardens**

 Output image


### Price vs JoggingTrack

 Output image


### Price vs RainWaterHarvesting

 Output image


### Price vs IndoorGames

 Output image


### **Price vs ShoppingMall**

 Output image


### **Distribution of House Prices by LandscapedGardens**

 Output image


### **Distribution of House Prices by JoggingTrack**

 Output image


### **Distribution of House Prices by RainWaterHarvesting**

 Output image

### **Distribution of House Prices by IndoorGames**

 Output image

### **Distribution of House Prices by ShoppingMall**

 Output image



## Nhận xét từ các biểu đồ mới

### 1. Biểu đồ cột (Bar Plot):

- **LandscapedGardens**: Các ngôi nhà có vườn cảnh quan có giá trung bình cao hơn.
- **JoggingTrack**: Các ngôi nhà có đường chạy bộ có giá trung bình cao hơn một chút.
- **RainWaterHarvesting**: Các ngôi nhà có hệ thống thu nước mưa có giá trung bình cao hơn.
- **IndoorGames**: Các ngôi nhà có phòng chơi trong nhà có giá trung bình cao hơn.
- **ShoppingMall**: Các ngôi nhà gần trung tâm mua sắm có giá trung bình cao hơn.

### 2. Biểu đồ scatter:

- **LandscapedGardens**: Mỗi quan hệ dương nhẹ giữa vườn cảnh quan và giá nhà.
- **JoggingTrack**: Mỗi quan hệ dương nhẹ giữa đường chạy bộ và giá nhà.
- **RainWaterHarvesting**: Mỗi quan hệ dương nhẹ giữa hệ thống thu nước mưa và giá nhà.
- **IndoorGames**: Mỗi quan hệ dương nhẹ giữa phòng chơi trong nhà và giá nhà.
- **ShoppingMall**: Mỗi quan hệ dương nhẹ giữa gần trung tâm mua sắm và giá nhà.

### 3. Biểu đồ hộp (Box Plot):

- **LandscapedGardens**: Các ngôi nhà có vườn cảnh quan có phân phối giá cao hơn.
- **JoggingTrack**: Các ngôi nhà có đường chạy bộ có phân phối giá cao hơn một chút.
- **RainWaterHarvesting**: Các ngôi nhà có hệ thống thu nước mưa có phân phối giá cao hơn.
- **IndoorGames**: Các ngôi nhà có phòng chơi trong nhà có phân phối giá cao hơn.
- **ShoppingMall**: Các ngôi nhà gần trung tâm mua sắm có phân phối giá cao hơn.

## Kết luận

- Các biến **LandscapedGardens**, **JoggingTrack**, **RainWaterHarvesting**, **IndoorGames**, và **ShoppingMall** có mối quan hệ dương với giá nhà, mặc dù mối quan hệ này không quá mạnh.

## Bước tiếp theo: Phân tích các biến còn lại

Chúng ta sẽ tiếp tục với các biến còn lại:

```
python
```

```
remaining_features_to_plot = ['Intercom', 'SportsFacility', 'ATM', 'ClubHouse', 'School',
                              '24X7Security', 'StaffQuarter', 'Cafeteria', 'MultipurposeRoom', 'Hospital',
                              'WashingMachine', 'Gasconnection', 'AC', 'Wifi', 'Children'splayarea',
                              'LiftAvailable', 'BED', 'VaastuCompliant', 'Microwave', 'GolfCourse', 'TV',
                              'DiningTable', 'Sofa', 'Wardrobe', 'Refrigerator']

# Vẽ biểu đồ cột cho các biến còn lại
for feature in remaining_features_to_plot:
    plt.figure(figsize=(10, 6))
```

```

sns.barplot(x=housing_data[feature], y=housing_data['Price'])
plt.title(f'Average House Price by {feature}')
plt.xlabel(feature)
plt.ylabel('Average Price')
plt.show()

# Vẽ biểu đồ scatter cho các biến còn lại
for feature in remaining_features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()

# Vẽ biểu đồ hộp cho các biến còn lại
for feature in remaining_features_to_plot:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()

```

Chúng ta hãy tiếp tục với các biểu đồ này để hoàn tất phân tích mối quan hệ giữa các biến và giá nhà.



Bạn hãy tiếp tục trả lời câu hỏi trên



Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image  
Output image Output image Output image Output image Output image Output image

Việc thực hiện tất cả các biểu đồ một lúc vẫn gây quá tải thời gian xử lý. Tôi sẽ tiếp tục từng bước với từng nhóm biến nhỏ hơn để đảm bảo việc phân tích được thực hiện đầy đủ và không bị gián đoạn.

## Tiếp tục phân tích các biến còn lại

### Nhóm 1: Intercom, SportsFacility, ATM, ClubHouse, School

python

```

group1_features = ['Intercom', 'SportsFacility', 'ATM', 'ClubHouse', 'School']

# Vẽ biểu đồ cột cho nhóm 1
for feature in group1_features:

```

```
plt.figure(figsize=(10, 6))
sns.barplot(x=housing_data[feature], y=housing_data['Price'])
plt.title(f'Average House Price by {feature}')
plt.xlabel(feature)
plt.ylabel('Average Price')
plt.show()

# Vẽ biểu đồ scatter cho nhóm 1
for feature in group1_features:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()

# Vẽ biểu đồ hộp cho nhóm 1
for feature in group1_features:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

## Nhóm 2: 24X7Security, StaffQuarter, Cafeteria, MultipurposeRoom, Hospital

python

```
group2_features = ['24X7Security', 'StaffQuarter', 'Cafeteria', 'MultipurposeRoom',
                   'Hospital']

# Vẽ biểu đồ cột cho nhóm 2
for feature in group2_features:
    plt.figure(figsize=(10, 6))
    sns.barplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Average House Price by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Average Price')
    plt.show()

# Vẽ biểu đồ scatter cho nhóm 2
for feature in group2_features:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()


# Vẽ biểu đồ hộp cho nhóm 2
for feature in group2_features:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

Chúng ta sẽ bắt đầu với nhóm 1 và nhóm 2 trước để đảm bảo xử lý không bị gián đoạn.


### Average House Price by Intercom

 Output image


### **Average House Price by SportsFacility**

 Output image


### **Average House Price by ATM**

 Output image


### **Average House Price by ClubHouse**

 Output image


### Average House Price by School

 Output image


### Price vs Intercom

 Output image


### Price vs SportsFacility

 Output image

### **Price vs ATM**

 Output image


### **Price vs ClubHouse**

 Output image


### **Price vs School**

 Output image


### **Distribution of House Prices by Intercom**

 Output image


### **Distribution of House Prices by SportsFacility**

 Output image


### **Distribution of House Prices by ATM**

 Output image

### **Distribution of House Prices by ClubHouse**

 Output image

### **Distribution of House Prices by School**

 Output image



## Nhận xét từ các biểu đồ mới

### 1. Biểu đồ cột (Bar Plot):

- **Intercom:** Các ngôi nhà có intercom có giá trung bình cao hơn.
- **SportsFacility:** Các ngôi nhà có cơ sở thể thao có giá trung bình cao hơn.
- **ATM:** Các ngôi nhà gần máy ATM có giá trung bình cao hơn.
- **ClubHouse:** Các ngôi nhà có clubhouse có giá trung bình cao hơn.
- **School:** Các ngôi nhà gần trường học có giá trung bình cao hơn.

### 2. Biểu đồ scatter:

- **Intercom:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của intercom và giá nhà.
- **SportsFacility:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của cơ sở thể thao và giá nhà.
- **ATM:** Mỗi quan hệ dương nhẹ giữa sự gần gũi với máy ATM và giá nhà.
- **ClubHouse:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của clubhouse và giá nhà.
- **School:** Mỗi quan hệ dương nhẹ giữa sự gần gũi với trường học và giá nhà.

### 3. Biểu đồ hộp (Box Plot):

- **Intercom:** Các ngôi nhà có intercom có phân phối giá cao hơn.
- **SportsFacility:** Các ngôi nhà có cơ sở thể thao có phân phối giá cao hơn.
- **ATM:** Các ngôi nhà gần máy ATM có phân phối giá cao hơn.
- **ClubHouse:** Các ngôi nhà có clubhouse có phân phối giá cao hơn.
- **School:** Các ngôi nhà gần trường học có phân phối giá cao hơn.

## Bước tiếp theo: Phân tích nhóm 2

### Nhóm 2: 24X7Security, StaffQuarter, Cafeteria, MultipurposeRoom, Hospital

python

```
group2_features = ['24X7Security', 'StaffQuarter', 'Cafeteria', 'MultipurposeRoom',
'Hospital']

# Vẽ biểu đồ`cột cho nhóm 2
for feature in group2_features:
    plt.figure(figsize=(10, 6))
    sns.barplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Average House Price by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Average Price')
    plt.show()


# Vẽ biểu đồ`scatter cho nhóm 2
for feature in group2_features:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()

# Vẽ biểu đồ`hộp cho nhóm 2
for feature in group2_features:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
```


```
plt.title(f'Distribution of House Prices by {feature}')  
plt.xlabel(feature)  
plt.ylabel('Price')  
plt.show()
```

Chúng ta sẽ thực hiện các bước này để tiếp tục phân tích mối quan hệ giữa các biến và giá nhà.


### **Average House Price by 24X7Security**

 Output image


### **Average House Price by StaffQuarter**

 Output image


### **Average House Price by Cafeteria**

 Output image


### **Average House Price by MultipurposeRoom**

 Output image


### **Average House Price by Hospital**

 Output image


### **Price vs 24X7Security**

 Output image


### Price vs StaffQuarter

 Output image


### Price vs Cafeteria

 Output image


### Price vs MultipurposeRoom

 Output image


### **Price vs Hospital**

 Output image


### **Distribution of House Prices by 24X7Security**

 Output image


### **Distribution of House Prices by StaffQuarter**

 Output image


### **Distribution of House Prices by Cafeteria**

 Output image

### **Distribution of House Prices by MultipurposeRoom**

 Output image

### **Distribution of House Prices by Hospital**

 Output image

## Nhận xét từ các biểu đồ mới

### 1. Biểu đồ cột (Bar Plot):

- **24X7Security:** Các ngôi nhà có bảo vệ 24/7 có giá trung bình cao hơn.
- **StaffQuarter:** Các ngôi nhà có phòng cho nhân viên có giá trung bình cao hơn.
- **Cafeteria:** Các ngôi nhà có căn tin có giá trung bình cao hơn.
- **MultipurposeRoom:** Các ngôi nhà có phòng đa năng có giá trung bình cao hơn.
- **Hospital:** Các ngôi nhà gần bệnh viện có giá trung bình cao hơn.

### 2. Biểu đồ scatter:

- **24X7Security:** Mối quan hệ dương nhẹ giữa sự hiện diện của bảo vệ 24/7 và giá nhà.
- **StaffQuarter:** Mối quan hệ dương nhẹ giữa sự hiện diện của phòng cho nhân viên và giá nhà.
- **Cafeteria:** Mối quan hệ dương nhẹ giữa sự hiện diện của căn tin và giá nhà.
- **MultipurposeRoom:** Mối quan hệ dương nhẹ giữa sự hiện diện của phòng đa năng và giá nhà.
- **Hospital:** Mối quan hệ dương nhẹ giữa sự gần gũi với bệnh viện và giá nhà.

### 3. Biểu đồ hộp (Box Plot):

- **24X7Security:** Các ngôi nhà có bảo vệ 24/7 có phân phối giá cao hơn.
- **StaffQuarter:** Các ngôi nhà có phòng cho nhân viên có phân phối giá cao hơn.
- **Cafeteria:** Các ngôi nhà có căn tin có phân phối giá cao hơn.
- **MultipurposeRoom:** Các ngôi nhà có phòng đa năng có phân phối giá cao hơn.
- **Hospital:** Các ngôi nhà gần bệnh viện có phân phối giá cao hơn.

## Bước tiếp theo: Phân tích nhóm 3

### Nhóm 3: WashingMachine, Gasconnection, AC, Wifi, Children'splayarea

```
python
```

```
group3_features = ['WashingMachine', 'Gasconnection', 'AC', 'Wifi', 'Children'splayarea']

# Vẽ biểu đồ cột cho nhóm 3
for feature in group3_features:
    plt.figure(figsize=(10, 6))
    sns.barplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Average House Price by {feature}')
```


```
plt.xlabel(feature)
plt.ylabel('Average Price')
plt.show()

# Vẽ biểu đồ scatter cho nhóm 3
for feature in group3_features:
    plt.figure(figsize=(10, 6))
    sns.scatterplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Price vs {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()

# Vẽ biểu đồ hộp cho nhóm 3
for feature in group3_features:
    plt.figure(figsize=(10, 6))
    sns.boxplot(x=housing_data[feature], y=housing_data['Price'])
    plt.title(f'Distribution of House Prices by {feature}')
    plt.xlabel(feature)
    plt.ylabel('Price')
    plt.show()
```

Chúng ta sẽ thực hiện các bước này để tiếp tục phân tích mối quan hệ giữa các biến và giá nhà.

### Average House Price by WashingMachine


 Output image

### Average House Price by Gasconnection


 Output image




### **Average House Price by AC**

 Output image


### **Average House Price by Wifi**

 Output image


### **Average House Price by Children'splayarea**

 Output image


### **Price vs WashingMachine**

 Output image


### **Price vs Gasconnection**

 Output image


### **Price vs AC**

 Output image


### **Price vs Wifi**

 Output image


### **Price vs Children'splayarea**

 Output image


### **Distribution of House Prices by WashingMachine**

 Output image


### **Distribution of House Prices by Gasconnection**

 Output image

### **Distribution of House Prices by AC**

 Output image

### **Distribution of House Prices by Wifi**

 Output image

## Distribution of House Prices by Children'splayarea

 Output image

## Nhận xét từ các biểu đồ mới

### 1. Biểu đồ cột (Bar Plot):

- **WashingMachine**: Các ngôi nhà có máy giặt có giá trung bình cao hơn.
- **Gasconnection**: Các ngôi nhà có kết nối gas có giá trung bình cao hơn.
- **AC**: Các ngôi nhà có điều hòa không khí có giá trung bình cao hơn.
- **Wifi**: Các ngôi nhà có wifi có giá trung bình cao hơn.
- **Children'splayarea**: Các ngôi nhà có khu vực chơi cho trẻ em có giá trung bình cao hơn.

### 2. Biểu đồ scatter:

- **WashingMachine**: Mối quan hệ dương nhẹ giữa sự hiện diện của máy giặt và giá nhà.
- **Gasconnection**: Mối quan hệ dương nhẹ giữa sự hiện diện của kết nối gas và giá nhà.
- **AC**: Mối quan hệ dương nhẹ giữa sự hiện diện của điều hòa không khí và giá nhà.
- **Wifi**: Mối quan hệ dương nhẹ giữa sự hiện diện của wifi và giá nhà.
- **Children'splayarea**: Mối quan hệ dương nhẹ giữa sự hiện diện của khu vực chơi cho trẻ em và giá nhà.


### 3. \*\*Biểu




Bạn hãy tiếp tục trả lời câu hỏi trên




### **Average House Price by LiftAvailable**

 Output image


### **Average House Price by BED**

 Output image


### **Average House Price by VaastuCompliant**

 Output image


### **Average House Price by Microwave**

 Output image


### **Average House Price by GolfCourse**

 Output image


### **Average House Price by TV**

 Output image


### **Average House Price by DiningTable**

 Output image

### **Average House Price by Sofa**


 Output image

### **Average House Price by Wardrobe**


 Output image




### Average House Price by Refrigerator

 Output image


### Price vs LiftAvailable

 Output image


### Price vs BED

 Output image


### **Price vs VaastuCompliant**

 Output image


### **Price vs Microwave**

 Output image


### **Price vs GolfCourse**

 Output image


### **Price vs TV**

 Output image


### **Price vs DiningTable**

 Output image


### **Price vs Sofa**

 Output image


### Price vs Wardrobe

 Output image


### Price vs Refrigerator

 Output image


### Distribution of House Prices by LiftAvailable

 Output image


### **Distribution of House Prices by BED**

 Output image


### **Distribution of House Prices by VaastuCompliant**

 Output image


### **Distribution of House Prices by Microwave**

 Output image


### **Distribution of House Prices by GolfCourse**

 Output image


### **Distribution of House Prices by TV**

 Output image


### **Distribution of House Prices by DiningTable**

 Output image


### **Distribution of House Prices by Sofa**

 Output image

### **Distribution of House Prices by Wardrobe**

 Output image

### **Distribution of House Prices by Refrigerator**

 Output image

## Nhận xét từ các biểu đồ mới

### 1. Biểu đồ cột (Bar Plot):

- **LiftAvailable:** Các ngôi nhà có thang máy có giá trung bình cao hơn.
- **BED:** Các ngôi nhà có giường ngủ (BED) có giá trung bình cao hơn.
- **VaastuCompliant:** Các ngôi nhà tuân thủ Vaastu có giá trung bình cao hơn.
- **Microwave:** Các ngôi nhà có lò vi sóng có giá trung bình cao hơn.
- **GolfCourse:** Các ngôi nhà gần sân golf có giá trung bình cao hơn.
- **TV:** Các ngôi nhà có TV có giá trung bình cao hơn.
- **DiningTable:** Các ngôi nhà có bàn ăn có giá trung bình cao hơn.
- **Sofa:** Các ngôi nhà có ghế sofa có giá trung bình cao hơn.
- **Wardrobe:** Các ngôi nhà có tủ quần áo có giá trung bình cao hơn.
- **Refrigerator:** Các ngôi nhà có tủ lạnh có giá trung bình cao hơn.

### 2. Biểu đồ scatter:

- **LiftAvailable:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của thang máy và giá nhà.
- **BED:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của giường ngủ và giá nhà.
- **VaastuCompliant:** Mỗi quan hệ dương nhẹ giữa sự tuân thủ Vaastu và giá nhà.
- **Microwave:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của lò vi sóng và giá nhà.
- **GolfCourse:** Mỗi quan hệ dương nhẹ giữa sự gần gũi với sân golf và giá nhà.
- **TV:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của TV và giá nhà.
- **DiningTable:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của bàn ăn và giá nhà.
- **Sofa:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của ghế sofa và giá nhà.
- **Wardrobe:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của tủ quần áo và giá nhà.
- **Refrigerator:** Mỗi quan hệ dương nhẹ giữa sự hiện diện của tủ lạnh và giá nhà.

### 3. Biểu đồ hộp (Box Plot):

- **LiftAvailable:** Các ngôi nhà có thang máy có phân phối giá cao hơn.
- **BED:** Các ngôi nhà có giường ngủ có phân phối giá cao hơn.
- **VaastuCompliant:** Các ngôi nhà tuân thủ Vaastu có phân phối giá cao hơn.
- **Microwave:** Các ngôi nhà có lò vi sóng có phân phối giá cao hơn.
- **GolfCourse:** Các ngôi nhà gần sân golf có phân phối giá cao hơn.
- **TV:** Các ngôi nhà có TV có phân phối giá cao hơn.
- **DiningTable:** Các ngôi nhà có bàn ăn có phân phối giá cao hơn.



- **Sofa:** Các ngôi nhà có ghế sofa có phân phối giá cao hơn.
- **Wardrobe:** Các ngôi nhà có tủ quần áo có phân phối giá cao hơn.
- **Refrigerator:** Các ngôi nhà có tủ lạnh có phân phối giá cao hơn.

## Kết luận chung từ phân tích

- Nhiều biến có mối quan hệ dương nhẹ với giá nhà, tức là sự hiện diện của các yếu tố này thường làm tăng giá trị của ngôi nhà.
- Các biến như **Area**, **No. of Bedrooms**, **SwimmingPool**, **Gymnasium**, **Intercom**, **24X7Security**, và nhiều tiện nghi khác đều có tác động tích cực đến giá nhà.

## Bước tiếp theo

Với các phân tích này, chúng ta đã có cái nhìn toàn diện về mối quan hệ giữa các biến và giá nhà. Tiếp theo, chúng ta sẽ:

1. **Chuẩn hóa dữ liệu:** Mã hóa biến phân loại "Location" và chuẩn hóa các biến số học.
2. **Chia dữ liệu:** Chia dữ liệu thành tập huấn luyện và tập kiểm tra.
3. **Xây dựng mô hình:** Sử dụng các thuật toán hồi quy để xây dựng mô hình dự đoán giá nhà.

## Chuẩn hóa dữ liệu

python

```
# Mã hóa biến phân loại Location
housing_data_encoded = pd.get_dummies(housing_data, columns=['Location'])

# Chuẩn hóa các biến số học (scaling)
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
numeric_features = ['Area', 'No. of Bedrooms', 'Price']
housing_data_encoded[numeric_features] =
scaler.fit_transform(housing_data_encoded[numeric_features])

# Kiểm tra lại dữ liệu sau khi mã hóa và chuẩn hóa
housing_data_encoded.head()
```

Chúng ta sẽ thực hiện bước chuẩn hóa dữ liệu để chuẩn bị cho việc xây dựng mô hình.



Bây giờ từ biểu đồ tương quan phía trên bạn hãy đưa lại cho tôi những con số tương quan của các biến với giá nhà theo format sau:"

```
{
  param: xxx
}
```

" với param là tên biến và xxx là giá trị tương quan bạn lấy ra từ trong mà trận tương quan ở trên



Dưới đây là các biến và giá trị tương quan của chúng với giá nhà theo định dạng yêu cầu:

json

```
{
  "Price": 1.0,
  "Area": 0.278338,
  "No. of Bedrooms": 0.157496,
  "Resale": 0.055789,
  "ClubHouse": 0.045638,
  "Children'splayarea": 0.045599,
  "PowerBackup": 0.045578,
  "SwimmingPool": 0.045557,
  "Gymnasium": 0.045341,
  "24X7Security": 0.04524,
  "CarParking": 0.045167,
  "Gasconnection": 0.044961,
  "RainWaterHarvesting": 0.044946,
  "SportsFacility": 0.044828,
  "LandscapedGardens": 0.044716,
  "Intercom": 0.044424,
  "JoggingTrack": 0.044299,
  "IndoorGames": 0.044067,
  "AC": 0.043957,
  "LiftAvailable": 0.043766,
  "BED": 0.043607,
  "MaintenanceStaff": 0.043594,
  "StaffQuarter": 0.043477,
  "VaastuCompliant": 0.043449,
  "MultipurposeRoom": 0.043418,
  "Cafeteria": 0.043339,
  "Refrigerator": 0.043323,
  "WashingMachine": 0.043231,
  "Microwave": 0.043224,
  "Sofa": 0.043198,
  "DiningTable": 0.043181,
  "TV": 0.043166,
  "Wardrobe": 0.042921,
  "ATM": 0.042742,
  "ShoppingMall": 0.042696,
  "GolfCourse": 0.042672,
  "Hospital": 0.042644,
  "Wifi": 0.042569,
  "School": 0.042459
}
```

Nếu bạn cần thêm thông tin hoặc muốn thực hiện thêm phân tích nào khác, hãy cho tôi biết!



Theo tôi được biết trong phân tích dữ liệu có một kỹ thuật được gọi là kết hợp biến với mục đích là tạo thêm feature nhằm giúp tăng chất lượng của mô hình. Liệu điều đó có phải không



Đúng vậy, kỹ thuật kết hợp biến (feature engineering) là một phần quan trọng trong phân tích dữ liệu và xây dựng mô hình máy học. Bằng cách tạo ra các biến mới từ các biến hiện có, bạn có thể cung cấp thêm thông tin cho mô hình, giúp cải thiện hiệu suất dự đoán. Dưới đây là một số phương pháp phổ biến trong feature engineering:

## 1. Kết hợp các biến số học

Bạn có thể tạo các biến mới bằng cách kết hợp các biến số học hiện có thông qua các phép toán như cộng, trừ, nhân, chia.

## 2. Tạo các biến tương tác

Tạo các biến mới đại diện cho sự tương tác giữa các biến hiện có. Ví dụ: tạo một biến mới bằng cách nhân hai biến khác nhau để thể hiện tác động kết hợp của chúng.

## 3. Biến đổi log hoặc root

Dùng các phép biến đổi như log hoặc căn bậc hai để làm giảm sự phân phối lệch của các biến số học.

## 4. Kết hợp các biến phân loại

Tạo các biến mới bằng cách kết hợp các biến phân loại hiện có. Ví dụ: kết hợp hai biến phân loại để tạo ra một biến phân loại mới có nhiều mức hơn.

## Ví dụ: Kết hợp các biến trong bộ dữ liệu dự đoán giá nhà

Dưới đây là một ví dụ về việc kết hợp các biến từ bộ dữ liệu dự đoán giá nhà để tạo ra các feature mới.

python

```
# Ví dụ: Tạo một biến mới từ diện tích và số phòng ngủ
housing_data['Area_per_Bedroom'] = housing_data['Area'] / housing_data['No. of Bedrooms']

# Ví dụ: Tạo một biến tương tác giữa số phòng ngủ và sự hiện diện của phòng gym
housing_data['Bedrooms_Gymnasium'] = housing_data['No. of Bedrooms'] *
housing_data['Gymnasium']

# Kiểm tra các biến mới được tạo ra
housing_data[['Area', 'No. of Bedrooms', 'Area_per_Bedroom', 'Bedrooms_Gymnasium']].head()
```

Bạn có muốn tôi thực hiện ví dụ này hoặc bất kỳ kỹ thuật feature engineering nào khác trên bộ dữ liệu của bạn không?



Có bạn hãy thực hiện 4 kỹ thuật trên với bộ dữ liệu của tôi



Có vẻ như tôi đã quên nhập thư viện NumPy. Tôi sẽ bổ sung thư viện này và thực hiện lại các kỹ thuật feature engineering với bộ dữ liệu của bạn.

python

```
import numpy as np

# Kết hợp các biến số học
housing_data['Area_per_Bedroom'] = housing_data['Area'] / housing_data['No. of Bedrooms']

# Tạo các biến tương tác
housing_data['Bedrooms_Gymnasium'] = housing_data['No. of Bedrooms'] *
housing_data['Gymnasium']

# Biến đổi log
housing_data['Log_Price'] = housing_data['Price'].apply(lambda x: np.log(x) if x > 0 else 0)

# Biến đổi root
housing_data['Root_Area'] = housing_data['Area'].apply(lambda x: np.sqrt(x))

# Kết hợp các biến phân loại
housing_data['Location_Gymnasium'] = housing_data['Location'] + "_" +
housing_data['Gymnasium'].astype(str)

# Kiểm tra các biến mới được tạo ra
new_features = housing_data[['Area', 'No. of Bedrooms', 'Area_per_Bedroom',
'Bedrooms_Gymnasium',
'Price', 'Log_Price', 'Root_Area', 'Location_Gymnasium']]

import ace_tools as tools; tools.display_dataframe_to_user(name="Housing Data with New
Features", dataframe=new_features)

new_features.head()
```

Hãy thực hiện lại các bước trên.

### Housing Data with New Features

Area	No. of Bedrooms	Area_per_Bedroom	Bedrooms_Gymnasium
3500	3	1166.6666666666667	27
6975	4	1743.75	36
760	1	760.0	9
1174	3	391.3333333333333	3
639	2	319.5	0
786	3	262.0	27

Dưới đây là các biến mới đã được tạo ra từ bộ dữ liệu dự đoán giá nhà của bạn:

- **Area\_per\_Bedroom:** Diện tích trên mỗi phòng ngủ.
- **Bedrooms\_Gymnasium:** Tương tác giữa số phòng ngủ và sự hiện diện của phòng gym.
- **Log\_Price:** Biến đổi log của giá nhà.
- **Root\_Area:** Biến đổi căn bậc hai của diện tích.
- **Location\_Gymnasium:** Kết hợp giữa vị trí và sự hiện diện của phòng gym.

Bảng dữ liệu mới với các biến này đã được hiển thị. Nếu bạn cần thực hiện thêm các phân tích hoặc xây dựng mô hình dựa trên các biến mới này, hãy cho tôi biết nhé!



Các biến khác thì sao có thể kết hợp được không nếu có bạn hãy kết hợp lại nếu không bạn hãy chỉ ra lý do tại sao không



Chúng ta có thể xem xét khả năng kết hợp các biến khác để tạo ra các tính năng mới. Một số biến có thể không phù hợp để kết hợp do không có ý nghĩa thực tiễn hoặc không tạo ra thông tin hữu ích cho mô hình. Dưới đây là các bước để xem xét việc kết hợp các biến khác:

## 1. Kết hợp các biến số học khác

Chúng ta có thể tạo ra các biến số học mới từ các biến hiện có.

## 2. Tạo các biến tương tác khác

Chúng ta có thể tạo ra các biến tương tác mới từ các biến hiện có.

## 3. Kết hợp các biến phân loại khác

Chúng ta có thể kết hợp các biến phân loại khác để tạo ra các biến phân loại mới.

## Thực hiện kết hợp các biến khác

### Kết hợp các biến số học

```
python
```

```
# Kết hợp các biến số học khác: MaintenanceStaff và CarParking
housing_data['Staff_per_CarParking'] = housing_data['MaintenanceStaff'] /
(housing_data['CarParking'] + 1)

# Kết hợp các biến số học khác: SwimmingPool và Gymnasium
housing_data['Pool_Gym'] = housing_data['SwimmingPool'] + housing_data['Gymnasium']

# Kết hợp các biến số học khác: PowerBackup và 24X7Security
housing_data['Power_Security'] = housing_data['PowerBackup'] + housing_data['24X7Security']
```

### Tạo các biến tương tác khác

```
python
```

```
# Tạo các biến tương tác khác: SportsFacility và ClubHouse
housing_data['Sports_Club'] = housing_data['SportsFacility'] * housing_data['ClubHouse']

# Tạo các biến tương tác khác: Intercom và VaastuCompliant
housing_data['Intercom_Vaastu'] = housing_data['Intercom'] * housing_data['VaastuCompliant']
```

## Kết hợp các biến phân loại khác

python

```
# Kết hợp các biến phân loại khác: Location và School
housing_data['Location_School'] = housing_data['Location'] + "_" +
housing_data['School'].astype(str)

# Kết hợp các biến phân loại khác: Location và ShoppingMall
housing_data['Location_ShoppingMall'] = housing_data['Location'] + "_" +
housing_data['ShoppingMall'].astype(str)
```

## Kiểm tra các biến mới được tạo ra

python

```
# Kiểm tra các biến mới được tạo ra
new_combined_features = housing_data[['Staff_per_CarParking', 'Pool_Gym', 'Power_Security',
'Sports_Club',
                                     'Intercom_Vaastu', 'Location_School',
'Location_ShoppingMall']]

import ace_tools as tools; tools.display_dataframe_to_user(name="Housing Data with Combined
Features", dataframe=new_combined_features)

new_combined_features.head()
```

Chúng ta sẽ thực hiện các bước trên để tạo ra các biến kết hợp mới và kiểm tra kết quả.

## Housing Data with Combined Features

Unable to display visualization

Dưới đây là các biến mới đã được tạo ra từ việc kết hợp các biến trong bộ dữ liệu của bạn:

1. **Staff\_per\_CarParking:** Số nhân viên bảo trì trên mỗi chỗ đỗ xe.
2. **Pool\_Gym:** Tổng hợp giữa sự hiện diện của hồ bơi và phòng gym.
3. **Power\_Security:** Tổng hợp giữa sự hiện diện của nguồn điện dự phòng và bảo vệ 24/7.
4. **Sports\_Club:** Tương tác giữa sự hiện diện của cơ sở thể thao và clubhouse.
5. **Intercom\_Vaastu:** Tương tác giữa sự hiện diện của intercom và tuân thủ Vaastu.
6. **Location\_School:** Kết hợp giữa vị trí và sự hiện diện của trường học.
7. **Location\_ShoppingMall:** Kết hợp giữa vị trí và sự hiện diện của trung tâm mua sắm.

Bảng dữ liệu mới với các biến này đã được hiển thị. Những biến mới này có thể giúp cung cấp thêm thông tin cho mô hình và cải thiện hiệu suất dự đoán. Nếu bạn cần thực hiện thêm các phân tích hoặc xây dựng mô hình dựa trên các biến mới này, hãy cho tôi biết nhé!



Còn gì nữa không hay hết rồi bạn hãy kiểm tra lại kỹ một lần nữa từng bước một



Chắc chắn rồi! Chúng ta sẽ kiểm tra kỹ lưỡng từng bước một để đảm bảo không bỏ sót bất kỳ khả năng kết hợp biến nào. Chúng ta sẽ xem xét từng nhóm biến và đánh giá khả năng kết hợp thêm biến.

## Kiểm tra các biến số học khác

- **Số nhân viên bảo trì (MaintenanceStaff) và số chỗ đỗ xe (CarParking):** Đã kết hợp thành **Staff\_per\_CarParking**.
- **Hồ bơi (SwimmingPool) và phòng gym (Gymnasium):** Đã kết hợp thành **Pool\_Gym**.
- **Nguồn điện dự phòng (PowerBackup) và bảo vệ 24/7 (24X7Security):** Đã kết hợp thành **Power\_Security**.

## Kiểm tra các biến tương tác khác

- **Cơ sở thể thao (SportsFacility) và clubhouse (ClubHouse):** Đã kết hợp thành **Sports\_Club**.
- **Intercom và tuân thủ Vaastu (VaastuCompliant):** Đã kết hợp thành **Intercom\_Vaastu**.

## Kiểm tra các biến phân loại khác

- **Vị trí (Location) và trường học (School):** Đã kết hợp thành **Location\_School**.
- **Vị trí (Location) và trung tâm mua sắm (ShoppingMall):** Đã kết hợp thành **Location\_ShoppingMall**.

## Kiểm tra các biến khác

Chúng ta sẽ kiểm tra các biến khác để đảm bảo không bỏ sót bất kỳ khả năng kết hợp nào.

## Kết hợp các biến số học khác

python

```
# Kết hợp các biến số học khác: ATM và CarParking
housing_data['ATM_per_CarParking'] = housing_data['ATM'] / (housing_data['CarParking'] + 1)

# Kết hợp các biến số học khác: SportsFacility và SwimmingPool
housing_data['Sports_Pool'] = housing_data['SportsFacility'] + housing_data['SwimmingPool']

# Kiểm tra các biến mới được tạo ra
additional_numeric_features = housing_data[['ATM_per_CarParking', 'Sports_Pool']]
```

## Tạo các biến tương tác khác

python

```
# Tạo các biến tương tác khác: ClubHouse và Gymnasium
housing_data['Club_Gym'] = housing_data['ClubHouse'] * housing_data['Gymnasium']

# Kiểm tra các biến mới được tạo ra
additional_interaction_features = housing_data[['Club_Gym']]
```

## Kết hợp các biến phân loại khác

python

```
# Kết hợp các biến phân loại khác: Location và ATM
housing_data['Location_ATM'] = housing_data['Location'] + "_" +
housing_data['ATM'].astype(str)

# Kết hợp các biến phân loại khác: Location và Cafeteria
housing_data['Location_Cafeteria'] = housing_data['Location'] + "_" +
housing_data['Cafeteria'].astype(str)

# Kiểm tra các biến mới được tạo ra
additional_categorical_features = housing_data[['Location_ATM', 'Location_Cafeteria']]
```

## Kiểm tra các biến mới được tạo ra

python

```
# Kiểm tra các biến mới được tạo ra
all_combined_features = housing_data[['Staff_per_CarParking', 'Pool_Gym', 'Power_Security',
'Sports_Club',
                                'Intercom_Vaastu', 'Location_School',
'Location_ShoppingMall',
                                'ATM_per_CarParking', 'Sports_Pool', 'Club_Gym',
                                'Location_ATM', 'Location_Cafeteria']]

import ace_tools as tools; tools.display_dataframe_to_user(name="All Combined Features",
dataframe=all_combined_features)

all_combined_features.head()
```

Chúng ta sẽ thực hiện các bước này để đảm bảo đã xem xét tất cả khả năng kết hợp biến.

### All Combined Features

Staff_per_CarParking	Pool_Gym	Power_Security	Sports_Club
0.0	2	2	0
0.9	18	18	81
0.9	18	18	81
0.9	18	18	81
0.9	18	18	81
1.0	2	2	1
0.0	1	1	1

Dưới đây là danh sách tất cả các biến kết hợp mới được tạo ra từ bộ dữ liệu của bạn:

1. **Staff\_per\_CarParking:** Số nhân viên bảo trì trên mỗi chỗ đỗ xe.
2. **Pool\_Gym:** Tổng hợp giữa sự hiện diện của hồ bơi và phòng gym.
3. **Power\_Security:** Tổng hợp giữa sự hiện diện của nguồn điện dự phòng và bảo vệ 24/7.
4. **Sports\_Club:** Tương tác giữa sự hiện diện của cơ sở thể thao và clubhouse.



5. **Intercom\_Vaastu**: Tương tác giữa sự hiện diện của intercom và tuân thủ Vaastu.
6. **Location\_School**: Kết hợp giữa vị trí và sự hiện diện của trường học.
7. **Location\_ShoppingMall**: Kết hợp giữa vị trí và sự hiện diện của trung tâm mua sắm.
8. **ATM\_per\_CarParking**: Số lượng máy ATM trên mỗi chỗ đỗ xe.
9. **Sports\_Pool**: Tổng hợp giữa sự hiện diện của cơ sở thể thao và hồ bơi.
10. **Club\_Gym**: Tương tác giữa sự hiện diện của clubhouse và phòng gym.
11. **Location\_ATM**: Kết hợp giữa vị trí và sự hiện diện của máy ATM.
12. **Location\_Cafeteria**: Kết hợp giữa vị trí và sự hiện diện của căn tin.

Như vậy, chúng ta đã kiểm tra kỹ và tạo ra tất cả các biến kết hợp có khả năng cải thiện mô hình. Nếu bạn cần thực hiện thêm các phân tích hoặc xây dựng mô hình dựa trên các biến mới này, hãy cho tôi biết nhé!



Tôi có một ý chưa được rõ lắm là khi tạo ra các biến này rồi thì chúng ta sẽ xóa các biến gốc đi hay chúng ta giữ lại các biến gốc và thêm biến mới vào



Việc giữ lại hay xóa bỏ các biến gốc sau khi tạo ra các biến mới phụ thuộc vào chiến lược cụ thể mà bạn sử dụng và bản chất của mô hình mà bạn đang xây dựng. Dưới đây là một số hướng dẫn chung:

## 1. Giữ lại cả biến gốc và biến mới

- **Ưu điểm**: Mô hình có thêm nhiều thông tin từ cả biến gốc và biến mới, có thể giúp cải thiện hiệu suất dự đoán.
- **Nhược điểm**: Có thể dẫn đến đa cộng tuyến (multicollinearity), nơi các biến độc lập có mối tương quan cao với nhau, gây khó khăn cho mô hình trong việc ước lượng chính xác.

## 2. Xóa bỏ biến gốc và chỉ giữ lại biến mới

- **Ưu điểm**: Giảm nguy cơ đa cộng tuyến, đơn giản hóa mô hình, và giảm số lượng biến, có thể giúp mô hình dễ hiểu hơn và nhanh hơn.
- **Nhược điểm**: Có thể mất thông tin giá trị nếu các biến gốc cung cấp thông tin quan trọng mà các biến mới không thể thay thế hoàn toàn.

## 3. Giữ lại một phần các biến gốc

- **Ưu điểm**: Kết hợp lợi ích của cả hai phương pháp trên, giữ lại những biến gốc thực sự quan trọng trong khi thêm vào các biến mới để tăng cường thông tin.
- **Nhược điểm**: Cần phải có sự lựa chọn cẩn thận về biến nào cần giữ lại và biến nào cần xóa bỏ.

## Đề xuất cụ thể

Một cách tiếp cận phổ biến là bắt đầu bằng cách giữ lại tất cả các biến, sau đó sử dụng các kỹ thuật lựa chọn biến (feature selection) để xác định những biến nào thực sự quan trọng đối với mô hình của bạn. Các kỹ thuật như Lasso Regression, Random Forest Feature Importance, và Principal Component Analysis (PCA) có thể giúp xác định biến nào nên giữ lại.

## Ví dụ: Giữ lại tất cả các biến và sử dụng Lasso Regression để lựa chọn biến

python

```
from sklearn.linear_model import LassoCV
from sklearn.model_selection import train_test_split

# Chuẩn bị dữ liệu
X = housing_data.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Huấn luyện mô hình LassoCV để lựa chọn biến
lasso = LassoCV(cv=5, random_state=42)
lasso.fit(X_train, y_train)

# Lựa chọn các biến quan trọng
importance = np.abs(lasso.coef_)
feature_names = X.columns
important_features = feature_names[importance > 0]

# Xem các biến quan trọng
important_features
```

Hãy cho tôi biết nếu bạn muốn thực hiện phương pháp cụ thể nào hoặc cần thêm thông tin chi tiết hơn.



Tương tự trên bạn hãy làm thêm 2 hàm lần lượt cho random forest feature importance và PCA cho việc chọn feature quan trọng



Tôi đã quên nhập các thư viện cần thiết cho việc chia dữ liệu và tính toán feature importance. Tôi sẽ bổ sung chúng và thực hiện lại.

python

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Chuẩn bị dữ liệu
X = housing_data.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
```

```

indices = np.argsort(importances[::-1])
return X_train.columns[indices], importances[indices]

# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng Random Forest
important_features_rf, importances_rf = select_features_random_forest(X_train, y_train)

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train, n_components=10)

# Hiện thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
import ace_tools as tools; tools.display_dataframe_to_user(name="Random Forest Feature
Importance", dataframe=rf_feature_importance)

# Hiện thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
tools.display_dataframe_to_user(name="PCA Feature Importance",
dataframe=pca_feature_importance)

# Hiện thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
tools.display_dataframe_to_user(name="PCA Explained Variance",
dataframe=explained_variance_df)

rf_feature_importance.head(), pca_feature_importance.head(), explained_variance_df.head()

```

Hãy thực hiện lại các bước trên.

Có vẻ như có vấn đề với dữ liệu phân loại (biến "Location" và các biến phân loại kết hợp). Chúng ta cần mã hóa các biến phân loại này trước khi sử dụng chúng trong các mô hình như Random Forest và PCA.

## Chuẩn bị dữ liệu

1. Mã hóa biến phân loại: Sử dụng one-hot encoding cho các biến phân loại.

## Thực hiện mã hóa và lại các bước lựa chọn biến quan trọng

python

```

from sklearn.preprocessing import OneHotEncoder

# Mã hóa các biến phân loại
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
X_encoded = pd.get_dummies(housing_data, columns=categorical_columns, drop_first=True)

# Chuẩn bị lại dữ liệu
X = X_encoded.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)

```

```

rf.fit(X_train, y_train)
importances = rf.feature_importances_
indices = np.argsort(importances)[::-1]
return X_train.columns[indices], importances[indices]

# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng Random Forest
important_features_rf, importances_rf = select_features_random_forest(X_train, y_train)

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train, n_components=10)

# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
import ace_tools as tools; tools.display_dataframe_to_user(name="Random Forest Feature
Importance", dataframe=rf_feature_importance)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
tools.display_dataframe_to_user(name="PCA Feature Importance",
dataframe=pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
tools.display_dataframe_to_user(name="PCA Explained Variance",
dataframe=explained_variance_df)

rf_feature_importance.head(), pca_feature_importance.head(), explained_variance_df.head()

```

Hãy thực hiện lại các bước này với việc mã hóa các biến phân loại trước.



Tương tự trên bạn hãy làm thêm 2 hàm lần lượt cho random forest feature importance và PCA cho việc chọn feature quan trọng



Dưới đây là các hàm để lựa chọn các biến quan trọng sử dụng Random Forest Feature Importance và PCA:

## Hàm chọn biến quan trọng bằng Random Forest Feature Importance

python

```

from sklearn.ensemble import RandomForestRegressor
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
import numpy as np
import pandas as pd

```

```
# Chuẩn bị lại dữ liệu
X_encoded = pd.get_dummies(housing_data, columns=categorical_columns, drop_first=True)
X = X_encoded.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
    indices = np.argsort(importances)[::-1]
    return X_train.columns[indices], importances[indices]

# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng Random Forest
important_features_rf, importances_rf = select_features_random_forest(X_train, y_train)

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train, n_components=10)

# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
import ace_tools as tools; tools.display_dataframe_to_user(name="Random Forest Feature
Importance", dataframe=rf_feature_importance)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
tools.display_dataframe_to_user(name="PCA Feature Importance",
dataframe=pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
tools.display_dataframe_to_user(name="PCA Explained Variance",
dataframe=explained_variance_df)

rf_feature_importance.head(), pca_feature_importance.head(), explained_variance_df.head()
```

Tôi sẽ thực hiện lại các bước này để đảm bảo hàm hoạt động chính xác.

Có vẻ như đã xảy ra lỗi trong quá trình thực hiện. Tôi sẽ tách riêng từng phần để kiểm tra và xử lý lỗi.

## Kiểm tra Random Forest Feature Importance

python

```
# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
    indices = np.argsort(importances)[::-1]
    return X_train.columns[indices], importances[indices]

# Lựa chọn biến quan trọng bằng Random Forest
important_features_rf, importances_rf = select_features_random_forest(X_train, y_train)
```

```
# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
tools.display_dataframe_to_user(name="Random Forest Feature Importance",
dataframe=rf_feature_importance)

rf_feature_importance.head()
```

## Kiểm tra PCA Feature Importance

python

```
# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train, n_components=10)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
tools.display_dataframe_to_user(name="PCA Feature Importance",
dataframe=pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
tools.display_dataframe_to_user(name="PCA Explained Variance",
dataframe=explained_variance_df)

pca_feature_importance.head(), explained_variance_df.head()
```

Hãy thử thực hiện lại từng phần này để xác định và xử lý lỗi.

Có vẻ như vấn đề vẫn chưa được giải quyết. Tôi sẽ tiếp tục kiểm tra từng bước để xác định nguồn gốc của lỗi.

## Kiểm tra bước mã hóa và chuẩn bị dữ liệu

python

```
# Mã hóa các biến phân loại
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
X_encoded = pd.get_dummies(housing_data, columns=categorical_columns, drop_first=True)

# Chuẩn bị lại dữ liệu
X = X_encoded.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

X_train.head(), y_train.head()
```

Hãy kiểm tra kết quả của bước này trước khi tiếp tục với các bước tiếp theo.

Có vẻ như vấn đề vẫn chưa được giải quyết. Tôi sẽ tiếp tục kiểm tra từng bước để xác định nguồn gốc của lỗi.

## Kiểm tra bước mã hóa và chuẩn bị dữ liệu

python

```
# Mã hóa các biến phân loại
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
                       'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
X_encoded = pd.get_dummies(housing_data, columns=categorical_columns, drop_first=True)

# Chuẩn bị lại dữ liệu
X = X_encoded.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

X_train.head(), y_train.head()
```

Hãy kiểm tra kết quả của bước này trước khi tiếp tục với các bước tiếp theo.

Có vẻ như vẫn có vấn đề trong việc mã hóa và chuẩn bị dữ liệu. Để tiếp tục phân tích và tránh các lỗi hệ thống, tôi sẽ thử tiếp cận bằng cách lưu trữ dữ liệu đã mã hóa và sau đó tiếp tục với các bước tiếp theo.

## Mã hóa dữ liệu và lưu trữ

Trước tiên, chúng ta sẽ mã hóa dữ liệu và lưu kết quả vào một tệp CSV, sau đó tiếp tục phân tích từ tệp này.

python

```
# Mã hóa các biến phân loại
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
                       'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
X_encoded = pd.get_dummies(housing_data, columns=categorical_columns, drop_first=True)

# Chuẩn bị lại dữ liệu
X_encoded['Price'] = housing_data['Price']

# Lưu dữ liệu đã mã hóa vào tệp CSV
encoded_file_path = '/mnt/data/housing_encoded.csv'
X_encoded.to_csv(encoded_file_path, index=False)

encoded_file_path
```

Chúng ta sẽ thực hiện bước này và tiếp tục từ tệp đã mã hóa.

Có vẻ như chúng ta vẫn gặp vấn đề khi xử lý dữ liệu trực tiếp trên nền tảng hiện tại. Để khắc phục điều này, tôi sẽ hướng dẫn bạn thực hiện các bước mã hóa và lựa chọn biến quan trọng trực tiếp trên môi trường của bạn.

# Hướng dẫn mã hóa và lựa chọn biến quan trọng

## Bước 1: Mã hóa dữ liệu

### 1. Mã hóa các biến phân loại sử dụng `pd.get\_dummies`:

python

```
import pandas as pd

# Giả sử `housing_data` là DataFrame chứa dữ liệu của bạn
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
                        'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
X_encoded = pd.get_dummies(housing_data, columns=categorical_columns, drop_first=True)

# Thêm cột giá nhà lại vào dữ liệu đã mã hóa
X_encoded['Price'] = housing_data['Price']
```

### 2. Lưu dữ liệu đã mã hóa vào tệp CSV (nếu cần thiết):

python

```
encoded_file_path = 'housing_encoded.csv'
X_encoded.to_csv(encoded_file_path, index=False)
```

## Bước 2: Lựa chọn biến quan trọng

### 1. Sử dụng Random Forest Feature Importance:

python

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
import numpy as np
import pandas as pd

# Chuẩn bị dữ liệu
X = X_encoded.drop(columns=['Price']) # Bỏ cột mục tiêu
y = X_encoded['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=42)

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
    indices = np.argsort(importances)[::-1]
    return X_train.columns[indices], importances[indices]

# Lựa chọn biến quan trọng bằng Random Forest
important_features_rf, importances_rf = select_features_random_forest(X_train, y_train)

# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
```



```
importances_rf})
print(rf_feature_importance)
```

## 2. Sử dụng PCA:

python

```
from sklearn.decomposition import PCA

# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train, n_components=10)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
print(pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
print(explained_variance_df)
```

Hãy thực hiện các bước này trong môi trường của bạn để tiếp tục phân tích và lựa chọn các biến quan trọng. Nếu bạn gặp bất kỳ vấn đề gì hoặc cần thêm trợ giúp, hãy cho tôi biết!



tôi chưa transform data ngay tại bước kết hợp xong trước khi đưa vào tìm important feature hãy giúp tôi viết một function transform data phù hợp



Chắc chắn rồi! Dưới đây là một hàm để chuẩn bị và chuyển đổi dữ liệu bao gồm các bước mã hóa, chia dữ liệu, và chuẩn hóa dữ liệu số học.

### Hàm chuẩn bị và chuyển đổi dữ liệu

python

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import pandas as pd

def prepare_transform_data(data, target_column, categorical_columns):
    # 1. Mã hóa các biến phân loại
    data_encoded = pd.get_dummies(data, columns=categorical_columns, drop_first=True)

    # 2. Chia dữ liệu thành tập đặc trưng (features) và mục tiêu (target)
```

```
X = data_encoded.drop(columns=[target_column])
y = data_encoded[target_column]

# 3. Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# 4. Chuẩn hóa các biến số học (scaling)
scaler = StandardScaler()
numeric_columns = X_train.select_dtypes(include=[np.number]).columns
X_train[numeric_columns] = scaler.fit_transform(X_train[numeric_columns])
X_test[numeric_columns] = scaler.transform(X_test[numeric_columns])

return X_train, X_test, y_train, y_test, scaler

# Sử dụng hàm này để chuẩn bị và chuyển đổi dữ liệu
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
target_column = 'Price'
X_train, X_test, y_train, y_test, scaler = prepare_transform_data(housing_data,
target_column, categorical_columns)
```

Hàm này thực hiện các bước sau:

1. **Mã hóa các biến phân loại:** Sử dụng `pd.get\_dummies` để mã hóa các biến phân loại.
2. **Chia dữ liệu thành tập đặc trưng và mục tiêu:** Chia dữ liệu thành `X` (features) và `y` (target).
3. **Chia dữ liệu thành tập huấn luyện và tập kiểm tra:** Sử dụng `train\_test\_split` để chia dữ liệu.
4. **Chuẩn hóa các biến số học:** Sử dụng `StandardScaler` để chuẩn hóa các biến số học.

Bây giờ, chúng ta có thể sử dụng hàm này để chuẩn bị dữ liệu trước khi tiến hành lựa chọn các biến quan trọng.

## Lựa chọn biến quan trọng

Sau khi dữ liệu đã được chuẩn bị và chuyển đổi, chúng ta có thể tiến hành lựa chọn các biến quan trọng.

### Lựa chọn biến quan trọng bằng Random Forest

python

```
from sklearn.ensemble import RandomForestRegressor

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
    indices = np.argsort(importances)[::-1]
    return X_train.columns[indices], importances[indices]

# Lựa chọn biến quan trọng bằng Random Forest
important_features_rf, importances_rf = select_features_random_forest(X_train, y_train)

# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
print(rf_feature_importance)
```

### Lựa chọn biến quan trọng bằng PCA

python

```
from sklearn.decomposition import PCA

# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train, n_components=10)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
print(pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
print(explained_variance_df)
```

Hãy thực hiện các bước này để chuẩn bị và chuyển đổi dữ liệu trước khi tiến hành lựa chọn các biến quan trọng.



tôi muốn transform cả cột non-numeric



Để chuẩn hóa cả các biến số học và các biến phân loại (non-numeric), chúng ta cần sử dụng các kỹ thuật mã hóa và chuẩn hóa phù hợp. Dưới đây là một hàm để chuẩn bị và chuyển đổi dữ liệu bao gồm mã hóa các biến phân loại và chuẩn hóa cả các biến số học và các biến đã được mã hóa.

## Hàm chuẩn bị và chuyển đổi dữ liệu

python

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
import pandas as pd
import numpy as np

def prepare_transform_data(data, target_column, categorical_columns):
    # 1. Tách các biến phân loại và số học
    numeric_columns = data.select_dtypes(include=[np.number]).columns.drop(target_column).tolist()
    all_columns = numeric_columns + categorical_columns

    # 2. Tạo bộ biến đổi (ColumnTransformer) cho cả biến số học và phân loại
    preprocessor = ColumnTransformer(
        transformers=[
            ('num', StandardScaler(), numeric_columns),
```

```

        ('cat', OneHotEncoder(drop='first'), categorical_columns)
    ])

    # 3. Chia dữ liệu thành tập đặc trưng (features) và mục tiêu (target)
    X = data[all_columns]
    y = data[target_column]

    # 4. Chia dữ liệu thành tập huấn luyện và tập kiểm tra
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

    # 5. Chuẩn hóa dữ liệu
    X_train_transformed = preprocessor.fit_transform(X_train)
    X_test_transformed = preprocessor.transform(X_test)

    return X_train_transformed, X_test_transformed, y_train, y_test, preprocessor

# Sử dụng hàm này để chuẩn bị và chuyển đổi dữ liệu
categorical_columns = ['Location', 'Location_Gymnasium', 'Location_School',
'Location_ShoppingMall', 'Location_ATM', 'Location_Cafeteria']
target_column = 'Price'
X_train_transformed, X_test_transformed, y_train, y_test, preprocessor =
prepare_transform_data(housing_data, target_column, categorical_columns)

```

## Lựa chọn biến quan trọng

Sau khi dữ liệu đã được chuẩn bị và chuyển đổi, chúng ta có thể tiến hành lựa chọn các biến quan trọng.

### Lựa chọn biến quan trọng bằng Random Forest

python

```

from sklearn.ensemble import RandomForestRegressor

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
    indices = np.argsort(importances)[::-1]
    return indices, importances[indices]

# Lựa chọn biến quan trọng bằng Random Forest
indices_rf, importances_rf = select_features_random_forest(X_train_transformed, y_train)

# Lấy tên các biến quan trọng
feature_names = preprocessor.named_transformers_['num'].get_feature_names_out().tolist() +
preprocessor.named_transformers_['cat'].get_feature_names_out(categorical_columns).tolist()
important_features_rf = [feature_names[i] for i in indices_rf]

# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
print(rf_feature_importance)

```

### Lựa chọn biến quan trọng bằng PCA

python

```

from sklearn.decomposition import PCA

# Hàm lựa chọn biến quan trọng bằng PCA

```

```
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
    return explained_variance, components

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train_transformed,
n_components=10)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=feature_names)
print(pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
print(explained_variance_df)
```

Hãy thực hiện các bước này để chuẩn bị và chuyển đổi dữ liệu trước khi tiến hành lựa chọn các biến quan trọng.



bạn hãy cho việc truyền categorical column một cách tự động giống như numeric columns



Chắc chắn rồi! Dưới đây là một phiên bản của hàm `prepare\_transform\_data` tự động xác định các cột phân loại và số học từ dữ liệu, thay vì truyền vào danh sách các cột phân loại.

## Hàm chuẩn bị và chuyển đổi dữ liệu tự động xác định các cột phân loại và số học

python

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer
import pandas as pd
import numpy as np

def prepare_transform_data(data, target_column):
    # Tự động xác định các cột số học và phân loại
    categorical_columns = data.select_dtypes(include=['object']).columns.tolist()
    numeric_columns = data.select_dtypes(include=[
np.number]).columns.drop(target_column).tolist()
    all_columns = numeric_columns + categorical_columns

    # Tạo bộ biến đổi (ColumnTransformer) cho cả biến số học và phân loại
    preprocessor = ColumnTransformer(
        transformers=[
            ('num', StandardScaler(), numeric_columns),
            ('cat', OneHotEncoder(drop='first'), categorical_columns)
        ])

    # Chia dữ liệu thành tập đặc trưng (features) và mục tiêu (target)
    X = data[all_columns]
    y = data[target_column]
```

```
# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Chuẩn hóa dữ liệu
X_train_transformed = preprocessor.fit_transform(X_train)
X_test_transformed = preprocessor.transform(X_test)

return X_train_transformed, X_test_transformed, y_train, y_test, preprocessor,
categorical_columns, numeric_columns

# Sử dụng hàm này để chuẩn bị và chuyển đổi dữ liệu
target_column = 'Price'
X_train_transformed, X_test_transformed, y_train, y_test, preprocessor, categorical_columns,
numeric_columns = prepare_transform_data(housing_data, target_column)
```

## Lựa chọn biến quan trọng

Sau khi dữ liệu đã được chuẩn bị và chuyển đổi, chúng ta có thể tiến hành lựa chọn các biến quan trọng.

### Lựa chọn biến quan trọng bằng Random Forest

python

```
from sklearn.ensemble import RandomForestRegressor

# Hàm lựa chọn biến quan trọng bằng Random Forest Feature Importance
def select_features_random_forest(X_train, y_train):
    rf = RandomForestRegressor(n_estimators=100, random_state=42)
    rf.fit(X_train, y_train)
    importances = rf.feature_importances_
    indices = np.argsort(importances)[::-1]
    return indices, importances[indices]

# Lựa chọn biến quan trọng bằng Random Forest
indices_rf, importances_rf = select_features_random_forest(X_train_transformed, y_train)

# Lấy tên các biến quan trọng
num_feature_names =
preprocessor.named_transformers_['num'].get_feature_names_out(numeric_columns).tolist()
cat_feature_names =
preprocessor.named_transformers_['cat'].get_feature_names_out(categorical_columns).tolist()
feature_names = num_feature_names + cat_feature_names
important_features_rf = [feature_names[i] for i in indices_rf]

# Hiển thị kết quả Random Forest Feature Importance
rf_feature_importance = pd.DataFrame({'Feature': important_features_rf, 'Importance':
importances_rf})
print(rf_feature_importance)
```

### Lựa chọn biến quan trọng bằng PCA

python

```
from sklearn.decomposition import PCA

# Hàm lựa chọn biến quan trọng bằng PCA
def select_features_pca(X_train, n_components=10):
    pca = PCA(n_components=n_components)
    pca.fit(X_train)
    explained_variance = pca.explained_variance_ratio_
    components = pca.components_
```

```

return explained_variance, components

# Lựa chọn biến quan trọng bằng PCA
explained_variance_pca, components_pca = select_features_pca(X_train_transformed,
n_components=10)

# Hiển thị kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=feature_names)
print(pca_feature_importance)

# Hiển thị phân trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
print(explained_variance_df)

```

Hãy thực hiện các bước này để chuẩn bị và chuyển đổi dữ liệu trước khi tiến hành lựa chọn các biến quan trọng.



bạn hãy làm từng bước chuyển đổi catecorical về numeric và sau đó scale dataframe



Chắc chắn rồi! Chúng ta sẽ thực hiện từng bước để chuyển đổi các cột phân loại thành số học và sau đó chuẩn hóa toàn bộ dataframe.

## Bước 1: Chuyển đổi các cột phân loại thành số học

Chúng ta sẽ sử dụng `OneHotEncoder` để mã hóa các cột phân loại.

python

```

import pandas as pd
import numpy as np
from sklearn.preprocessing import OneHotEncoder

# Tự động xác định các cột phân loại và số học
categorical_columns = housing_data.select_dtypes(include=['object']).columns.tolist()
numeric_columns = housing_data.select_dtypes(include=[np.number]).columns.tolist()

# Mã hóa các cột phân loại
encoder = OneHotEncoder(drop='first', sparse=False)
encoded_categorical_data = encoder.fit_transform(housing_data[categorical_columns])

# Chuyển đổi các cột mã hóa thành DataFrame
encoded_categorical_df = pd.DataFrame(encoded_categorical_data,
columns=encoder.get_feature_names_out(categorical_columns))

# Kết hợp các cột số học và các cột đã mã hóa
housing_data_encoded = pd.concat([housing_data[numeric_columns], encoded_categorical_df],
axis=1)

# Hiển thị DataFrame đã được mã hóa
import ace_tools as tools; tools.display_dataframe_to_user(name="Housing Data Encoded",
dataframe=housing_data_encoded)
housing_data_encoded.head()

```

## Bước 2: Chuẩn hóa các cột số học

Chúng ta sẽ sử dụng `StandardScaler` để chuẩn hóa các cột số học.

python

```
from sklearn.preprocessing import StandardScaler

# Tạo scaler
scaler = StandardScaler()

# Chuẩn hóa các cột số học
scaled_data = scaler.fit_transform(housing_data_encoded)

# Chuyển đổi dữ liệu đã chuẩn hóa thành DataFrame
housing_data_scaled = pd.DataFrame(scaled_data, columns=housing_data_encoded.columns)

# Hiển thị DataFrame đã được chuẩn hóa
tools.display_dataframe_to_user(name="Housing Data Scaled", dataframe=housing_data_scaled)
housing_data_scaled.head()
```

## Bước 3: Chia dữ liệu thành tập huấn luyện và tập kiểm tra

Chúng ta sẽ sử dụng `train\_test\_split` để chia dữ liệu.

python

```
from sklearn.model_selection import train_test_split

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X = housing_data_scaled.drop(columns=['Price']) # Bỏ cột mục tiêu
y = housing_data_scaled['Price']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

X_train.head(), y_train.head()
```

Chúng ta sẽ thực hiện từng bước này để đảm bảo quá trình chuyển đổi và chuẩn hóa dữ liệu được thực hiện chính xác.



giờ tôi có thêm data\_test là tập test với cùng số lượng cột tương tự bạn hãy giúp tôi one hot và scale cho cả train và data\_test



Chắc chắn rồi! Chúng ta sẽ thực hiện mã hóa và chuẩn hóa cho cả tập huấn luyện và tập kiểm tra bằng cách áp dụng cùng các phép biến đổi. Dưới đây là cách thực hiện:

### Chuẩn bị dữ liệu



Giả sử bạn có `data\_train` và `data\_test` là các DataFrame chứa dữ liệu huấn luyện và kiểm tra tương ứng.

## Bước 1: Chuyển đổi các cột phân loại thành số học

### 1. Mã hóa các cột phân loại sử dụng `OneHotEncoder`:

python

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import OneHotEncoder

# Giả sử data_train và data_test là DataFrame chứa dữ liệu của bạn
data_train = housing_data.copy() # Giả sử housing_data là data_train
data_test = data_train.copy() # Thay đổi bằng tập data_test thực tế của bạn

# Tự động xác định các cột phân loại và số học
categorical_columns = data_train.select_dtypes(include=['object']).columns.tolist()
numeric_columns = data_train.select_dtypes(include=[np.number]).columns.tolist()

# Mã hóa các cột phân loại
encoder = OneHotEncoder(drop='first', sparse=False)
encoded_categorical_data_train = encoder.fit_transform(data_train[categorical_columns])
encoded_categorical_data_test = encoder.transform(data_test[categorical_columns])

# Chuyển đổi các cột mã hóa thành DataFrame
encoded_categorical_df_train = pd.DataFrame(encoded_categorical_data_train,
columns=encoder.get_feature_names_out(categorical_columns))
encoded_categorical_df_test = pd.DataFrame(encoded_categorical_data_test,
columns=encoder.get_feature_names_out(categorical_columns))

# Kết hợp các cột số học và các cột đã mã hóa
data_train_encoded = pd.concat([data_train[numeric_columns], encoded_categorical_df_train],
axis=1)
data_test_encoded = pd.concat([data_test[numeric_columns], encoded_categorical_df_test],
axis=1)
```

## Bước 2: Chuẩn hóa các cột số học

### 2. Chuẩn hóa các cột số học sử dụng `StandardScaler`:

python

```
from sklearn.preprocessing import StandardScaler

# Tạo scaler
scaler = StandardScaler()

# Chuẩn hóa các cột số học cho tập huấn luyện
scaled_data_train = scaler.fit_transform(data_train_encoded)

# Áp dụng scaler lên tập kiểm tra
scaled_data_test = scaler.transform(data_test_encoded)

# Chuyển đổi dữ liệu đã chuẩn hóa thành DataFrame
data_train_scaled = pd.DataFrame(scaled_data_train, columns=data_train_encoded.columns)
data_test_scaled = pd.DataFrame(scaled_data_test, columns=data_test_encoded.columns)

# Hiển thị DataFrame đã được chuẩn hóa
import acetoools as tools; tools.display_dataframe_to_user(name="Train Data Scaled",
dataframe=data_train_scaled)
```

```
tools.display_dataframe_to_user(name="Test Data Scaled", dataframe=data_test_scaled)
data_train_scaled.head(), data_test_scaled.head()
```

## Chia dữ liệu thành tập huấn luyện và tập kiểm tra

### 3. Chia dữ liệu thành tập đặc trưng (features) và mục tiêu (target):

python

```
# Giả sử cột mục tiêu là 'Price'
target_column = 'Price'

# Chia dữ liệu thành tập đặc trưng và mục tiêu cho tập huấn luyện
X_train = data_train_scaled.drop(columns=[target_column])
y_train = data_train_scaled[target_column]

# Chia dữ liệu thành tập đặc trưng và mục tiêu cho tập kiểm tra
X_test = data_test_scaled.drop(columns=[target_column])
y_test = data_test_scaled[target_column]

X_train.head(), y_train.head(), X_test.head(), y_test.head()
```

Bây giờ, dữ liệu của bạn đã được mã hóa và chuẩn hóa cho cả tập huấn luyện và tập kiểm tra. Bạn có thể sử dụng các tập dữ liệu này để xây dựng và đánh giá mô hình.



tôi làm tương tự nhưng bị gặp lỗi sau với data test



Lỗi này xảy ra vì các giá trị trong cột phân loại của tập dữ liệu kiểm tra (`data\_test`) không khớp hoàn toàn với các giá trị của tập huấn luyện (`data\_train`). Để khắc phục vấn đề này, chúng ta có thể sử dụng tham số `handle\_unknown='ignore'` trong `OneHotEncoder` để bỏ qua các giá trị không xác định.

## Bước 1: Chuyển đổi các cột phân loại thành số học

python

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import OneHotEncoder

# Giả sử data_train và data_test là DataFrame chứa dữ liệu của bạn
data_train = housing_data.copy() # Giả sử housing_data là data_train
data_test = data_train.copy() # Thay đổi bằng tập data_test thực tế của bạn

# Tự động xác định các cột phân loại và số học
categorical_columns = data_train.select_dtypes(include=['object']).columns.tolist()
numeric_columns = data_train.select_dtypes(include=[np.number]).columns.tolist()

# Mã hóa các cột phân loại
encoder = OneHotEncoder(drop='first', sparse=False, handle_unknown='ignore')
```

```
encoded_categorical_data_train = encoder.fit_transform(data_train[categorical_columns])
encoded_categorical_data_test = encoder.transform(data_test[categorical_columns])

# Chuyển đổi các cột mã hóa thành DataFrame
encoded_categorical_df_train = pd.DataFrame(encoded_categorical_data_train,
columns=encoder.get_feature_names_out(categorical_columns))
encoded_categorical_df_test = pd.DataFrame(encoded_categorical_data_test,
columns=encoder.get_feature_names_out(categorical_columns))

# Kết hợp các cột số học và các cột đã mã hóa
data_train_encoded = pd.concat([data_train[numeric_columns], encoded_categorical_df_train],
axis=1)
data_test_encoded = pd.concat([data_test[numeric_columns], encoded_categorical_df_test],
axis=1)
```

## Bước 2: Chuẩn hóa các cột số học

python

```
from sklearn.preprocessing import StandardScaler

# Tạo scaler
scaler = StandardScaler()

# Chuẩn hóa các cột số học cho tập huấn luyện
scaled_data_train = scaler.fit_transform(data_train_encoded)

# Áp dụng scaler lên tập kiểm tra
scaled_data_test = scaler.transform(data_test_encoded)

# Chuyển đổi dữ liệu đã chuẩn hóa thành DataFrame
data_train_scaled = pd.DataFrame(scaled_data_train, columns=data_train_encoded.columns)
data_test_scaled = pd.DataFrame(scaled_data_test, columns=data_test_encoded.columns)

# Hiển thị DataFrame đã được chuẩn hóa
import ace_tools as tools; tools.display_dataframe_to_user(name="Train Data Scaled",
dataframe=data_train_scaled)
tools.display_dataframe_to_user(name="Test Data Scaled", dataframe=data_test_scaled)
data_train_scaled.head(), data_test_scaled.head()
```

## Chia dữ liệu thành tập huấn luyện và tập kiểm tra

python

```
# Giả sử cột mục tiêu là 'Price'
target_column = 'Price'

# Chia dữ liệu thành tập đặc trưng và mục tiêu cho tập huấn luyện
X_train = data_train_scaled.drop(columns=[target_column])
y_train = data_train_scaled[target_column]

# Chia dữ liệu thành tập đặc trưng và mục tiêu cho tập kiểm tra
X_test = data_test_scaled.drop(columns=[target_column])
y_test = data_test_scaled[target_column]

X_train.head(), y_train.head(), X_test.head(), y_test.head()
```

Bằng cách sử dụng `handle\_unknown='ignore'`, chúng ta có thể bỏ qua các giá trị không xác định trong tập kiểm tra, giúp giải quyết vấn đề gặp phải. Hãy thử thực hiện lại các bước này để đảm bảo quá trình chuyển đổi và chuẩn hóa dữ liệu được thực hiện chính xác.



tôi đã có feature của PCR và random forest vậy bây giờ tôi muốn train với linear regression và metric RMSE thì tôi sử dụng feature như thế nào để training



Chúng ta có thể thực hiện các bước để sử dụng các đặc trưng quan trọng từ PCA và Random Forest để huấn luyện một mô hình Linear Regression và tính toán chỉ số RMSE (Root Mean Squared Error). Dưới đây là các bước chi tiết để thực hiện điều này:

## Chuẩn bị Dữ liệu

1. Mã hóa và Chuẩn hóa Dữ liệu (như đã làm ở bước trước)
2. Chọn các đặc trưng quan trọng từ PCA và Random Forest

### Bước 1: Chọn các đặc trưng quan trọng từ PCA và Random Forest

python

```
# Giả sử explained_variance_pca và components_pca đã được tính toán trước đó
# Chọn các đặc trưng từ PCA
pca_num_components = 10
pca_features = [f'PC{i+1}' for i in range(pca_num_components)]
pca_feature_importance = pd.DataFrame(components_pca.T, columns=pca_features,
index=feature_names)

# Giả sử important_features_rf và importances_rf đã được tính toán trước đó
# Chọn các đặc trưng từ Random Forest
num_rf_features = 10 # Số lượng đặc trưng quan trọng cần chọn
important_features_rf_selected = important_features_rf[:num_rf_features]

# Kết hợp các đặc trưng quan trọng từ PCA và Random Forest
selected_features = list(set(pca_features + important_features_rf_selected))
```

### Bước 2: Huấn luyện mô hình Linear Regression với các đặc trưng quan trọng và tính toán RMSE

python

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import numpy as np

# Lấy các cột tương ứng với các đặc trưng quan trọng từ tập dữ liệu đã chuẩn hóa
X_train_selected = X_train[selected_features]
X_test_selected = X_test[selected_features]

# Huấn luyện mô hình Linear Regression
linear_regressor = LinearRegression()
linear_regressor.fit(X_train_selected, y_train)

# Dự đoán trên tập kiểm tra
y_pred = linear_regressor.predict(X_test_selected)
```

```
# Tính toán RMSE
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f'RMSE: {rmse}')
```

## Chi tiết từng bước

### 1. Chọn các đặc trưng từ PCA

python

```
pca_num_components = 10
pca_features = [f'PC{i+1}' for i in range(pca_num_components)]
pca_feature_importance = pd.DataFrame(components_pca.T, columns=pca_features,
index=feature_names)

# Chọn 10 đặc trưng quan trọng nhất từ PCA
top_pca_features = pca_feature_importance.columns[:pca_num_components].tolist()
```

### 2. Chọn các đặc trưng từ Random Forest

python

```
num_rf_features = 10 # Số lượng đặc trưng quan trọng cần chọn
important_features_rf_selected = important_features_rf[:num_rf_features]
```

### 3. Kết hợp các đặc trưng từ PCA và Random Forest

python

```
selected_features = list(set(top_pca_features + important_features_rf_selected))
```

### 4. Huấn luyện mô hình Linear Regression và tính toán RMSE

python

```
# Lấy các cột tương ứng với các đặc trưng quan trọng từ tập dữ liệu đã chuẩn hóa
X_train_selected = pd.DataFrame(X_train_transformed, columns=feature_names)
[selected_features]
X_test_selected = pd.DataFrame(X_test_transformed, columns=feature_names)[selected_features]

# Huấn luyện mô hình Linear Regression
linear_regressor = LinearRegression()
linear_regressor.fit(X_train_selected, y_train)

# Dự đoán trên tập kiểm tra
y_pred = linear_regressor.predict(X_test_selected)

# Tính toán RMSE
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f'RMSE: {rmse}')
```

Hãy thực hiện các bước này để huấn luyện mô hình Linear Regression với các đặc trưng quan trọng từ PCA và Random Forest và tính toán chỉ số RMSE.



hãy giúp tôi lấy feature PCA để huấn luyện linear regression với metric RMSE



Chúng ta sẽ lấy các thành phần chính từ PCA và sử dụng chúng để huấn luyện mô hình Linear Regression và tính toán chỉ số RMSE. Dưới đây là các bước chi tiết:

## Bước 1: Huấn luyện PCA và lấy các thành phần chính

python

```
from sklearn.decomposition import PCA

# Giả sử X_train_transformed và X_test_transformed đã được chuẩn bị từ trước
# Huấn luyện PCA
n_components = 10
pca = PCA(n_components=n_components)
X_train_pca = pca.fit_transform(X_train_transformed)
X_test_pca = pca.transform(X_test_transformed)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance = pca.explained_variance_ratio_
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance))], 'Explained Variance': explained_variance})
print(explained_variance_df)
```

## Bước 2: Huấn luyện mô hình Linear Regression và tính toán RMSE

python

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import numpy as np

# Huấn luyện mô hình Linear Regression với các thành phần chính từ PCA
linear_regressor = LinearRegression()
linear_regressor.fit(X_train_pca, y_train)

# Dự đoán trên tập kiểm tra
y_pred = linear_regressor.predict(X_test_pca)

# Tính toán RMSE
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f'RMSE: {rmse}')
```

## Tổng hợp mã nguồn đầy đủ

python

```
# Bước 1: Huấn luyện PCA và lấy các thành phần chính
from sklearn.decomposition import PCA
import pandas as pd

n_components = 10
pca = PCA(n_components=n_components)
X_train_pca = pca.fit_transform(X_train_transformed)
X_test_pca = pca.transform(X_test_transformed)

explained_variance = pca.explained_variance_ratio_
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance))], 'Explained Variance': explained_variance})
print(explained_variance_df)

# Bước 2: Huấn luyện mô hình Linear Regression và tính toán RMSE
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import numpy as np

linear_regressor = LinearRegression()
linear_regressor.fit(X_train_pca, y_train)

y_pred = linear_regressor.predict(X_test_pca)

rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f'RMSE: {rmse}')
```

Hãy thực hiện các bước này để lấy các thành phần chính từ PCA và sử dụng chúng để huấn luyện mô hình Linear Regression và tính toán chỉ số RMSE.



làm sao để tôi có thể lấy được feature từ kết quả pca



Để lấy các đặc trưng từ kết quả PCA và sử dụng chúng để huấn luyện mô hình Linear Regression, bạn có thể thực hiện các bước sau:

## Bước 1: Huấn luyện PCA và lấy các thành phần chính

python

```
from sklearn.decomposition import PCA
import pandas as pd

# Huấn luyện PCA
n_components = 10
pca = PCA(n_components=n_components)
X_train_pca = pca.fit_transform(X_train_transformed)
X_test_pca = pca.transform(X_test_transformed)

# Hiển thị phân trăm phương sai được giải thích bởi các thành phần chính
explained_variance = pca.explained_variance_ratio_
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
```

```
range(len(explained_variance))], 'Explained Variance': explained_variance})
print(explained_variance_df)
```

## Bước 2: Huấn luyện mô hình Linear Regression và tính toán RMSE

python

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import numpy as np

# Huấn luyện mô hình Linear Regression với các thành phần chính từ PCA
linear_regressor = LinearRegression()
linear_regressor.fit(X_train_pca, y_train)

# Dự đoán trên tập kiểm tra
y_pred = linear_regressor.predict(X_test_pca)

# Tính toán RMSE
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f'RMSE: {rmse}')
```

## Tổng hợp mã nguồn đầy đủ

python

```
# Bước 1: Huấn luyện PCA và lấy các thành phần chính
from sklearn.decomposition import PCA
import pandas as pd

n_components = 10
pca = PCA(n_components=n_components)
X_train_pca = pca.fit_transform(X_train_transformed)
X_test_pca = pca.transform(X_test_transformed)

explained_variance = pca.explained_variance_ratio_
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance))], 'Explained Variance': explained_variance})
print(explained_variance_df)

# Bước 2: Huấn luyện mô hình Linear Regression và tính toán RMSE
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import numpy as np

linear_regressor = LinearRegression()
linear_regressor.fit(X_train_pca, y_train)

y_pred = linear_regressor.predict(X_test_pca)

rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f'RMSE: {rmse}')
```

Nếu bạn muốn lấy tên của các đặc trưng quan trọng từ PCA, bạn có thể sử dụng các thành phần chính (principal components) đã được tạo ra trong bước huấn luyện PCA.

## Hiển thị các đặc trưng quan trọng từ PCA



python

```
# Giả sử bạn đã có `components_pca` từ kết quả PCA
pca_feature_importance = pd.DataFrame(components_pca.T, columns=[f'PC{i+1}' for i in
range(components_pca.shape[0])], index=X_train.columns)
print(pca_feature_importance)

# Hiển thị phần trăm phương sai được giải thích bởi các thành phần chính
explained_variance_df = pd.DataFrame({'Principal Component': [f'PC{i+1}' for i in
range(len(explained_variance_pca))], 'Explained Variance': explained_variance_pca})
print(explained_variance_df)
```

Bằng cách này, bạn có thể lấy các đặc trưng từ kết quả PCA và sử dụng chúng để huấn luyện mô hình Linear Regression và tính toán chỉ số RMSE.



bạn hãy giúp tôi viết một hàm covert lại ban đầu sau khi dùng log housing\_data['Log\_Price'] = housing\_data['Price'].apply(lambda x: np.log(x) if x > 0 else 0)



Chắc chắn rồi! Để chuyển đổi lại giá trị `Log\_Price` về giá trị ban đầu của `Price`, bạn có thể sử dụng hàm mũ (`exp`) của numpy. Dưới đây là cách thực hiện việc này.

### Hàm chuyển đổi giá trị `Log\_Price` về `Price`

python

```
import numpy as np

def revert_log_price(log_prices):
    """
    Chuyển đổi giá trị Log_Price về giá trị ban đầu của Price.

    Args:
        log_prices (pd.Series): Series chứa các giá trị Log_Price.

    Returns:
        pd.Series: Series chứa các giá trị Price đã được chuyển đổi lại.
    """
    return log_prices.apply(lambda x: np.exp(x) if x != 0 else 0)

# Sử dụng hàm để chuyển đổi Log_Price về Price
housing_data['Reverted_Price'] = revert_log_price(housing_data['Log_Price'])

# Kiểm tra kết quả
housing_data[['Price', 'Log_Price', 'Reverted_Price']].head()
```

### Tổng hợp mã nguồn

Dưới đây là mã nguồn đầy đủ để chuyển đổi `Log\_Price` về `Price`:

python

```
import numpy as np
import pandas as pd

# Giả sử bạn đã có housing_data với cột Log_Price
housing_data = pd.DataFrame({
    'Price': [100000, 200000, 300000, 400000, 500000],
    'Log_Price': [np.log(100000), np.log(200000), np.log(300000), np.log(400000),
np.log(500000)]
})

def revert_log_price(log_prices):
    """
    Chuyển đổi giá trị Log_Price về giá trị ban đầu của Price.

    Args:
    log_prices (pd.Series): Series chứa các giá trị Log_Price.

    Returns:
    pd.Series: Series chứa các giá trị Price đã được chuyển đổi lại.
    """
    return log_prices.apply(lambda x: np.exp(x) if x != 0 else 0)

# Sử dụng hàm để chuyển đổi Log_Price về Price
housing_data['Reverted_Price'] = revert_log_price(housing_data['Log_Price'])

# Kiểm tra kết quả
print(housing_data[['Price', 'Log_Price', 'Reverted_Price']])
```

Hãy chạy mã nguồn này để đảm bảo rằng giá trị `Log\_Price` được chuyển đổi lại thành công về giá trị ban đầu của `Price`. Nếu bạn cần bất kỳ điều chỉnh nào khác, hãy cho tôi biết!



Tôi đã hoàn thành nhưng bây giờ cần tốt hơn nữa. Chúng ta hãy bắt đầu với việc xử lý các điểm nhiễu trong dữ liệu sau khi kết hợp feature



Để xử lý các điểm nhiễu (outliers) trong dữ liệu sau khi đã kết hợp các đặc trưng, chúng ta có thể sử dụng một số phương pháp phổ biến như:

1. **Sử dụng IQR (Interquartile Range):** Loại bỏ các giá trị nằm ngoài khoảng  $Q1 - 1.5 * IQR$  và  $Q3 + 1.5 * IQR$ .
2. **Sử dụng Z-score:** Loại bỏ các giá trị có z-score lớn hơn một ngưỡng nhất định (ví dụ: 3 hoặc -3).
3. **Visualization:** Sử dụng các biểu đồ như Box Plot để trực quan hóa và loại bỏ các điểm nhiễu.

## Bước 1: Kết hợp các đặc trưng

Giả sử bạn đã thực hiện các bước kết hợp đặc trưng như đã đề cập ở các bước trước.

## Bước 2: Xử lý các điểm nhiễu bằng IQR

python

```
import pandas as pd
import numpy as np

# Hàm xử lý điểm nhiễu bằng IQR
def remove_outliers_iqr(data, column):
    Q1 = data[column].quantile(0.25)
    Q3 = data[column].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    return data[(data[column] >= lower_bound) & (data[column] <= upper_bound)]

# Giả sử bạn đã có DataFrame housing_data_encoded sau khi mã hóa và chuẩn hóa
housing_data_cleaned = housing_data_encoded.copy()

# Áp dụng hàm xử lý điểm nhiễu cho các cột số học
for column in housing_data_cleaned.select_dtypes(include=[np.number]).columns:
    housing_data_cleaned = remove_outliers_iqr(housing_data_cleaned, column)

# Kiểm tra kết quả
housing_data_cleaned.describe()
```

### Bước 3: Sử dụng Z-score để xử lý điểm nhiễu

python

```
from scipy import stats

# Hàm xử lý điểm nhiễu bằng Z-score
def remove_outliers_zscore(data, threshold=3):
    z_scores = stats.zscore(data.select_dtypes(include=[np.number]))
    abs_z_scores = np.abs(z_scores)
    return data[(abs_z_scores < threshold).all(axis=1)]

# Áp dụng hàm xử lý điểm nhiễu cho toàn bộ dữ liệu
housing_data_cleaned_zscore = remove_outliers_zscore(housing_data_encoded)

# Kiểm tra kết quả
housing_data_cleaned_zscore.describe()
```

### Bước 4: Visualization với Box Plot

python

```
import seaborn as sns
import matplotlib.pyplot as plt

# Vẽ Box Plot để trực quan hóa các điểm nhiễu
plt.figure(figsize=(15, 10))
sns.boxplot(data=housing_data_encoded.select_dtypes(include=[np.number]))
plt.xticks(rotation=90)
plt.title('Box Plot for Numerical Features')
plt.show()
```

### Kết hợp và Sử dụng dữ liệu đã xử lý

Sau khi loại bỏ các điểm nhiễu, bạn có thể tiếp tục sử dụng dữ liệu đã làm sạch để huấn luyện mô hình Linear Regression và các bước tiếp theo.

python

```
# Sử dụng dữ liệu đã làm sạch để tiếp tục phân tích và huấn luyện mô hình
X_cleaned = housing_data_cleaned.drop(columns=['Price']) # Bỏ cột mục tiêu
y_cleaned = housing_data_cleaned['Price']

# Chia dữ liệu thành tập huấn luyện và tập kiểm tra
X_train_cleaned, X_test_cleaned, y_train_cleaned, y_test_cleaned =
train_test_split(X_cleaned, y_cleaned, test_size=0.2, random_state=42)

# Huấn luyện mô hình Linear Regression với dữ liệu đã làm sạch
linear_regressor_cleaned = LinearRegression()
linear_regressor_cleaned.fit(X_train_cleaned, y_train_cleaned)

# Dự đoán trên tập kiểm tra đã làm sạch
y_pred_cleaned = linear_regressor_cleaned.predict(X_test_cleaned)

# Tính toán RMSE với dữ liệu đã làm sạch
rmse_cleaned = np.sqrt(mean_squared_error(y_test_cleaned, y_pred_cleaned))
print(f'RMSE after cleaning: {rmse_cleaned}')
```

Hãy thực hiện các bước này để loại bỏ các điểm nhiễu và cải thiện mô hình của bạn. Nếu bạn cần bất kỳ điều chỉnh nào khác, hãy cho tôi biết!

ChatGPT can make mistakes. Check important info.