```
!pip install xgboost seaborn
```

```
Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packages (2.1.4)
Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.0.2)
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.21.5)
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from xgboost) (1.14.1)
Requirement already satisfied: pandas>=1.2 in /usr/local/lib/python3.11/dist-packages (from seaborn) (2.2.2)
Requirement already satisfied: matplotlib!=3.6.1,>=3.4 in /usr/local/lib/python3.11/dist-packages (from seaborn) (3.10.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (0.12
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (2
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (11.1.0
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seaborn)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.11/dist-packages (from matplotlib!=3.6.1,>=3.4->seabor
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.2->seaborn) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.2->seaborn) (2025.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3
```

```python
#  Import libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, MinMaxScaler
from sklearn.metrics import mean_squared_error
from xgboost import XGBRegressor
```

```python
import pandas as pd
path="/content/AmesHousing (1).csv"
data=pd.read_csv(path)
print(data)
```

```
      Order        PID  MS SubClass MS Zoning  Lot Frontage  Lot Area Street  \
0         1  526301100           20        RL         141.0     31770   Pave
1         2  526350040           20        RH          80.0     11622   Pave
2         3  526351010           20        RL          81.0     14267   Pave
3         4  526353030           20        RL          93.0     11160   Pave
4         5  527105010           60        RL          74.0     13830   Pave
...     ...        ...          ...       ...           ...       ...    ...
2925   2926  923275080           80        RL          37.0      7937   Pave
2926   2927  923276100           20        RL           NaN      8885   Pave
2927   2928  923400125           85        RL          62.0     10441   Pave
2928   2929  924100070           20        RL          77.0     10010   Pave
2929   2930  924151050           60        RL          74.0      9627   Pave

      Alley Lot Shape Land Contour  ... Pool Area Pool QC  Fence Misc Feature  \
0       NaN       IR1          Lvl  ...         0     NaN    NaN          NaN
1       NaN       Reg          Lvl  ...         0     NaN  MnPrv          NaN
2       NaN       IR1          Lvl  ...         0     NaN    NaN         Gar2
3       NaN       Reg          Lvl  ...         0     NaN    NaN          NaN
4       NaN       IR1          Lvl  ...         0     NaN  MnPrv          NaN
...     ...       ...          ...  ...       ...     ...    ...          ...
2925    NaN       IR1          Lvl  ...         0     NaN  GdPrv          NaN
2926    NaN       IR1          Low  ...         0     NaN  MnPrv          NaN
2927    NaN       Reg          Lvl  ...         0     NaN  MnPrv         Shed
2928    NaN       Reg          Lvl  ...         0     NaN    NaN          NaN
2929    NaN       Reg          Lvl  ...         0     NaN    NaN          NaN

      Misc Val Mo Sold Yr Sold Sale Type  Sale Condition  SalePrice
0            0       5    2010        WD          Normal     215000
1            0       6    2010        WD          Normal     105000
2        12500       6    2010        WD          Normal     172000
3            0       4    2010        WD          Normal     244000
4            0       3    2010        WD          Normal     189900
...        ...     ...     ...       ...             ...        ...
2925         0       3    2006        WD          Normal     142500
2926         0       6    2006        WD          Normal     131000
2927       700       7    2006        WD          Normal     132000
2928         0       4    2006        WD          Normal     170000
2929         0      11    2006        WD          Normal     188000

[2930 rows x 82 columns]
```

```python
#  Basic info
print("Dataset shape:", data.shape)
print("Columns:", data.columns.tolist())
```

```
Dataset shape: (2930, 82)
Columns: ['Order', 'PID', 'MS SubClass', 'MS Zoning', 'Lot Frontage', 'Lot Area', 'Street', 'Alley', 'Lot Shape', 'Land Contour', 'L
```

```python
data.head()
```

| | Order | PID | MS SubClass | MS Zoning | Lot Frontage | Lot Area | Street | Alley | Lot Shape | Land Contour | ... | Pool Area | Pool QC | Fence | Misc Feature | Misc Val | Mo Sold |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 526301100 | 20 | RL | 141.0 | 31770 | Pave | NaN | IR1 | Lvl | ... | 0 | NaN | NaN | NaN | 0 | 5 |
| 1 | 2 | 526350040 | 20 | RH | 80.0 | 11622 | Pave | NaN | Reg | Lvl | ... | 0 | NaN | MnPrv | NaN | 0 | 6 |
| 2 | 3 | 526351010 | 20 | RL | 81.0 | 14267 | Pave | NaN | IR1 | Lvl | ... | 0 | NaN | NaN | Gar2 | 12500 | 6 |
| 3 | 4 | 526353030 | 20 | RL | 93.0 | 11160 | Pave | NaN | Reg | Lvl | ... | 0 | NaN | NaN | NaN | 0 | 4 |
| 4 | 5 | 527105010 | 60 | RL | 74.0 | 13830 | Pave | NaN | IR1 | Lvl | ... | 0 | NaN | MnPrv | NaN | 0 | 3 |

5 rows × 82 columns

```python
data.isnull().sum()
```

| | 0 |
|---|---|
| Order | 0 |
| PID | 0 |
| MS SubClass | 0 |
| MS Zoning | 0 |
| Lot Frontage | 490 |
| ... | ... |
| Mo Sold | 0 |
| Yr Sold | 0 |
| Sale Type | 0 |
| Sale Condition | 0 |
| SalePrice | 0 |

82 rows × 1 columns

dtype: int64

```python
data['Lot Frontage'].fillna(data['Lot Frontage'].mode()[0],inplace=True)
```

<ipython-input-13-951fcd9818e3>:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

  data['Lot Frontage'].fillna(data['Lot Frontage'].mode()[0],inplace=True)

```python
data.isnull().sum()
```

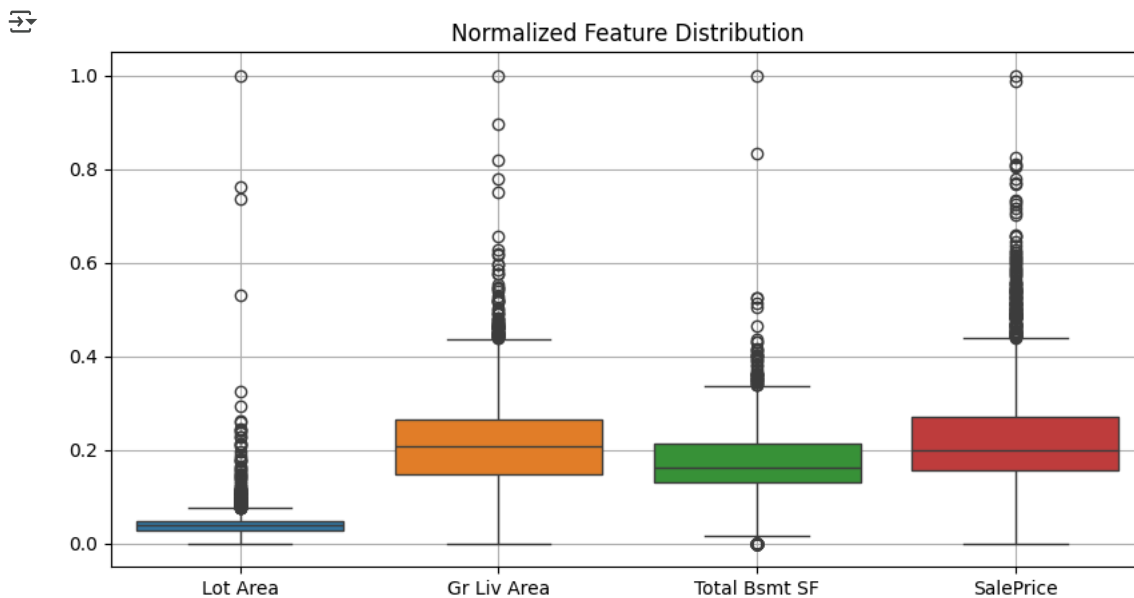|  | 0 |
| --- | --- |
| **Order** | 0 |
| **PID** | 0 |
| **MS SubClass** | 0 |
| **MS Zoning** | 0 |
| **Lot Frontage** | 0 |
| ... | ... |
| **Mo Sold** | 0 |
| **Yr Sold** | 0 |
| **Sale Type** | 0 |
| **Sale Condition** | 0 |
| **SalePrice** | 0 |

82 rows × 1 columns

dtype: int64

```python
# Encode categorical variables
le = LabelEncoder()
for col in data.columns:
    if data[col].dtype == "object":
        data[col] = le.fit_transform(data[col])
```

```python
#  Visualize normalized features (for learning)
scaler = MinMaxScaler()
scaled_data = pd.DataFrame(scaler.fit_transform(data), columns=data.columns)

plot_cols = ['Lot Area', 'Gr Liv Area', 'Total Bsmt SF', 'SalePrice']
plot_cols = [col for col in plot_cols if col in scaled_data.columns]

plt.figure(figsize=(10, 5))
sns.boxplot(data=scaled_data[plot_cols])
plt.title("Normalized Feature Distribution")
plt.grid()
plt.show()
```



```python
#  Split features and target
target_col = 'SalePrice'
X = data.drop(target_col, axis=1)
y = data[target_col]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train XGBoost Model with Tuned Parameters
xgb_model = XGBRegressor(
```

```
    n_estimators=250,
    learning_rate=0.07,
    max_depth=6,
    subsample=0.8,
    colsample_bytree=0.8,
    random_state=42
)

xgb_model.fit(X_train, y_train)
```

```
                                    XGBRegressor                                    ⓘ

XGBRegressor(base_score=None, booster=None, callbacks=None,
             colsample_bylevel=None, colsample_bynode=None,
             colsample_bytree=0.8, device=None, early_stopping_rounds=None,
             enable_categorical=False, eval_metric=None, feature_types=None,
             gamma=None, grow_policy=None, importance_type=None,
             interaction_constraints=None, learning_rate=0.07, max_bin=None,
             max_cat_threshold=None, max_cat_to_onehot=None,
             max_delta_step=None, max_depth=6, max_leaves=None,
             min_child_weight=None, missing=nan, monotone_constraints=None,
             multi_strategy=None, n_estimators=250, n_jobs=None,
             num_parallel_tree=None, random_state=42, ...)
```

```
# Evaluate the model
y_pred = xgb_model.predict(X_test)
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print(f" RMSE of the XGBoost model: {rmse:.2f}")
```
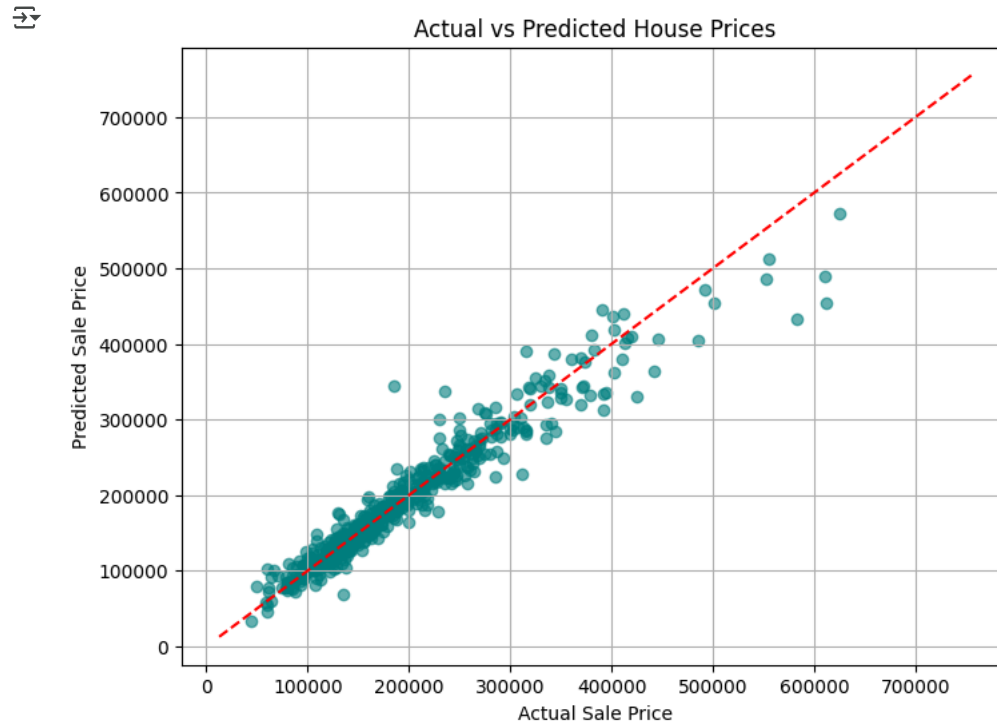
```
    RMSE of the XGBoost model: 23183.27
```

```
# Plot predicted vs actual prices
plt.figure(figsize=(8, 6))
plt.scatter(y_test, y_pred, alpha=0.6, color='teal')
plt.xlabel("Actual Sale Price")
plt.ylabel("Predicted Sale Price")
plt.title("Actual vs Predicted House Prices")
plt.grid(True)
plt.plot([y.min(), y.max()], [y.min(), y.max()], color='red', linestyle='--')
plt.show()
```
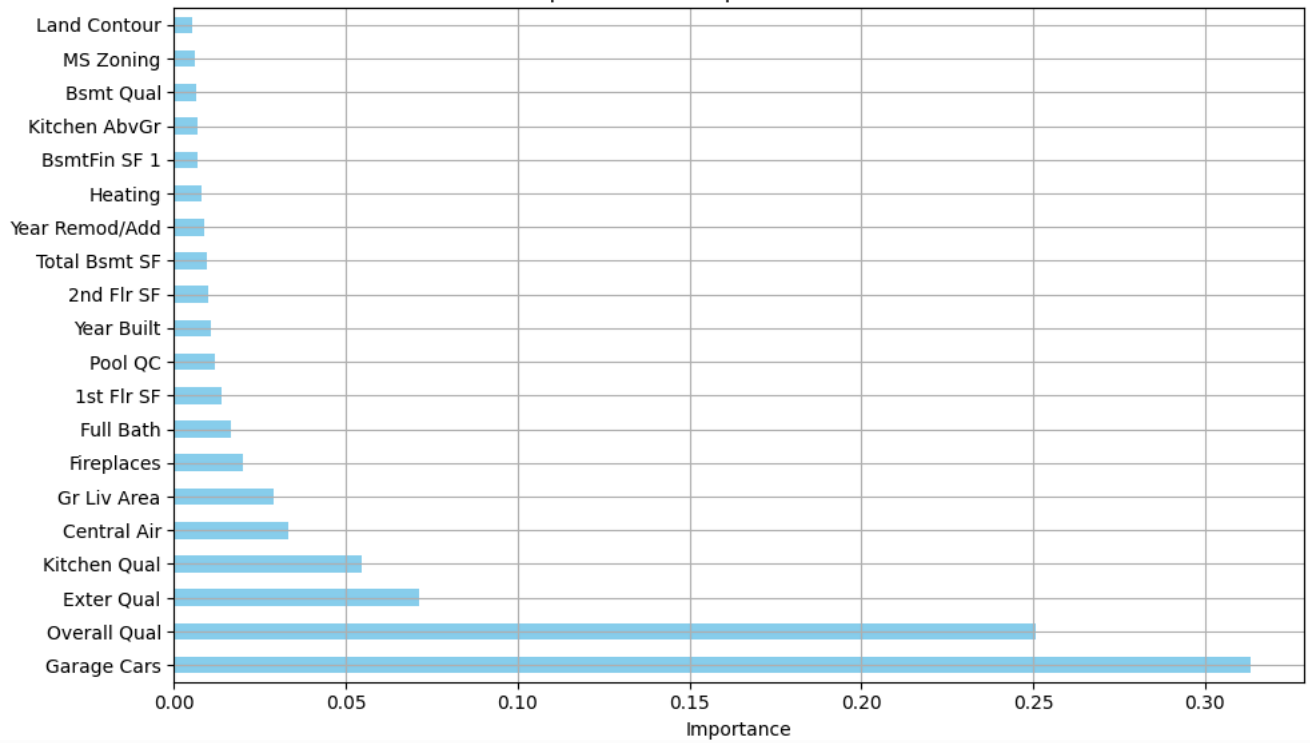


```
# : Feature Importance
plt.figure(figsize=(10, 6))
xgb_model.feature_importances_.argsort()
feat_imp = pd.Series(xgb_model.feature_importances_, index=X.columns)
feat_imp.nlargest(20).plot(kind='barh', color='skyblue')
plt.title("Top 20 Feature Importances - XGBoost")
plt.xlabel("Importance")
plt.grid()
plt.tight_layout()
plt.show()
```
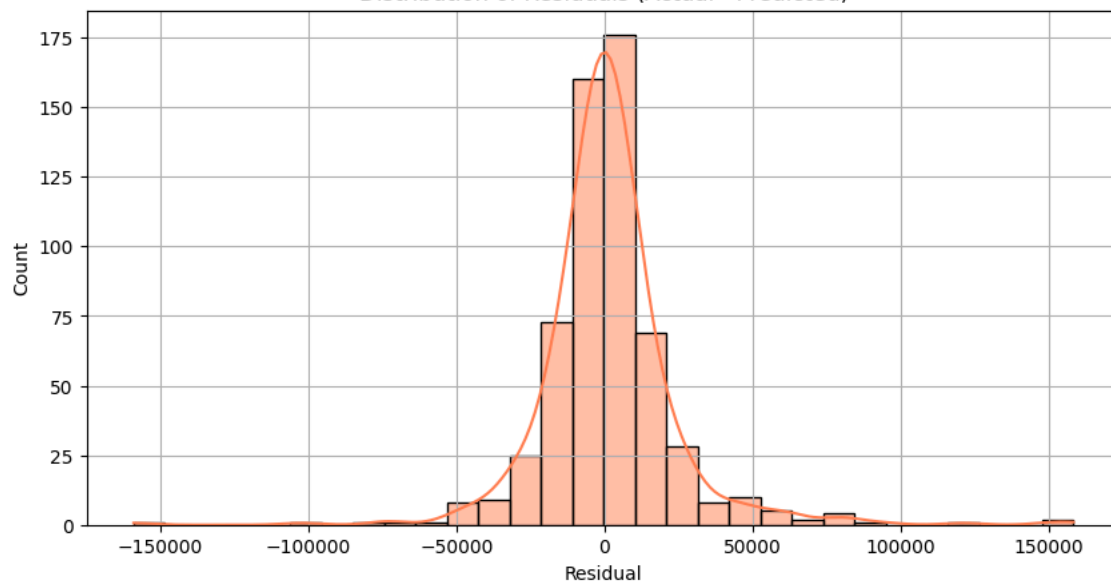
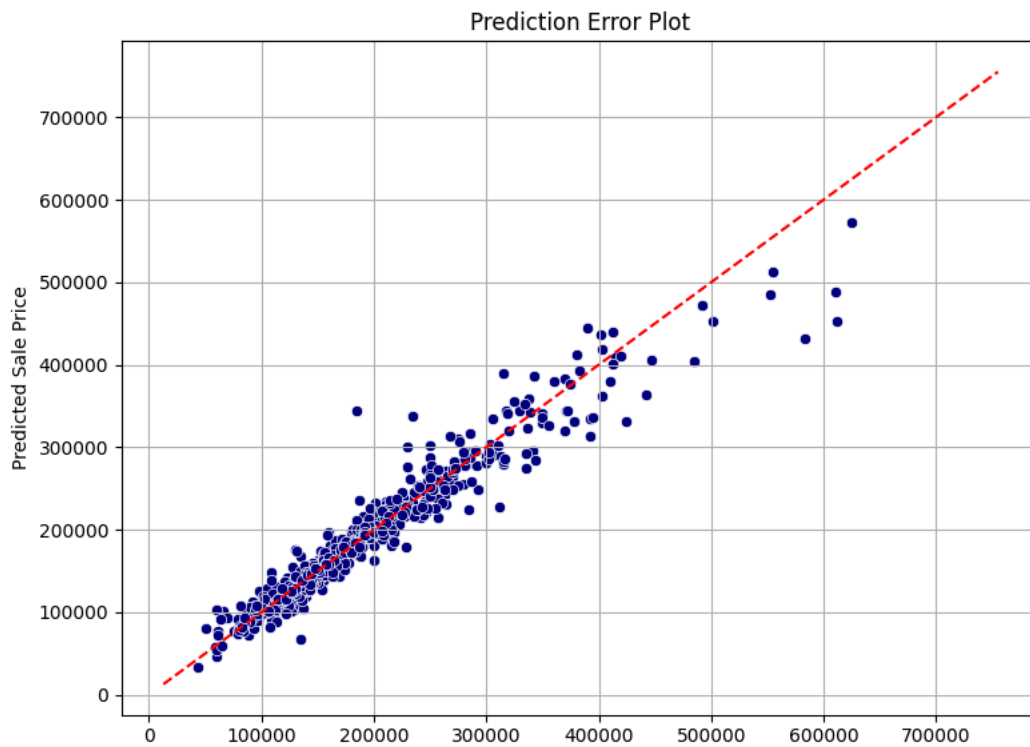## Top 20 Feature Importances - XGBoost



```
#  Residuals Plot
residuals = y_test - y_pred
plt.figure(figsize=(10, 5))
sns.histplot(residuals, bins=30, kde=True, color='coral')
plt.title("Distribution of Residuals (Actual - Predicted)")
plt.xlabel("Residual")
plt.grid()
plt.show()
```
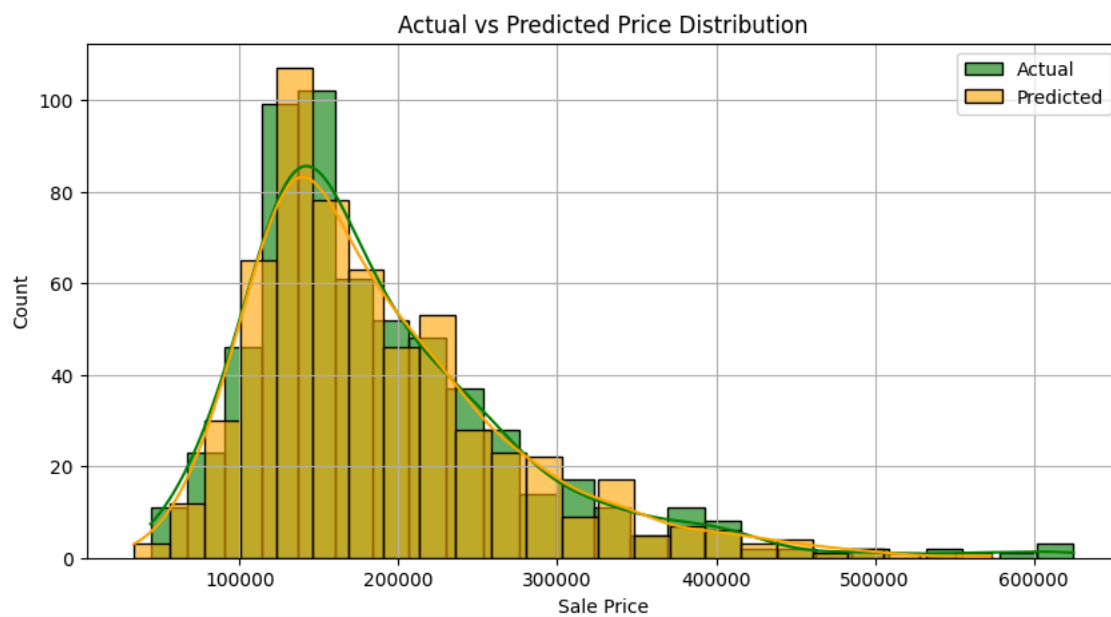
## Distribution of Residuals (Actual - Predicted)



```
#  Prediction Error (Actual vs Predicted)
plt.figure(figsize=(8, 6))
sns.scatterplot(x=y_test, y=y_pred, color='navy')
plt.plot([y.min(), y.max()], [y.min(), y.max()], color='red', linestyle='--')
plt.xlabel("Actual Sale Price")
plt.ylabel("Predicted Sale Price")
plt.title("Prediction Error Plot")
plt.grid(True)
plt.tight_layout()
plt.show()
```

## Prediction Error Plot



```python
# Overlayed Histogram for Actual vs Predicted
plt.figure(figsize=(10, 5))
sns.histplot(y_test, label="Actual", kde=True, color='green', alpha=0.6)
sns.histplot(y_pred, label="Predicted", kde=True, color='orange', alpha=0.6)
plt.title("Actual vs Predicted Price Distribution")
plt.xlabel("Sale Price")
plt.legend()
plt.grid()
plt.show()
```

## Actual vs Predicted Price Distribution



```python
from sklearn.metrics import r2_score, mean_squared_error
import numpy as np

# Evaluate the model
r2 = r2_score(y_test, y_pred)
rmse = np.sqrt(mean_squared_error(y_test, y_pred))

print(" Final Model Evaluation:")
print(f" R² Score (Accuracy): {r2:.4f}")
print(f" RMSE (Error in Price): {rmse:.2f}")
```

```
 Final Model Evaluation:
  R² Score (Accuracy): 0.9330
  RMSE (Error in Price): 23183.27
```