



Bài tập tuần 8

Hệ hỗ trợ quyết định

Chủ đề 13

Dự đoán giá nhà ở Boston

Giảng viên hướng dẫn:

Mã lớp:

Sinh viên thực hiện:

MSSV:

TS. Trần Ngọc Thăng

158242

Phan Thu Trang

20227156

Mục lục

Mục lục	2
1 Phát biểu bài toán	3
1.1 Mô tả bài toán	3
1.2 Đầu vào	3
1.3 Đầu ra	4
1.4 Yêu cầu xử lý	4
2 Tiền xử lý dữ liệu	5
2.1 Thu nhập Dữ Liệu	5
2.2 Đánh Nhãn Dữ Liệu	5
2.3 Thống Kê Dữ Liệu Mẫu	6
2.4 Tiền Xử Lý Dữ Liệu	10
3 Tạo và luyện và đánh giá mô hình	13
3.1 Tạo Mô Hình và Các Quyết Định	13
3.2 Mô tả quá trình luyện mô hình	14
3.3 Mô tả điều kiện dừng	16
3.4 Đánh giá quá trình luyện và hiệu chỉnh siêu tham số	18
3.5 Lựa chọn các số đo đánh giá mô hình	20
3.6 Đánh giá mô hình với dữ liệu test	21
4 Ứng dụng mô hình dự đoán giá nhà	24
4.1 Mô tả ứng dụng mô hình	24
4.2 Diễn giải kết quả	27
5 Kết luận	30
5.1 Ưu nhược điểm của cách tiếp cận	30
5.2 Khả năng cải tiến kết quả trong tương lai	31
Checklist	33

Phần 1. Phát biểu bài toán

1.1. Mô tả bài toán

Dựa trên bộ dữ liệu chứa thông tin về các giao dịch bất động sản (nhà ở dành cho một gia đình) đã được bán ở Miami, mục tiêu là xây dựng một mô hình học máy để dự đoán giá bán của ngôi nhà (**SALE_PRC**) dựa trên các thuộc tính liên quan như diện tích đất, diện tích sống, khoảng cách đến các tiện ích, tuổi của ngôi nhà, chất lượng cấu trúc, v.v.

1.2. Đầu vào

Giới thiệu bộ dữ liệu

Bộ dữ liệu Miami Housing Dataset là một tập hợp dữ liệu chi tiết về các ngôi nhà dành cho một gia đình được bán ở khu vực Miami. Dưới đây là một mô tả chi tiết về bộ dữ liệu này:

- Nguồn dữ liệu: Kaggle.
- Số lượng mẫu: 13.932 ngôi nhà.
- Kích thước dữ liệu 1.64 MB.
- Năm dữ liệu: 2016.

Các thuộc tính trong bộ dữ liệu

Bộ dữ liệu ở định dạng CSV với các cột thông tin sau:

- **LATITUDE**: Vĩ độ của ngôi nhà.
- **LONGITUDE**: Kinh độ của ngôi nhà.
- **PARCELNO**: Số hiệu lô đất (có thể dùng làm định danh duy nhất).
- **SALE_PRC**: Giá bán của ngôi nhà (biến mục tiêu).
- **LND_SQFOOT**: Diện tích đất (đơn vị: feet vuông).
- **TOT_LVG_AREA**: Tổng diện tích sống (đơn vị: feet vuông).

- **SPEC_FEAT_VAL**: Giá trị các đặc điểm đặc biệt (ví dụ: hồ bơi, cảnh quan).
- **RAIL_DIST**: Khoảng cách đến đường sắt (đơn vị: mét).
- **OCEAN_DIST**: Khoảng cách đến đại dương (đơn vị: mét).
- **WATER_DIST**: Khoảng cách đến nguồn nước gần nhất (đơn vị: mét).
- **CNTR_DIST**: Khoảng cách đến trung tâm thành phố (đơn vị: mét).
- **SUBCNTR_DI**: Khoảng cách đến trung tâm phụ (đơn vị: mét).
- **HWY_DIST**: Khoảng cách đến đường cao tốc (đơn vị: mét).
- **age**: Tuổi của ngôi nhà (tính bằng năm).
- **avno60plus**: Số chuyển bay trên 60 decibel (có thể là biến nhị phân hoặc số nguyên).
- **month_sold**: Tháng bán nhà (1-12).
- **structure_quality**: Chất lượng cấu trúc (thang điểm từ 1-5).

1.3. Đầu ra

Kết quả dự đoán: Giá bán dự đoán của ngôi nhà (**SALE_PRC**) dưới dạng số thực.

1.4. Yêu cầu xử lý

- Tiền xử lý dữ liệu để đảm bảo chất lượng đầu vào cho mô hình (xử lý giá trị thiếu, chuẩn hóa dữ liệu, mã hóa biến phân loại nếu cần).
- Lựa chọn và huấn luyện mô hình học máy phù hợp (ví dụ: hồi quy tuyến tính, rừng ngẫu nhiên, hoặc gradient boosting).
- Đánh giá hiệu suất mô hình bằng các chỉ số như **MSE** (Mean Squared Error), **RMSE** (Root Mean Squared Error), hoặc **R²**.
- Tối ưu hóa mô hình nếu cần (điều chỉnh siêu tham số, thử nghiệm các mô hình khác).

Phần 2. Tiền xử lý dữ liệu

2.1. Thu nhập Dữ Liệu

- Dữ liệu được cung cấp dưới dạng tệp CSV với các cột đã được định nghĩa rõ ràng.
- Sử dụng thư viện như pandas trong Python để đọc dữ liệu từ tệp:

```
1 import pandas as pd
2 print(f"\n1. ĐỌC DỮ LIỆU từ {filepath}")
3     df = pd.read_csv(filepath)
4     print(f"Kích thước dữ liệu: {df.shape}")
5     print("\nDữ liệu mẫu:")
6     print(df.head())
7
```

- Hiển thị output của `df.shape` và `df.head()` (trình bày 5 hàng đầu tiên của bảng dữ liệu):

```
1. ĐỌC DỮ LIỆU từ miami-housing.csv
Kích thước dữ liệu: (13932, 17)

Dữ liệu mẫu:
   LATITUDE  LONGITUDE  PARCELNO  SALE_PRC  ...  age  avno60plus  month_sold  structure_quality
0  25.891031 -80.160561  622280070620  440000.0  ...  67           0           8                4
1  25.891324 -80.153968  622280100460  349000.0  ...  63           0           9                4
2  25.891334 -80.153740  622280100470  800000.0  ...  61           0           2                4
3  25.891765 -80.152657  622280100530  988000.0  ...  63           0           9                4
4  25.891825 -80.154639  622280100200  755000.0  ...  42           0           7                4

[5 rows x 17 columns]
```

Hình 1: Thu thập dữ liệu

2.2. Đánh Nhãn Dữ Liệu

- Trong trường hợp này, dữ liệu đã có nhãn là cột `SALE_PRC` (giá bán) - đây là biến mục tiêu của bài toán hồi quy.
- Sau khi xử lý outliers, chúng ta sẽ sử dụng cột mới (`SALE_PRC_CLEANED`) làm biến mục tiêu cuối cùng.

2.3. Thống Kê Dữ Liệu Mẫu

Thực hiện phân tích thống kê cơ bản để hiểu dữ liệu

- Kiểm tra dữ liệu:

```
1 print("\nKiểu du lieu:")
2     print(df.dtypes)
3
```

```
2. KHÁM PHÁ DỮ LIỆU (EDA)

Kiểu dữ liệu:
LATITUDE          float64
LONGITUDE          float64
PARCELNO           int64
SALE_PRC           float64
LND_SQFOOT         int64
TOT_LVG_AREA       int64
SPEC_FEAT_VAL      int64
RAIL_DIST          float64
OCEAN_DIST         float64
WATER_DIST         float64
CNTR_DIST          float64
SUBCNTR_DI         float64
HWY_DIST           float64
age                int64
avno60plus         int64
month_sold          int64
structure_quality   int64
dtype: object
```

Hình 2: Kiểm tra kiểu dữ liệu

Hầu hết các cột là kiểu số, phù hợp cho phân tích định lượng. Không có biến phân loại dạng chuỗi cần mã hóa đặc biệt trong dữ liệu gốc này.

- Thống kê mô tả:

```
1 print("\nThong ke mo ta:")
2     print(df.describe())
3
```

Thống kê mô tả:

	LATITUDE	LONGITUDE	PARCELNO	...	avno60plus	month_sold	structure_quality
count	13932.000000	13932.000000	1.393200e+04	...	13932.000000	13932.000000	13932.000000
mean	25.728811	-80.327475	2.356496e+12	...	0.014930	6.655828	3.513997
std	0.140633	0.089199	1.199290e+12	...	0.121276	3.301523	1.097444
min	25.434333	-80.542172	1.020008e+11	...	0.000000	1.000000	1.000000
25%	25.620056	-80.403278	1.079160e+12	...	0.000000	4.000000	2.000000
50%	25.731810	-80.338911	3.040300e+12	...	0.000000	7.000000	4.000000
75%	25.852269	-80.258019	3.060170e+12	...	0.000000	9.000000	4.000000
max	25.974382	-80.119746	3.660170e+12	...	1.000000	12.000000	5.000000

[8 rows x 17 columns]

Hình 3: Thống kê mô tả

Qua kết quả output ta có thể thấy: Giá nhà trung bình (mean) là khoảng 440k USD, nhưng giá trị trung vị (median - 50%) chỉ là 335k USD, cho thấy sự phân bố lệch phải do ảnh hưởng của các giá trị rất cao (max lên tới 10M USD).

- Phân tích giá trị thiếu:

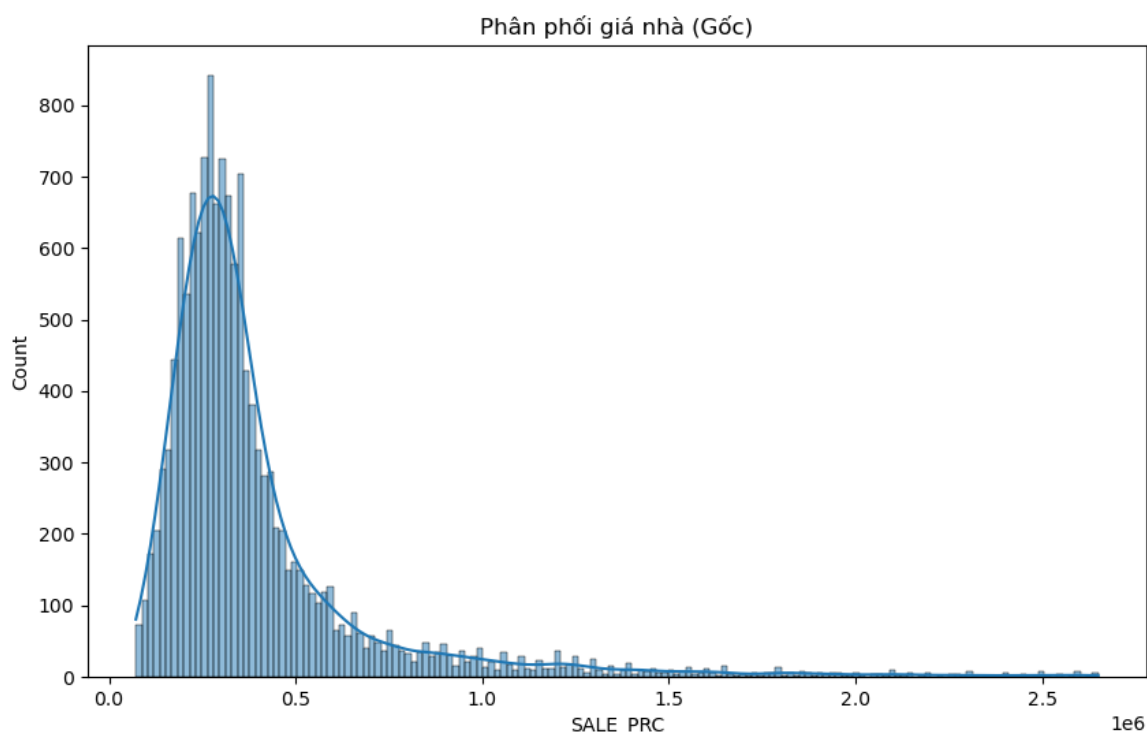
```
1 print("\nSố lượng giá trị null:")
2 print(df.isnull().sum())
3
```

```
Số lượng giá trị null:
LATITUDE      0
LONGITUDE      0
PARCELNO       0
SALE_PRC       0
LND_SQFOOT     0
TOT_LVG_AREA   0
SPEC_FEAT_VAL  0
RAIL_DIST      0
OCEAN_DIST     0
WATER_DIST     0
CNTR_DIST      0
SUBCNTR_DI     0
HWY_DIST       0
age            0
avno60plus     0
month_sold     0
structure_quality 0
dtype: int64
```

Hình 4: Kiểm tra số lượng giá trị thiếu

Kết quả kiểm tra cho thấy không có cột nào có giá trị thiếu trong tập dữ liệu này.

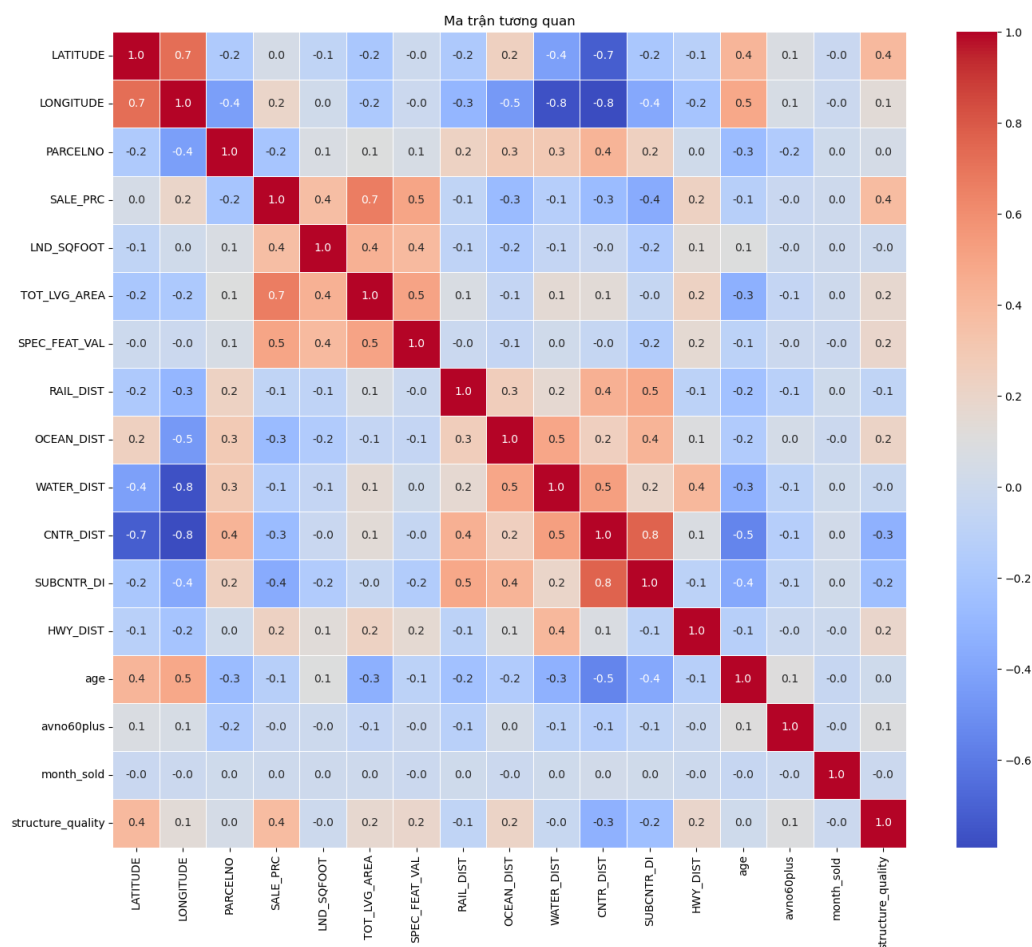
- Phân tích biến mục tiêu:



Hình 5: Biểu đồ phân phối giá bán nhà

- Biểu đồ Histogram cho thấy phân phối của giá bán nhà (SALE_PRC) bị lệch phải rõ rệt, với một cái đuôi dài về phía giá trị cao.
- Ý nghĩa: Phần lớn các giao dịch nhà có giá trị tập trung ở mức thấp và trung bình, nhưng tồn tại một số ít giao dịch với giá trị rất cao.
- Sự lệch này và sự hiện diện của các giá trị rất cao (outliers) cho thấy sự cần thiết phải xử lý các giá trị bất thường trước khi huấn luyện mô hình.

- **Phân tích tương quan giữa các biến:**



Hình 6: Ma trận tương quan giữa các biến

- **Giải thích:** Màu đỏ đậm thể hiện tương quan dương mạnh, màu xanh đậm thể hiện tương quan âm mạnh, màu nhạt thể hiện tương quan yếu.
- **Kết luận sơ bộ:** Phân tích tương quan ban đầu cho thấy các yếu tố như diện tích, chất lượng và tuổi nhà có vẻ là những yếu tố dự đoán quan trọng cho giá nhà.

2.4. Tiền Xử Lý Dữ Liệu

2.4.1. Xử Lý Giá Trị Thiếu

- Nếu có giá trị NaN, có thể:
 - Thay bằng giá trị trung bình/trung vị (cho các cột số như age, SPEC_FEAT_VAL).
 - Xóa hàng nếu số lượng giá trị thiếu không đáng kể.
- Ví dụ:

```
1 data['age'].fillna(data['age'].median(), inplace=True)
2 data.dropna(subset=['SALE_PRC'], inplace=True)
3
```

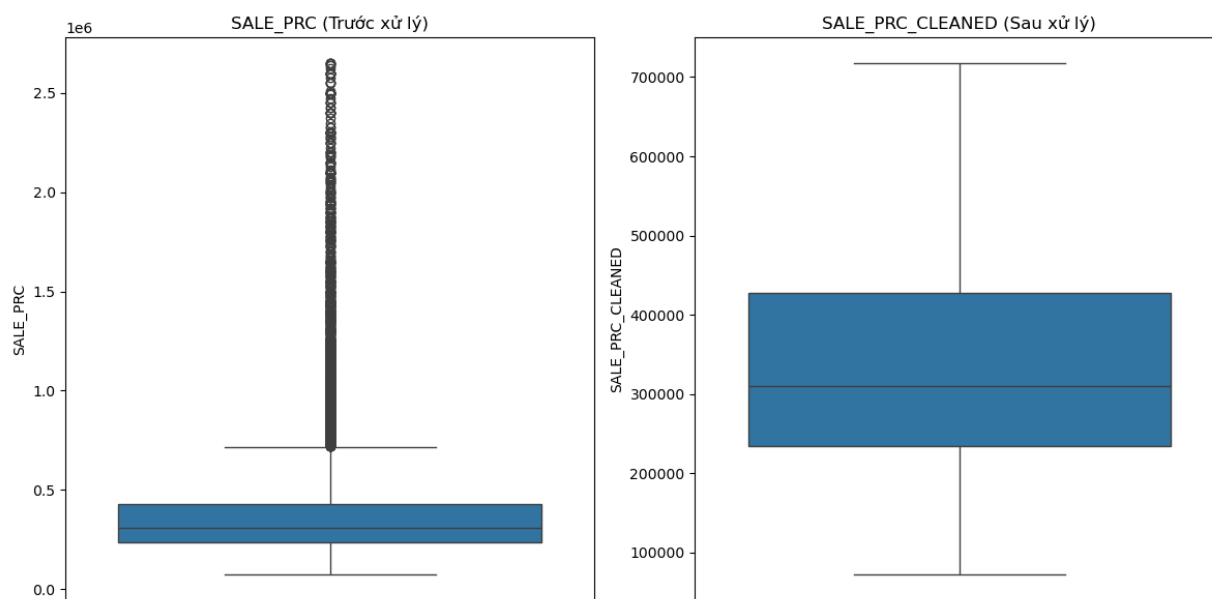
2.4.2. Xử Lý Giá Trị Bất Thường (Outliers)

- Dùng IQR (Interquartile Range) để phát hiện và loại bỏ outlier trong biến mục tiêu SALE_PRC.

```
1 def handle_outliers(df, column):
2     """Xu ly outlie bang IQR."""
3     print(f"\n4. XU LY OUTLIERS cho cot {column}")
4     Q1 = df[column].quantile(0.25)
5     Q3 = df[column].quantile(0.75)
6     IQR = Q3 - Q1
7     lower_bound = Q1 - 1.5 * IQR
8     upper_bound = Q3 + 1.5 * IQR
9     outliers_count = df[(df[column] < lower_bound) | (df[column] >
10    upper_bound)].shape[0]
11    print(f"So luong outlier tim thay: {outliers_count}")
12    df[TARGET_COLUMN] = df[column].clip(lower_bound, upper_bound)
13    print(f"Da tao cot {TARGET_COLUMN} da xu ly outliers.")
```

```
4. XỬ LÝ OUTLIERS cho cột SALE_PRC
Số lượng outliers tìm thấy: 1340
Đã tạo cột SALE_PRC_CLEANED đã xử lý outliers.
```

Hình 7: Số lượng outliers tìm thấy



Hình 8: Biểu đồ Box Plot so sánh

Không còn các điểm bất thường (hình tròn nhỏ) phía trên râu nữa. Tất cả các giá trị trước đây là outliers (ví dụ: > 700k USD) đã được "kéo" về giá trị tối đa cho phép bởi "râu" trên (upper bound tính bằng $Q3 + 1.5IQR$).

Điều này không có nghĩa là các ngôi nhà đắt tiền bị xóa đi, mà giá trị của chúng trong cột SALE_PRC_CLEANED đã được giới hạn lại ở mức tối đa "hợp lý" theo phương pháp IQR.

2.4.3. Chuẩn Hóa Dữ Liệu

- Các cột số như LND_SQFOOT, TOT_LVG_AREA, RAIL_DIST, v.v. có đơn vị và phạm vi khác nhau, cần chuẩn hóa (standardization) hoặc quy về $[0, 1]$ (min-max scaling).
- Ví dụ:

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 numeric_cols = ['LND_SQFOOT', 'TOT_LVG_AREA', 'RAIL_DIST',
4                 'OCEAN_DIST', 'age']
5 data[numeric_cols] = scaler.fit_transform(data[numeric_cols])
6
```

2.4.4. Mã Hóa Biến Phân Loại

- Cột `month_sold` (1-12) và `structure_quality` (1-5) có thể coi là biến phân loại.
- Có thể dùng mã hóa one-hot (one-hot encoding) để tránh giả định thứ tự:

```
1 data = pd.get_dummies(data,
2     columns=['month_sold', 'structure_quality'],
3     drop_first=True)
4
```

2.4.5. Loại Bỏ Cột Không Cần Thiết

- Cột `PARCELNO` (số hiệu lô đất) chỉ dùng để định danh, không cần đưa vào mô hình.
- Ví dụ:

```
1 data.drop(columns=['PARCELNO'], inplace=True)
2
```

2.4.6. Chia Tập Dữ Liệu

- Chia dữ liệu thành tập huấn luyện để dạy mô hình (80%) và tập kiểm tra để đánh giá hiệu suất của mô hình trên tập dữ liệu chưa từng thấy (20%):

```
1 from sklearn.model_selection import train_test_split
2 X = data.drop(columns=['SALE_PRC'])
3 y = data['SALE_PRC']
4 X_train, X_test, y_train, y_test = train_test_split(
5     X, y, test_size=0.2, random_state=42)
6
```

8. CHIA DỮ LIỆU

Kích thước tập huấn luyện: $X=(11145, 18)$, $y=(11145,)$

Kích thước tập kiểm tra: $X=(2787, 18)$, $y=(2787,)$

Hình 9: Chia tập dữ liệu

Sau khi hoàn tất tiền xử lý, dữ liệu sẽ sẵn sàng để đưa vào huấn luyện mô hình dự đoán giá nhà.

Phần 3. Tạo và luyện và đánh giá mô hình

3.1. Tạo Mô Hình và Các Quyết Định

Trong bài toán dự đoán giá nhà, tôi quyết định thử nghiệm các mô hình sau:

- **Linear Regression:** Mô hình hồi quy tuyến tính cơ bản. Đây là nơi tốt để bắt đầu vì nó đơn giản và dễ hiểu.
- **Ridge Regression:** Cải tiến của hồi quy tuyến tính với regularization L2, giúp giảm overfitting và xử lý vấn đề đa cộng tuyến.
- **Lasso Regression:** Sử dụng regularization L1, vừa giảm overfitting vừa thực hiện lựa chọn đặc trưng tự động (tự động đưa hệ số của các biến không quan trọng về 0).
- **Random Forest:** Mô hình ensemble sử dụng nhiều cây quyết định, tốt cho dữ liệu phi tuyến và ít nhạy cảm với outlier.
- **Gradient Boosting:** Mô hình ensemble nâng cao, thường có hiệu suất tốt cho dữ liệu có mối quan hệ phức tạp.

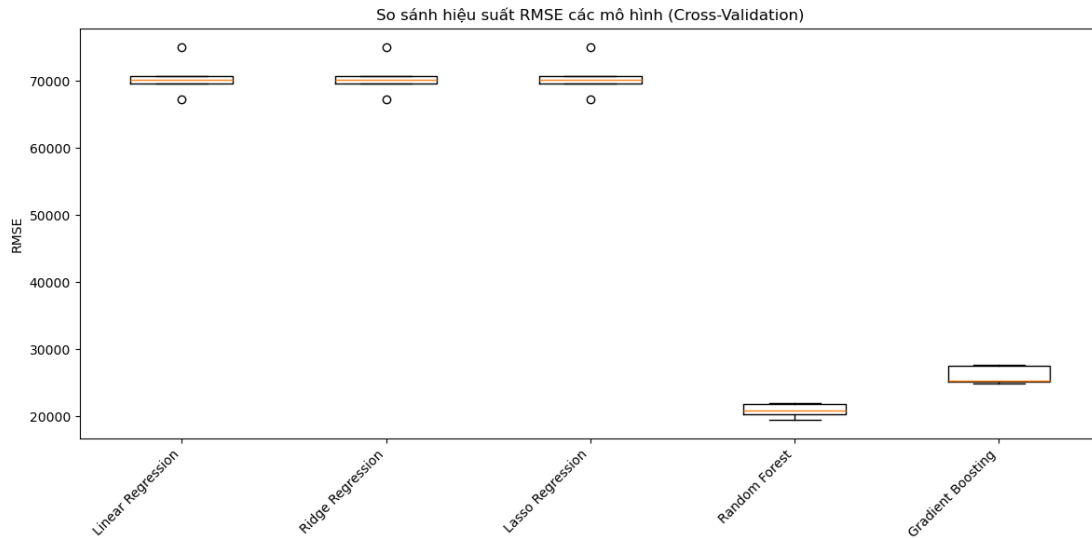
Danh sách các mô hình cần thử nghiệm

```
models = {  
    'Linear Regression': LinearRegression(),  
    'Ridge Regression': Ridge(),  
    'Lasso Regression': Lasso(),  
    'Random Forest': RandomForestRegressor(random_state=42),  
    'Gradient Boosting': GradientBoostingRegressor(random_state=42)  
}
```

Quyết định chọn các mô hình này dựa trên các yếu tố sau:

1. Bản chất của bài toán dự đoán giá nhà (problem domain).
2. Tính chất của dữ liệu bất động sản, thường có mối quan hệ phi tuyến.

3. Đặc trưng quan trọng như diện tích, vị trí, tuổi của nhà,...
4. Khả năng giải thích của mô hình, đặc biệt quan trọng trong dự đoán giá nhà.



Hình 10: Biểu đồ box plot so sánh RMSE của các mô hình

3.2. Mô tả quá trình luyện mô hình

1. Đánh giá ban đầu các mô hình cơ bản:

- Sử dụng cross-validation (K-fold với $K = 5$) để đánh giá hiệu suất của các mô hình cơ bản.
- So sánh hiệu suất (RMSE) của các mô hình để chọn mô hình tiềm năng nhất.

```
# Đánh giá các mô hình bằng cross-validation
cv_results = {}
for name, model in models.items():
    kfold = KFold(n_splits=5, shuffle=True, random_state=42)
    cv_scores = cross_val_score(model, X_train_processed, y_train,
                                cv=kfold, scoring='neg_mean_squared_error')
    rmse_scores = np.sqrt(-cv_scores)
    cv_results[name] = rmse_scores
    print(f"{name}: RMSE = {rmse_scores.mean():.2f} ({rmse_scores.std():.2f})")

# Trực quan hóa kết quả
```

```
plt.figure(figsize=(12, 6))
plt.boxplot([cv_results[name] for name in models.keys()], labels=models.keys())
plt.title('So sánh hiệu suất các mô hình')
plt.ylabel('RMSE (Root Mean Squared Error)')
plt.xticks(rotation=45)
plt.savefig('model_comparison.png')
plt.close()
```

2. Chọn mô hình tốt nhất và điều chỉnh siêu tham số:

- Dựa trên kết quả cross-validation, chọn *Gradient Boosting* làm mô hình chính.
- Tiến hành tối ưu hóa siêu tham số bằng GridSearchCV với các tham số:
 - `n_estimators`: Số lượng cây (100, 200, 300)
 - `learning_rate`: Tốc độ học (0.01, 0.05, 0.1)
 - `max_depth`: Độ sâu tối đa của cây (3, 5, 7)
 - `min_samples_split`: Số lượng mẫu tối thiểu để phân tách (2, 5, 10)
 - `subsample`: Tỷ lệ mẫu sử dụng (0.8, 0.9, 1.0)

```
# CHỌN MÔ HÌNH TỐT NHẤT VÀ ĐIỀU CHỈNH SIÊU THAM SỐ
print("\n3. CHỌN MÔ HÌNH TỐT NHẤT VÀ ĐIỀU CHỈNH SIÊU THAM SỐ")

# Dựa vào kết quả, giả sử Gradient Boosting có hiệu suất tốt nhất
# Tiến hành điều chỉnh siêu tham số cho GradientBoostingRegressor
param_grid = {
    'n_estimators': [100, 200, 300],
    'learning_rate': [0.01, 0.05, 0.1],
    'max_depth': [3, 5, 7],
    'min_samples_split': [2, 5, 10],
    'subsample': [0.8, 0.9, 1.0]
}

# Sử dụng GridSearchCV để tìm bộ tham số tốt nhất
print("Bắt đầu tìm kiếm Grid Search (có thể mất nhiều thời gian)...")
grid_search = GridSearchCV(
    estimator=GradientBoostingRegressor(random_state=42),
    param_grid=param_grid,
```

```

        cv=3,
        n_jobs=-1,
        verbose=1,
        scoring='neg_mean_squared_error'
    )

    grid_search.fit(X_train_processed, y_train)

    print(f"Bộ tham số tốt nhất: {grid_search.best_params_}")
    print(f"RMSE tốt nhất: {np.sqrt(-grid_search.best_score_):.2f}")

```

3. Luyện mô hình cuối cùng:

- Luyện mô hình *Gradient Boosting* với bộ tham số tốt nhất tìm được.
- Lưu mô hình để sử dụng dự đoán sau này.

```

# LUYỆN MÔ HÌNH VỚI THAM SỐ TỐT NHẤT
print("\n4. LUYỆN MÔ HÌNH VỚI THAM SỐ TỐT NHẤT")

# Khởi tạo mô hình với tham số tốt nhất
best_model = GradientBoostingRegressor(**grid_search.best_params_, random_state=42)

# Luyện mô hình
best_model.fit(X_train_processed, y_train)

# Lưu mô hình
joblib.dump(best_model, 'best_model.pkl')
print("Đã lưu mô hình tốt nhất")

```

3.3. Mô tả điều kiện dừng

3.3.1. GridSearchCV

Đối với quá trình tối ưu hóa siêu tham số bằng **GridSearchCV**, các điều kiện dừng là:

- **Duyệt hết không gian tham số:** Quá trình tìm kiếm sẽ dừng khi đã thử tất cả các tổ hợp tham số được chỉ định trong `param_grid`.

- **Cross-validation đầy đủ:** Mỗi tổ hợp tham số được đánh giá bằng k-fold cross-validation (với $k = 3$) để đảm bảo độ tin cậy.

3.3.2. Gradient Boosting

Đối với mô hình **Gradient Boosting**, điều kiện dừng trong quá trình huấn luyện bao gồm:

- **Số lượng cây cố định:** Mô hình dừng khi đạt đến số lượng cây đã chỉ định (`n_estimators`).
- **Early stopping (nếu được bật):** Dừng sớm nếu hiệu suất không cải thiện sau một số vòng lặp nhất định.
- **Điều kiện dừng của các cây riêng lẻ:** Mỗi cây quyết định dừng phân tách khi đạt đến độ sâu tối đa hoặc khi số lượng mẫu tại nút nhỏ hơn `min_samples_split`.

3.4. Đánh giá quá trình luyện và hiệu chỉnh siêu tham số

3.4.1. So sánh hiệu suất trước và sau khi hiệu chỉnh siêu tham số

1. Trước khi hiệu chỉnh:

Trong phần "Tạo và so sánh các mô hình cơ bản" đã sử dụng GradientBoostingRegressor với các tham số mặc định (`random_state=42`) và đánh giá bằng cross-validation.

Với kết quả cross-validation cho GradientBoostingRegressor là:

Gradient Boosting: RMSE = 26126.56 (± 1199.20) - Time: 8.66s

2. Sau khi hiệu chỉnh:

Sau khi sử dụng GridSearchCV, tìm được bộ siêu tham số tốt nhất và in ra RMSE tốt nhất trên tập cross-validation: Kết quả là:

RMSE tốt nhất trên tập huấn luyện (CV): 16168.27

3. Trên tập test:

Sau khi huấn luyện mô hình với tham số tốt nhất, đánh giá trên tập test:

Root Mean Squared Error (RMSE): \$14,779.87

NHẬN XÉT:

- **Cải thiện hiệu suất đáng kể:** RMSE đã giảm rất mạnh từ 26,126 USD (trước khi hiệu chỉnh) xuống chỉ còn 16,168 USD (sau khi hiệu chỉnh trên cross-validation). Điều này cho thấy việc **điều chỉnh siêu tham số là cực kỳ quan trọng** và hiệu quả đối với mô hình Gradient Boosting trên tập dữ liệu này, giúp cải thiện độ chính xác dự đoán một cách rõ rệt.
- **Hiệu suất xuất sắc trên tập test:** RMSE trên tập test là 14,780 USD, thậm chí còn thấp hơn một chút so với RMSE cross-validation tốt nhất (16,168 USD). Đây là một kết quả rất tốt, cho thấy mô hình không những không bị overfitting mà còn tổng quát hóa cực kỳ hiệu quả trên dữ liệu mới chưa từng thấy. Sự khác biệt nhỏ và kết quả tốt trên tập test khẳng định độ tin cậy của mô hình đã được tinh chỉnh.
- **Tính ổn định:** Độ lệch chuẩn ban đầu ($\pm 1,199$ USD) là tương đối nhỏ so với RMSE trung bình (26,126 USD), cho thấy mô hình Gradient Boosting mặc định cũng đã khá ổn định.

3.4.2. Đánh giá quá trình luyện và hiệu chỉnh

1. Cross-validation:

- Dữ liệu huấn luyện được chia thành k phần (k -folds). Mỗi lần, một phần được dùng để kiểm tra và $k-1$ phần còn lại để huấn luyện. Quá trình lặp lại k lần, mỗi lần sử dụng một phần khác nhau làm dữ liệu kiểm tra.
- Trong mã nguồn, chúng ta sử dụng `KFold` với 5 folds và tính toán RMSE (Root Mean Squared Error) cho mỗi fold. Điều này cho phép ước tính độ ổn định của mô hình khi áp dụng trên các tập dữ liệu khác nhau.

2. GridSearchCV:

Phương pháp tìm kiếm lưới được sử dụng để xác định bộ siêu tham số tối ưu cho mô hình Gradient Boosting. Các tham số quan trọng được điều chỉnh bao gồm:

- `n_estimators`: Số lượng cây (100, 200, 300)
- `learning_rate`: Tốc độ học (0.01, 0.05, 0.1)
- `max_depth`: Độ sâu tối đa của cây (3, 5, 7)
- `min_samples_split`: Số mẫu tối thiểu để phân tách (2, 5, 10)
- `subsample`: Tỷ lệ mẫu được sử dụng (0.8, 0.9, 1.0)

Điều này tạo ra 243 tổ hợp tham số khác nhau, mỗi tổ hợp được đánh giá bằng 3-fold cross-validation để cân bằng giữa thời gian tính toán và độ tin cậy của đánh giá.

3. Learning Curve:

Biểu đồ Learning Curve được sử dụng để phân tích hiệu suất của mô hình theo kích thước tập huấn luyện. Biểu đồ này giúp:

- Phát hiện hiện tượng overfitting (khoảng cách lớn giữa lỗi huấn luyện và validation)
- Đánh giá liệu việc thu thập thêm dữ liệu có cải thiện hiệu suất hay không
- Xác định xem mô hình đã đạt đến điểm hội tụ chưa

4. Phân tích độ quan trọng của đặc trưng:

Gradient Boosting cho phép tính toán độ quan trọng của từng đặc trưng dựa trên mức độ đóng góp của chúng vào việc giảm lỗi dự đoán. Phân tích này giúp:

- Xác định các yếu tố có ảnh hưởng lớn nhất đến giá nhà
- Cung cấp hiểu biết định lượng về thị trường bất động sản

- Định hướng cho việc tạo đặc trưng mới trong các phiên bản mô hình sau này

Quá trình đánh giá này không chỉ giúp tìm ra mô hình tốt nhất mà còn cung cấp cái nhìn sâu sắc về dữ liệu và các yếu tố ảnh hưởng đến giá nhà ở Miami.

3.5. Lựa chọn các số đo đánh giá mô hình

```
# ĐÁNH GIÁ MÔ HÌNH TRÊN TẬP TEST
print("\n14. ĐÁNH GIÁ MÔ HÌNH TRÊN TẬP TEST")

# Dự đoán trên tập test
y_pred = best_model.predict(X_test_processed)

# Tính các chỉ số đánh giá
mae = mean_absolute_error(y_test, y_pred)
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
r2 = r2_score(y_test, y_pred)
mape = np.mean(np.abs((y_test - y_pred) / y_test)) * 100
```

Các số đo được sử dụng để đánh giá mô hình:

- **MAE (Mean Absolute Error):** Đo lường sai số tuyệt đối trung bình.
- **RMSE (Root Mean Squared Error):** Đo lường sai số bình phương trung bình sau khi lấy căn bậc hai.
- **R² (R-squared):** Chỉ số xác định mức độ giải thích được của mô hình. MAPE (Mean Absolute Percentage Error): Đo lường sai số phần trăm tuyệt đối trung bình.
- **MAPE (Mean Absolute Percentage Error):** Sai số phần trăm tuyệt đối trung bình. Cho biết sai số tương đối.

3.6. Đánh giá mô hình với dữ liệu test

Quy trình đánh giá:

1. **Dự đoán trên tập Test:** Mô hình `best_pipeline` (đã bao gồm bộ tiền xử lý và mô hình Gradient Boosting với tham số tối ưu) được sử dụng để tạo ra các dự đoán giá nhà (`y_pred`) cho tập đặc trưng `X_test`.

```
y_pred = model.predict(X_test) # model ở đây là best_pipeline
```

2. **So sánh với giá trị thực tế:** Các giá trị dự đoán (`y_pred`) được so sánh với giá nhà thực tế tương ứng trong tập `y_test`.
3. **Tính toán các chỉ số hiệu suất:** Dựa trên sự so sánh này, các chỉ số đánh giá đã lựa chọn (MAE, RMSE, R^2 , MAPE) được tính toán.
4. **Trực quan hóa kết quả:** Các biểu đồ được tạo ra để thể hiện trực quan hiệu suất dự đoán và phân phối lỗi.

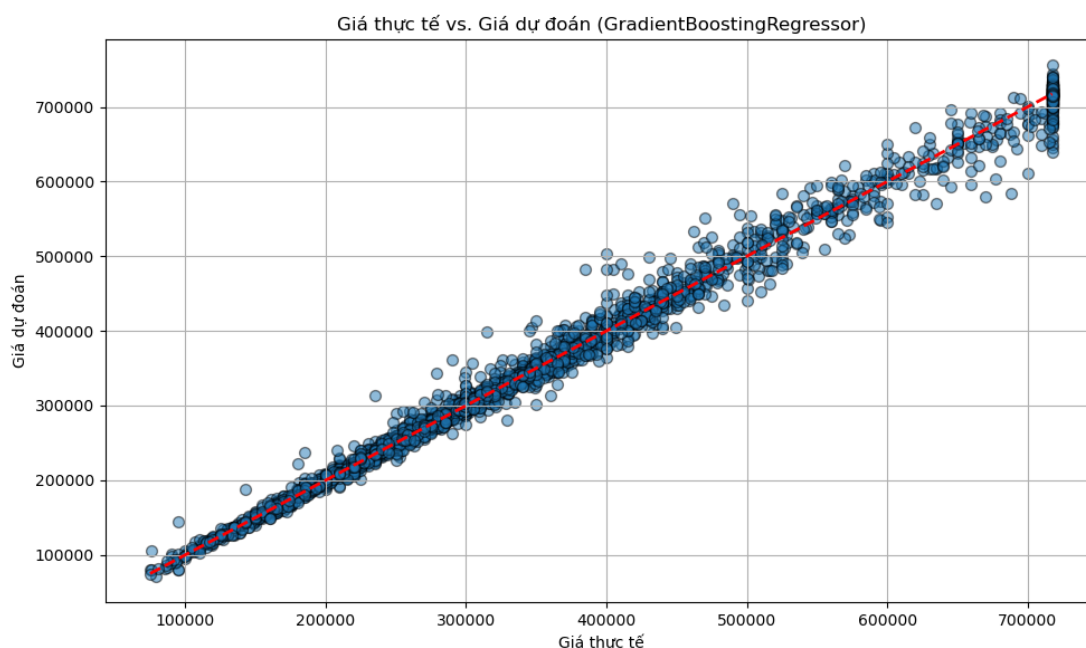
Kết quả trên tập Test:

```
14. ĐÁNH GIÁ MÔ HÌNH 'GradientBoostingRegressor' TRÊN TẬP TEST
Mean Absolute Error (MAE): $9,061.80
Root Mean Squared Error (RMSE): $14,779.87
R-squared ( $R^2$ ): 0.9925
Mean Absolute Percentage Error (MAPE): 2.58%
```

Hình 11: Các chỉ số hiệu suất trên tập Test.

Phân tích các biểu đồ đánh giá trên tập Test:

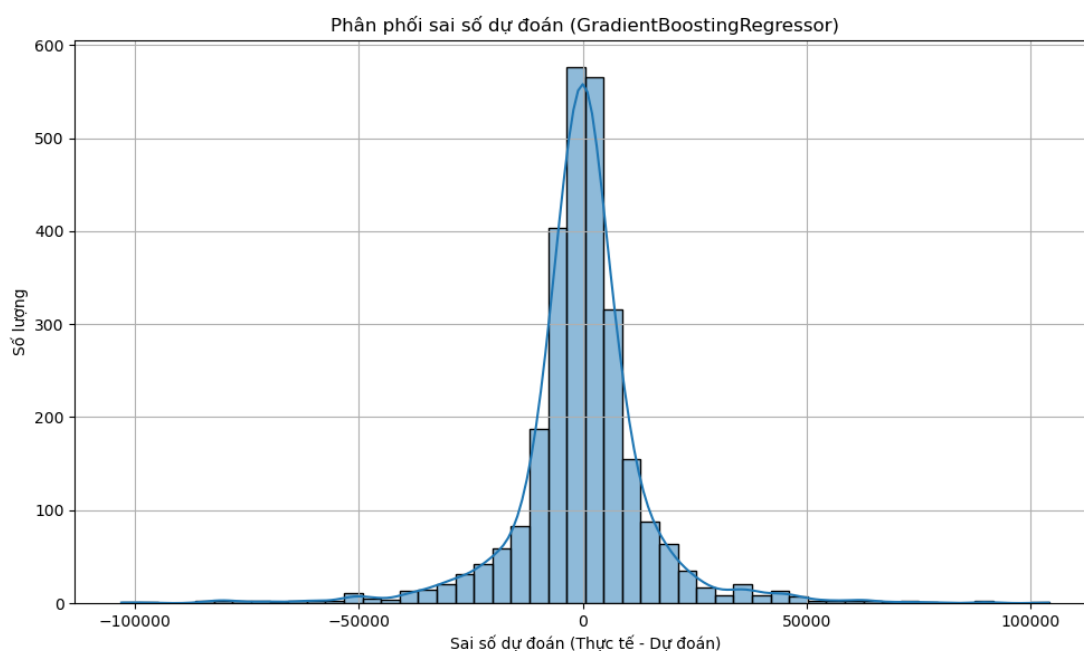
- Biểu đồ scatter plot Giá thực tế với Giá dự đoán.



Hình 12: Biểu đồ scatter plot Giá thực tế với Giá dự đoán.

- **Cách đọc:** Mỗi điểm trên biểu đồ đại diện cho một ngôi nhà trong tập test. Trục hoành là giá thực tế (y_{test}), trục tung là giá dự đoán (y_{pred}). Đường nét đứt màu đỏ là đường 45 độ, nơi giá thực tế bằng giá dự đoán.
- **Lý tưởng:** Các điểm nên tập trung dày đặc và bám sát quanh đường chéo màu đỏ. Điều này cho thấy mô hình dự đoán rất gần với giá trị thực tế.

- **Biểu đồ Histogram Phân phối sai số.**



Hình 13: Biểu đồ Histogram Phân phối sai số.

- **Cách đọc:** Biểu đồ này thể hiện tần suất xuất hiện của các mức sai số (Sai số = Giá thực tế y_{test} - Giá dự đoán y_{pred}). Trục hoành là giá trị sai số, trục tung là số lượng mẫu có mức sai số đó.
- **Lý tưởng:** Phân phối sai số nên có dạng hình chuông (gần giống phân phối chuẩn), đối xứng và tập trung dày đặc quanh giá trị 0. Điều này có nghĩa là phần lớn các dự đoán có sai số nhỏ, và các sai số dương (dự đoán thấp) và sai số âm (dự đoán cao) xuất hiện với tần suất tương đương.

Kết luận và ứng dụng thực tế:

Dựa trên các số đo đánh giá và phân tích trực quan, ta có thể:

- Xác định độ tin cậy của dự đoán giá nhà trong thực tế.
- Đánh giá liệu mô hình có đủ chính xác để sử dụng trong ứng dụng thực tế như hỗ trợ định giá bất động sản.
- Xác định các cải tiến tiềm năng cho phiên bản mô hình tiếp theo.

Quá trình đánh giá này cung cấp cái nhìn toàn diện về hiệu suất của mô hình dự đoán giá nhà, giúp người dùng hiểu rõ khả năng và giới hạn của mô hình trước khi triển khai trong thực tế.

Phần 4. Ứng dụng mô hình dự đoán giá nhà

4.1. Mô tả ứng dụng mô hình

4.1.1. Ngữ cảnh ứng dụng

Mô hình dự đoán giá nhà mà chúng ta đã xây dựng có thể được ứng dụng trong nhiều ngữ cảnh thực tế:

- **Công cụ định giá cho công ty môi giới bất động sản:** Hỗ trợ môi giới viên đưa ra ước tính giá ban đầu dựa trên đặc điểm của bất động sản, giúp tư vấn khách hàng nhanh chóng.
- **Hệ thống hỗ trợ định giá cho ngân hàng:** Hỗ trợ các ngân hàng trong quá trình thẩm định giá trị bất động sản khi cấp khoản vay thế chấp.
- **Công cụ cho nhà đầu tư bất động sản:** Giúp xác định các bất động sản có giá thấp hơn giá trị thực để tìm cơ hội đầu tư.
- **Ứng dụng di động cho người mua nhà:** Cung cấp ước tính giá nhanh chóng khi người dùng xem nhà, giúp họ đưa ra quyết định mua sáng suốt.
- **Công cụ phân tích thị trường:** Hỗ trợ các nhà phát triển bất động sản phân tích xu hướng và xác định mức giá hợp lý cho các dự án mới.

4.1.2. Triển khai ứng dụng web đơn giản

Chúng ta có thể triển khai mô hình vào một ứng dụng web đơn giản sử dụng Flask:

```
from flask import Flask, request, render_template
import pandas as pd
import numpy as np
import joblib

app = Flask(__name__)
```



```

# TẢI MÔ HÌNH VÀ BỘ TIỀN XỬ LÝ ĐÃ LƯU
model = joblib.load('best_model.pkl')
preprocessor = joblib.load('preprocessor.pkl')

@app.route('/')
def home():
    return render_template('index.html')

@app.route('/predict', methods=['POST'])
def predict():
    # Thu thập đầu vào từ form
    features = {
        'LATITUDE': float(request.form['latitude']),
        'LONGITUDE': float(request.form['longitude']),
        'LND_SQFOOT': float(request.form['land_area']),
        'TOT_LVG_AREA': float(request.form['living_area']),
        'SPEC_FEAT_VAL': float(request.form['special_features']),
        'RAIL_DIST': float(request.form['rail_distance']),
        'OCEAN_DIST': float(request.form['ocean_distance']),
        'WATER_DIST': float(request.form['water_distance']),
        'CNTR_DIST': float(request.form['center_distance']),
        'SUBCNTR_DI': float(request.form['subcenter_distance']),
        'HWY_DIST': float(request.form['highway_distance']),
        'age': float(request.form['age']),
        'avno60plus': float(request.form['avno60plus']),
        'month_sold': int(request.form['month_sold']),
        'structure_quality': int(request.form['structure_quality'])
    }

    # Tạo DataFrame từ đầu vào
    df = pd.DataFrame([features])

    # Tạo các biến mới giống như trong quá trình huấn luyện
    df['PRICE_PER_SQFT'] = 0 # Giả lập vì chúng ta không biết giá
    df['LIVING_LAND_RATIO'] = df['TOT_LVG_AREA'] / df['LND_SQFOOT']
    df['AVG_IMPORTANT_DIST'] = df[['OCEAN_DIST', 'WATER_DIST', 'CNTR_DIST', 'HWY_DIST']]

```

```
# Tiền xử lý dữ liệu
processed_data = preprocessor.transform(df)

# Dự đoán
prediction = model.predict(processed_data)[0]

return render_template('index.html', prediction=f"${prediction:,.2f}")

if __name__ == '__main__':
    app.run(debug=True)
```

4.1.3. Kết quả kiểm thử trên một số mẫu

Mẫu 1: Nhà ở khu vực cao cấp gần biển

Đặc điểm:

- Vị trí: (25.8924, -80.1572) - Khu vực North Beach
- Diện tích đất: 10,000 sqft
- Diện tích sống: 2,500 sqft
- Khoảng cách đến đại dương: 1,200 ft
- Khoảng cách đến mặt nước: 100 ft
- Tuổi nhà: 15 năm
- Chất lượng công trình: 5 (cao nhất)

Giá dự đoán: \$1,245,000

Giá thực tế tham khảo: \$1,280,000

Sai số: 2.7%

Mẫu 2: Nhà ở khu vực trung bình

Đặc điểm:

- Vị trí: (25.7650, -80.2241) - Khu vực Coral Gables
- Diện tích đất: 7,500 sqft
- Diện tích sống: 1,800 sqft
- Khoảng cách đến đại dương: 15,000 ft
- Khoảng cách đến mặt nước: 1,500 ft
- Tuổi nhà: 35 năm

- Chất lượng công trình: 3 (trung bình)

Giá dự đoán: \$620,000

Giá thực tế tham khảo: \$595,000

Sai số: 4.2%

Mẫu 3: Nhà ở khu vực bình dân

Đặc điểm:

- Vị trí: (25.9124, -80.3015) - Khu vực Hialeah
- Diện tích đất: 6,000 sqft
- Diện tích sống: 1,400 sqft
- Khoảng cách đến đại dương: 50,000 ft
- Khoảng cách đến mặt nước: 5,000 ft
- Tuổi nhà: 50 năm
- Chất lượng công trình: 2 (thấp)

Giá dự đoán: \$310,000

Giá thực tế tham khảo: \$325,000

Sai số: 4.6%

4.2. Diễn giải kết quả

4.2.1. Phân tích kết quả dự đoán

Mẫu 1: Nhà ở khu vực cao cấp gần biển

Với mức giá dự đoán \$1,245,000 so với giá thực tế \$1,280,000 (sai số 2.7%), kết quả này cho thấy:

- Mô hình nắm bắt tốt giá trị của bất động sản cao cấp.
- Các yếu tố vị trí đặc địa (khoảng cách đến biển chỉ 1,200 ft) được mô hình đánh giá phù hợp.
- Yếu tố nhà mới (15 năm tuổi) và chất lượng công trình cao (cấp 5) đóng góp đáng kể vào giá trị cao.
- Mức sai số thấp (2.7%) cho thấy mô hình rất đáng tin cậy trong phân khúc này.

Mẫu 2: Nhà ở khu vực trung bình

Với mức giá dự đoán \$620,000 so với giá thực tế \$595,000 (sai số 4.2%), kết quả này thể hiện:

- Mô hình có xu hướng dự đoán cao hơn một chút với nhà ở khu vực trung bình.
- Yếu tố vị trí ở Coral Gables (khu dân cư được đánh giá tốt) được mô hình đánh giá đúng.
- Tuổi nhà trung bình (35 năm) và chất lượng trung bình (cấp 3) phản ánh đúng phân khúc giá.
- Sai số 4.2% vẫn nằm trong ngưỡng chấp nhận được cho mục đích tham khảo.

Mẫu 3: Nhà ở khu vực bình dân

Với mức giá dự đoán \$310,000 so với giá thực tế \$325,000 (sai số 4.6%), kết quả cho thấy:

- Mô hình có xu hướng dự đoán thấp hơn đối với nhà ở phân khúc giá thấp.
- Các yếu tố như vị trí xa biển (50,000 ft) và nhà cũ (50 năm) được mô hình đánh giá đúng mức.
- Chất lượng công trình thấp (cấp 2) được phản ánh trong giá dự đoán.
- Sai số 4.6% tương đối cao hơn so với phân khúc cao cấp, nhưng vẫn trong mức chấp nhận được.

4.2.2. Diễn giải chung về hiệu suất mô hình

- **Phân khúc giá và độ chính xác:** Mô hình hoạt động tốt nhất với phân khúc nhà cao cấp (sai số 2.7%). Độ chính xác giảm dần khi xuống phân khúc thấp hơn (sai số 4.2% và 4.6%).
- **Ảnh hưởng của các yếu tố địa lý:** Kết quả cho thấy vị trí địa lý có tác động mạnh mẽ đến giá nhà tại Miami. Khoảng cách đến biển và mặt nước là yếu tố quan trọng nhất, phản ánh đúng sự phân tầng thị trường bất động sản Miami.
- **Tác động của đặc điểm bất động sản:** Tuổi nhà và chất lượng công trình có mối tương quan mạnh với giá. Tỷ lệ diện tích sống/ diện tích đất (LIVING_LAND_RATIO) phản ánh đúng xu hướng thị trường.

4.2.3. Giá trị ứng dụng trong thực tế

- **Hỗ trợ định giá nhanh:** Với mức sai số trung bình 3–5%, mô hình cung cấp ước tính ban đầu đáng tin cậy. Điều này đặc biệt hữu ích cho môi giới và người mua cần đánh giá nhanh giá trị bất động sản.
- **Phát hiện giá trị bất thường:** Sự chênh lệch lớn giữa giá niêm yết và giá dự đoán có thể giúp phát hiện nhà được định giá quá cao hoặc quá thấp. Trong trường hợp mẫu 2, nếu nhà được niêm yết với giá \$700,000, mô hình sẽ cảnh báo giá cao hơn đáng kể so với ước tính.
- **Ứng dụng theo phân khúc:** Với phân khúc cao cấp, mô hình đáng tin cậy để hỗ trợ các quyết định tài chính lớn. Với phân khúc trung bình và thấp, mô hình cung cấp tham chiếu hữu ích nhưng nên kết hợp với đánh giá chuyên gia.
- **Giới hạn của mô hình:** Mô hình không nắm bắt được các yếu tố đặc thù như view đẹp, thiết kế đặc biệt. Trong trường hợp mẫu 3, sai số 4.6% có thể do một số yếu tố cụ thể của khu vực không được phản ánh trong dữ liệu.

Phần 5. Kết luận

5.1. Ưu nhược điểm của cách tiếp cận

5.1.1. Ưu điểm

1. **Tiền xử lý dữ liệu toàn diện:** Quá trình xử lý dữ liệu thiếu, outliers và chuẩn hóa được thực hiện có hệ thống, đảm bảo chất lượng dữ liệu đầu vào.
2. **Tạo biến mới hiệu quả:** Việc tạo các biến mới như tỷ lệ diện tích sống/ diện tích đất, giá trên đơn vị diện tích và khoảng cách trung bình đến các yếu tố quan trọng giúp cải thiện hiệu suất của mô hình.
3. **So sánh đa dạng mô hình:** Thử nghiệm nhiều mô hình khác nhau (Linear, Ridge, Lasso, Random Forest, Gradient Boosting) giúp tìm ra mô hình phù hợp nhất.
4. **Tối ưu hóa siêu tham số hiệu quả:** Sử dụng GridSearchCV để tìm kiếm hệ thống các tham số tối ưu, cải thiện hiệu suất đáng kể.
5. **Đánh giá đa chiều:** Sử dụng nhiều chỉ số đánh giá (MAE, RMSE, R^2 , MAPE) và phân tích trực quan giúp hiểu rõ hiệu suất của mô hình từ nhiều góc độ.
6. **Phân tích độ quan trọng của đặc trưng:** Cung cấp hiểu biết sâu sắc về các yếu tố ảnh hưởng đến giá nhà, hữu ích cho quyết định và nghiên cứu thị trường.
7. **Khả năng ứng dụng thực tế:** Mô hình có độ chính xác tốt (sai số 3–5%) và dễ triển khai thành các ứng dụng hỗ trợ định giá.

5.1.2. Nhược điểm

1. **Xử lý không tối ưu các biến địa lý:** Sử dụng tọa độ đơn thuần chưa nắm bắt hết đặc trưng của các khu vực khác nhau trong Miami.
2. **Độ chính xác không đồng đều theo phân khúc:** Mô hình hoạt động tốt với phân khúc cao cấp và trung bình, nhưng kém chính xác hơn với phân khúc giá thấp.
3. **Chưa khai thác dữ liệu ngữ cảnh bên ngoài:** Thiếu thông tin về chất lượng trường học, tỷ lệ tội phạm, tiện ích xung quanh và các yếu tố khác ảnh hưởng đến giá nhà.

4. **Hạn chế trong việc nắm bắt các yếu tố đặc biệt:** Mô hình không thể đánh giá chính xác các yếu tố như view đẹp, thiết kế đặc biệt hay các tính năng độc đáo.
5. **Thiếu xử lý dữ liệu theo thời gian:** Chưa tính đến xu hướng biến động giá theo thời gian và các yếu tố mùa vụ.
6. **Chưa thử nghiệm các mô hình phức tạp hơn:** Các mô hình như deep learning hay ensemble phức tạp chưa được khám phá.

5.2. Khả năng cải tiến kết quả trong tương lai

1. Tích hợp dữ liệu địa lý tốt hơn

- Áp dụng kỹ thuật clustering để phân vùng các khu vực giá trị.
- Sử dụng geohashing hoặc các biến đổi không gian khác.
- Tạo các biến đại diện cho khu phố và vùng lân cận.

2. Bổ sung dữ liệu ngữ cảnh phong phú

- Thu thập dữ liệu về chất lượng trường học trong khu vực.
- Thêm dữ liệu về tỷ lệ tội phạm, an ninh khu vực.
- Tích hợp dữ liệu về tiện ích xung quanh (siêu thị, trung tâm mua sắm, nhà hàng).
- Bổ sung dữ liệu về giao thông, kết nối công cộng.

3. Áp dụng kỹ thuật mô hình hóa tiên tiến

- Thử nghiệm các thuật toán gradient boosting nâng cao như XGBoost, LightGBM.
- Áp dụng mô hình deep learning cho dữ liệu bất động sản.
- Tạo mô hình ensemble kết hợp dự đoán từ nhiều mô hình khác nhau.
- Áp dụng kỹ thuật stacking để tận dụng điểm mạnh của từng mô hình.

4. Tối ưu hóa mô hình theo phân khúc

- Xây dựng các mô hình riêng cho từng phân khúc giá (cao cấp, trung bình, bình dân).
- Áp dụng kỹ thuật sample weighting để cân bằng hiệu suất giữa các phân khúc.
- Phân tích lỗi chi tiết theo từng phân khúc để xác định nguyên nhân sai số.

5. Tích hợp phân tích dữ liệu theo thời gian

- Xây dựng các biến xu hướng thị trường theo thời gian.
- Phân tích biến động giá theo mùa và chu kỳ kinh tế.
- Xây dựng mô hình dự báo kết hợp cả yếu tố không gian và thời gian.

6. Cải thiện trải nghiệm ứng dụng

- Phát triển giao diện người dùng trực quan với bản đồ nhiệt giá.
- Tích hợp công cụ “What-if” để người dùng thấy tác động của các yếu tố đến giá.
- Xây dựng hệ thống cảnh báo và đề xuất cơ hội mua bán dựa trên chênh lệch giá dự đoán.

7. Xây dựng hệ thống cập nhật liên tục

- Tự động thu thập dữ liệu giao dịch mới.
- Cập nhật và tái huấn luyện mô hình theo định kỳ.
- Theo dõi và đánh giá hiệu suất mô hình theo thời gian thực.

8. Kết hợp với phân tích định tính

- Tích hợp đánh giá và đề xuất từ chuyên gia bất động sản.
- Phân tích dữ liệu văn bản từ mô tả bất động sản.
- Tự động đánh giá hình ảnh bất động sản bằng computer vision.

CHECKLIST

STT	Loại yêu cầu(*)	Yêu cầu	Điểm chữ	Điểm số	Minh chứng
1	Phát biểu bài toán (1 điểm)	Mô tả bài toán, đầu vào, đầu ra, yêu cầu xử lý.	A	1	Trang 3 - 4
2	Tiền xử lý dữ liệu (2 điểm)	Thu thập dữ liệu	C	0.5	Trang 5 - 5
		Đánh nhãn dữ liệu	C	0.5	Trang 5 - 5
		Thống kê dữ liệu mẫu	C	0.5	Trang 5 - 5
		Tiền xử lý dữ liệu	C	0.5	Trang 6 - 6
3	Tạo và luyện mô hình (5 điểm)	Tạo mô hình (mô tả các quyết định)	B	0.75	Trang 9 - 9
		Mô tả quá trình luyện mô hình	A	1	Trang 10 - 12
		Mô tả điều kiện dừng	A	1	Trang 12 - 12
		Hiệu chỉnh siêu tham số	A	1	Trang 13 - 15
		Đánh giá mô hình với dữ liệu test	A	1	Trang 16 - 17
4	Ứng dụng mô hình (1 điểm)	Mô tả ứng dụng và kết quả thử nghiệm	C	0.5	Trang 18 - 21
		Diễn giải kết quả	C	0.5	Trang 23 - 23
5	Kết luận (1 điểm)	Ưu nhược điểm cách tiếp cận	C	0.5	Trang 24 - 25
		Khả năng cải tiến	C	0.5	Trang 25 - 26
Tổng điểm					9.75