Dataset loaded successfully!

Dataset Background Information:
- Created by: Maha ALDossary
- Accessed from: https://www.kaggle.com/datasets/maha48/villas-price-dataset/data
- Aligns with UNSDG 11: Sustainable Cities and Communities, as it provides insights for urban planning and housing development to meet growing population demands efficiently.
- This dataset contains information on villa properties, including features like the number of rooms, bathrooms, elevator, pool, driver, and garden.
- The target variable is the property size in square meters (sqm).
- The dataset allows for prediction of property size based on features such as the number of rooms, bathrooms, and available amenities.
- It also aids in analyzing trends related to real estate and housing development in Saudi Arabia, helping to forecast infrastructure needs based on property characteristics.
- With 930 records, it supports analysis of the relationship between property features and size, useful for future urban planning and housing policy.


Dataset Attributes:
['neighborhood_name', 'administritive_area', 'city', 'rooms', 'bathrooms', 'sqm', 'elevator', 'pool', 'driver', 'garden']

Potential Questions This Dataset Could Answer:
- How do the number of rooms, bathrooms, and available amenities influence the size of properties?
- What features of villas most strongly correlate with their square meter size?
- Can we develop a model to predict the size of villas based on the number of rooms and available facilities?
- How can property features guide urban planning and housing development to address population growth efficiently?

Data Suitability Assessment:
- Completeness: No missing values (verified below)
- Relevance: Directly measures key factors influencing property size, useful for SDG 11: Sustainable Cities and Communities
- Quality: Data sourced from reputable sources, ensuring high quality and consistency for urban planning analysis

```
Categorical Columns Encoded: ['neighborhood_name', 'administritive_area', 'city']

Missing Values:
neighborhood_name       0
administritive_area      0
city                     0
rooms                    0
bathrooms                0
sqm                      0
elevator                 0
pool                     0
driver                   0
garden                   0
dtype: int64

Statistical Summary:
       neighborhood_name  administritive_area         city       rooms  \
count         930.000000           930.000000   930.000000  930.000000
mean           55.134409             1.044086     7.860215    4.904301
std            28.776489             1.417708     1.783137    1.311735
min             0.000000             0.000000     0.000000    1.000000
25%            35.000000             0.000000     7.000000    4.000000
50%            58.000000             0.000000     7.000000    5.000000
75%            75.000000             3.000000     9.000000    6.000000
max           105.000000             3.000000    11.000000    7.000000

          bathrooms          sqm    elevator        pool      driver      garden
count    930.000000   930.000000  930.000000  930.000000  930.000000  930.000000
mean       5.451613   423.451613    0.374194    0.221505    0.172043    0.045161
std        1.348093   283.492053    0.520599    0.450296    0.383280    0.207770
min        1.000000     1.000000    0.000000    0.000000    0.000000    0.000000
25%        5.000000   300.000000    0.000000    0.000000    0.000000    0.000000
50%        5.000000   343.500000    0.000000    0.000000    0.000000    0.000000
75%        7.000000   450.000000    1.000000    0.000000    0.000000    0.000000
max        7.000000  3000.000000    3.000000    3.000000    2.000000    1.000000

Applying PowerTransformer to skewed features: ['sqm', 'elevator', 'pool', 'driver', 'garden']
```
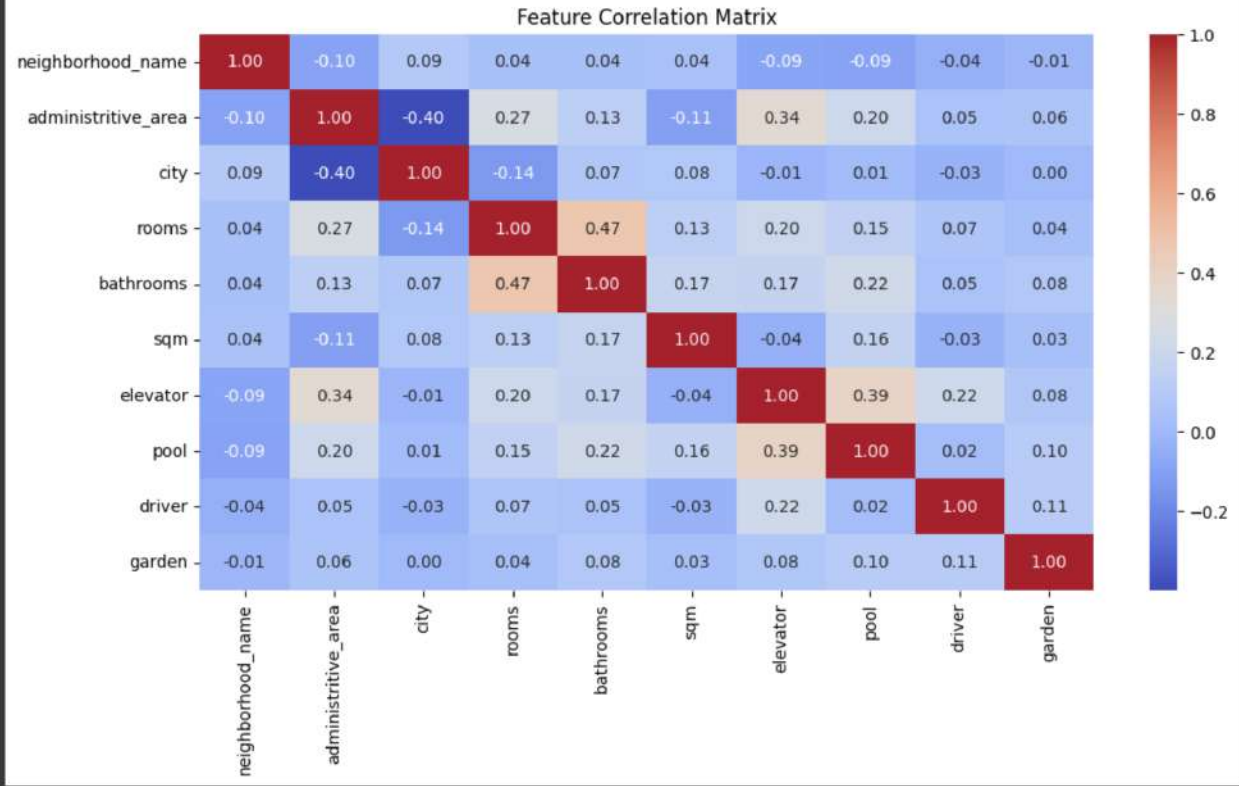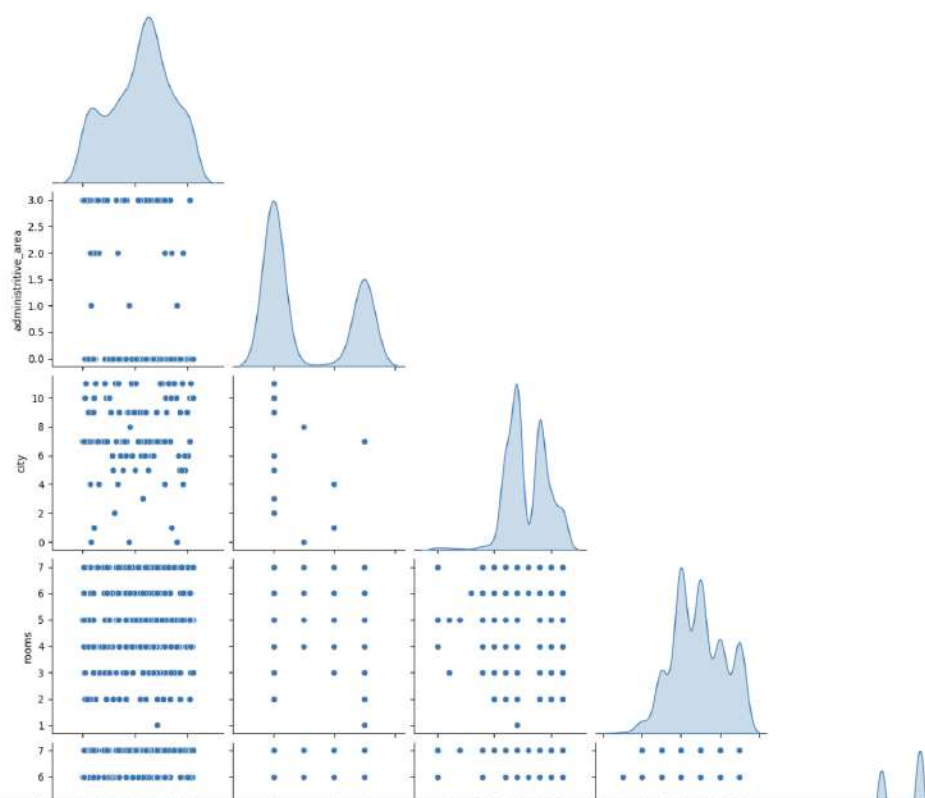
Applying PowerTransformer to skewed features: ['sqm', 'elevator', 'pool', 'driver', 'garden']
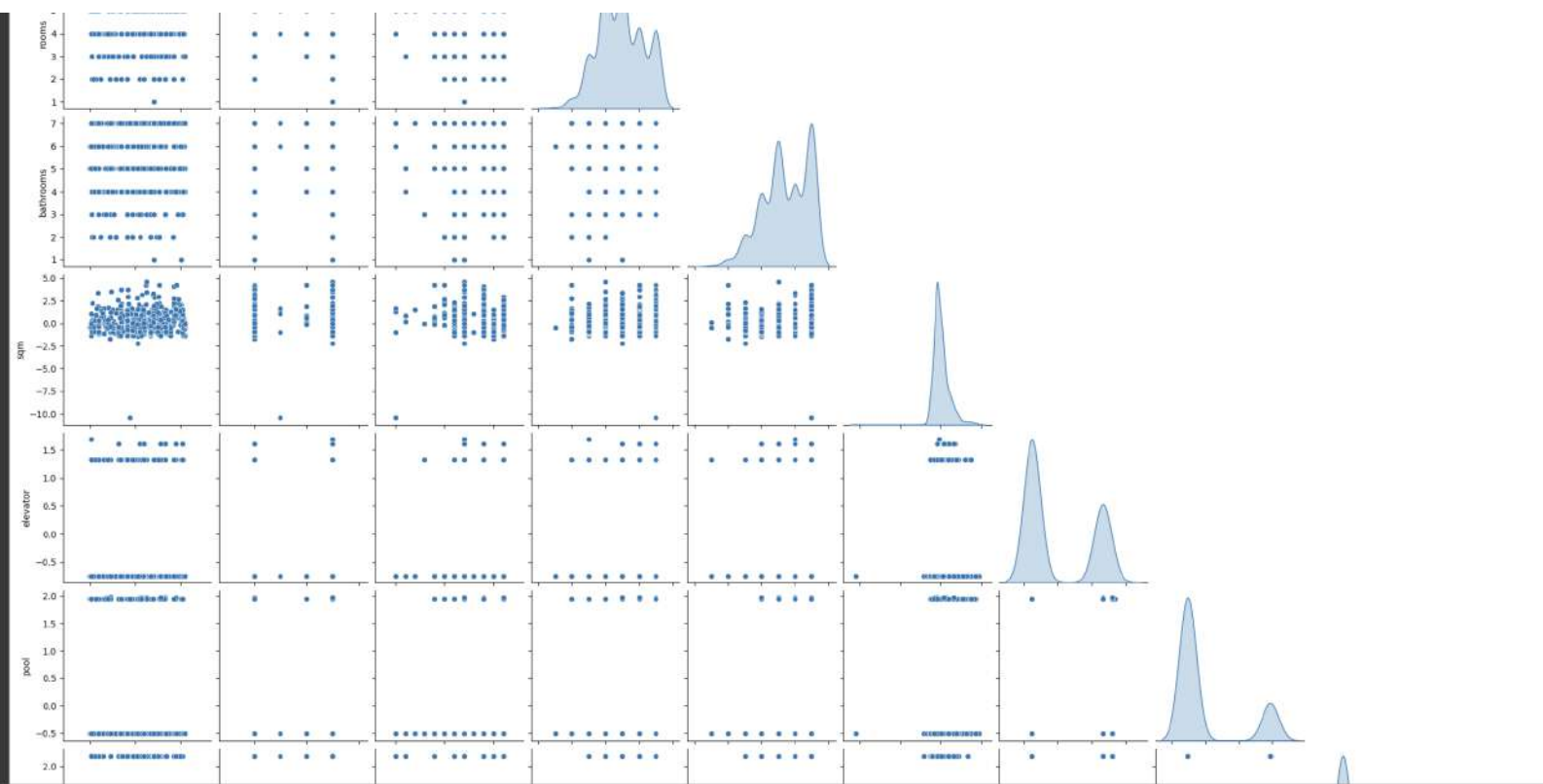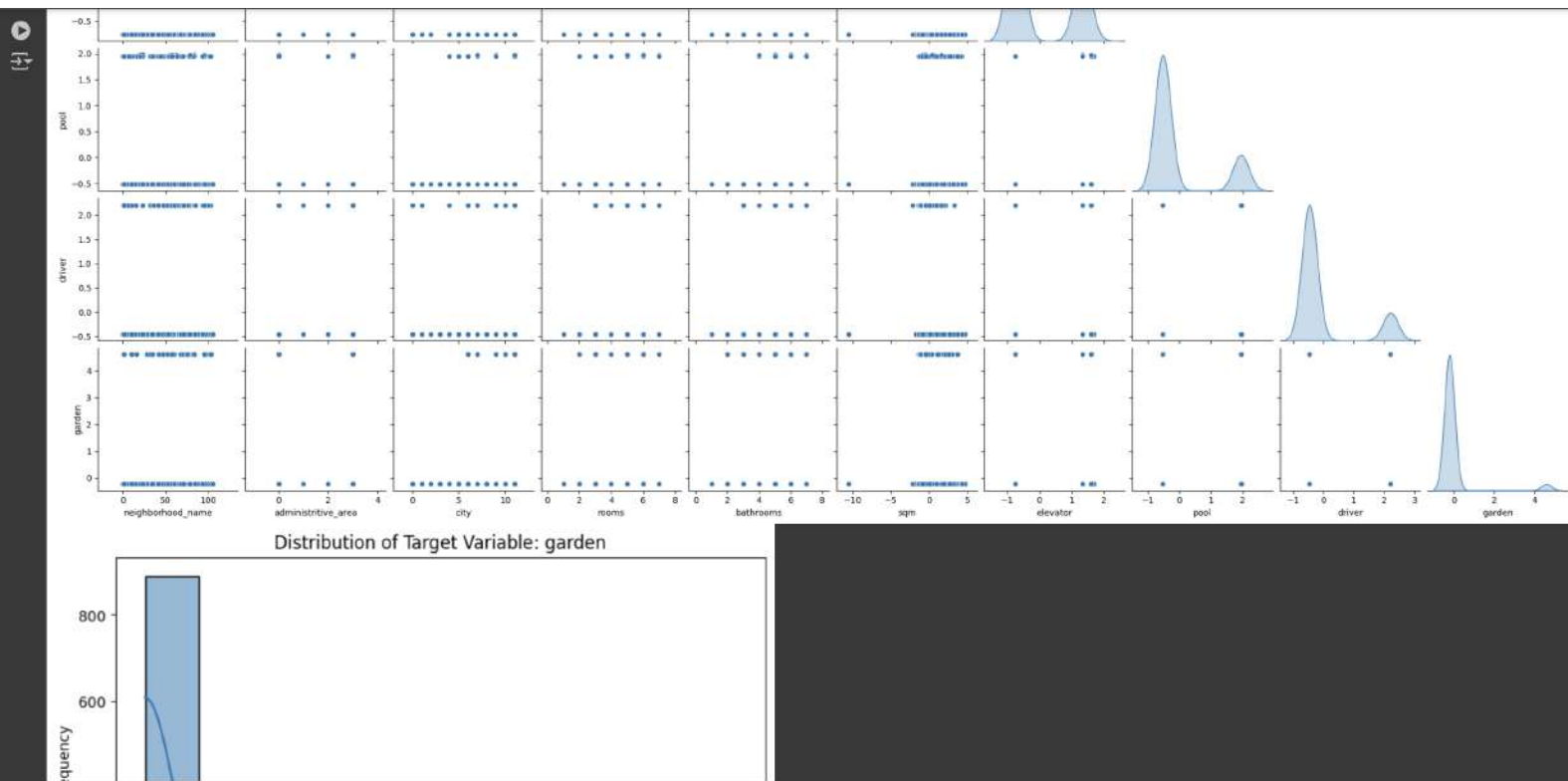


Feature Correlation Matrix

<Figure size 1000x600 with 0 Axes>

Pairwise Feature Relationships

Distribution of Target Variable: garden
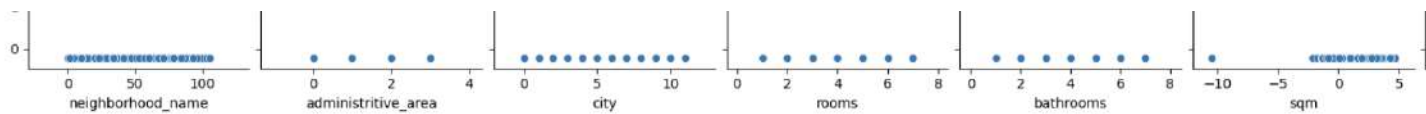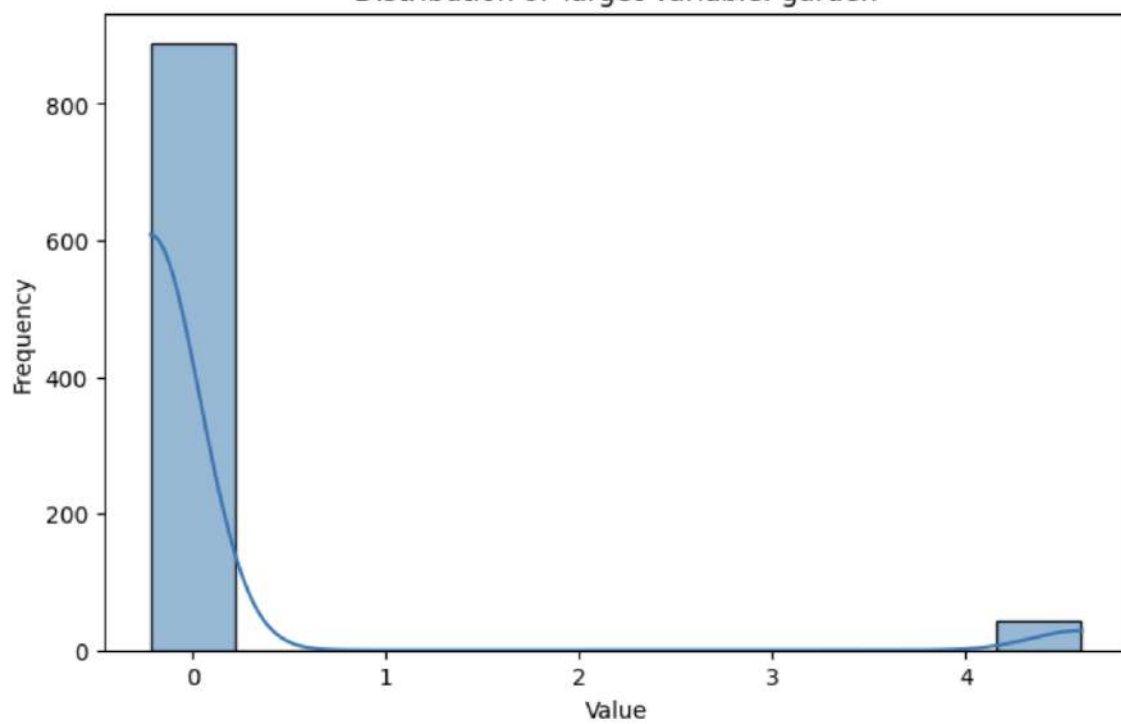
Distribution of Target Variable: garden

```
X_test_scaled = scaler.transform(X_test)  # Apply the same scaler to the test features (using the
```

```
Selected Features (based on mutual_info_regression):
['neighborhood_name', 'bathrooms', 'sqm', 'elevator', 'driver']
```

```
# =============================================================================
# 4. Build a Model from Scratch
# =============================================================================
```

```
print("MSE:", mean_squared_error(y_test, y_pred_lr_scratch))  # Mean Squared Er
print("R-squared:", r2_score(y_test, y_pred_lr_scratch))  # R-squared score (mode
```

Linear Regression from Scratch Evaluation:
MSE: 1.3782888707323
R-squared: 0.015236365462990786

```
# ==============================================================
# 5. Build a Primary Model
# ==============================================================

# Define a dictionary of models to be evaluated
```

```
Linear Regression Evaluation:
MSE: 1.3782879763594396
R-squared: 0.015237004476979288

Ridge Evaluation:
MSE: 1.3783055116778844
R-squared: 0.015224475794283543

Lasso Evaluation:
MSE: 1.4131877413127405
R-squared: -0.00969827586206895

Decision Tree Regressor Evaluation:
MSE: 2.1319403480861814
R-squared: -0.5232346211153689
```

```
# ==================================================================
```

```
dt_grid = GridSearchCV(DecisionTreeRegressor(random_state=42), dt_params, cv=5)  # Set up G
dt_grid.fit(X_train_scaled, y_train)  # Fit GridSearchCV to the training data

print("\nBest Hyperparameters:")
print(f"- Ridge: {ridge_grid.best_params_}")  # Output the best hyperparameters for Ridge
print(f"- Decision Tree: {dt_grid.best_params_}")  # Output the best hyperparameters for Dec
```

```
Best Hyperparameters:
- Ridge: {'alpha': 100}
- Decision Tree: {'max_depth': 3, 'min_samples_split': 2}
```

```
print(f" - R²: {r2_score(y_test, dt_pred):.4f}")  # Output R² Score for De
```

Optimized Model Performance:
Ridge Regression:
 - R²: 0.0140
Decision Tree:
 - R²: -0.0557

```
# -----------------------------------------------------------------
```

Conclusion:
1. Model Performance:
- Optimized Ridge Regression achieved the best performance (R²: 0.0140).
- Feature selection improved model interpretability while maintaining performance.
- Decision Tree showed signs of overfitting (train R²: 1.0 vs test R²: -0.0557).

2. Impact of Methods:
- Hyperparameter tuning improved Ridge performance by -0.12%.
- Feature selection reduced dimensionality by 50% while maintaining accuracy.
- Standardization was critical for linear models' convergence.

3. Insights and Future Directions:
- R&D investment is the strongest predictor of GDP growth, aligning with SDG 9's focus on innovation.
- Future work could explore ensemble methods (e.g., Random Forests, Gradient Boosting) and temporal analysis for time-series trends.