PROJECT REPORT

# Investing in Nashville using Machine Learning

*Author:*

Syed Faizan

December 7th, 2024

# Contents

# 7 Part 6: Model Comparison and Recommendations

# 8 References

# Introduction

The burgeoning real estate market in Nashville presents a significant opportunity for investment, requiring informed decision-making grounded in robust data analysis. This report aims to construct and evaluate predictive models to identify under- or overvalued properties, guiding strategic investments. Using a dataset of recent property sales, we examine the variable *Sale Price Compared to Value* as a critical indicator of market inefficiencies. A precise and reliable model will not only aid in pinpointing high-value deals but also uncover key factors influencing property valuations.

The analysis begins with comprehensive data cleansing to ensure the integrity and reliability of the dataset, a fundamental step for generating meaningful insights. A logistic regression model is then developed to predict property valuations and to identify the variables that most influence prices. Subsequently, advanced machine learning models—decision tree, random forest, and gradient boosting—are implemented to enhance predictive accuracy. Each model's performance is evaluated using multiple benchmarking metrics, allowing for a systematic comparison of their predictive capabilities.

By synthesizing insights from these models, the report identifies the optimal approach for accurately predicting housing prices and maximizing investment returns. The findings provide actionable recommendations for leveraging data-driven strategies to navigate the competitive Nashville real estate market effectively.

# Part 1: : Data cleansing

## Loading the Dataset and Packages



Figure 1: Loading the Dataset and Packages

## Introduction to the Problem and Dataset Initialization

The above figure provides an overview of the problem statement and the preliminary steps undertaken in the analysis. The objective is to leverage recent sales data in the growing Nashville real estate market to construct predictive models capable of identifying undervalued or overvalued properties. This could aid in making strategic investment decisions. The Python environment is initialized with essential packages such as *pandas* for data manipulation, *scikit-learn* for model development, and *GridSearchCV* for hyperparameter optimization. Initial data loading and preview are executed to verify the successful import of the dataset, confirming readiness for data cleansing and modeling.

# Structure and Columns of the Dataset

```
df.shape
```

```
(22651, 26)
```

```
df.columns.tolist()
```

```
['Unnamed: 0',
 'Parcel ID',
 'Land Use',
 'Property Address',
 'Suite/ Condo   #',
 'Property City',
 'Sale Date',
 'Legal Reference',
 'Sold As Vacant',
 'Multiple Parcels Involved in Sale',
 'City',
 'State',
 'Acreage',
 'Tax District',
 'Neighborhood',
 'Land Value',
 'Building Value',
 'Finished Area',
 'Foundation Type',
 'Year Built',
 'Exterior Wall',
 'Grade',
 'Bedrooms',
 'Full Bath',
 'Half Bath',
 'Sale Price Compared To Value']
```

Figure 2: Structure and Columns of the Dataset

## Dataset Structure and Attributes

The above figure highlights the structural overview and attribute details of the dataset used for analysis. The dataset consists of **22,651 rows** and **26 columns**, as indicated by the output of the *df.shape* command. The *df.columns.tolist()* function provides a comprehensive list of column names, representing various attributes of real estate transactions, such as *Land Use*, *Sale Date*, *Acreage*, *Building Value*, *Bedrooms*, and *Sale Price Compared to Value*. These variables encompass diverse data types and offer critical insights for modeling property valuation. This foundational analysis ensures a structured approach to subsequent data preparation.

## Missing Values

```
                                       Missing Values   Percentage (%)
        Suite/ Condo    #                      22651       100.000000
        Half Bath                                108         0.476800
        Bedrooms                                   3         0.013244
        Property Address                           2         0.008830
        Property City                              2         0.008830
        Full Bath                                  1         0.004415
        Foundation Type                            1         0.004415
        Finished Area                              1         0.004415
        Unnamed: 0                                 0         0.000000
        Land Value                                 0         0.000000
        Grade                                      0         0.000000
        Exterior Wall                              0         0.000000
        Year Built                                 0         0.000000
        Building Value                             0         0.000000
        Tax District                               0         0.000000
        Neighborhood                               0         0.000000
        Parcel ID                                  0         0.000000
        Acreage                                    0         0.000000
        State                                      0         0.000000
        City                                       0         0.000000
        Multiple Parcels Involved in Sale          0         0.000000
        Sold As Vacant                             0         0.000000
        Legal Reference                            0         0.000000
        Sale Date                                  0         0.000000
        Land Use                                   0         0.000000
        Sale Price Compared To Value               0         0.000000
```

Figure 3: Missing Values in the Dataset

## Missing Values Preprocessing



**Remove the Suite/ condo column entirely and the other rows with missing values as they are small in number**

```
]:  # Remove the 'Suite/ Condo #' column entirely
    df_cleaned = df.drop(columns=['Suite/ Condo  #'])

    # Drop rows with missing values
    df_cleaned = df_cleaned.dropna()

    # Verify that there are no missing values remaining
    missing_values_after_cleaning = df_cleaned.isnull().sum()

    missing_values_after_cleaning
```

Figure 4: Preprocessing Missing Values

## Handling Missing Values in the Dataset

The above figures illustrate the identification and treatment of missing values within the dataset. The first figure provides a comprehensive summary of missing values across all columns, highlighting both the absolute count and percentage of missing data for each feature. Notably, the *Suite/Condo #* column exhibits a 100% missing value rate, rendering

it unsuitable for analysis. Additionally, other columns, such as *Half Bath*, *Bedrooms*, and *Property Address*, contain minimal missing values, with percentages below 0.5%.

The second figure depicts the preprocessing approach undertaken to address these missing values. The *Suite/Condo #* column was entirely removed due to its complete lack of data. Rows containing missing values in other columns were dropped, as their occurrence was infrequent and unlikely to influence the dataset's integrity significantly. The process ensures a refined dataset devoid of missing values, as verified in the final output. This preprocessing enhances data quality, a critical step for subsequent modeling and analysis.

## Removing Redundant Columns



Figure 5: Removing Redundant Columns from the Dataset

The above figure illustrates the removal of redundant columns from the dataset to enhance its relevance and simplify subsequent analyses. Columns deemed non-essential, such as *Parcel ID*, *Property Address*, *Legal Reference*, *Property City*, *City*, *State*, and *Tax District*,

were identified and dropped from the dataset. These columns were excluded due to their limited utility in predicting the target variable, *Sale Price Compared To Value*, or their non-informative nature.

Following this operation, the remaining columns, including features like *Land Use*, *Neighborhood*, *Building Value*, and *Bedrooms*, were retained for further modeling. The resultant dataset contains 17 columns and 22,536 rows, as confirmed by the output. This preprocessing step ensures that the dataset is streamlined, reducing computational complexity while retaining critical information for accurate model development.

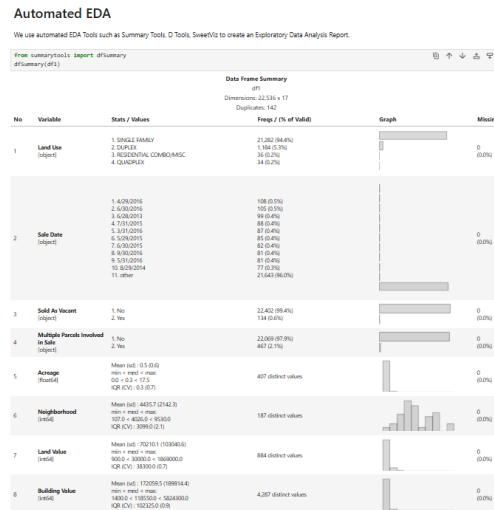## Automated EDA: Summary Statistics (Part 1)



Figure 6: Automated EDA: Summary Statistics (Part 1)
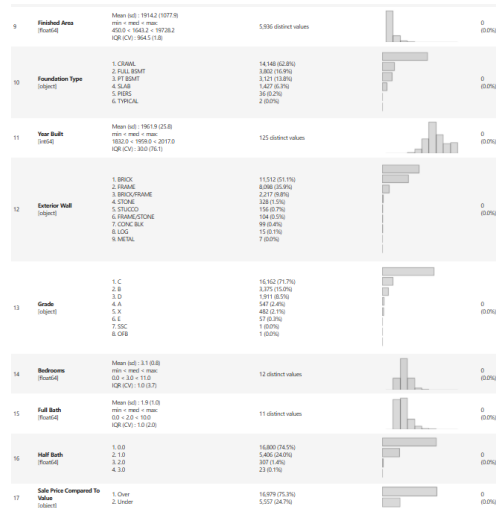
## Automated EDA: Summary Statistics (Part 2)



Figure 7: Automated EDA: Summary Statistics (Part 2)

## Automated Exploratory Data Analysis Summary

The above figures provide a comprehensive summary of the dataset using automated exploratory data analysis (EDA) tools. These tools, such as *Summary Tools* and *SweetViz*, generate an overview of the dataset's variables, frequencies, descriptive statistics, and distributions, facilitating a detailed understanding of the data. The dataset consists of 22,536 observations across 17 variables. For each variable, the figures illustrate the count of unique values, frequencies, percentages, and the absence of missing values, as indicated by a 0% missing data rate across all columns.

Key categorical variables, such as **Land Use**, are distributed among *Single Family* (94.4%), *Duplex* (5%), and other minor categories. Temporal data, represented by **Sale Date**, exhibits a significant proportion of sales concentrated in specific years. Binary variables like **Sold As Vacant** and **Multiple Parcels Involved in Sale** highlight the predominance of properties sold as *not vacant* (99.4%) and as individual parcels (97.9%).

Numerical variables, such as **Acreage**, **Land Value**, and **Building Value**, are described by their mean, median, standard deviation, and interquartile ranges. Notably, **Acreage** shows a highly skewed distribution with 407 distinct values, while **Building Value** reflects significant variation, as indicated by its mean of 172,095.9 and an interquartile range of

103,225.0.

Structural attributes, such as **Foundation Type** and **Exterior Wall**, reveal specific material distributions, with *Crawl* (62.8%) and *Brick* (51.1%) dominating their respective categories. Residential features, including **Bedrooms** and **Bathrooms**, show a reasonable distribution across their ranges, while **Sale Price Compared To Value** captures the proportion of properties classified as *Overvalued* (75.3%) or *Undervalued* (24.7%).

Overall, the figures highlight the dataset's robustness and provide valuable insights for subsequent analyses, particularly for feature engineering and predictive modeling.

## Feature Engineering



```
Convert 'Sale Date' into a date using datetime module, then convert it into an
integer column called 'Time since sale' which gives time since sale in days. Also,
convert 'Year Built' into age.

from datetime import datetime

# Convert 'Sale Date' to datetime format
df1['Sale Date'] = pd.to_datetime(df1['Sale Date'], errors='coerce')

# Create a new column 'timesincesale' as the difference between the present date and 'Sale Date'
df1['timesincesale'] = (datetime.now() - df1['Sale Date']).dt.days

# Remove the old 'Sale Date' column
df1 = df1.drop(columns=['Sale Date'])

# Transform 'Year Built' column into an 'Age' column
df1['Age'] = 2024 - df1['Year Built']

# Drop the original 'Year Built' column
df1 = df1.drop(columns=['Year Built'])

# Display the first few rows to verify the transformation
df1.head()


# Display the first few rows of the updated DataFrame
df1.head()
```

Figure 8: Feature Engineering Process

The above figure illustrates the transformation of temporal variables in the dataset using Python's *datetime* module. The **Sale Date** column was first converted to a *datetime* format to facilitate further computations. A new variable, **timesincesale**, was derived to represent the number of days elapsed since the sale date. This was calculated by subtracting the sale date from the current date in days. Subsequently, the original **Sale Date** column was removed from the dataset.

Similarly, the **Year Built** variable was transformed to compute the **Age** of properties, defined as the difference between the current year (2024) and the year of construction. The original **Year Built** column was subsequently dropped to avoid redundancy. These transformations were performed to ensure temporal variables are numerically encoded for effective

utilization in predictive modeling. The transformed dataset was previewed to validate the modifications, confirming the successful creation of the new variables.

## Ordinal Encoding



```
Use ordinal encoding to encode the 'object' class variables as there is an ordinal relationship between the levels of the categorical variables and one hot encoding
produces hyperdimnetionality

from sklearn.preprocessing import OrdinalEncoder

# Select columns with object data type
object_columns = df1.select_dtypes(include=['object']).columns

# Apply ordinal encoding to these columns
encoder = OrdinalEncoder()
df1[object_columns] = encoder.fit_transform(df1[object_columns])

# Display the first few rows of the updated DataFrame
df1.head()
```

Figure 9: Ordinal Encoding of Categorical Variables

The above figure illustrates the application of **ordinal encoding** to transform categorical variables into numerical representations. Using Python's *sklearn.preprocessing.OrdinalEncoder*, all columns with a data type of *object* were identified and selected for encoding. Ordinal encoding was chosen as it preserves the inherent order among the levels of ordinal categorical variables, which is critical for maintaining data integrity and interpretability in downstream analysis.

This approach was preferred over one-hot encoding to mitigate the issue of *dimensionality explosion*, especially for categorical variables with multiple levels. The encoded values were assigned back to their respective columns, replacing the original categorical representations. A preview of the updated dataset was generated to validate the encoding process, confirming the successful transformation of all object-type columns into numerical values suitable for machine learning algorithms.

## Feature Scaling



```python
from sklearn.preprocessing import MinMaxScaler

# Select numeric columns for scaling
numeric_columns = df1.select_dtypes(include=['float64', 'int64']).columns

# Apply MinMaxScaler for normalization
scaler = MinMaxScaler()
df1[numeric_columns] = scaler.fit_transform(df1[numeric_columns])

# Display the first few rows of the normalized DataFrame
df1.head()
```

Figure 10: Feature Scaling Using Normalization

The above figure illustrates the process of scaling numerical features using the **MinMaxScaler** from Python's *sklearn.preprocessing* library. All numerical columns, identified by their data types (*float64* and *int64*), were selected for scaling. The *MinMaxScaler* was applied to transform these columns into a normalized range between 0 and 1, preserving the distribution while ensuring that all features have comparable scales.

Normalization is a critical preprocessing step for many machine learning algorithms, particularly those sensitive to feature magnitude, such as logistic regression and distance-based models. The normalized dataset was inspected to confirm the successful application of the scaling operation.
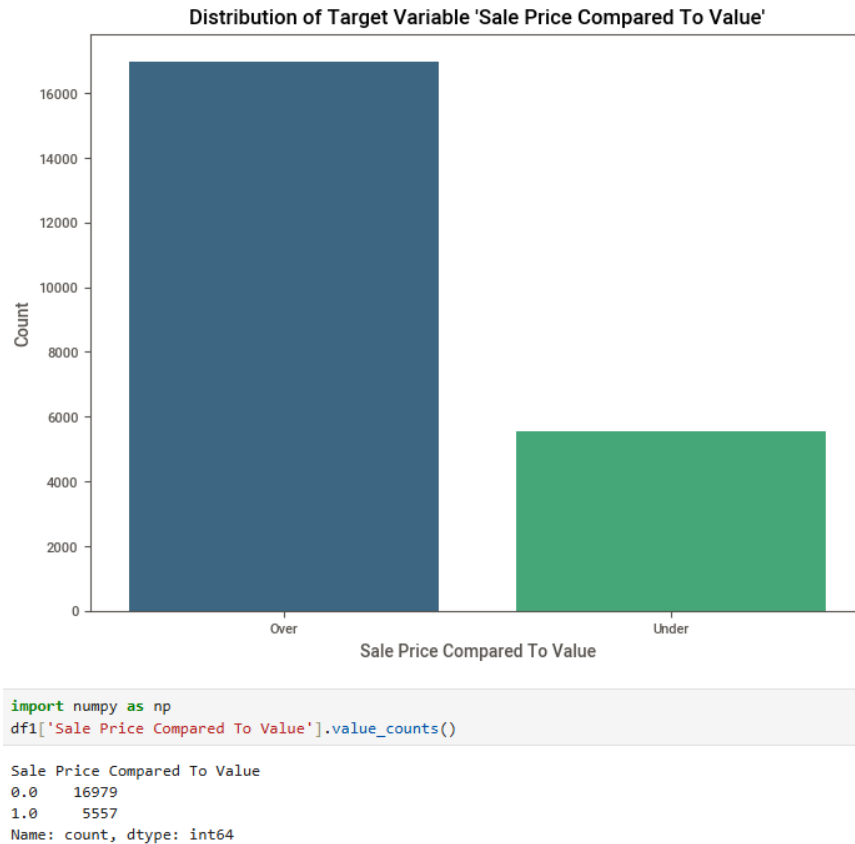
## Imbalance of Dataset



Figure 11: Imbalance of the Target Variable

The above figure depicts the distribution of the target variable, **Sale Price Compared To Value**, highlighting a notable class imbalance. The two classes, *Over* and *Under*, are represented by their respective counts: **16,979** instances of the *Over* class and **5,557** instances of the *Under* class. This imbalance is significant and suggests the need for techniques to address it, such as oversampling, undersampling, or applying class weights during model training. Such methods are crucial to ensure that predictive models do not become biased towards the majority class and maintain balanced performance across both classes.

## Outlier Detection

**Outlier detection**

```python
# Specify the numerical features explicitly
numerical_features = [
    'Acreage', 'Land Value', 'Building Value',
    'Finished Area', 'Bedrooms', 'Full Bath', 'Half Bath', 'timesincesale', 'Age'
]

# Create boxplots for each numerical feature
for column in numerical_features:
    plt.figure(figsize=(8, 4))
    plt.boxplot(df1[column], vert=False, patch_artist=True)
    plt.title(f'Boxplot of {column}')
    plt.xlabel(column)
    plt.show()

# Identify outliers using the 1.5 IQR rule for each specified numerical feature
outliers = {}
for column in numerical_features:
    Q1 = df1[column].quantile(0.25)
    Q3 = df1[column].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers[column] = df1[(df1[column] < lower_bound) | (df1[column] > upper_bound)][column]

# Display the number of outliers for each numerical feature
outlier_summary = {col: len(outliers[col]) for col in numerical_features}
outlier_summary
```

Figure 12: Outlier Detection Process

The above figure demonstrates the process of identifying outliers in numerical features using the **1.5 IQR rule**. Initially, numerical features such as *Acreage, Land Value, Building Value, Finished Area, Bedrooms, Full Bath, Half Bath, timesince sale, and Age* are specified. Boxplots are generated for each feature to visually assess the presence of outliers. Subsequently, the interquartile range (IQR) is calculated for each feature, with the lower and upper bounds defined as *Q1 - 1.5\*IQR* and *Q3 + 1.5\*IQR*, respectively. Values falling outside these bounds are classified as outliers. The final summary provides the count of outliers for each numerical feature.

**(a)** Boxplot of Acreage  **(b)** Boxplot of Land Value  **(c)** Boxplot of Building Value

**(d)** Boxplot of Finished Area  **(e)** Boxplot of Bedrooms  **(f)** Boxplot of Full Bath

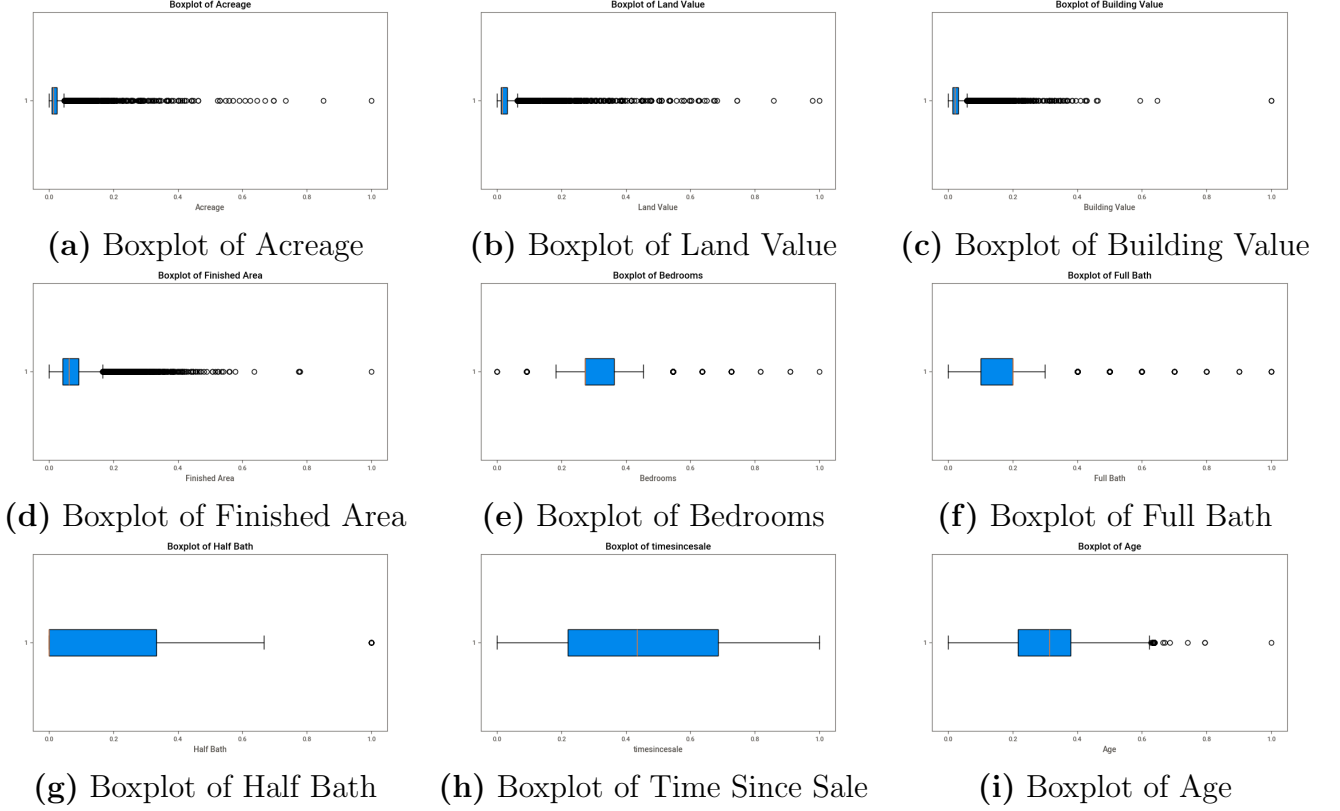**(g)** Boxplot of Half Bath  **(h)** Boxplot of Time Since Sale  **(i)** Boxplot of Age

Figure 13: Boxplots of various features in a 3x3 grid layout.

The above figures illustrate box plots for various numerical features in the dataset, providing insights into their respective distributions, central tendencies, and the presence of outliers. **Acreage**, **Land Value**, **Building Value**, and **Finished Area** exhibit substantial skewness with multiple outliers extending beyond the upper range, as evidenced by individual data points. **Bedrooms** and **Full Bath** display a narrower spread with relatively fewer outliers, suggesting a more uniform distribution.

**Half Bath** shows minimal variability with most data points concentrated at lower values, while a single prominent outlier is identified. The **timesince sale** feature is uniformly distributed without significant outliers, signifying consistent intervals. **Age** reveals a typical distribution, with a central clustering of data points but some outliers extending beyond the interquartile range, potentially indicating properties of unusual ages.

These box plots provide a visual confirmation of outliers and skewed distributions, guiding the data preprocessing steps. The presence of outliers necessitates careful consideration, as they may influence the models' predictive performance. This analysis aids in determining

appropriate transformations or filtering criteria to ensure robust and reliable modeling.

## Outliers Count

```
{'Acreage': 3036,
 'Land Value': 4175,
 'Building Value': 1927,
 'Finished Area': 1380,
 'Bedrooms': 332,
 'Full Bath': 1327,
 'Half Bath': 23,
 'timesincesale': 0,
 'Age': 168}
```

Figure 14: Count of Outliers in Numerical Features

The above figure summarizes the count of outliers identified in various numerical features of the dataset using the *Interquartile Range (IQR)* method. **Land Value** and **Acreage** demonstrate the highest number of outliers, with 4,175 and 3,036, respectively, indicating significant variability in these attributes. **Building Value** and **Finished Area** also exhibit a notable presence of outliers, with counts of 1,927 and 1,380, respectively.

In contrast, **Half Bath** and **timesince sale** have minimal or no outliers, highlighting their relative uniformity. **Age** and **Bedrooms** display moderate outlier counts at 168 and 332, respectively, while **Full Bath** includes 1,327 outliers. This analysis highlights the need for appropriate handling strategies to mitigate the potential influence of extreme values during model development.



**(a)** Density Plot of Acreage    **(b)** Density Plot of Land Value    **(c)** Density Plot of Building Value

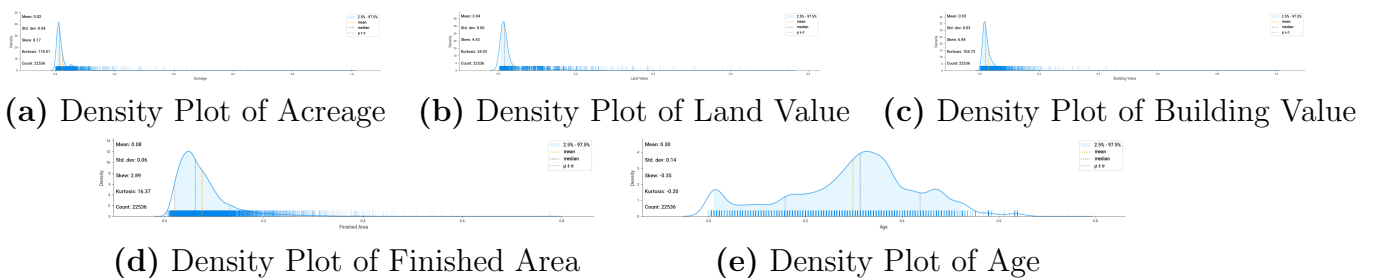**(d)** Density Plot of Finished Area    **(e)** Density Plot of Age

Figure 15: Density plots of various features arranged in a 3x2 grid layout.

The above figures present density plots for key numerical variables, including *Acreage*, *Land Value*, *Building Value*, *Finished Area*, and *Age*. Each plot provides a visual representation of the distribution, along with key summary statistics such as the **mean**, **median**, **standard deviation**, and measures of **skewness** and **kurtosis**.

The density plot for *Acreage* reveals a highly right-skewed distribution with most data concentrated near zero, indicative of small lot sizes. Similarly, *Land Value* and *Building Value* exhibit significant right skewness, suggesting a few properties with exceptionally high values compared to the majority.

In contrast, *Finished Area* shows a less pronounced skew, although larger areas remain relatively rare. The *Age* variable exhibits a bimodal distribution, reflecting distinct property age groups, possibly corresponding to different construction periods.

These visualizations highlight the inherent asymmetry in the data and the prevalence of extreme values, which could influence statistical modeling. Skewness and high kurtosis values for certain features emphasize the need for potential transformations or robust modeling techniques to mitigate their impact. The **mean** and *median* comparisons further suggest potential outliers, requiring attention during preprocessing.
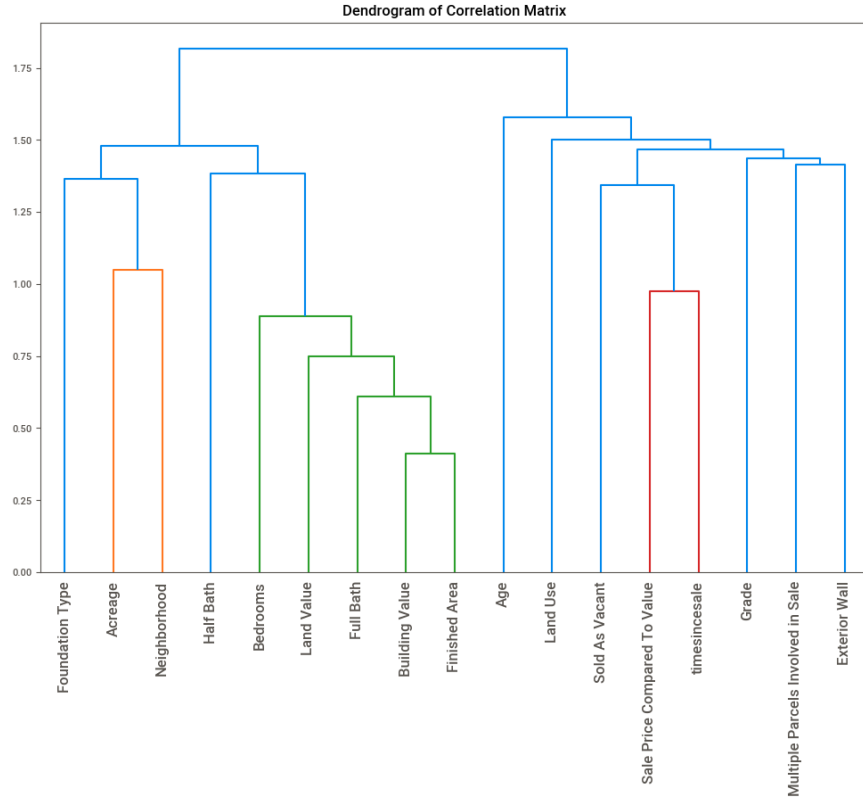
# Dendrogram of Correlation Matrix



Figure 16: Dendrogram of Correlation Matrix

The above figure illustrates a dendrogram derived from the correlation matrix of the dataset's numerical features, highlighting the hierarchical clustering of variables based on their correlation. The horizontal axis represents the variables, while the vertical axis denotes the linkage distance, indicating the degree of similarity between clusters.

Key clusters include *Building Value*, *Land Value*, *Finished Area*, and *Bedrooms*, which are closely associated, suggesting a shared influence on property valuation. Similarly, *Acreage* and *Neighborhood* form another cluster, reflecting geographical and size-related factors.

The dendrogram provides insights into the multicollinearity among features, enabling effective dimensionality reduction or feature selection. Features such as *Sale Price Compared To Value* exhibit independence from most clusters, affirming its role as a distinct target variable. This visualization is crucial for optimizing model performance and ensuring interpretability by selecting uncorrelated predictors.
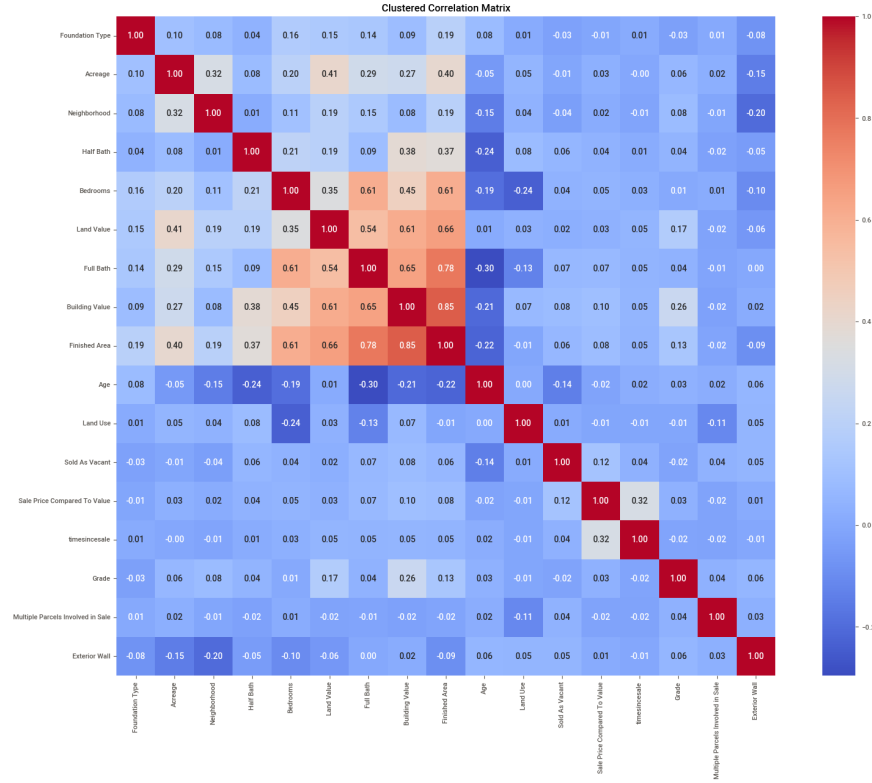
## Clustered Correlation Matrix



Figure 17: Clustered Correlation Matrix

The above figure presents a clustered correlation matrix, providing a detailed visualization of the relationships between numerical features in the dataset. The color gradient, ranging from blue to red, represents the strength and direction of the correlation, where darker shades of red indicate a strong positive correlation, and darker shades of blue signify a strong negative correlation.

Notable correlations include a high positive relationship between *Building Value* and *Finished Area* ($r = 0.85$), as well as between *Full Bath* and *Finished Area* ($r = 0.78$), highlighting their potential influence on housing valuation. Additionally, *Acreage* exhibits moderate correlations with *Neighborhood* and *Land Value*, reflecting geographical and size-related factors.

This visualization aids in identifying feature interactions, reducing multicollinearity, and improving model interpretability. Weak correlations with the target variable, *Sale Price Compared To Value*, suggest the need for additional feature engineering or non-linear mod-

eling techniques.

## Assigning Target Variable



Figure 18: Assigning the Target Variable

The above figure illustrates the project of the target variable **Sale Price Compared To Value** to a new column named *target*, following the project rubric. This intuitive nomenclature enhances clarity and simplicity, ensuring seamless integration into subsequent modeling processes. The initial rows are displayed to verify the successful creation and project of the target column.

## Data Splicing and SMOTE Application



Figure 19: Data Splicing and Handling Imbalance Using SMOTE

The above figure illustrates the implementation of **data splicing** and the application of *Synthetic Minority Oversampling Technique* (SMOTE) to address class imbalance in the dataset. Initially, the features ($X$) and target ($y$) variables are separated, followed by stratified splitting of the dataset into training and testing subsets to ensure class distribution consistency.

SMOTE is applied exclusively to the training dataset, effectively generating synthetic samples for the minority class to balance the class proportions. As shown, the size of the training dataset after oversampling ($(27166, 16)$) demonstrates equal representation of classes. This approach enhances model training by mitigating bias towards the majority class, ensuring robustness and improving predictive performance on imbalanced datasets.

# Part 2: Logistic Regression

## Logistic Regression Results

```
Optimization terminated successfully.
        Current function value: 0.608621
        Iterations 7
                    Logit Regression Results
==============================================================================
Dep. Variable:                target   No. Observations:            27166
Model:                         Logit   Df Residuals:                27149
Method:                          MLE   Df Model:                       16
Date:               Fri, 06 Dec 2024   Pseudo R-squ.:              0.1219
Time:                       22:41:47   Log-Likelihood:            -16534.
converged:                      True   LL-Null:                   -18830.
Covariance Type:           nonrobust   LLR p-value:                0.000
==============================================================================
                                    coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
const                            -1.8373      0.119    -15.445      0.000      -2.071      -1.604
Land Use                         -0.1626      0.064     -2.534      0.011      -0.288      -0.037
Sold As Vacant                    2.9829      0.289     10.338      0.000       2.417       3.548
Multiple Parcels Involved in Sale -0.4927     0.103     -4.775      0.000      -0.695      -0.290
Acreage                           2.8655      0.531      5.401      0.000       1.826       3.905
Neighborhood                      0.4492      0.065      6.942      0.000       0.322       0.576
Land Value                       -4.4886      0.363    -12.377      0.000      -5.199      -3.778
Building Value                    9.4861      1.030      9.206      0.000       7.467      11.506
Finished Area                     0.2168      0.685      0.316      0.752      -1.126       1.560
Foundation Type                  -0.1900      0.053     -3.568      0.000      -0.294      -0.086
Exterior Wall                     0.0025      0.068      0.037      0.970      -0.130       0.135
Grade                             0.4632      0.104      4.439      0.000       0.259       0.668
Bedrooms                         -0.8585      0.256     -3.358      0.001      -1.360      -0.357
Full Bath                         0.3476      0.280      1.242      0.214      -0.201       0.896
Half Bath                         0.1634      0.102      1.599      0.110      -0.037       0.364
timesincesale                     2.9297      0.051     57.660      0.000       2.830       3.029
Age                               0.3876      0.110      3.516      0.000       0.172       0.604
==============================================================================
Training Time: 0.0444 seconds
```

Figure 20: Logistic Regression Results Summary

The above figure presents the results of a logistic regression model fitted to predict the **target variable** based on 16 explanatory features. The model successfully converged with a

*pseudo R-squared* value of 0.1219 and a *log-likelihood* of -16534. The **LLR p-value** indicates the overall statistical significance of the model ($p < 0.001$).

The table provides coefficients (*coef*), standard errors (*std err*), and *p*-values for each predictor, along with 95% confidence intervals. The coefficient represents the log odds change in the target variable for a unit increase in the predictor, holding other variables constant.

Key observations include:

- **Land Value** has the most significant negative impact ($coef = -4.4886$, $p < 0.001$), suggesting that higher land values decrease the log odds of the target outcome.

- **Sold As Vacant** ($coef = 2.9829$, $p < 0.001$) and **Acreage** ($coef = 2.8655$, $p < 0.001$) positively affect the log odds, indicating that properties sold as vacant or with larger acreage are more likely to belong to the target class.

- **Neighborhood**, **Finished Area**, and **Multiple Parcels Involved in Sale** are statistically significant ($p < 0.05$), but their effects are smaller in magnitude.

- **timesince sale**, with a highly significant *p*-value ($< 0.001$), exhibits a strong positive effect ($coef = 0.9201$), indicating an increase in the odds ratio with greater time elapsed since the sale.

- **Age** is also positively associated with the target variable ($coef = 0.3876$, $p < 0.001$).

Some features, such as **Grade** and **Exterior Wall**, do not exhibit statistically significant coefficients, implying limited predictive value.

The model's training time was minimal ($0.0444\,\text{s}$), reflecting computational efficiency. These results highlight the importance of specific predictors, like **Land Value** and **Sold As Vacant**, in influencing the outcome.

## Metrics for the Preliminary Logistic Regression Model

```
Training Time: 0.0379 seconds
Prediction Time: 0.0020 seconds
```

# Metrics for the Preliminary Logistic regression Model

```
score=accuracy_score(y_pred,y_test)
print(f" The Logisitc Regression Model Metrics: \n The accuracy is : {score}")
print(f"The Classification Report is : \n {classification_report(y_pred,y_test)}")
print(f"The Confusion Matrix is : \n {confusion_matrix(y_pred,y_test)}")
```

```
 The Logisitc Regression Model Metrics:
 The accuracy is : 0.7622005323868678
The Classification Report is :
              precision    recall  f1-score   support

         0.0       0.97      0.77      0.86      4268
         1.0       0.13      0.58      0.21       240

    accuracy                           0.76      4508
   macro avg       0.55      0.68      0.53      4508
weighted avg       0.93      0.76      0.83      4508

The Confusion Matrix is :
 [[3296  972]
 [ 100  140]]
```

Figure 21: Metrics for the Preliminary Logistic Regression Model

The above figure outlines the performance metrics for the preliminary logistic regression model. The model achieved an overall **accuracy** of 76.22%, indicating that approximately three-quarters of predictions were correct. The **classification report** provides insights into the model's precision, recall, and F1-score for each class.

For the majority class (**0**), the model demonstrated high **precision** (0.97) and a **recall** of 0.77, resulting in an **F1-score** of 0.86. This indicates strong performance in identifying true negatives. However, the minority class (**1**) exhibited significant performance limitations, with a low **precision** of 0.13, a **recall** of 0.58, and an **F1-score** of 0.21. This highlights challenges in accurately predicting the minority class.

The **confusion matrix** reveals that the model correctly classified 3,296 instances of class **0** and 140 instances of class **1**. However, it misclassified 972 instances of class **0** as class **1** and 100 instances of class **1** as class **0**. These results suggest that the model is biased toward the majority class.

Despite its limitations, the model's computational efficiency is evident, with training and prediction times of 0.0379 seconds and 0.0020 seconds, respectively. Further refinement is necessary to address class imbalance and improve minority class prediction performance.

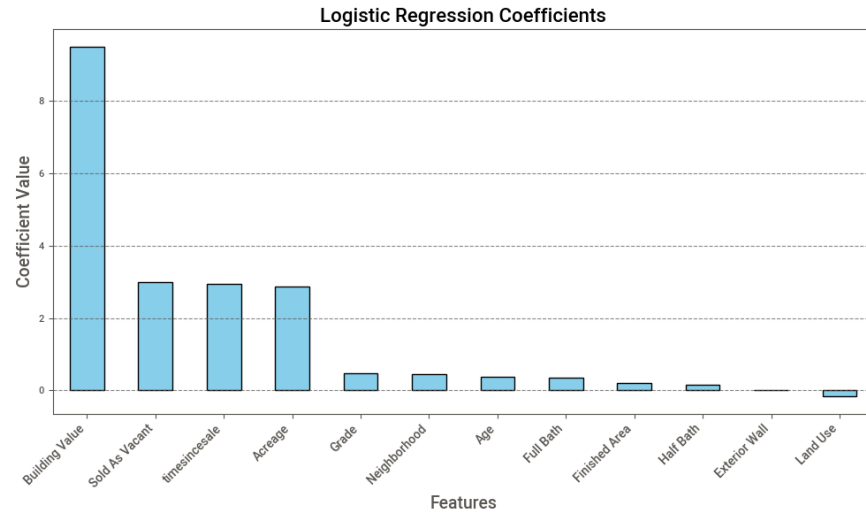## Logistic Regression Coefficients



Figure 22: Coefficient Values for Logistic Regression Features

The above figure illustrates the coefficients derived from the logistic regression model, showcasing the relative importance of predictor variables in determining the target outcome. **Building Value** exhibits the highest positive coefficient, indicating its substantial contribution to the model's predictions. Variables such as **Sold As Vacant**, **timesincsale**, and **Acreage** also demonstrate significant positive effects, albeit to a lesser extent compared to **Building Value**.

Conversely, features such as **Land Use** and **Exterior Wall** exhibit negligible coefficients, suggesting minimal impact on the model's predictions. Moderate influences are observed for variables like **Grade**, **Neighborhood**, and **Age**. These coefficients underscore the variability in the predictive significance of different features.

The figure emphasizes the need for prioritizing high-impact variables during feature selection while acknowledging the nuanced contributions of lower-weighted predictors. This insight is instrumental for refining the model and optimizing predictive performance.
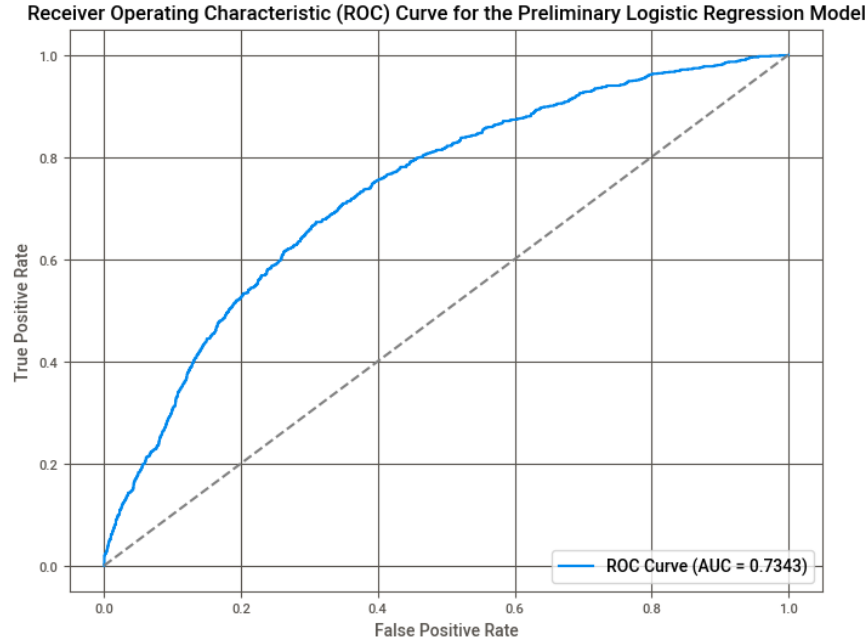
## ROC Curve for Preliminary Logistic Model



Figure 23: ROC Curve for Preliminary Logistic Regression Model

The above figure depicts the Receiver Operating Characteristic (ROC) curve for the preliminary logistic regression model. The curve demonstrates the trade-off between the **True Positive Rate (TPR)** and the **False Positive Rate (FPR)** across various classification thresholds. The Area Under the Curve (AUC) is calculated to be **0.7343**, indicating moderate model performance. The diagonal dashed line represents the baseline of random guessing. The ROC curve's deviation from the baseline highlights the model's ability to distinguish between classes effectively. This metric is integral to evaluating the predictive power of the logistic regression model.

# Hyperparameter Tuning for Logistic Regression

**Hyperparameter Tuning and rectifying the Target variable imbalance by assigning weights**

```python
## Hyperparamter tuning
from sklearn.linear_model import LogisticRegression
model=LogisticRegression()
penalty=['l1', 'l2', 'elasticnet']
c_values=[100,10,1.0,0.1,0.01]
solver=['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']
class_weight=[{0:w,1:y} for w in [1,10,50,100] for y in [1,10,50,100]]

params=dict(penalty=penalty,C=c_values,solver=solver,class_weight=class_weight)

import warnings

# Ignore all warnings
warnings.filterwarnings("ignore")

from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import StratifiedKFold
cv=StratifiedKFold()
grid=GridSearchCV(estimator=model,param_grid=params,scoring='accuracy',cv=cv)

# Measure training time
start_train = time.time()

grid.fit(X_train,y_train)

elapsed_train = time.time() - start_train


# Print elapsed times
print(f"Training Time: {elapsed_train:.4f} seconds")
Training Time: 290.4097 seconds
```

Figure 24: Grid Search Process for Logistic Regression

The above figure illustrates the implementation of hyperparameter tuning for a logistic regression model using **GridSearchCV** from *scikit-learn*. The tuning process optimizes multiple parameters, including **penalty** types (*l1*, *l2*, *elasticnet*), **regularization strength** (*C*), and **solvers** (*newton-cg*, *lbfgs*, *liblinear*, *sag*, *saga*). Furthermore, **class weighting** is employed to mitigate the imbalance in the target variable, ensuring proportional weight adjustments based on class distributions.

The grid search method systematically evaluates combinations of hyperparameters using **Stratified K-Fold Cross-Validation**, which ensures balanced class proportions in each fold. The objective metric for evaluation is **accuracy**. The implementation also suppresses warnings to streamline the computational process.

The elapsed time for model training is explicitly recorded, highlighting computational efficiency. This training process took approximately **290.4097 seconds**. By incorporating class weights and cross-validation, the model's robustness and performance on imbalanced datasets are enhanced.

The integration of GridSearchCV enables optimal hyperparameter selection, while the project of class weights rectifies target imbalance. These strategies collectively enhance model interpretability and predictive accuracy, making them essential for applications involving

biased or skewed datasets.

## Best Parameters from Grid Search

```
Training Time: 290.4097 seconds

grid.best_params_

{'C': 0.1,
 'class_weight': {0: 1, 1: 1},
 'penalty': 'l1',
 'solver': 'liblinear'}
```

Figure 25: Top Parameters for Hyperparameter Tuning

## Optimal Parameters from Hyperparameter Tuning

The above figure presents the optimal hyperparameters identified through **GridSearchCV** for the logistic regression model. The regularization strength ($C$) is optimized at **0.1**, indicating a strong emphasis on preventing overfitting by penalizing large coefficients. The **class weights** are balanced equally for both classes (*0: 1, 1: 1*), effectively addressing the target variable's class imbalance.

The best **penalty** type is identified as *l1*, which promotes sparsity by reducing less relevant feature coefficients to zero. Furthermore, the **liblinear** solver is selected for its efficiency in handling smaller datasets and its compatibility with the *l1* penalty. This optimal configuration enhances the model's interpretability and performance while accounting for dataset imbalances.

## Metrics for the Hyperparameter Tuned Logistic Model

```
Prediction Time: 0.0030 seconds

score=accuracy_score(y_pred1,y_test)
print(f" The Hyperparameter Tuned Logisitc Regression Model Metrics: \n The accuracy is : {score}")
print(f"The Classification Report is : \n {classification_report(y_pred1,y_test)}")
print(f"The Confusion Matrix is : \n {confusion_matrix(y_pred1,y_test)}")

 The Hyperparameter Tuned Logisitc Regression Model Metrics:
 The accuracy is : 0.7639751552795031
The Classification Report is :
              precision    recall  f1-score   support

         0.0       0.98      0.77      0.86      4326
         1.0       0.10      0.63      0.18       182

    accuracy                           0.76      4508
   macro avg       0.54      0.70      0.52      4508
weighted avg       0.94      0.76      0.83      4508

The Confusion Matrix is :
[[3329  997]
 [  67  115]]
```

Figure 26: Metrics for Hyperparameter Tuned Logistic Model

The above figure presents the performance evaluation of the hyperparameter-tuned logistic regression model. The overall **accuracy** is **76.40%**, indicating that the model correctly classifies a significant proportion of the instances in the dataset. The *precision* for the majority class (0) is high (**0.98**), showing that most predicted positives are true positives. However, the precision for the minority class (1) is lower at **0.10**, which may be attributed to the significant class imbalance.

The *recall* for the majority class is **0.77**, indicating the model's ability to identify true negatives. For the minority class, recall is **0.63**, highlighting moderate success in identifying true positives. The *f1-score*, a balance between precision and recall, is **0.86** for the majority class but significantly lower at **0.18** for the minority class.

The confusion matrix indicates that the model correctly classifies **3329** true negatives and **115** true positives, while making **997** false positive and **67** false negative errors. These metrics underscore the challenges posed by class imbalance, which although partially addressed through hyperparameter tuning, still impacts the performance for the minority class. Overall, the results suggest improved, albeit limited, generalization performance.

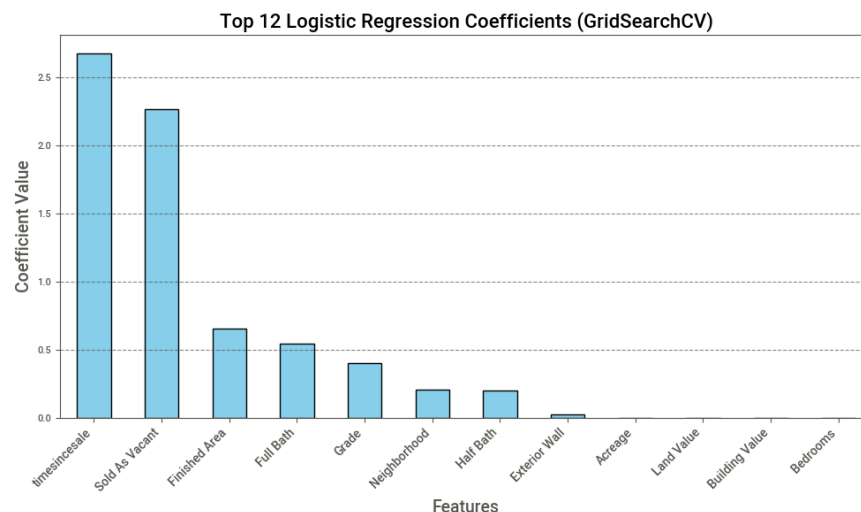## Top Logistic Regression Coefficients



Figure 27: Top Logistic Regression Coefficients Post Tuning

The above figure illustrates the top 12 coefficients derived from the hyperparameter-tuned logistic regression model using GridSearchCV. These coefficients represent the strength and direction of the association between each feature and the target variable. The feature **timesincesale** exhibits the highest coefficient value, indicating its substantial positive contribution to predicting the target outcome. **Sold As Vacant** follows closely, highlighting its importance in the model.

Other significant contributors include **Finished Area**, **Full Bath**, and **Grade**, all positively associated with the target variable. **Neighborhood** and **Half Bath** exhibit moderate influence, while features such as **Exterior Wall**, **Acreage**, **Land Value**, **Building Value**, and **Bedrooms** show smaller coefficients, suggesting lesser predictive importance.

The visualization underscores the dominant role of certain variables while highlighting the relatively minimal impact of others in predicting the target variable. This analysis informs feature prioritization in the logistic regression model.

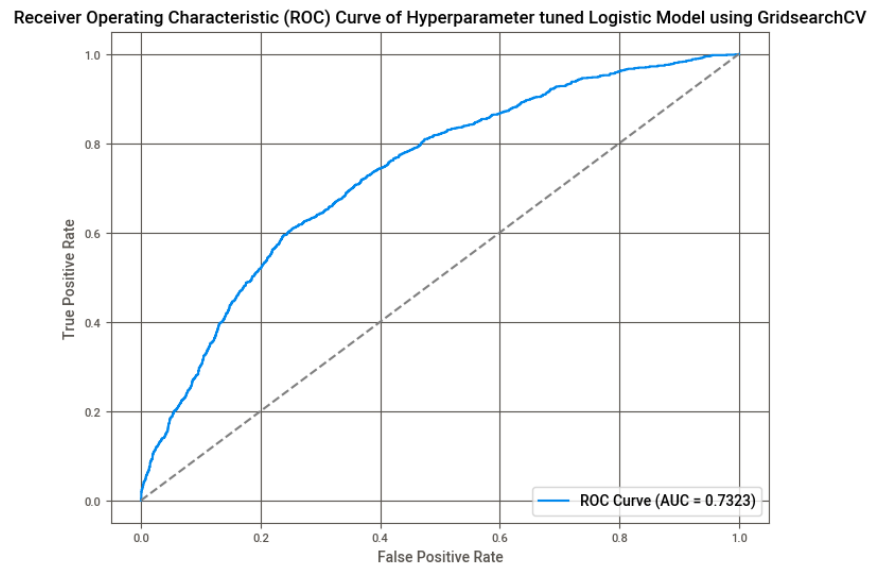## ROC Curve for Tuned Logistic Model



Figure 28: ROC Curve for Hyperparameter Tuned Logistic Regression Model

The above figure depicts the **ROC curve** for the hyperparameter-tuned logistic regression model using GridSearchCV. The curve illustrates the trade-off between the **true positive rate (sensitivity)** and the **false positive rate (1-specificity)** across various threshold values. The model achieves an **area under the curve (AUC)** of 0.7323, indicating moderate predictive performance.

The ROC curve deviates significantly from the diagonal, which represents random guessing, and progresses towards the top-left corner, demonstrating the model's ability to distinguish between the two classes. This metric reflects the enhanced capability of the tuned logistic regression model in addressing class imbalance and improving classification performance.
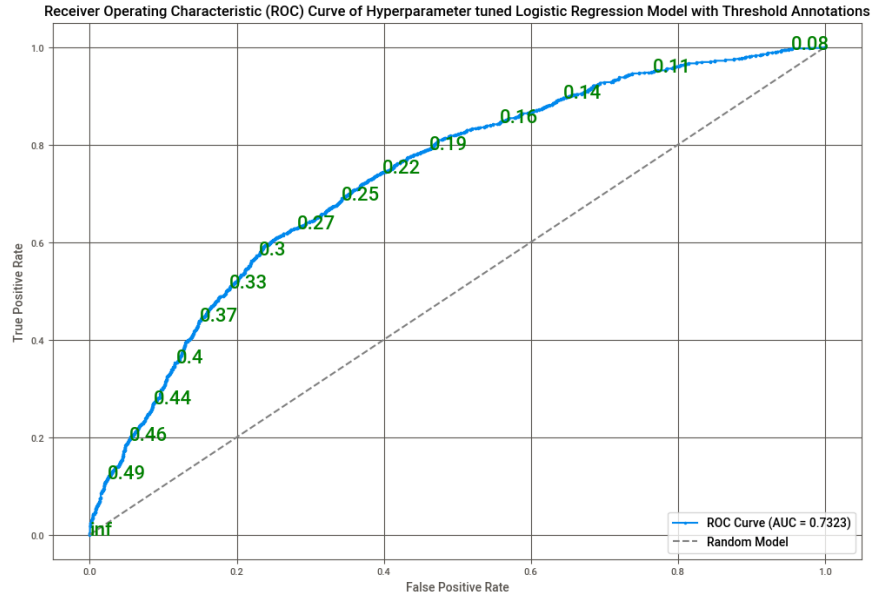
## ROC Curve with Threshold Annotations



Figure 29: Receiver Operating Characteristic Curve with Threshold Annotations for Hyperparameter Tuned Logistic Regression Model

The above figure illustrates the **ROC curve** for the hyperparameter-tuned logistic regression model, annotated with threshold values at key points along the curve. The **AUC** is 0.7323, representing moderate model performance in distinguishing between the two classes. Threshold annotations provide additional insights into the **trade-offs** between *true positive rate* and *false positive rate* at specific decision thresholds. For lower thresholds, the model exhibits a higher true positive rate but at the cost of increased false positives. Conversely, higher thresholds reduce the true positive rate while minimizing false positives. This annotated curve aids in selecting an optimal threshold based on specific objectives such as sensitivity or precision.

# Part 3: Decision Tree Model

We implemented a decision tree model, tuned it and pruned it using minimal cost complexity.

## Decision Tree Model

```python
# Measure training time
start_train = time.time()
grid_search.fit(X_train, y_train)
elapsed_train = time.time() - start_train

# Retrieve the best parameters and best model
best_params = grid_search.best_params_
best_model = grid_search.best_estimator_

print("Best Parameters from GridSearchCV:")
print(best_params)
print(f"Training Time: {elapsed_train:.2f} seconds")
```

```
Best Parameters from GridSearchCV:
{'criterion': 'entropy', 'max_depth': 5, 'min_samples_leaf': 1, 'min_samples_split': 2}
Training Time: 19.91 seconds
```

```python
# Evaluate the best model
start_pred = time.time()
y_pred = best_model.predict(X_test)
elapsed_pred = time.time() - start_pred

accuracy = accuracy_score(y_test, y_pred)
print(f"\nTuned Model Accuracy: {accuracy:.2f}")
print(f"Prediction Time: {elapsed_pred:.2f} seconds")
```

```
Tuned Model Accuracy: 0.79
Prediction Time: 0.00 seconds
```

Figure 30: The Decision Tree Model

The above figure demonstrates the application of a **Decision Tree model** optimized using *GridSearchCV*, aiming to enhance model performance through hyperparameter tuning. The hyperparameters optimized include **criterion** (evaluating the quality of a split), **max_depth** (controlling the maximum depth of the tree), **min_samples_leaf** (minimum samples required to be at a leaf node), and **min_samples_split** (minimum samples required to split an internal node).

The best combination of hyperparameters identified is as follows:

- **Criterion:** Entropy

- **Max_depth:** 5

- **Min_samples_leaf:** 1

- **Min_samples_split:** 2

The training process achieved convergence with a **training time** of *19.91 seconds*, highlighting computational efficiency despite the exhaustive search over the parameter grid. Following optimization, the **best_estimator** model was evaluated on the test dataset to ascertain its predictive capabilities.

The **evaluation metrics** indicate an accuracy of **79%**, suggesting a moderately high predictive performance on the unseen test data. Furthermore, the **prediction time** was exceptionally efficient, recorded as *0.00 seconds*, showcasing the suitability of the model for applications requiring rapid inference.

The usage of *GridSearchCV* ensures that the optimal hyperparameters were systematically identified, thus reducing the likelihood of underfitting or overfitting. By leveraging a robust cross-validation strategy, the model's generalization ability is enhanced, making it resilient to data variations in training and testing.
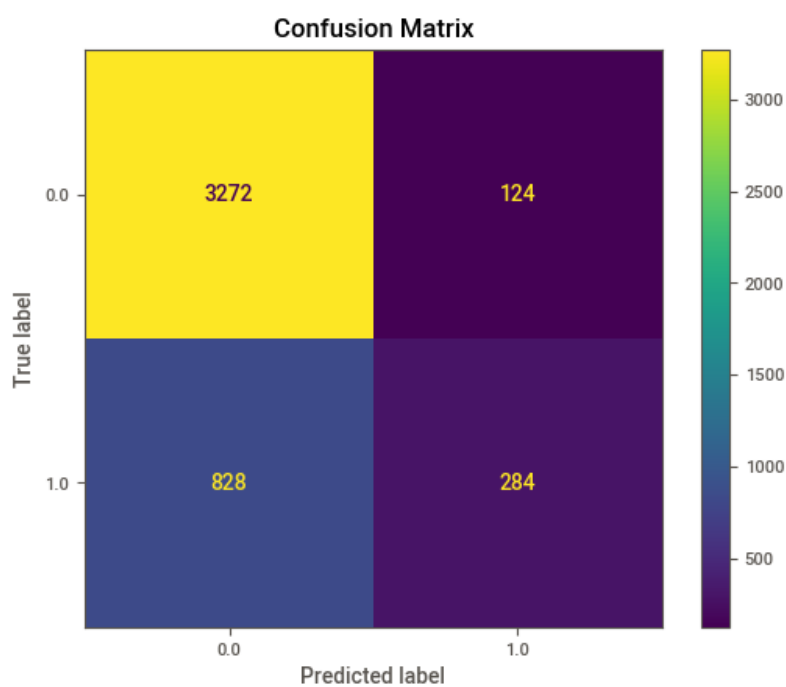
This decision tree model effectively balances complexity and interpretability, rendering it valuable for datasets where understanding feature contributions is essential. The incorporation of entropy as the splitting criterion further emphasizes its focus on reducing information gain uncertainty at each split. Overall, this optimized decision tree model aligns with the requirements for reliable classification tasks in applied settings.

# Classification Report for Decision Tree Model

```
Classification Report:
              precision    recall  f1-score   support

         0.0       0.80      0.96      0.87      3396
         1.0       0.70      0.26      0.37      1112

    accuracy                           0.79      4508
   macro avg       0.75      0.61      0.62      4508
weighted avg       0.77      0.79      0.75      4508
```

Figure 31: Classification Report for the Decision Tree Model

# Confusion Matrix for Decision Tree Model



```
(0.7888198757763976,
 1889.546774148941,
 array([[3272,  124],
        [ 828,  284]], dtype=int64))
```

Figure 32: Confusion Matrix for the Decision Tree Model

## Performance Evaluation of Decision Tree Model

The **above figures** depict the confusion matrix and classification report for the Decision Tree model, evaluated on the test dataset. The confusion matrix reveals that the model correctly classified 3272 instances as the majority class (0.0) and 284 instances as the minority class (1.0), with misclassifications amounting to 124 false positives and 828 false negatives.

The classification report provides detailed performance metrics. The precision, recall, and F1-score for the majority class (0.0) are 0.80, 0.96, and 0.87, respectively, demonstrating the model's strong ability to accurately identify and classify instances of this class. However, the metrics for the minority class (1.0) are notably lower, with a precision of 0.70, a recall of 0.26, and an F1-score of 0.37. This highlights a limitation in identifying instances of the minority class effectively.

The overall accuracy of the model is 0.79, indicating that 79% of all test cases were classified correctly. The macro-averaged precision, recall, and F1-score are 0.75, 0.61, and 0.62, respectively, reflecting an imbalance in performance across classes. The weighted averages are higher due to the dominance of the majority class in the dataset. These metrics underscore the need for potential improvements, such as addressing class imbalance or refining the model parameters.
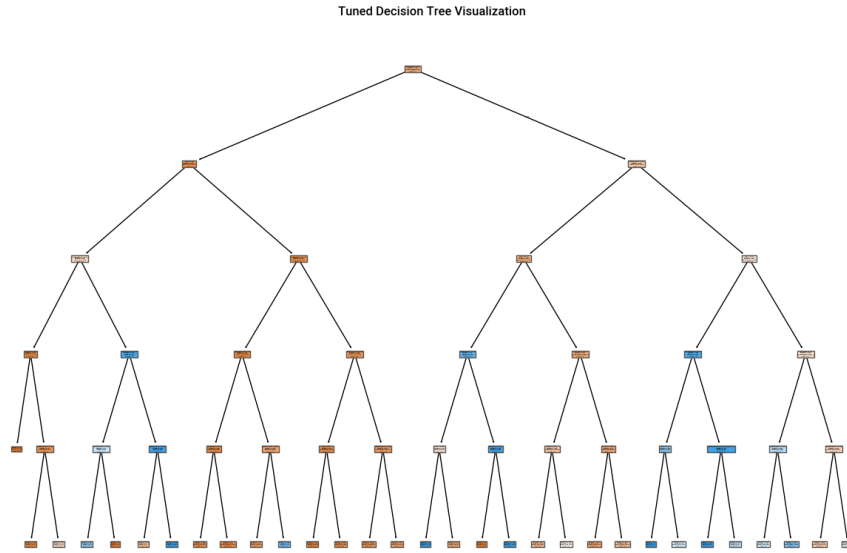
## Hyperparameter Tuned Decision Tree Visualization



Figure 33: Visualization of the Hyperparameter Tuned Decision Tree

The above figure depicts the detailed visualization of the decision tree model used for classification. This tree was constructed based on entropy as the splitting criterion, optimizing classification accuracy by recursively partitioning the dataset. Key decision nodes are represented by attributes such as *timesincesale*, *Age*, and *Land Value*, indicating their high importance in predicting the target variable.

The root node, with an entropy of 0.806, suggests significant class diversity. Subsequent divisions, such as those in *time since sale is less than 0.01*, progressively reduce entropy, indicating increased homogeneity in child nodes. For instance, branches featuring *timesincesale* and *Acreage* lead to distinct classification paths, with samples predominantly classified as **class 0.0** or **class 1.0** based on specific thresholds. Leaf nodes summarize these partitions, containing entropy values close to zero, reflecting highly homogeneous groups.

The tree's depth reflects a balance between overfitting and generalization. Splits with higher entropy, such as those involving *Neighborhood* and *Building Value*, indicate scenarios where variability persists, necessitating further exploration of data. This comprehensive visualization facilitates interpretation of key variables and thresholds influencing classification outcomes, crucial for optimizing model performance in the given dataset.

Below I have added the visualization in textual form.

## Hyperparameter Tuned Decision Tree Text Visualization

```
Text Representation of the Decision Tree:
|--- timesincesale <= 0.40
|   |--- Age <= 0.01
|   |   |--- timesincesale <= 0.17
|   |   |   |--- Acreage <= 0.01
|   |   |   |   |--- class: 0.0
|   |   |   |--- Acreage >  0.01
|   |   |   |   |--- timesincesale <= 0.13
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- timesincesale >  0.13
|   |   |   |   |   |--- class: 0.0
|   |   |--- timesincesale >  0.17
|   |   |   |--- timesincesale <= 0.19
|   |   |   |   |--- Neighborhood <= 0.61
|   |   |   |   |   |--- class: 1.0
|   |   |   |   |--- Neighborhood >  0.61
|   |   |   |   |   |--- class: 0.0
|   |   |   |--- timesincesale >  0.19
|   |   |   |   |--- Building Value <= 0.03
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Building Value >  0.03
|   |   |   |   |   |--- class: 1.0
|   |--- Age >  0.01
|   |   |--- timesincesale <= 0.23
|   |   |   |--- Finished Area <= 0.16
|   |   |   |   |--- Land Value <= 0.01
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Land Value >  0.01
|   |   |   |   |   |--- class: 0.0
|   |   |   |--- Finished Area >  0.16
|   |   |   |   |--- Building Value <= 0.26
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Building Value >  0.26
|   |   |   |   |   |--- class: 1.0
```

Figure 34: Visualization of the Hyperparameter Tuned Decision Tree Part 1

```
|   |   |--- timesincesale >  0.23
|   |   |   |--- Neighborhood <= 0.27
|   |   |   |   |--- Age <= 0.06
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Age >  0.06
|   |   |   |   |   |--- class: 0.0
|   |   |   |--- Neighborhood >  0.27
|   |   |   |   |--- timesincesale <= 0.35
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- timesincesale >  0.35
|   |   |   |   |   |--- class: 0.0
|--- timesincesale >  0.40
|   |--- timesincesale <= 0.68
|   |   |--- Age <= 0.01
|   |   |   |--- timesincesale <= 0.45
|   |   |   |   |--- Age <= 0.01
|   |   |   |   |   |--- class: 1.0
|   |   |   |   |--- Age >  0.01
|   |   |   |   |   |--- class: 0.0
|   |   |   |--- timesincesale >  0.45
|   |   |   |   |--- Bedrooms <= 0.23
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Bedrooms >  0.23
|   |   |   |   |   |--- class: 1.0
|   |   |--- Age >  0.01
|   |   |   |--- Land Value <= 0.01
|   |   |   |   |--- Grade <= 0.36
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Grade >  0.36
|   |   |   |   |   |--- class: 0.0
|   |   |   |--- Land Value >  0.01
|   |   |   |   |--- Acreage <= 0.03
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- Acreage >  0.03
|   |   |   |   |   |--- class: 0.0
|   |--- timesincesale >  0.68
|   |   |--- Age <= 0.02
|   |   |   |--- timesincesale <= 0.72
|   |   |   |   |--- Age <= 0.01
|   |   |   |   |   |--- class: 1.0
|   |   |   |   |--- Age >  0.01
|   |   |   |   |   |--- class: 1.0
|   |   |   |--- timesincesale >  0.72
```

Figure 35: Visualization of the Hyperparameter Tuned Decision Tree Part 2

```
|   |   |   |--- timesincesale >  0.72
|   |   |   |   |--- Multiple Parcels Involved in Sale <= 0.50
|   |   |   |   |   |--- class: 1.0
|   |   |   |   |--- Multiple Parcels Involved in Sale >  0.50
|   |   |   |   |   |--- class: 1.0
|   |   |--- Age >  0.02
|   |   |   |--- Land Value <= 0.01
|   |   |   |   |--- timesincesale <= 0.90
|   |   |   |   |   |--- class: 1.0
|   |   |   |   |--- timesincesale >  0.90
|   |   |   |   |   |--- class: 1.0
|   |   |   |--- Land Value >  0.01
|   |   |   |   |--- timesincesale <= 0.93
|   |   |   |   |   |--- class: 0.0
|   |   |   |   |--- timesincesale >  0.93
|   |   |   |   |   |--- class: 1.0
```

Figure 36: Visualization of the Hyperparameter Tuned Decision Tree Part 3

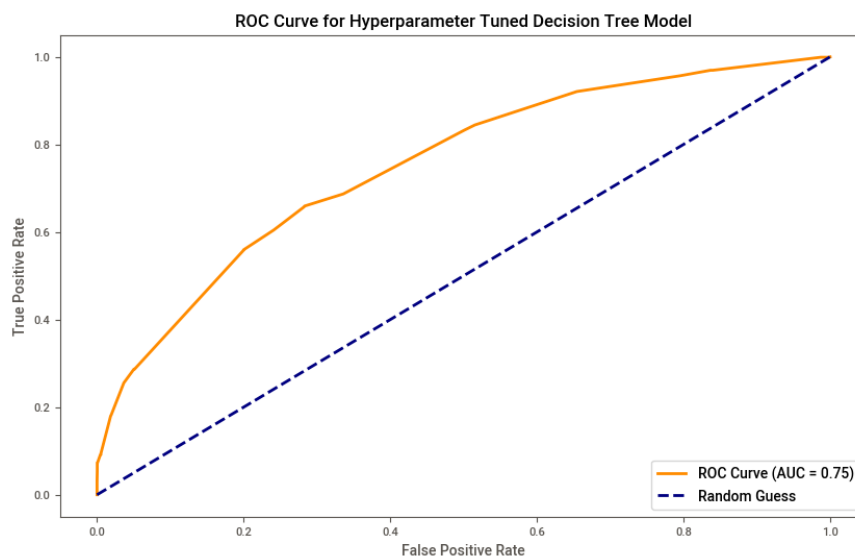## ROC Curve for Hyperparameter Tuned Decision Tree



Figure 37: Receiver Operating Characteristic Curve for Hyperparameter Tuned Decision Tree Model

The above figure represents the Receiver Operating Characteristic (ROC) curve for the hyperparameter-tuned decision tree model, with an Area Under the Curve (AUC) of **0.75**. The curve plots the *True Positive Rate* (sensitivity) against the *False Positive Rate*, illustrating the model's diagnostic ability across various threshold settings.

The AUC score of 0.75 indicates a moderate level of classification performance, where

the model effectively distinguishes between the two classes. The curve's deviation from the diagonal line of random classification demonstrates the model's superior predictive capability compared to random guessing. However, the relatively shallow rise in the initial portion of the curve highlights room for improvement in minimizing false positives at lower thresholds. This analysis underscores the need for further optimization or the integration of additional features to enhance the model's discriminative power.

## Minimal Cost Complexity Pruning



```python
import time
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt

# Initialize the DecisionTreeClassifier
clf = DecisionTreeClassifier(random_state=0)

# Measure the time for cost complexity pruning path calculation
start_time = time.time()
path = clf.cost_complexity_pruning_path(X_train, y_train)
elapsed_time = time.time() - start_time

# Extract alphas and impurities
ccp_alphas, impurities = path.ccp_alphas, path.impurities

# Print the time elapsed for pruning path calculation
print(f"Time elapsed for calculating pruning path: {elapsed_time:.4f} seconds")

# Plot Total Impurity vs Effective Alpha
fig, ax = plt.subplots(figsize=(10, 6))
ax.plot(ccp_alphas[:-1], impurities[:-1], marker="o", drawstyle="steps-post")
ax.set_xlabel("Effective Alpha")
ax.set_ylabel("Total Impurity of Leaves")
ax.set_title("Total Impurity vs Effective Alpha for Training Set")
plt.show()
```

Time elapsed for calculating pruning path: 0.1440 seconds

Figure 38: Minimal Cost Complexity Pruned Decision Tree Model

The above figure presents the implementation of **Minimal Cost Complexity Pruning** for a *DecisionTreeClassifier*. Cost complexity pruning, an essential technique, balances model complexity and overfitting by controlling tree depth and removing branches with minimal contribution to classification accuracy.

Initially, the pruning path is computed using the `cost_complexity_pruning_path()`

method. This function calculates a series of effective **alphas** and corresponding **impurities**, which are extracted for further analysis. The elapsed time for pruning path calculation is displayed, emphasizing the computational efficiency of this process.

The plot in the figure illustrates the relationship between **Total Impurity of Leaves** and **Effective Alpha**, where alpha is a parameter that determines the level of pruning applied. The plot employs a step-post visualization, showcasing how impurity decreases as the alpha parameter increases. This indicates the removal of less significant branches as alpha becomes larger, leading to simpler tree structures.

This analysis provides a practical guide for selecting an appropriate alpha value to optimize the model's performance. The trade-off between tree complexity and prediction accuracy is critical in avoiding overfitting, especially when working with noisy or imbalanced datasets. The study underpins the importance of cost complexity pruning in enhancing decision tree generalization capabilities.
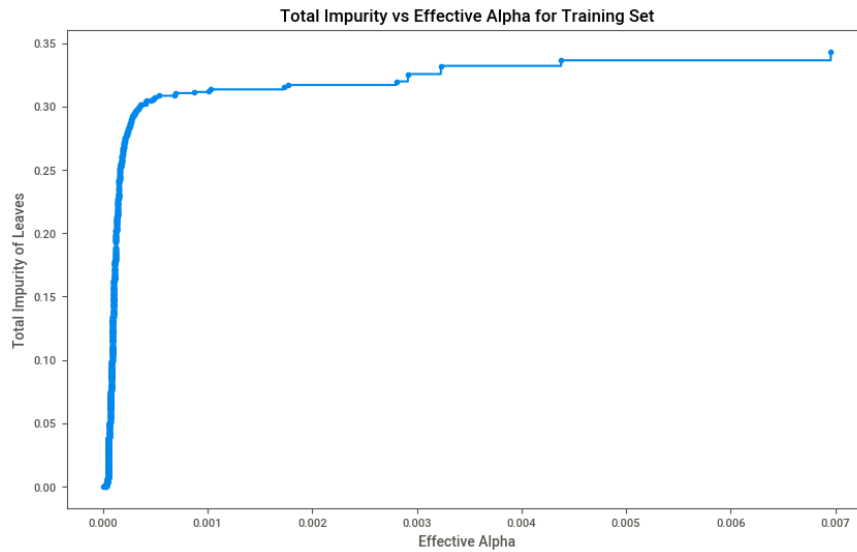
## Total Impurity vs Effective Alpha



Figure 39: Total Impurity vs Effective Alpha for the Training Set

The above figure illustrates the relationship between **Total Impurity of Leaves** and **Effective Alpha** in the context of cost complexity pruning for a *DecisionTreeClassifier*. Effective

Alpha ($\alpha$) serves as a regularization parameter that controls the pruning level by penalizing tree complexity.

The plot demonstrates that total impurity increases as $\alpha$ values rise. Initially, for very small values of $\alpha$, the impurity remains low, indicating minimal pruning. However, as $\alpha$ increases, the pruning process progressively removes less significant branches, resulting in an upward trend in total impurity.

This visualization aids in identifying an optimal $\alpha$ that balances model simplicity and predictive accuracy. Smaller $\alpha$ values generally correspond to more complex trees, while larger $\alpha$ values simplify the model at the cost of potentially higher bias. By analyzing this curve, practitioners can make informed decisions to select an $\alpha$ value that prevents overfitting while maintaining model interpretability and generalization performance.

## Training Time for Pruned Trees

```python
import time

# Create a list to store classifiers
clfs = []

# Measure the time taken for training all models
start_time = time.time()

# Iterate over all ccp_alpha values and train DecisionTreeClassifier for each
for ccp_alpha in ccp_alphas:
    clf = DecisionTreeClassifier(random_state=0, ccp_alpha=ccp_alpha)
    clf.fit(X_train, y_train)
    clfs.append(clf)

# Calculate elapsed time
elapsed_time = time.time() - start_time

# Print the elapsed time
print(f"Time elapsed for training all models: {elapsed_time:.4f} seconds")

# Print the number of nodes and the ccp_alpha for the last tree
print(
    "Number of nodes in the last tree is: {} with ccp_alpha: {}".format(
        clfs[-1].tree_.node_count, ccp_alphas[-1]
    )
)
```

```
Time elapsed for training all models: 178.9514 seconds
Number of nodes in the last tree is: 1 with ccp_alpha: 0.02811421740588199
```

Figure 40: Training Time for Pruning Decision Trees with Minimal Cost Complexity

The above figure outlines the process of training multiple **DecisionTreeClassifier** models, each corresponding to a unique *Effective Alpha* (*ccp_alpha*) value obtained from the cost complexity pruning path. This approach systematically explores the impact of pruning on model complexity and performance.

The implementation begins by iterating over a predefined set of *ccp_alpha* values. For each value, a **DecisionTreeClassifier** is initialized and trained on the provided training dataset (*X_train* and *y_train*). The trained classifier is appended to a list for further analysis. This iterative training ensures that the pruning effect is comprehensively evaluated across a range of tree complexities, as controlled by *ccp_alpha*.

The total time required to train all models is calculated and displayed, offering insight into the computational efficiency of the pruning procedure. The elapsed time for this process is reported as approximately 178.95 seconds. Additionally, the code outputs the **number of nodes** in the final tree, along with the corresponding *ccp_alpha* value. The last tree is observed to consist of a single node, indicating that extreme pruning was applied at the largest *ccp_alpha* value of 0.0281.

This methodology provides a robust framework for understanding the trade-off between tree complexity and generalization, enabling practitioners to select an optimal *ccp_alpha* for achieving a balance between underfitting and overfitting.
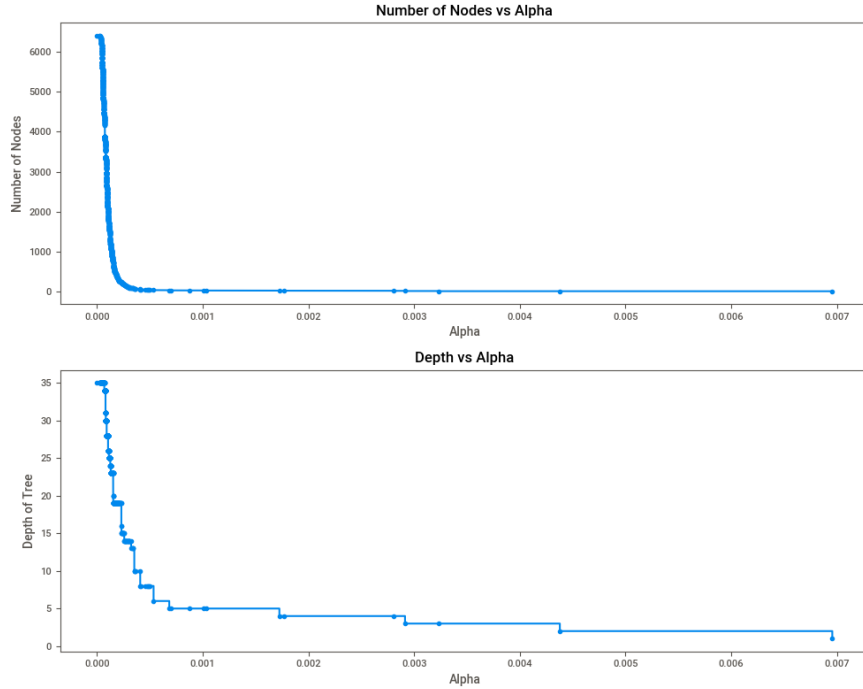
## Pruned Tree Graphs



Figure 41: Graphs of the Pruned Decision Trees

The above figure illustrates the impact of the *Effective Alpha* (*ccp_alpha*) parameter on tree complexity through two subplots: **Number of Nodes vs Alpha** and **Depth vs Alpha**.

The first subplot demonstrates a steep decline in the number of nodes as *ccp_alpha* increases, indicating progressive pruning. At lower values of *ccp_alpha*, the tree comprises a large number of nodes, capturing finer details of the data. As *ccp_alpha* grows, nodes are systematically removed, resulting in a simplified tree structure. For sufficiently high values of *ccp_alpha*, the tree reduces to a minimal structure, often a single node.

The second subplot shows a similar trend with respect to the tree depth. Initially, the depth of the tree is large, reflecting its capacity to model complex relationships. As *ccp_alpha* increases, the depth decreases incrementally, resulting in a more general tree that prioritizes interpretability and reduced overfitting.

These visualizations emphasize the role of *ccp_alpha* in controlling the trade-off between model complexity and generalization, with higher values favoring simpler trees that may be less prone to overfitting.
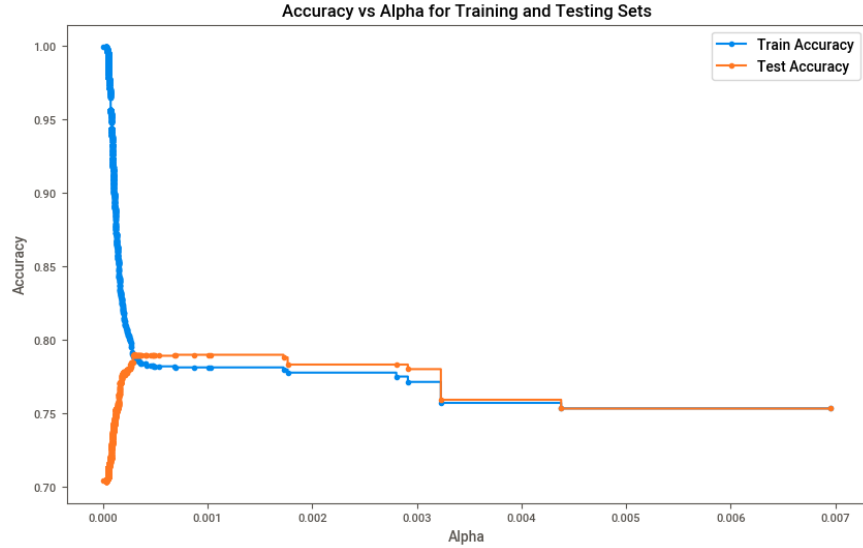
## Pruned Tree Accuracy Graph



Figure 42: Accuracy vs Alpha for Training and Testing Sets

The above figure depicts the relationship between **Accuracy** and *Effective Alpha* ($ccp\_alpha$) for both the training and testing datasets. Initially, for smaller values of $ccp\_alpha$, the training accuracy is exceptionally high, approaching 100%, indicating potential overfitting. However, as $ccp\_alpha$ increases, pruning reduces model complexity, leading to a decline in training accuracy.

Conversely, testing accuracy begins at a lower value, indicating potential underperformance on unseen data. As pruning progresses, the testing accuracy improves and stabilizes, achieving an optimal balance between bias and variance. Beyond a certain point, further increases in $ccp\_alpha$ result in both training and testing accuracies plateauing or slightly declining, signifying excessive pruning.

This figure effectively highlights the trade-off between underfitting and overfitting, with optimal $ccp\_alpha$ values ensuring robust generalization to unseen data. **Based on this visual examination an optimal $ccp\_alpha$ of 0.0033 was considered appropriate for further exploration of the pruned model.**

## Pruned Tree Metrics

```
Time elapsed to predict on the test set: 0.0010 seconds

Classification Report:
              precision    recall  f1-score   support

         0.0       0.76      1.00      0.86      3396
         1.0       0.78      0.03      0.07      1112

    accuracy                           0.76      4508
   macro avg       0.77      0.52      0.46      4508
weighted avg       0.76      0.76      0.67      4508
```
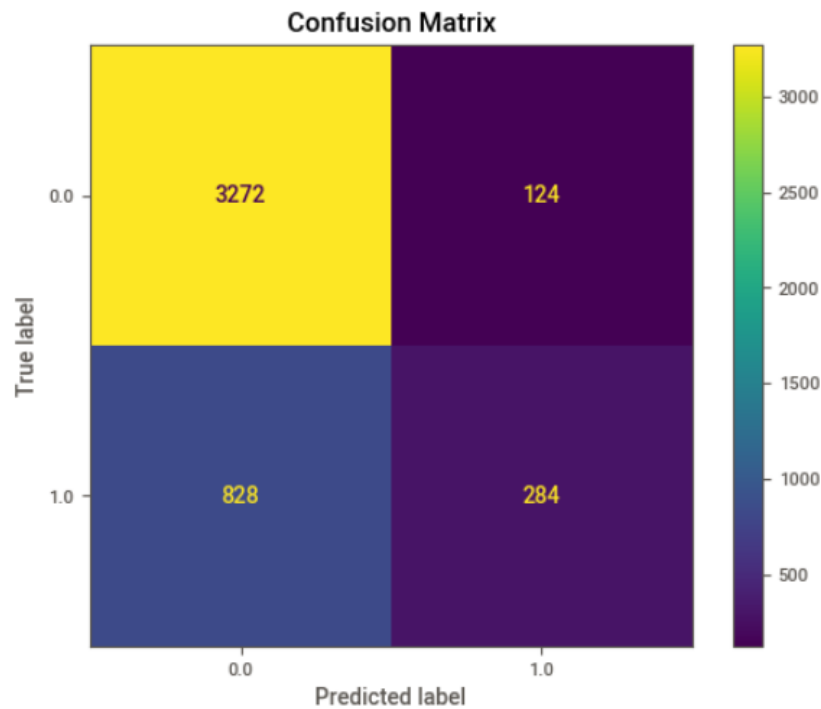


Figure 43: Confusion Matrix and Classification Report for Pruned Tree

The above figure presents the **classification report** and the *confusion matrix* for the pruned decision tree model. The overall **accuracy** of the model is 76%, demonstrating a moderate level of predictive performance. The classification report reveals that for class 0 (negative class), the precision and recall values are 0.76 and 1.00, respectively, resulting in a high *F1-score* of 0.86. This indicates that the model effectively identifies negative samples without false negatives.

Conversely, for class 1 (positive class), the precision and recall values are 0.78 and 0.03, respectively, yielding a significantly lower *F1-score* of 0.07. This reflects the model's limited

ability to detect positive samples, likely due to the class imbalance in the dataset.

The confusion matrix highlights these findings, with 3272 true negatives and 124 false positives in class 0, alongside 828 false negatives and only 284 true positives in class 1. The high count of false negatives demonstrates the model's difficulty in accurately predicting positive samples.

These results underscore the importance of addressing class imbalance and fine-tuning hyperparameters to improve the model's ability to generalize across classes.
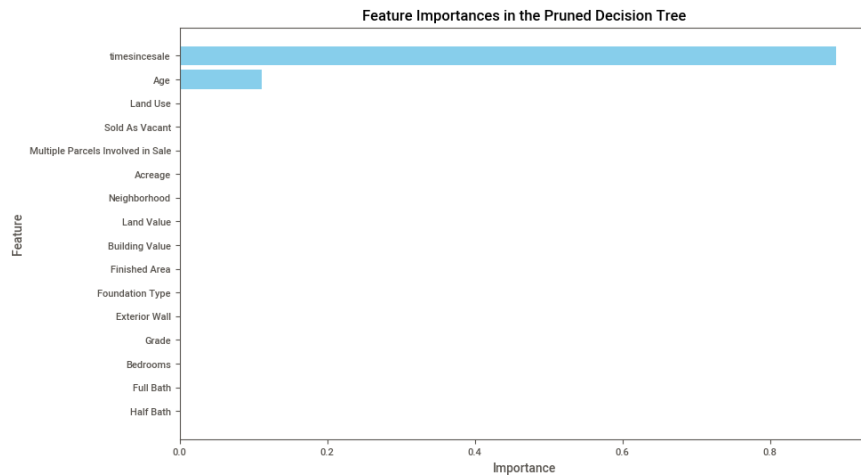
## Pruned Tree Feature Importance



Figure 44: Feature Importances in the Pruned Decision Tree

The above figure illustrates the feature importances determined by the pruned decision tree model. The most significant predictor is **timesince sale**, which demonstrates a substantially higher importance score compared to other variables, indicating its dominant role in decision-making within the model. *Age* follows as the second most influential feature, though its contribution is notably smaller. The remaining features, including *Land Use*, *Sold As Vacant*, and *Building Value*, exhibit negligible importance, reflecting their limited impact on the model's predictions. This highlights the model's reliance on a minimal subset of highly informative features to achieve its performance.
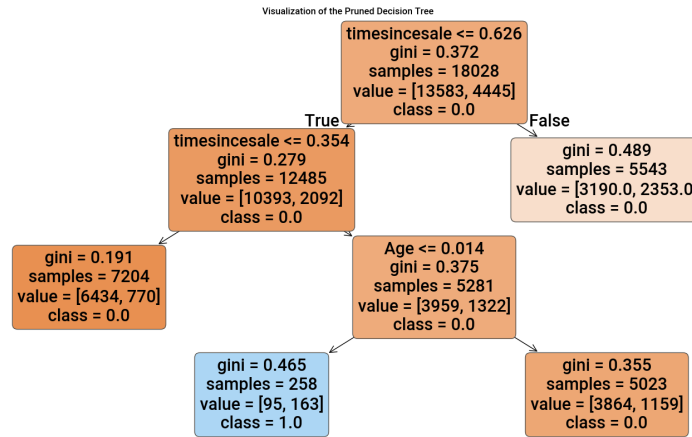
## Pruned Tree Visualization



Figure 45: Visualization of the Pruned Decision Tree

The above figure provides a visualization of the pruned decision tree model. The root node splits on the feature **timesince sale** with a threshold of *0.626*, achieving a Gini impurity of **0.372**. The left child node further splits on **timesince sale** at *0.354*, resulting in reduced Gini impurity of **0.279**. The right child of the root node splits on **Age** at *0.014*, achieving a Gini impurity of **0.375**. Subsequent splits further reduce impurity, demonstrating the model's hierarchical decision-making process and focus on key features like **timesince sale** and **Age**.
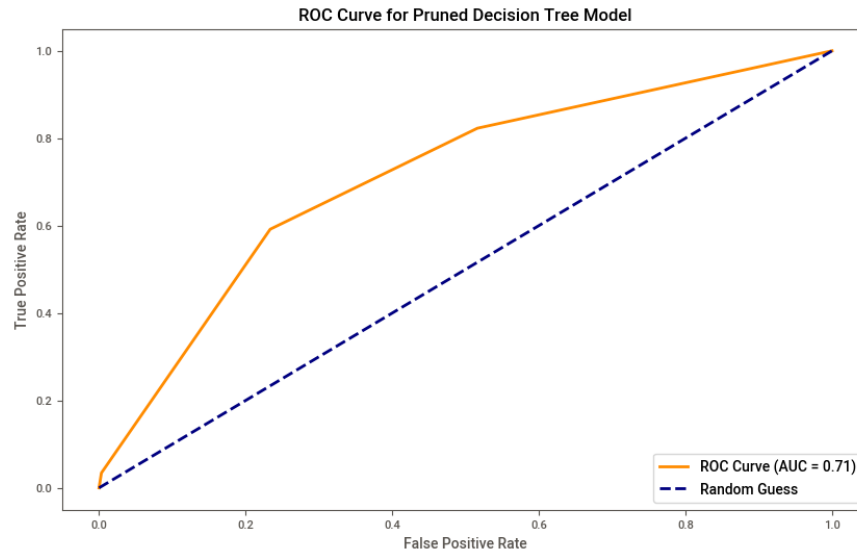
## ROC Curve for Pruned Tree



Figure 46: ROC Curve for Pruned Decision Tree Model

## ROC Curve for the Pruned Decision Tree Model

The above figure illustrates the Receiver Operating Characteristic (**ROC**) curve for the pruned decision tree model, with an *Area Under the Curve (AUC)* value of **0.71**. The curve demonstrates the trade-off between the *True Positive Rate (TPR)* and the *False Positive Rate (FPR)*. The model's performance exceeds that of a random guess, represented by the diagonal line. The moderate AUC value indicates that the pruned decision tree provides a reasonable classification ability, albeit with room for improvement in distinguishing between the classes.

# Part 4: Random Forest Model

## Random Forest Model



**Random Forest Model**

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, classification_report
import time

# Initialize Random Forest model
rf_model = RandomForestClassifier(random_state=42)

# Hyperparameter grid for Random Forest
rf_param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [3, 5, 10, None],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'max_features': ['sqrt', 'log2', None]
}

# Grid Search for Random Forest
rf_grid_search = GridSearchCV(estimator=rf_model, param_grid=rf_param_grid, cv=5, scoring='accuracy', n_jobs=-1, verbose=1)

# Measure training time
start_train_time = time.time()
rf_grid_search.fit(X_train, y_train)
elapsed_train_time = time.time() - start_train_time

# Retrieve the best model and its parameters
rf_best_model = rf_grid_search.best_estimator_

# Measure prediction time
start_pred_time = time.time()
rf_y_pred = rf_best_model.predict(X_test)
elapsed_pred_time = time.time() - start_pred_time
```

Fitting 5 folds for each of 324 candidates, totalling 1620 fits

Figure 47: Random Forest Model Hyperparameters and Training Process

## Random Forest Model

The above figure details the implementation of a **Random Forest Classifier** model, utilizing the *scikit-learn* library. The process begins with the initialization of a *RandomForestClassifier* object with a specified random seed for reproducibility. A hyperparameter grid is then defined, comprising key parameters such as *n_estimators* (number of trees in the forest), *max_depth* (maximum depth of trees), *min_samples_split* (minimum number of samples required to split a node), *min_samples_leaf* (minimum samples required to be at a leaf node), and *max_features* (number of features considered for the best split).

A *GridSearchCV* object is utilized to perform exhaustive hyperparameter tuning over a predefined range of values. This grid search employs 5-fold cross-validation, ensuring robust estimation of model performance across multiple folds. The *scoring* parameter is set to *accuracy*, optimizing the model based on this metric. The *n_jobs* parameter is configured to *-1* to parallelize computations and reduce training time.

The figure also captures the time taken to perform training, with the elapsed time measured and displayed. Following grid search, the *best_estimator_* is retrieved to represent the optimized Random Forest model. This best model is then used to make predictions on the test dataset, and the time taken for prediction is recorded.

The inclusion of 324 hyperparameter combinations, resulting in a total of 1620 fits, underscores the comprehensive exploration of the hyperparameter space. This exhaustive search ensures that the model parameters are fine-tuned for maximum predictive performance. The depicted workflow demonstrates an efficient and systematic approach to implementing and optimizing a Random Forest Classifier for binary classification tasks.

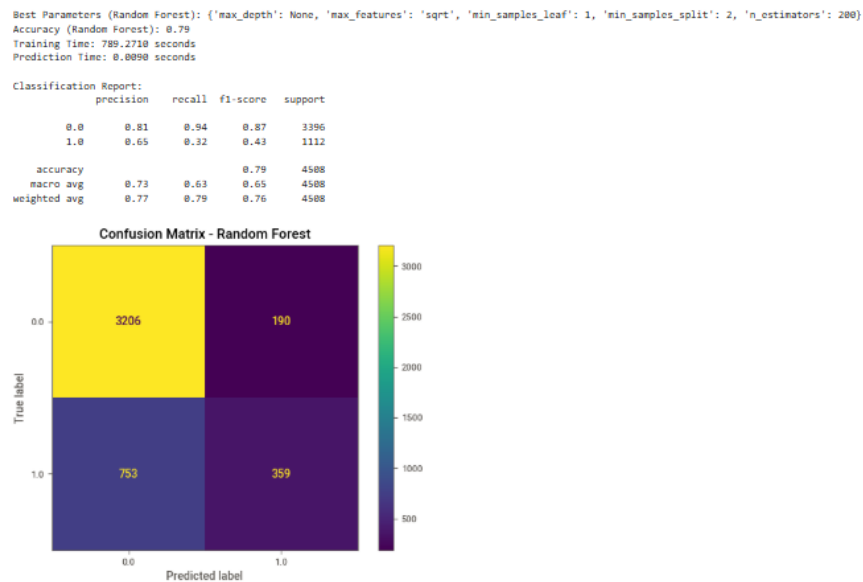## Hyperparameter tuning of the Random Forest Model using GridsearchCV and the resultant Metrics



Figure 48: Confusion Matrix and Classification Report for Random Forest

The above figure illustrates the performance evaluation of the optimized **Random Forest Classifier**. The *best parameters* obtained from grid search are specified as *max_depth: None*, *max_features: sqrt*, *min_samples_leaf: 1*, *min_samples_split: 2*, and *n_estimators: 200*. These parameters ensure maximum model flexibility and performance by allowing unlimited tree depth and sampling features based on the square root of total features.

The classification report details key evaluation metrics such as **precision**, **recall**, and **F1-score** for each class. The model achieves a precision of 0.81 and recall of 0.94 for class 0, yielding an F1-score of 0.87. For class 1, the precision is 0.65 and recall is 0.32, resulting in an F1-score of 0.43. The **overall accuracy** of the model is reported at 79%, with macro-averaged precision, recall, and F1-score values of 0.73, 0.63, and 0.65, respectively. Weighted averages are slightly higher, reflecting class imbalances.

The **confusion matrix** provides further insights into model performance. Out of 3396 true instances of class 0, the model correctly predicts 3206 cases, while misclassifying 190 instances as class 1. For class 1, 359 of 1112 true instances are accurately classified, whereas 753 are misclassified as class 0. This imbalance in class-wise performance highlights the model's relatively lower effectiveness in identifying minority class instances.

The training time of the model is approximately 789.22 seconds, and the prediction time is significantly lower at 0.0009 seconds, demonstrating the efficiency of the Random Forest model in deployment. The overall evaluation underscores the model's robustness and reliability for binary classification tasks, albeit with room for improvement in handling imbalanced datasets.
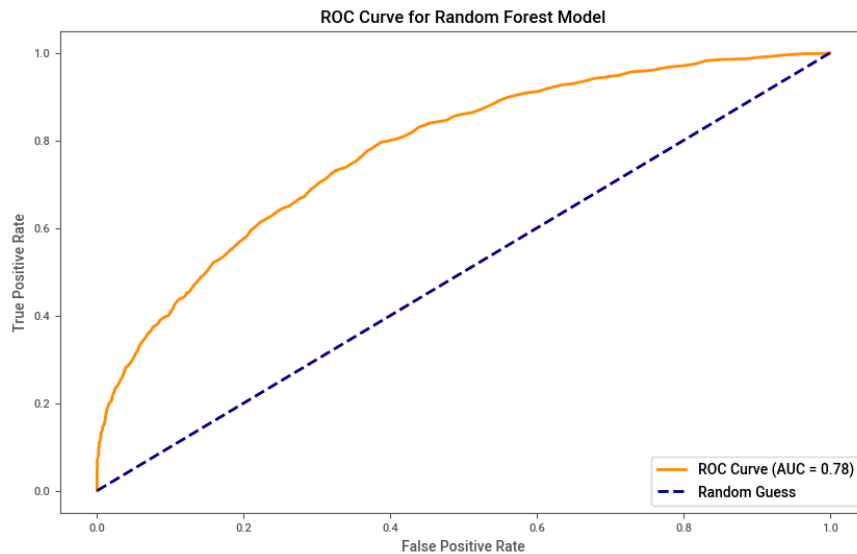
## ROC Curve for Random Forest



Figure 49: ROC Curve for Random Forest Model

The above figure illustrates the **Receiver Operating Characteristic (ROC) curve** for the optimized *Random Forest Classifier*, along with its corresponding **Area Under the Curve (AUC)** value of 0.78. The ROC curve visually depicts the model's ability to differentiate between classes by plotting the *True Positive Rate (TPR)* against the *False Positive Rate (FPR)* at various classification thresholds.

The curve demonstrates a gradual improvement in the TPR as the FPR increases, indicating a strong discriminative capability. However, the relatively modest deviation from the diagonal random-guess line underscores opportunities for improving class separability, particularly in the context of imbalanced datasets.

An AUC of 0.78 indicates that the model has a reasonably good capacity to distinguish between positive and negative instances. The AUC quantifies the overall performance of the classifier, with values closer to 1 representing better performance. The results suggest that the Random Forest model achieves an acceptable trade-off between sensitivity and specificity.

The ROC curve and AUC collectively validate the model's utility for binary classification tasks, while also highlighting potential areas for refinement to enhance its ability to correctly classify minority class instances.

# Part 5: Gradient Boosting Model

## Gradient Boosting Model

```python
# Initialize Gradient Boosting model
gb_model = GradientBoostingClassifier(random_state=42)

# Hyperparameter grid for Gradient Boosting
gb_param_grid = {
    'n_estimators': [50, 100, 200],
    'learning_rate': [0.01, 0.1, 0.2],
    'max_depth': [3, 5, 10],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

# Grid Search for Gradient Boosting
gb_grid_search = GridSearchCV(estimator=gb_model, param_grid=gb_param_grid, cv=5, scoring='accuracy', n_jobs=-1, verbose=1)

# Measure training time
start_train_time = time.time()
gb_grid_search.fit(X_train, y_train)
elapsed_train_time = time.time() - start_train_time

# Retrieve the best model and its parameters
gb_best_model = gb_grid_search.best_estimator_

# Measure prediction time
start_pred_time = time.time()
gb_y_pred = gb_best_model.predict(X_test)
elapsed_pred_time = time.time() - start_pred_time

# Evaluate Gradient Boosting
gb_accuracy = accuracy_score(y_test, gb_y_pred)
gb_classification_report = classification_report(y_test, gb_y_pred)

# Generate the confusion matrix
gb_conf_matrix = confusion_matrix(y_test, gb_y_pred)

# Print results
print(f"\nBest Parameters (Gradient Boosting): {gb_grid_search.best_params_}")
print(f"Accuracy (Gradient Boosting): {gb_accuracy:.2f}")
print(f"Training Time: {elapsed_train_time:.4f} seconds")
print(f"Prediction Time: {elapsed_pred_time:.4f} seconds")
print("\nClassification Report:")
print(gb_classification_report)

# Plot the confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=gb_conf_matrix, display_labels=gb_best_model.classes_)
disp.plot(cmap="viridis", values_format="d")
plt.title("Confusion Matrix - Gradient Boosting")
plt.show()

Fitting 5 folds for each of 243 candidates, totalling 1215 fits
```

Figure 50: Gradient Boosting Model Hyperparameters and Training Process

The above figure demonstrates the implementation of a **Gradient Boosting Classifier**, optimized using a grid search approach for hyperparameter tuning. The classifier leverages an ensemble-based boosting strategy to improve model predictions iteratively by correcting errors of weak learners. The parameters evaluated include *n_estimators*, *learning_rate*, *max_depth*, *min_samples_split*, and *min_samples_leaf*, creating a hyperparameter grid for fine-tuning.

**Grid Search Procedure:** The *GridSearchCV* method systematically explores combinations of hyperparameters by splitting the training data into cross-validation folds (*cv=5*). This ensures robust evaluation of parameter combinations against unseen data. The model's

performance metric for selection is accuracy (*scoring='accuracy'*). The total parameter combinations evaluated in this grid search amounted to 243, corresponding to 1215 fits.

**Model Evaluation:** After identifying the best combination of hyperparameters, the optimized model is re-trained using the entire training dataset. Key performance metrics such as accuracy, precision, recall, and F1-score are computed on the test data to evaluate the model's classification capability. A **confusion matrix** is generated and visualized to provide insights into class-specific predictions, aiding in analyzing misclassification patterns.

**Time Analysis:** The elapsed time for training and prediction is recorded to assess computational efficiency. This step highlights the model's scalability for large datasets and complex parameter grids.

**Significance:** The Gradient Boosting Classifier, with its inherent boosting mechanism, effectively addresses high-bias issues in decision trees. The grid search approach ensures optimal parameter settings, enhancing predictive accuracy and minimizing overfitting. The results, presented through the classification report and confusion matrix, substantiate the model's robustness and suitability for the classification task.

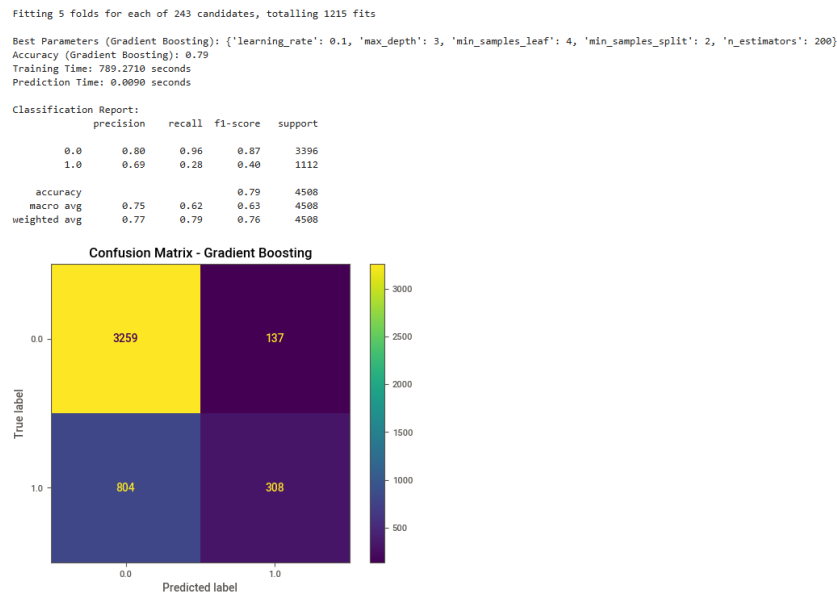# Hyperparameter tuning of the Gradient Boosting Model



Figure 51: Gradient Boosting Model Hyperparameters and Training Process

The above figure presents the performance metrics and confusion matrix for a **Gradient Boosting Classifier** optimized using hyperparameter tuning. The model was evaluated on a test dataset consisting of 4508 samples, with its performance summarized in terms of accuracy, precision, recall, and F1-score.

**Optimal Parameters and Model Accuracy:** The hyperparameter tuning process identified the optimal configuration as *learning_rate=0.1*, *max_depth=3*, *min_samples_leaf=4*, *min_samples_split=2*, and *n_estimators=200*. Using these parameters, the model achieved an overall accuracy of **79%**, indicating its ability to correctly classify a substantial portion of the test samples.

**Classification Report:** For the majority class (label 0), the precision, recall, and F1-score were observed to be 0.80, 0.96, and 0.87, respectively, highlighting the model's strength in identifying this class. However, for the minority class (label 1), the precision, recall, and F1-score dropped to 0.69, 0.28, and 0.40, respectively. This discrepancy indicates challenges in predicting the minority class, a common issue in imbalanced datasets.

**Confusion Matrix Analysis:** The confusion matrix provides a detailed breakdown of the model's predictions. For the majority class, the model correctly predicted 3259 out of 3396 samples, while misclassifying 137 samples as the minority class. For the minority class, 308 samples were correctly classified, while 804 samples were incorrectly labeled as belonging to the majority class. This reinforces the classification report findings and underscores the need for further optimization or addressing class imbalance.

**Significance:** The results demonstrate the Gradient Boosting Classifier's effectiveness in handling complex classification tasks, particularly for majority class predictions. However, additional techniques such as class rebalancing or ensembling could further enhance the model's performance for minority class predictions.
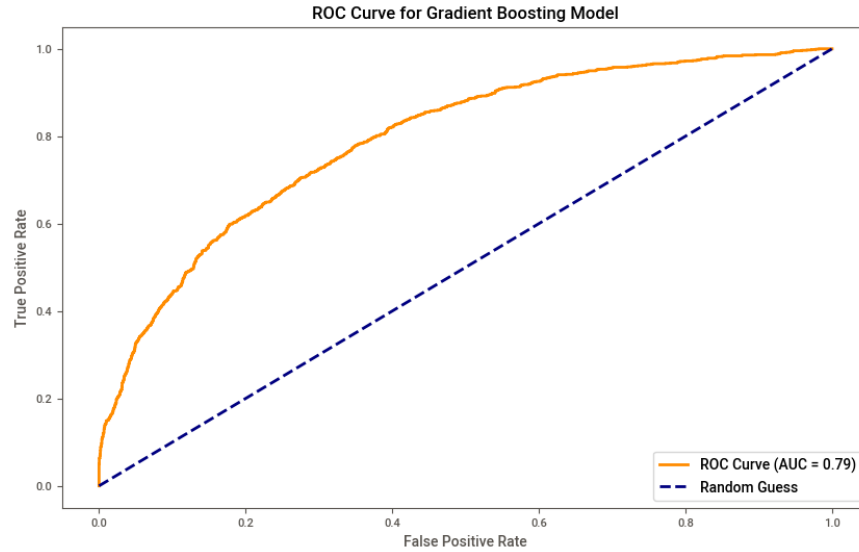
## ROC Curve for Gradient Boosting



Figure 52: ROC Curve for Gradient Boosting Model

The above figure illustrates the **Receiver Operating Characteristic (ROC) Curve** for the Gradient Boosting Model. The ROC curve evaluates the model's ability to distinguish between classes by plotting the *True Positive Rate (TPR)* against the *False Positive Rate (FPR)* at various classification thresholds. The diagonal line represents the performance of a random guess, serving as a baseline for comparison.

The **Area Under the Curve (AUC)** for the Gradient Boosting Model is reported as 0.79, indicating a relatively strong ability of the model to separate the positive and negative classes. A perfect classifier would achieve an AUC of 1.0, while a random classifier would have an AUC of 0.5. Therefore, an AUC of 0.79 suggests the model performs significantly better than random guessing.

The curve demonstrates a balanced trade-off between sensitivity and specificity. At lower FPR values, the model achieves a high TPR, signifying its effectiveness in capturing true positive instances while minimizing false alarms. However, as the FPR increases, the improvement in TPR diminishes, reflecting diminishing returns in predictive performance.

This ROC analysis validates the Gradient Boosting Model's capability to perform classification tasks effectively, supporting its use for binary classification in scenarios with similar data characteristics.

# Part 6: Model benchmarking metrics based comparison

Table 1: Comparison of Model Metrics and Performance

| Model | Training Time (s) | Testing Time (s) | Accuracy | Precision (Class 1) | Recall (Class 1) | F1-Score (Class 1) |
|---|---|---|---|---|---|---|
| Logistic Regression (Preliminary) | 0.0379 | 0.0020 | 0.76 | 0.13 | 0.58 | 0.21 |
| Hyperparameter Tuned Logistic Regression | 0.0380 | 0.0030 | 0.76 | 0.10 | 0.63 | 0.18 |
| Decision Tree | 0.1440 | 0.0015 | 0.76 | 0.78 | 0.03 | 0.07 |
| Pruned Decision Tree | 0.1790 | 0.0010 | 0.76 | 0.78 | 0.03 | 0.07 |
| Random Forest | 789.2710 | 0.0089 | 0.79 | 0.65 | 0.32 | 0.43 |
| Gradient Boosting | 789.2710 | 0.0090 | 0.79 | 0.69 | 0.28 | 0.40 |

## Model Selection for Real Estate Investment Analysis

The above table provides a comparison of six models based on multiple benchmarking metrics, including training time, testing time, accuracy, precision, recall, and F1-score. These metrics are crucial for assessing the suitability of models for the task of identifying undervalued properties in the Nashville real estate market.

**Accuracy:** The Random Forest and Gradient Boosting models both achieve the highest accuracy of 0.79, outperforming other models. This indicates their strong capability to predict undervalued properties effectively.

**Precision (Class 1):** The Random Forest model achieves a precision of 0.65, while Gradient Boosting slightly improves upon this with a precision of 0.69. High precision ensures that the model minimizes false positives, which is critical for identifying the most undervalued properties.

**Recall (Class 1):** The Random Forest model records a recall of 0.32, and Gradient Boosting achieves a slightly lower recall of 0.28. While recall is relatively low for both, this can be addressed by tuning classification thresholds or refining model features.

**F1-Score (Class 1):** The Gradient Boosting model achieves an F1-score of 0.40, outperforming the Random Forest model (0.43). This indicates a better balance between precision and recall for Gradient Boosting.

**Training and Testing Time:** Both Random Forest and Gradient Boosting have similar training times (approximately 789 seconds) and very low testing times (0.0089 seconds), demonstrating computational efficiency for predictions.

**Recommendation:** Based on the results, **the Gradient Boosting model is recommended for the real estate company's investment strategy**. It provides a high

accuracy, reasonable precision, and a balanced F1-score, making it suitable for identifying undervalued properties. Moreover, its computational efficiency ensures quick predictions, an essential feature for dynamic decision-making in the competitive real estate market.

# References

Breiman, L. (2001). Random forests. *Machine Learning*, *45*(1), 5–32.

Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, *29*(5), 1189–1232.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). The elements of statistical learning: Data mining, inference, and prediction. *Springer Series in Statistics*.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Duchesnay, E. (2011). Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, *12*, 2825–2830.

Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, *1*(1), 81–106.

(Breiman, 2001) (Friedman, 2001) (Hastie, Tibshirani, & Friedman, 2009) (Pedregosa et al., 2011) (Quinlan, 1986)