# College of Professional Studies

# Northeastern University San Jose

**MPS Analytics**

**Course: ALY6020: Predictive Analytics**

**Assignment:**

Module 4 Project - Investing in Nashville

**Submitted on:**

Feb 8, 2024

**Submitted to:**                          **Submitted by:**

Prof: BEHZAD AHMADI          NIKSHITA RANGANATHAN

# INTRODUCTION

**Understanding the dataset:**

The Housing Data dataset, is a collection of housing transactions in Nashville. This dataset is valuable for real estate market analysis, predictive modeling, and urban planning. It serves as a rich source for understanding housing trends, pricing strategies, and the impact of location on property values.

Below are the data descriptions of each variable of the data that briefly describe the contents of the data set:

- **Unnamed: 0**: An index or unique identifier for each record.
- **Parcel ID**: A unique identifier for a parcel of land in the real estate registry.
- **Land Use**: The designated use of the property, such as "SINGLE FAMILY."
- **Property Address**: The street address of the property.
- **Suite/Condo #**: If applicable, the suite or condo number of the property.
- **Property City**: The city in which the property is located, here it's "NASHVILLE."
- **Sale Date**: The date on which the property was sold.
- **Legal Reference**: A reference number for the legal document of the sale.
- **Sold As Vacant**: Indicates whether the property was vacant at the time of sale ("Yes" or "No").
- **Multiple Parcels Involved in Sale**: Indicates if the sale involved multiple parcels ("Yes" or "No").
- **Owner Name**: The name of the property's owner.
- **Address**: Likely a duplication or variation of property address details.
- **Acreage**: The size of the property in acres.
- **Tax District**: The tax district in which the property is located.
- **Land Value**: The assessed value of the land.
- **Building Value**: The assessed value of the building on the property.
- **Total Value**: The total assessed value of the property (land + building).
- **Finished Area**: The total finished area of the property in square feet.
- **Foundation Type**: The type of foundation the building has (e.g., slab, crawl space).
- **Year Built**: The year in which the building was constructed.
- **Exterior Wall**: The material of the building's exterior wall (e.g., brick, frame).
- **Grade**: An assessment of the property's overall quality or grade.
- **Bedrooms**: The number of bedrooms in the property.
- **Full Bath**: The number of full bathrooms in the property.
- **Half Bath**: The number of half bathrooms in the property.
- **Sale Price Compared To Value**: A categorical variable indicating whether the sale price was over, under, or at the assessed value ("Over", "Under").

The dataset consists of 22,651 rows and 26 columns. The dataset includes a mix of data types, with 15 object (string) columns, 6 float columns, and 5 int columns.

# DATA CLEANING

- **Handling missing values**
  'Suite/ Condo #' column has 22,651 missing entries, indicating the information is either consistently unavailable or not applicable. Minor missing data in 'Property Address', 'Property City', and a few other variables suggest targeted cleaning or imputation could easily remedy these gaps.

  The missing values in the 'Half Bath', 'Full Bath', 'Bedrooms', and 'Finished Area' columns were substituted with their respective column medians.

- **Dropping irrelevant columns**
  Columns 'Unnamed: 0', 'Suite/ Condo    #', 'Legal Reference', 'Parcel ID', 'Property Address', and 'State' were removed because they do not contribute to the analysis or the insights.

- **Removing duplicate rows**
  Duplicate rows in the dataset can lead to inconsistencies and affect the accuracy of data analysis. It is essential to find and remove any duplicate rows from the dataset before starting the analysis.
  144 rows were identified as duplicate rows and they were eliminated for enhancing data quality.

- **Simplifying Categorical Data**
  The categorical columns were mapped to binary representations (0 or 1) in three columns ('Sold As Vacant', 'Multiple Parcels Involved in Sale', and 'Sale Price Compared To Value').

- **Combining columns**
  The 'Total Bath' and 'Total Value' columns were created by combining data from the from existing columns ('Half Bath' + 'Full Bath') and ('Land Value' + 'Building Value') respectively.

- **Changing datatype**
  Column 'Sale Date' was converted to datetime format.

# EDA

**Observations from the EDA:**

**Figure 1:**

Key observations from the histograms:

- Acreage varies significantly with a minimum of 0.04 acres to a maximum of 17.5 acres, indicating a wide range of property sizes.
- Year Built ranges from 1832 to 2017, showcasing the diversity in property ages.
- Finished Area and Bedrooms suggest a variety of property sizes and capacities, from smaller homes to larger estates.

**Figure 2 &3 :**

- The correlation analysis indicates that there is a strong positive correlation between 'Finished Area' and 'Total Value', highlighting a significant relationship where larger areas are typically associated with higher property values.
- Variables such as 'Sold As Vacant,' 'Finished Area,' and 'Total Value' are positively correlated to 'Sale Price Compared To Value'.

# LOGISTIC REGRESSION

We started with encoding the categorical values present in the data. Encoding is a crucial data preprocessing step in machine learning and statistical analysis, particularly when working with categorical variables. Without encoding, models may not consider categorical data, leading to errors. This process also ensures consistency in representing categories and plays a vital role in improving model performance by allowing models to learn from categorical data.

Then, the dataset is split into two parts: one containing the features and the other with the target variable, excluding specific columns ('Sale Price Compared To Value', 'Sale Date', 'City', 'Property City') from the features.
Subsequently, the dataset is further divided into training and testing sets using a 80-20 split ratio.

A **logistic regression model (Figure 4)** is constructed. The model's accuracy is calculated, and a confusion matrix and classification report are generated to assess its performance in detail.

The **Logit** function (Figure 5) is then applied to the data, and a summary of the logistic regression analysis is produced, providing coefficients, standard errors, z-values, p-values, and confidence intervals for each predictor. This summary helps to understand the significance of each predictor and the model's overall fit to the data.

**Accuracy :** The model attained an accuracy rate of 75.6%, indicating it correctly classified about 75.6% of the instances.

**Confusion Matrix:**
- True Positives (TP): Only 6 instances of accurately identified cases where housing prices are inaccurately represented.
- True Negatives (TN): High number of instances correctly identified as accurately represented prices (3396).
- False Positives (FP): High number of instances incorrectly classified as inaccurately represented prices (1098).
- False Negatives (FN): Only 1 instance incorrectly classified as accurately represented price.

Features such as 'Sold As Vacant' and 'Multiple Parcels Involved in Sale' appear to be significant drivers of housing prices based on their coefficients.

# DECISION TREE

**Accuracy**: The accuracy of the current model is 76% (Figure 6).

**Confusion Matrix:**

- True Positives (TP): The model correctly identifies 32 instances of inaccurately represented prices.
- False Positives (FP): There are 1072 false positives, indicating instances where the model incorrectly predicts inaccurately represented prices.
- True Negatives (TN): The model correctly identifies 3388 instances of accurately represented prices.
- False Negatives (FN): Only 9 instances are incorrectly classified as accurately represented prices.

Both models have a similar accuracy, with the decision tree model achieving a slightly higher accuracy. The decision tree model identifies more true positives and fewer false negatives compared to the logistic regression model. Therefore, the decision tree model appears to be more effective in this scenario.

This decision tree model would be used to predict whether a data sample belongs to class 0 or class 1 based on features such as 'Total Bath', 'Foundation Type', 'Land Use', and 'Bedrooms'.

The decision tree (Figure 7) provides a clear and interpretable model, showing the decision rules derived from the training data. This is useful for understanding feature importance and the decision-making process in classification tasks. The first split is made based on 'Land Use', dividing the dataset into two major branches.

# RANDOM FOREST

**Accuracy**: The accuracy of the Random Forest model is 72.3% (Figure 8).

**Confusion Matrix:**

- True Positives (TP): The model correctly identifies 164 instances of inaccurately represented prices.
- False Positives (FP): There are 307 false positives, indicating instances where the model incorrectly predicts inaccurately represented prices.
- True Negatives (TN): The model correctly identifies 3090 instances of accurately represented prices.
- False Negatives (FN): There are 940 false negatives, indicating instances where the model incorrectly classifies accurately represented prices.

The random forest model has a lower accuracy and significantly higher F1 score compared to the logistic regression and decision tree models. The random forest model has a higher number of false negatives and false positives compared to the logistic regression and decision tree models. It also has a lower number of true positives and true negatives.

# GRADIENT BOOST

**Accuracy**: The Gradient Boost model achieves an accuracy of 76.1%, indicating that it correctly predicts housing prices in approximately 76.1% of cases (Figure 9).

**Confusion matrix:**
- True Positives (TP): The model correctly predicts 49 instances as inaccurately represented prices.
- False Positives (FP): There are 21 instances where the model incorrectly predicts inaccurately represented prices.
- True Negatives (TN): The model correctly predicts 3376 instances as accurately represented prices.
- False Negatives (FN): There are 1055 instances where the model incorrectly predicts accurately represented prices.

The Gradient Boost model outperforms all the previous models but has a lower F1 score compared to the random forest model.

# CONCLUSION

From this comparison (Figure 10 & 11), we can observe that Gradient Boost has the highest accuracy among the models, followed closely by Decision Tree. However, Random Forest achieves the highest F1 score, indicating better balance between precision and recall. Gradient Boost also shows a relatively higher F1 score compared to Logistic Regression and Decision Tree, indicating better performance in capturing both precision and recall.

The Logistic Regression model is simple and interpretable, making it easier to understand the relationships between features and housing prices.

# REFERENCES

Saini, A. (2024, January 10). *Gradient Boosting Algorithm: A complete guide for beginners*. Analytics Vidhya.

https://www.analyticsvidhya.com/blog/2021/09/gradient-boosting-algorithm-a-complete-guide-for-beginners/

*Figure 1 – Histograms for understanding distribution*

*Figure 2 – Correlation matrix*

```
Sale Price Compared To Value          1.000000
Sold As Vacant                        0.114710
Finished Area                         0.088250
Total Value                           0.086470
Total Bath                            0.079010
Bedrooms                              0.048628
Acreage                               0.032082
Year Built                            0.025391
Neighborhood                          0.017786
Multiple Parcels Involved in Sale    -0.017709
Name: Sale Price Compared To Value, dtype: float64
```

*Figure 3 – Correlation of other features with target variable*

*Figure 4 – Logistic Regression Model*



*Figure 5 – Logistic model result summary (statsmodels)*



*Figure 6 – Decision tree Model*

Predictive Analytics

*Figure 7 – Decision tree*



*Figure 8 – Random forest Model*



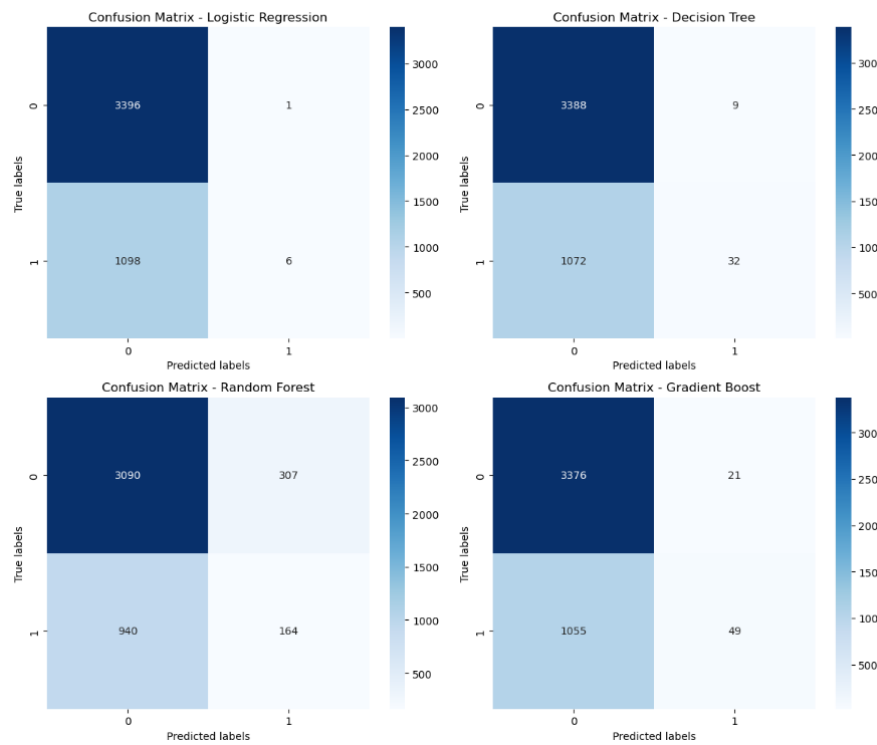*Figure 9 – Gradient boost Model*

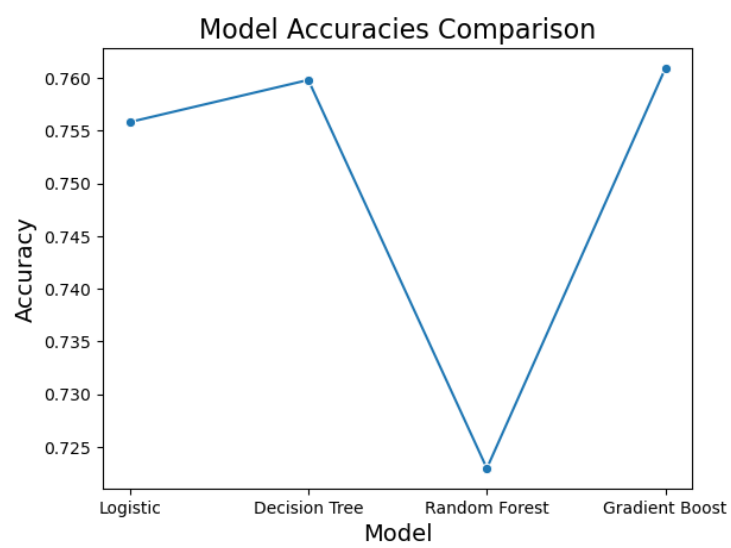Predictive Analytics

*Figure 10 – Confusion matrix for all models*



*Figure 11 – Graph comparing accuracies of models*