# REPORT

**Prepared by - Agathiyan K (045005)** PGDM - BDA04

Project Title: Demographic Patterns and Home Loan Applicant Profiling

Project Objectives | Problem Statements 1.1. PO1 | PS1: Classification of Loan Dataset into Segments | Clusters | Classes using Supervised Learning Classification Algorithms 1.2. PO2 | PS2: Determination of an Appropriate Classification Model 1.3. PO3 | PS3: Identification of Important | Contributing | Significant Variables or Features and their Thresholds for Classification

The Pre-processing Report will be in the end (after the managerial insights)

**Input Variables:** ['application_date_indicator', 'msamd_name', 'purchaser_type_name', 'loan_type_name', 'loan_purpose_name', 'hud_median_family_income', 'loan_amount_000s']

**Output Variable:** action_taken_name

### 3.2. Data Analysis 3.2.1.1. PO1 | PS1:: Supervised Machine Learning Classification Algorithm: Decision Tree (Base Model) | Metrics Used - Gini Coefficient, Entropy

Decision Tree

- Entropy: 0.15079860160668734
- Gini Impurity: 0.08122381021613079

Entropy: With an entropy value of 0.1508, the Decision Tree model using entropy as the impurity measure achieves a moderate level of disorder or unpredictability in the dataset. This indicates that the Decision Tree splits the data into nodes based on features in a way that minimizes the disorder within each node.

Gini Impurity: The Gini impurity value of 0.0812 suggests a relatively low impurity or uncertainty within the Decision Tree model. Gini impurity measures the probability of incorrectly classifying a randomly chosen element in the dataset, with lower values indicating better purity of the dataset.

### 3.2.1.2. PO1 | PS1:: Supervised Machine Learning Classification Algorithms: {Random Forest | XGBoost} (Comparison Models) | Metrics Used

**Random Forest:**

- Entropy for Random Forest: 0.4448700986425233
- Gini Impurity for Random Forest: 0.03587907150255809

Entropy: The entropy value for the Random Forest model is approximately 0.445. Entropy measures the uncertainty or disorder in the dataset. A higher entropy value indicates higher disorder, suggesting that the Random Forest model's decision boundaries are less clear or well-defined.

Gini Impurity: The Gini impurity value for the Random Forest model is approximately 0.036. Gini impurity measures the probability of misclassifying an instance in a dataset. A lower Gini impurity value indicates that the Random Forest model's decision boundaries are more precise and less prone to misclassification.

**XG Boost:**

- Entropy for XGBoost: 0.11704169720141437
- Gini Impurity for XGBoost: 0.04694253206253052

Entropy: The entropy value for XGBoost is approximately 0.117. Entropy measures the impurity or disorder in a set of data. A lower entropy value indicates that the data is more pure or well-separated into distinct classes.

Gini Impurity: The Gini impurity for XGBoost is approximately 0.047. Gini impurity measures the probability of incorrectly classifying a randomly chosen element if it were randomly labeled according to the distribution of labels in the subset. Similar to entropy, a lower Gini impurity value suggests a more pure or well-separated dataset.

### 3.2.2.1.1. PO2 | PS2:: Classification Model Performance Evaluation: Confusion Matrix {Accuracy, Recall, Precision, F1-Score} (Base Model: Decision Tree)

|   | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.00 | 0.00 | | 63 |
| 1 | 0.00 | 0.00 | | 121 |
| 2 | 0.00 | 0.00 | | 2000 |
| 3 | 0.00 | 0.00 | | 334 |
| 4 | 0.94 | 1.00 | | 0.97 |
| 5 | 1.00 | 1.00 | | 1.00 |
| 6 | 0.00 | 0.00 | | 13 |
| 7 | 0.00 | 0.00 | | 26 |

| accuracy | macro avg | weighted avg |
|---|---|---|
| 0.94 | 0.24 | 0.89 |
| | 0.25 | 0.94 |
| | 0.25 | 0.92 |

Precision: Precision measures the ratio of correctly predicted instances of a class to the total instances predicted as that class. Classes 4 and 5 have high precision scores of 0.94 and 1.00, respectively, indicating that the model correctly identifies most instances of these classes. However, other classes (0, 1, 2, 3, 6, 7) have precision scores of 0.00, suggesting that the model fails to correctly predict instances for these classes.

Recall: Recall, also known as sensitivity, measures the ratio of correctly predicted instances of a class to the total actual instances of that class. Classes 4 and 5 have high recall scores of 1.00, indicating that the model correctly identifies almost all actual instances of these classes. Similar to precision, other classes (0, 1, 2, 3, 6, 7) have recall scores of 0.00, indicating that the model fails to capture most actual instances of these classes.

F1-score: The F1-score is the harmonic mean of precision and recall, providing a balance between the two metrics. Classes 4 and 5 have high F1-scores, reflecting the balance between precision and recall for these classes. For classes with a precision or recall of 0.00, the F1-score is also 0.00, indicating poor performance.

Accuracy: The overall accuracy of the model is 0.94, indicating that it correctly predicts the class labels for 94% of the instances in the dataset.

- Observations: The model performs well in classifying instances of classes 4 and 5, achieving high precision, recall, and F1-scores. However, for other classes, the model's performance is poor, as indicated by precision, recall, and F1-scores of 0.00.

Overall, while the model demonstrates high accuracy, its effectiveness varies significantly across different classes, suggesting potential imbalance or bias in the dataset or model training process. Further investigation and possibly model refinement may be necessary to address these issues.

### 3.2.2.1.2. PO2 | PS2:: Classification Model Performance Evaluation: Time Statistics | Memory Statistics (Base Model: Decision Tree)

Decision Tree:

- Training Time (s): 0.0698244571685791
- Memory Used (MB): 713.6015625

3.2.2.2.1. PO2 | PS2:: Classification Model Performance Evaluation: Confusion Matrix {Accuracy, Recall, Precision, F1-Score} (Comparison Models: Random Forest | XGBoost)

Random Forest

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.89 | 0.47 | 0.61 | |
| 1 | 0.77 | 0.51 | 0.61 | |
| 2 | 0.84 | 0.85 | 0.84 | |
| 3 | 0.80 | 0.64 | 0.71 | |
| 4 | 0.99 | 0.99 | 0.99 | |
| 5 | 1.00 | 1.00 | 1.00 | |
| 6 | 1.00 | 0.73 | 0.85 | |
| 7 | 0.96 | 0.89 | 0.92 | |

| accuracy | macro avg | weighted avg |
|---|---|---|
| 0.98 | 0.91 | 0.98 |
| | 0.76 | 0.98 |
| | 0.82 | 0.98 |

Precision: Precision measures the ratio of correctly predicted positive observations to the total predicted positives. For the Random Forest model:

- Across all classes, the precision values indicate a generally high level of accuracy in the model's predictions. Classes 4 and 5 stand out with precision scores of 99% and 100%, respectively, indicating near-perfect prediction accuracy for these classes. Additionally, classes 0, 1, 2, 3, and 7 exhibit precision scores ranging from 77% to 96%, indicating strong prediction accuracy across a diverse range of classes. Overall, the model demonstrates the ability to make accurate predictions across multiple classes, with precision scores consistently above 77%.

Recall: Recall measures the ratio of correctly predicted positive observations to the actual positives in the dataset.

F1-score: The F1-score is the harmonic mean of precision and recall and provides a balance between the two metrics.

Support: Support represents the number of actual occurrences of each class in the specified dataset.

- Observations: The Random Forest model demonstrates high precision and recall across most classes, indicating its effectiveness in correctly classifying instances across different categories. The model achieves an overall accuracy of 98%, suggesting strong performance in classifying instances correctly. The macro average and weighted average F1-scores are high, indicating a good balance between precision and recall across all classes. The model performs exceptionally well in predicting classes 4 and 5, as indicated by their high precision, recall, and F1-scores. Class 0 and Class 1 have relatively lower precision and recall compared to other classes, suggesting room for improvement in predicting these classes.

Overall, the Random Forest model demonstrates robust performance across multiple evaluation metrics, highlighting its effectiveness in classification tasks.

**XG Boost**

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.88 | 0.21 | 0.34 |
| 1 | 0.84 | 0.20 | 0.32 |
| 2 | 0.73 | 0.74 | 0.74 |
| 3 | 0.75 | 0.38 | 0.50 |
| 4 | 0.98 | 0.99 | 0.99 |
| 5 | 1.00 | 1.00 | 1.00 |
| 6 | 0.88 | 0.47 | 0.61 |
| 7 | 0.96 | 0.89 | 0.92 |

| accuracy | macro avg | weighted avg |
|---|---|---|
| 0.97 | 0.88 | 0.97 |
| | 0.61 | 0.97 |
| | 0.68 | 0.97 |

Precision: Precision measures the proportion of true positive predictions among all positive predictions made by the model. It indicates the model's ability to avoid false positives. Here:

- The precision values for the different classes vary, ranging from 73% to 100%. Overall, the model demonstrates strong prediction accuracy, with most classes having precision scores above 80%. Classes 4 and 5 stand out with precision scores of 98% and 100%, respectively, indicating very high accuracy in predicting these classes. However, classes 0, 1, 6, and 7 also show respectable precision scores ranging from 84% to 96%. Class 2 has the lowest precision at 73%, suggesting that predictions for this class may be less accurate compared to others. Nonetheless, the model's performance across most classes indicates reliable predictive capability.

Recall: Recall, also known as sensitivity, measures the proportion of true positives that were correctly identified by the model among all actual positives in the dataset. Here:

- The recall values for the different classes show considerable variation, ranging from 20% to 100%. Classes 4 and 5 have exceptionally high recall scores of 99% and 100%, respectively, indicating that the model effectively captures the majority of instances belonging to these classes. Conversely, classes 0 and 1 have low recall values of 21% and 20%, respectively, suggesting that the model misses a significant portion of instances belonging to these classes. Classes 2, 3, 6, and 7 exhibit moderate recall scores ranging from 38% to 74%, indicating varying degrees of effectiveness in capturing instances from these classes. Overall, the recall values provide insight into how well the model identifies instances from each class, with higher values indicating better performance in correctly identifying instances.

F1-score: The F1-score is the harmonic mean of precision and recall. It provides a balance between precision and recall. Here:

- The F1-scores for the different classes exhibit a wide range of values, indicating varying levels of model performance across classes. Classes 4 and 5 have exceptionally high F1-scores of 0.99 and 1.00, respectively, suggesting excellent precision and recall, resulting in high overall model performance for these classes. Conversely, classes 0 and 1 have low F1-scores of 0.34 and 0.32, respectively, indicating poorer model performance in terms of both precision and recall for these classes. Classes 2, 3, and 6

exhibit moderate F1-scores, ranging from 0.50 to 0.74, indicating relatively balanced performance in terms of precision and recall for these classes. Class 7 also shows a high F1-score of 0.92, indicating a good balance between precision and recall for this class. Overall, the F1-scores provide a comprehensive measure of model performance, considering both precision and recall, and highlight areas where the model may need improvement for certain classes.

Support: Support indicates the number of actual occurrences of each class in the dataset.

Accuracy: The overall accuracy of the model is 97%. Accuracy measures the proportion of correctly classified instances among all instances in the dataset.

- Observations: The XGBoost model demonstrates high precision and recall for most classes, particularly for classes 4 and 5, indicating its effectiveness in correctly classifying instances across multiple classes. Classes 0, 1, and 3 have relatively lower precision and recall, suggesting challenges in correctly identifying instances belonging to these classes. The weighted average F1-score of 0.97 indicates overall good performance of the XGBoost model in terms of balancing precision and recall across all classes.

### 3.2.2.2.2. PO2 | PS2:: Classification Model Performance Evaluation: Time Statistics | Memory Statistics (Comparison Models: Random Forest | XGBoost)

**Random Forest**

- Training Time (s): 5.672376871109009
- Memory Used (MB): 885.76953125

**XG Boost**

- Training Time (s): 45.33557963371277
- Memory Used (MB): 897.80859375

### 3.2.3.1. PO3 | PS3:: Variable or Feature Analysis: Base Model (Decision Tree)

| feature | importance |
| --- | --- |
| application_date_indicator | 0.479 |
| msamd_name | 0.292 |
| purchaser_type_name | 0.229 |
| loan_type_name | 0.000 |
| loan_purpose_name | 0.000 |
| hud_median_family_income | 0.000 |
| loan_amount_000s | 0.000 |

The most important feature is application_date_indicator with an importance value of 0.479. The second most important feature is msamd_name with an importance value of 0.292. The third most important feature is purchaser_type_name with an importance value of 0.229. Features loan_type_name, loan_purpose_name, hud_median_family_income, and loan_amount_000s have importance values of 0.000, indicating they are not contributing significantly to the model's predictions.

### 3.2.3.2. PO3 | PS3:: Variable or Feature Analysis: Comparison Models (Random Forest | XGBoost) Random Forest

Feature Importances:

| feature | Importance |
|---|---|
| loan_amount_000s | 0.333727 |
| purchaser_type_name | 0.201634 |
| application_date_indicator | 0.195447 |
| hud_median_family_income | 0.120389 |
| msamd_name | 0.096998 |
| loan_type_name | 0.027990 |
| loan_purpose_name | 0.023816 |

The most important feature is loan_amount_000s with an importance value of 0.334. The second most important feature is purchaser_type_name with an importance value of 0.202. The third most important feature is application_date_indicator with an importance value of 0.195. Features hud_median_family_income and msamd_name also have relatively significant importance values of 0.120 and 0.097, respectively. Features loan_type_name and loan_purpose_name have lower importance values, indicating they contribute less to the model's predictions.

**XG Boost**

| feature | Importance |
|---|---|
| application_date_indicator | 0.765484 |
| purchaser_type_name | 0.199578 |
| hud_median_family_income | 0.014517 |
| msamd_name | 0.012562 |
| loan_type_name | 0.003441 |
| loan_purpose_name | 0.003015 |
| loan_amount_000s | 0.001402 |

The most important feature is application_date_indicator with a high importance value of 0.765. The second most important feature is purchaser_type_name with a relatively lower importance value of 0.200. Features hud_median_family_income and msamd_name have very low importance values of 0.015 and 0.013, respectively. Features loan_type_name, loan_purpose_name, and loan_amount_000s contribute even less to the model's predictions, with importance values below 0.01.

**Results | Observations 4.1. Classification Model Parameters: Base Model (Decision Tree) | Comparison Models (Random Forest | XGBoost)**

| Algorithm | Entropy | Gini Impurity |
|---|---|---|
| Decision Tree | 0.1508 | 0.0812 |
| Random Forest | 0.4449 | 0.0359 |
| XGBoost | 0.1170 | 0.0469 |

The entropy values indicate the average uncertainty in classifying a sample in the dataset. Lower entropy values suggest better decision tree purity, indicating more accurate predictions. Random Forest shows the highest entropy, indicating higher uncertainty compared to Decision Tree and XGBoost. Gini impurity measures the probability of misclassifying a sample. Lower Gini impurity values suggest better decision tree purity and thus, more accurate predictions. Random Forest exhibits the lowest Gini impurity, indicating better purity and classification performance compared to Decision Tree and XGBoost.

## 4.2. Classification Model Performance: Time & Memory Statistics [Base Model (Decision Tree) | Comparison Models (Random Forest | XGBoost)]

| Algorithm | Training Time (s) | Memory Used (MB) |
|---|---|---|
| Decision Tree | 0.0698 | 713.60 |
| Random Forest | 5.6724 | 885.77 |
| XGBoost | 45.3356 | 897.81 |

Decision Tree has the shortest training time, indicating fast model training compared to Random Forest and XGBoost. Random Forest exhibits longer training time than Decision Tree but shorter than XGBoost, suggesting a trade-off between training time and model complexity. XGBoost requires the longest training time, indicating more computationally intensive model training compared to Decision Tree and Random Forest. Memory usage is highest for XGBoost, followed by Random Forest and then Decision Tree, indicating the memory requirements of each algorithm.

## 4.3. Variable or Feature Analysis: Base Model (Decision Tree) | Comparison Models (Random Forest | XGBoost)

| Algorithm | Precision | Recall | F1-score |
|---|---|---|---|
| Decision Tree | 0.89 | 0.94 | 0.92 |
| Random Forest | 0.98 | 0.98 | 0.98 |
| XGBoost | 0.97 | 0.97 | 0.97 |

Random Forest achieves the highest precision, recall, and F1-score among the three algorithms, indicating superior overall performance in terms of classification accuracy. Decision Tree exhibits the lowest precision, recall, and F1-score, suggesting comparatively lower classification accuracy compared to Random Forest and XGBoost. XGBoost achieves slightly lower precision, recall, and F1-score than Random Forest but higher than Decision Tree, indicating its performance falls between the other two algorithms.

### 4.3.1. List of Relevant or Important Variables or Features and their Thresholds

| Feature | Importance (Decision Tree) | Importance (Random Forest) | Importance (XGBoost) |
|---|---|---|---|
| application_date_indicator | 0.479 | 0.195 | 0.765 |
| purchaser_type_name | 0.229 | 0.202 | 0.200 |
| msamd_name | 0.292 | 0.097 | 0.013 |
| loan_type_name | 0.000 | 0.028 | 0.003 |
| loan_purpose_name | 0.000 | 0.024 | 0.003 |
| hud_median_family_income | 0.000 | 0.120 | 0.015 |
| loan_amount_000s | 0.000 | 0.334 | 0.001 |

Analysis:

Application Date Indicator: This feature appears to be crucial across all models, especially in XGBoost, where it has the highest importance. This suggests that the application date might have a significant impact on the outcome.

Purchaser Type Name: While it has a considerable importance in all models, it ranks lower than the application date indicator. It indicates that the type of purchaser involved might also play a role in the outcome.

MSAMD Name: Despite being significant in the Decision Tree model, its importance diminishes in Random Forest and XGBoost. This suggests that its predictive power might not be as strong as other features.

Loan Type Name and Loan Purpose Name: These features have negligible importance across all models, indicating they might not contribute much to the predictions.

HUD Median Family Income: While it has some importance in Random Forest, it's relatively low in XGBoost. This might indicate that income levels of families in certain areas have varying impacts depending on the model.

Loan Amount (000s): It's significant only in the Random Forest model. This could mean that loan amount plays a more critical role in its decision-making process compared to other models.

**Best Method:** Based on the provided data and analysis, it's challenging to determine the absolute "best" method as it heavily depends on various factors such as dataset characteristics, interpretability requirements, computational resources, and specific objectives of the analysis.

- However, based on the insights gained: Decision Trees are easy to interpret and computationally efficient but might suffer from overfitting. Random Forests tend to provide better generalization by averaging multiple decision trees and handle overfitting well. XGBoost often delivers superior performance due to its boosting technique, which iteratively improves model performance, but it may require more computational resources and tuning effort. Therefore, if interpretability and computational efficiency are crucial, Decision Trees might be preferred. If high predictive accuracy and generalization are the primary objectives, Random Forests or XGBoost could be more suitable depending on the specific requirements and computational resources available.

### 4.3.2. List of Non-Relevant or Non-Important Variables or Features

| Feature | Importance (Decision Tree) | Importance (Random Forest) | Importance (XGBoost) |
|---|---|---|---|
| loan_type_name | 0.000 | 0.028 | 0.003 |
| loan_purpose_name | 0.000 | 0.024 | 0.003 |
| hud_median_family_income | 0.000 | 0.120 | 0.015 |
| loan_amount_000s | 0.000 | 0.334 | 0.001 |

Analysis:

Loan Type Name and Loan Purpose Name: These variables are consistently deemed unimportant across all three models. This suggests that the type of loan and its purpose might not significantly influence the outcome in this context.

HUD Median Family Income: While this feature has some importance in Random Forest, it ranks low in XGBoost and is considered unimportant in Decision Trees. This disparity could indicate that its predictive power varies between models or that it interacts with other variables differently.

Loan Amount (000s): This feature is deemed unimportant in both Decision Trees and XGBoost but holds significant importance in Random Forest. This discrepancy suggests that the importance of loan amount might vary depending on the modeling technique used.

The importance of variables can vary significantly between different modeling techniques, highlighting the importance of exploring multiple approaches to gain a comprehensive understanding of the data. Features that are deemed unimportant in one model may still have predictive power in others, underscoring the need for model comparison and interpretation. Non-important variables can be removed from the model to simplify it without sacrificing predictive performance, potentially improving model efficiency and interpretability.

Overall, understanding the importance of variables across different models provides valuable insights into the underlying relationships within the data and aids in making informed decisions about feature selection and model optimization.

**Managerial Insights**

**5.1. Appropriate Model: Compare and Contrast {Decision Tree | Random Forest | XGBoost}**

| Algorithm | Precision (weighted avg) | Recall (weighted avg) | F1-score (weighted avg) | Training Time (s) | Memory Used (MB) |
|---|---|---|---|---|---|
| Decision Tree | 0.89 | 0.94 | 0.92 | 0.0698 | 713.60 |
| Random Forest | 0.98 | 0.98 | 0.98 | 5.6724 | 885.77 |
| XGBoost | 0.97 | 0.97 | 0.97 | 45.3356 | 897.81 |

**Insights:** Performance Metrics: Precision, Recall, and F1-score: Random Forest outperforms both Decision Tree and XGBoost in terms of precision, recall, and F1-score, indicating superior overall predictive performance. Consistency: Decision Tree and XGBoost exhibit similar performance across precision, recall, and F1-score metrics, suggesting comparable effectiveness in classification tasks. Resource Utilization: Training Time: Decision Tree has the lowest training time, indicating quick model development. Random Forest and XGBoost require significantly more time due to their ensemble and gradient boosting nature, respectively. Memory Usage: Decision Tree consumes the least memory, making it suitable for resource-constrained environments. However, Random Forest and XGBoost utilize more memory due to their complex ensemble structures and the need to store multiple trees.

**Trade-offs and Considerations:** Accuracy vs. Resource Consumption: Random Forest offers the highest accuracy but requires more computational resources. Managers should weigh the trade-off between model accuracy and resource utilization based on their specific requirements and constraints. Time Sensitivity: Decision Tree's quick training time makes it ideal for scenarios where rapid model deployment is critical. However, if accuracy is paramount and time is not a constraint, Random Forest or XGBoost may be preferable despite longer training times.

**Model Selection Strategy:** Use Case Suitability: Decision Tree may be preferred for quick exploratory analysis or when interpretability is essential. Random Forest and XGBoost are better suited for scenarios where high predictive accuracy is required, such as fraud detection or customer churn prediction. Iterative Improvement: If initial results with Decision Tree are promising but lack the desired accuracy, managers can consider transitioning to Random Forest or XGBoost for iterative improvement without sacrificing interpretability entirely.

Overall, managers should carefully evaluate the trade-offs between model performance, training time, and resource consumption to select the most suitable algorithm for their specific business objectives and constraints.

**5.2. Relevant or Important Variables or Features (Given the Appropriate Model)**

| Feature | Importance (Decision Tree) | Importance (Random Forest) | Importance (XGBoost) |
|---|---|---|---|
| application_date_indicator | 0.479 | 0.196 | 0.765 |
| msamd_name | 0.292 | 0.097 | 0.013 |
| purchaser_type_name | 0.229 | 0.202 | 0.200 |
| hud_median_family_income | 0.000 | 0.120 | 0.015 |
| loan_type_name | 0.000 | 0.028 | 0.003 |
| loan_purpose_name | 0.000 | 0.024 | 0.003 |

| Feature | Importance (Decision Tree) | Importance (Random Forest) | Importance (XGBoost) |
|---|---|---|---|
| loan_amount_000s | 0.000 | 0.335 | 0.001 |

Application Date Indicator: This variable consistently shows high importance across all three methods, particularly in Decision Trees and XGBoost. Its importance suggests that the timing of loan applications has a significant impact on the outcome, implying potential seasonality or temporal trends that affect decision-making.

MSAMD Name: While moderately important in Decision Trees, MSAMD Name holds relatively low importance in Random Forest and XGBoost. This disparity indicates that the geographic area, represented by MSAMD Name, might have varying levels of influence on the outcome depending on the modeling technique used.

Purchaser Type Name: Purchaser Type Name ranks consistently high in importance across all methods, suggesting its significant impact on the outcome. Understanding the different categories of purchasers and their behaviors could provide valuable insights for targeted marketing or risk assessment strategies.

HUD Median Family Income: Despite being deemed unimportant in Decision Trees, HUD Median Family Income holds substantial importance in Random Forest. This discrepancy warrants further investigation to understand its varying impact on the outcome and its interaction with other variables.

Loan Type Name and Loan Purpose Name: These variables are consistently deemed unimportant across all methods. Simplifying the model by removing these features may streamline the decision-making process without sacrificing predictive performance.

Loan Amount (000s): While considered unimportant in Decision Trees and XGBoost, Loan Amount is highly significant in Random Forest. The differing importance of this variable suggests that its impact on the outcome may be model-dependent, emphasizing the need for model comparison and interpretation.

- Decision Trees: Interpretability: Decision Trees offer straightforward interpretation, as they represent decision rules in a tree-like structure, making them easy to understand for non-technical stakeholders. Variable Importance: Decision Trees provide feature importance metrics, allowing managers to identify key drivers of the outcome. Simplicity: Decision Trees are simple and intuitive, making them suitable for quick decision-making and hypothesis generation. Overfitting: They are prone to overfitting, especially with complex datasets, which may lead to poor generalization performance on unseen data. Single Tree Limitation: Decision Trees may lack predictive accuracy compared to ensemble methods like Random Forest and XGBoost, particularly when dealing with high-dimensional or noisy data.
- Random Forest: Ensemble Learning: Random Forest is an ensemble learning technique that combines multiple Decision Trees to improve predictive performance and reduce overfitting. Robustness: Random Forest is robust to overfitting and noise, making it suitable for a wide range of datasets and predictive tasks. Variable Importance: It provides feature importance scores, enabling managers to prioritize variables based on their contribution to the model. Computational Efficiency: Random Forest can handle large datasets efficiently, thanks to its parallelized training process. Black Box Nature: While Random Forest improves prediction accuracy, its complex ensemble nature makes it less interpretable compared to individual Decision Trees.
- XGBoost: Gradient Boosting: XGBoost is an advanced implementation of gradient boosting, which builds models sequentially, iteratively improving upon the residuals of the previous models. High Accuracy: XGBoost often achieves state-of-the-art performance on structured/tabular data, making it well-suited for predictive modeling tasks where accuracy is paramount. Feature Importance: Like Random Forest, XGBoost provides feature importance scores, aiding managers in understanding which variables drive predictions. Regularization: XGBoost offers various regularization techniques to prevent overfitting, such as shrinkage (learning rate) and tree depth control. Resource Intensive: Training XGBoost models can be computationally expensive and may require tuning hyperparameters to achieve optimal performance.

**Considerations:**

Model Selection: Managers should select the appropriate method based on the trade-off between interpretability, predictive accuracy, and computational resources.

Data Complexity: Consider the complexity of the dataset, including the number of features, the presence of interactions, and the volume of data, when choosing the modeling approach.

Interpretability vs. Accuracy: Balance the need for model interpretability with the desire for high predictive accuracy, depending on the specific business requirements.

Validation and Monitoring: Continuously validate and monitor model performance to ensure that it remains effective over time and in response to changing data patterns.

# Preprocessing Report

Project Title: Demographic Patterns and Home Loan Applicant Profiling

1. Project Objectives | Problem Statements 1.1. PO1 | PS1: Classification of Consumer Data into Segments | Clusters | Classes using Supervised Learning Classification Algorithms 1.2. PO2 | PS2: Determination of an Appropriate Classification Model 1.3. PO3 | PS3: Identification of Important | Contributing | Significant Variables or Features and their Thresholds for Classification

2. Description of Data 2.1. Data Source, Size, Shape 2.1.1. Data Source (Website Link)
   https://www.kaggle.com/datasets/miker400/washington-state-home-mortgage-hdma2016
   (https://www.kaggle.com/datasets/miker400/washington-state-home-mortgage-hdma2016)
   https://drive.google.com/file/d/13fp1-YgAuSiR_bWZwetJiEcJZOeDeAtu/view?usp=sharing
   (https://drive.google.com/file/d/13fp1-YgAuSiR_bWZwetJiEcJZOeDeAtu/view?usp=sharing)

   2.1.2. Data Size (in KB | MB | GB …) 30.1 MB

   2.1.3. Data Shape (Dimension: Number of Variables | Number of Records) 39 Variables (39 columns) Maximum number of rows 60,000 Total number of records: 60000 Total number of filled cells: 2038780 Missed cells: 301220

2.2. Description of Variables

2.2.1. Index Variable(s): I1, I2, … The index variable is S.no

2.2.2. Outcome Variable or Feature: OV The outcome variable, labeled as 'action_taken_name', represents the target variable in our analysis. It is the variable of interest that we aim to predict or classify using our machine learning models.

2.2.3. Input Variables or Features having Categories | Input Categorical Variables or Features (ICV)

['msamd_name', 'loan_type_name', 'loan_purpose_name']

2.2.3.1. Input Variables or Features having Nominal Categories | Categorical Variables or Features - Nominal Type: ICNV1, ICNV2, … All the categorical variables available in the dataset for nominal variables

2.2.3.2. Input Variables or Features having Ordinal Categories | Categorical Variables or Features - Ordinal Type: ICOV1, ICOV2, … No ordinal data available in the dataset

2.2.3. Input Non-Categorical Variables or Features: INCV1, INCV2, …

['hud_median_family_income', 'loan_amount_000s']

2.3. Descriptive Statistics 2.3.1. Descriptive Statistics: Outcome Variable or Feature (Categorical) 2.3.1.1. Count | Frequency Statistics count unique

| Cluster_Label | Count |
|---:|---:|
| 3 | 14162 |
| 1 | 14049 |
| 0 | 11435 |
| 4 | 10743 |
| 2 | 9611 |

2.3.1.2. Proportion (Relative Frequency) Statistics

| Cluster_Label | Proportion |
|---:|---:|
| 3 | 0.236033 |
| 1 | 0.234150 |
| 0 | 0.190583 |
| 4 | 0.179050 |
| 2 | 0.160183 |

## 2.3.2. Descriptive Statistics: Input Categorical Variables or Features 2.3.2.1. Count | Frequency Statistics

| Variable | Count Unique |
|---:|---:|
| state_name | 1 |
| state_abbr | 1 |
| respondent_id | 593 |
| purchaser_type_name | 10 |
| property_type_name | 3 |
| preapproval_name | 3 |
| owner_occupancy_name | 3 |
| msamd_name | 14 |
| loan_type_name | 4 |
| loan_purpose_name | 3 |
| lien_status_name | 4 |
| hoepa_status_name | 2 |
| county_name | 39 |
| co_applicant_sex_name | 5 |
| co_applicant_ethnicity_name | 5 |
| applicant_sex_name | 4 |
| applicant_ethnicity_name | 4 |
| agency_name | 6 |
| agency_abbr | 6 |
| action_taken_name | 8 |

| Variable | Value |
|---:|---:|
| state_name | Washington |
| state_abbr | WA |
| respondent_id | 32489 |
| purchaser_type_name | Loan was not originated or was not sold in calendar year covered by the loan/application register |
| property_type_name | One-to-four family dwelling (other than manufactured housing) |
| preapproval_name | Not applicable |
| owner_occupancy_name | Owner-occupied as a principal dwelling |
| msamd_name | Seattle, Bellevue, Everett - WA |

| Variable | Value |
|---|---|
| loan_type_name | Conventional |
| loan_purpose_name | Refinancing |
| lien_status_name | Secured by a first lien |
| hoepa_status_name | Not a HOEPA loan |
| county_name | King County |
| co_applicant_sex_name | No co-applicant |
| co_applicant_ethnicity_name | No co-applicant |
| applicant_sex_name | Male |
| applicant_ethnicity_name | Not Hispanic or Latino |
| agency_name | Department of Housing and Urban Development |
| agency_abbr | HUD |
| action_taken_name | Loan originated |

| Variable | Frequency |
|---|---|
| state_name | 60000 |
| state_abbr | 60000 |
| respondent_id | 5006 |
| purchaser_type_name | 16112 |
| property_type_name | 57630 |
| preapproval_name | 47832 |
| owner_occupancy_name | 53940 |
| msamd_name | 17965 |
| loan_type_name | 42917 |
| loan_purpose_name | 28576 |
| lien_status_name | 57046 |
| hoepa_status_name | 59996 |
| county_name | 12915 |
| co_applicant_sex_name | 26987 |
| co_applicant_ethnicity_name | 26987 |
| applicant_sex_name | 37070 |
| applicant_ethnicity_name | 44014 |
| agency_name | 27514 |
| agency_abbr | 27514 |
| action_taken_name | 55815 |

In the context of catdf dataset: state_name and state_abbr: These columns have only one unique value, which is "Washington" for state_name and "WA" for state_abbr. This suggests that these columns may not provide much information for analysis as they have constant values for all rows. respondent_id: This column has 593 unique values, and the most frequent respondent_id is "32489" with a frequency of 5006. This column likely identifies different respondents. Other categorical columns: Each column represents a categorical variable, and the summary provides information about the number of unique categories, the most

frequent category (top), and its frequency. action_taken_name: This column represents the action taken for the loan application. It has 8 unique values, and "Loan originated" is the most frequent action with a frequency of 55815.

### 2.3.2.2. Proportion (Relative Frequency) Statistics

| Variable | Frequency |
| --- | --- |
| state_name | Washington: 100.00% |
| state_abbr | WA: 100.00% |
| respondent_id | 32489: 8.34% |
| purchaser_type_name | Loan was not originated or was not sold in cal... |
| property_type_name | One-to-four family dwelling (other than manufa... |
| preapproval_name | Not applicable: 79.72% |
| owner_occupancy_name | Owner-occupied as a principal dwelling: 89.90% |
| msamd_name | Seattle, Bellevue, Everett - WA: 34.80% |
| loan_type_name | Conventional: 71.53% |
| loan_purpose_name | Refinancing: 47.63% |
| lien_status_name | Secured by a first lien: 95.08% |
| hoepa_status_name | Not a HOEPA loan: 99.99% |
| county_name | King County: 21.56% |
| co_applicant_sex_name | No co-applicant: 44.98% |
| co_applicant_ethnicity_name | No co-applicant: 44.98% |
| applicant_sex_name | Male: 61.78% |
| applicant_ethnicity_name | Not Hispanic or Latino: 73.36% |
| agency_name | Department of Housing and Urban Development: 4... |
| agency_abbr | HUD: 45.86% |
| action_taken_name | Loan originated: 93.03% |

### 2.3.3. Descriptive Statistics: Input Non-Categorical Variables or Features 2.3.3.1. Measures of Central Tendency

| Variable | Count | Mean/Std/Min/25%/50%/75%/Max |
| --- | --- | --- |
| tract_to_msamd_income | 59878 | Mean: 107.62, Std: 28.23, Min: 14.05, 25%: 88.97, 50%: 105.55, 75%: 123.33, Max: 257.14 |
| population | 59878 | Mean: 5278.78, Std: 1716.10, Min: 98.00, 25%: 4070.00, 50%: 5145.00, 75%: 6382.00, Max: 13025.00 |
| minority_population | 59878 | Mean: 23.24, Std: 14.42, Min: 2.04, 25%: 12.95, 50%: 19.42, 75%: 29.68, Max: 94.79 |
| number_of_owner_occupied_units | 59876 | Mean: 1399.04, Std: 518.33, Min: 15.00, 25%: 1034.00, 50%: 1359.00, 75%: 1722.00, Max: 2997.00 |
| number_of_1_to_4_family_units | 59878 | Mean: 1873.28, Std: 738.51, Min: 27.00, 25%: 1414.00, 50%: 1770.00, 75%: 2249.00, Max: 5893.00 |
| loan_amount_000s | 60000 | Mean: 291.36, Std: 604.96, Min: 1.00, 25%: 170.00, 50%: 242.00, 75%: 337.00, Max: 55000.00 |
| hud_median_family_income | 59878 | Mean: 73869.41, Std: 12811.24, Min: 48700.00, 25%: 63100.00, 50%: 73300.00, 75%: 90300.00, Max: 90300.00 |

| Variable | Count | Mean/Std/Min/25%/50%/75%/Max |
|---|---|---|
| applicant_income_000s | 53630 | Mean: 112.82, Std: 122.86, Min: 1.00, 25%: 61.00, 50%: 89.00, 75%: 132.00, Max: 6161.00 |
| sequence_number | 60000 | Mean: 77526.47, Std: 150515.70, Min: 1.00, 25%: 3231.50, 50%: 16481.00, 75%: 72762.25, Max: 1241590.00 |
| census_tract_number | 59878 | Mean: 1750.60, Std: 3359.68, Min: 1.00, 25%: 114.02, 50%: 403.02, 75%: 713.10, Max: 9757.00 |
| as_of_year | 60000 | Mean: 2016.00, Std: 0.00, Min: 2016.00, 25%: 2016.00, 50%: 2016.00, 75%: 2016.00, Max: 2016.00 |
| application_date_indicator | 60000 | Mean: 0.03, Std: 0.23, Min: 0.00, 25%: 0.00, 50%: 0.00, 75%: 0.00, Max: 2.00 |

2.3.3.3. Correlation Statistics (with Test of Correlation)

1. Analysis of Data 3.1. Data Pre-Processing 3.1.1. Missing Data Statistics and Treatment 3.1.1.1.1. Missing Data Statistics: Records Number of rows with missing data: 60000 Number of rows with more than 50% missing data: 0

3.1.1.1.2. Missing Data Treatment: Records 3.1.1.1.2.1. Removal of Records with More Than 50% Missing Data: None | R1, R2, … No rows with more than 50% missing values

3.1.1.2.1. Missing Data Statistics: Categorical Variables or Features

| Variable | Missing Records | Percentage Missing |
|---|---|---|
| denial_reason_name_3 | 59999 | 99.998333 |
| denial_reason_name_2 | 59957 | 99.928333 |
| denial_reason_name_1 | 59882 | 99.803333 |
| rate_spread | 58134 | 96.890000 |
| edit_status_name | 47549 | 79.248333 |
| msamd_name | 8381 | 13.968333 |
| applicant_income_000s | 6370 | 10.616667 |
| number_of_owner_occupied_units | 124 | 0.206667 |
| census_tract_number | 122 | 0.203333 |
| tract_to_msamd_income | 122 | 0.203333 |
| hud_median_family_income | 122 | 0.203333 |
| number_of_1_to_4_family_units | 122 | 0.203333 |
| minority_population | 122 | 0.203333 |
| population | 122 | 0.203333 |
| county_name | 92 | 0.153333 |

3.1.1.2.2. Missing Data Treatment: Categorical Variables or Features

3.1.1.2.2.1. Removal of Variables or Features with More Than 50% Missing Data: None | CV1, CV2, … Removed the below columns as they have more than 50% data missing • denial_reason_name_3 • denial_reason_name_2 • denial_reason_name_1 • rate_spread (non-cat) • edit_status_name

3.1.1.2.2.2. Imputation of Missing Data using Descriptive Statistics: Mode

3.1.1.3.1. Missing Data Statistics: Non-Categorical Variables or Features

| Feature | Missing Records |
|---|---|
| tract_to_msamd_income | 122 |
| population | 122 |
| minority_population | 122 |
| number_of_owner_occupied_units | 124 |
| number_of_1_to_4_family_units | 122 |
| loan_amount_000s | 0 |
| hud_median_family_income | 122 |
| applicant_income_000s | 6370 |
| sequence_number | 0 |
| census_tract_number | 122 |
| as_of_year | 0 |
| application_date_indicator | 0 |

3.1.1.3.2. Missing Data Treatment: Non-Categorical Variables or Features 3.1.1.3.2.1. Removal of Variables or Features with More Than 50% Missing Data: None | NCV1, NCV2, … • rate_spread

3.1.1.3.2.2. Imputation of Missing Data using Descriptive Statistics: Mean | Median • Imputing the missing values using mean

3.1.2. Numerical Encoding of Categorical Variables or Features (Encoding Schema - Alphanumeric Order) (Encoding Schema - Alphanumeric Order)

| Feature | Number of Unique Values |
|---|---|
| state_name | 1 |
| state_abbr | 1 |
| respondent_id | 593 |
| purchaser_type_name | 10 |
| property_type_name | 3 |
| preapproval_name | 3 |
| owner_occupancy_name | 3 |
| msamd_name | 14 |
| loan_type_name | 4 |
| loan_purpose_name | 3 |
| lien_status_name | 4 |
| hoepa_status_name | 2 |
| county_name | 39 |
| co_applicant_sex_name | 5 |
| co_applicant_ethnicity_name | 5 |
| applicant_sex_name | 4 |
| applicant_ethnicity_name | 4 |
| agency_name | 6 |
| agency_abbr | 6 |

| Feature | Number of Unique Values |
|---|---|
| action_taken_name | 8 |

We are converting the above variables into numeric format in the alpha numeric order

3.1.3. Outlier Statistics and Treatment (Scaling | Transformation) 3.1.3.1.1. Outlier Statistics: Non-Categorical Variables or Features

Outliers count for the Non-Categorical Variables

| Feature | Number of Unique Values |
|---|---|
| tract_to_msamd_income | 1309 |
| population | 553 |
| minority_population | 2641 |
| number_of_owner_occupied_units | 478 |
| number_of_1_to_4_family_units | 2150 |
| loan_amount_000s | 2467 |
| hud_median_family_income | 0 |
| applicant_income_000s | 3765 |
| sequence_number | 7898 |
| census_tract_number | 9391 |
| as_of_year | 0 |
| application_date_indicator | 776 |

3.1.3.1.2. Outlier Treatment: Non-Categorical Variables or Features 3.1.3.1.2.1. Standardization: OV1, OV2, … 3.1.3.1.2.2. Normalization using Min-Max Scaler: OV3, OV4, … 3.1.3.1.2.3. Log Transformation: OV5, OV6, … I performed scaling using normalization using min-max scaler. But post the scaling, bubbles were still visible in the box plot. This signifies that the outliers present in the non categorical datset are not heavily influenced by the scaling method. The count of outliers seems consistent across different scaling methods.

### 3.1.4. Data Bifurcation: Training & Testing Sets [Bifurcation Schema: Random Sampling or Stratified Sampling (Based on Outcome Variable or Feature) with {70% | 75% | 80%} Data in Training Set and {30% | 25% | 20%} Data in Testing Set]

The dataset was systematically divided into two distinct subsets: a training set and a testing set. This division is crucial for evaluating the performance and generalization of machine learning models.

Bifurcation Schema Sampling Technique: Stratified Sampling Ratio: 80% of the data in the training set and 20% in the testing set Stratified Sampling Stratified sampling was employed to ensure that each subset (training and testing) maintains the same proportion of classes as the original dataset. This approach is particularly beneficial when dealing with imbalanced datasets or when preserving class representation is essential for model training and evaluation.

**The 80-20 split ratio was chosen to allocate a significant portion of the data to the training set (80%), allowing models to learn from a substantial amount of information while retaining a separate testing set (20%) for unbiased evaluation and validation.**

### 3.2. Data Analysis

3.2.1.1. PO1 | PS1:: Supervised Machine Learning Classification Algorithm: Decision Tree (Base Model)

3.2.1.2. PO1 | PS1:: Supervised Machine Learning Classification Algorithms: {Logistic Regression | Support Vector Machine | K Nearest Neighbour} (Comparison Models)

## *Defining the library*

In [2]:
```python
# Required Libraries
import pandas as pd, numpy as np # For Data Manipulation
from sklearn.preprocessing import LabelEncoder, OrdinalEncoder # For Encodi
ng Categorical Data [Nominal | Ordinal]
from sklearn.preprocessing import OneHotEncoder # For Creating Dummy Variab
les of Categorical Data [Nominal]
from sklearn.impute import SimpleImputer, KNNImputer # For Imputation of Mi
ssing Data
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScale
r # For Rescaling Data
from sklearn.model_selection import train_test_split # For Splitting Data i
nto Training & Testing Sets
import matplotlib.pyplot as plt
import numpy as np
from scipy.stats import pearsonr
from scipy import stats

# Required Libraries
import pandas as pd, numpy as np # For Data Manipulation
import matplotlib.pyplot as plt, seaborn as sns # For Data Visualization
import scipy.cluster.hierarchy as sch # For Hierarchical Clustering
from sklearn.cluster import AgglomerativeClustering as agclus, KMeans as km
clus # For Agglomerative & K-Means Clustering
from sklearn.metrics import silhouette_score as sscore, davies_bouldin_scor
e as dbscore # For Clustering Model Evaluation

# @title load library { display-mode: "form" }
# Load IPython extension for measuring time
!pip install ipython-autotime
%reload_ext autotime

# Load IPython extension for memory profiling
!pip install memory-profiler
%reload_ext memory_profiler

# Your imports
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.cluster.hierarchy as sch
from sklearn.cluster import AgglomerativeClustering as agclus, KMeans as km
clus
from sklearn.metrics import silhouette_score as sscore, davies_bouldin_scor
e as dbscore
from scipy.cluster.hierarchy import dendrogram, linkage
import plotly.graph_objects as go

# Load preprocessing libraries
from sklearn.preprocessing import LabelEncoder, OrdinalEncoder, OneHotEncod
er
from sklearn.impute import SimpleImputer, KNNImputer
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScale
r
from sklearn.model_selection import train_test_split

from scipy.stats import f_oneway


# Import
```

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, StratifiedShuffleSplit
from sklearn.tree import DecisionTreeClassifier, export_text, plot_tree  # For Decision Tree Model
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.metrics import confusion_matrix, classification_report # For Decision Tree Model Evaluation
from sklearn.neighbors import KNeighborsClassifier
from sklearn.decomposition import PCA
from matplotlib.colors import ListedColormap
from sklearn.svm import SVC
from sklearn.metrics import confusion_matrix, accuracy_score
from matplotlib.colors import ListedColormap


import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, StratifiedShuffleSplit
from sklearn.tree import DecisionTreeClassifier, export_text, plot_tree  # For Decision Tree Model
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.metrics import confusion_matrix, classification_report # For Decision Tree Model Evaluation
from sklearn.neighbors import KNeighborsClassifier
from sklearn.decomposition import PCA
from matplotlib.colors import ListedColormap
from sklearn.svm import SVC
from sklearn.metrics import confusion_matrix, accuracy_score
from matplotlib.colors import ListedColormap


# Load preprocessing libraries
from sklearn.preprocessing import LabelEncoder, OrdinalEncoder, OneHotEncoder
from sklearn.impute import SimpleImputer, KNNImputer
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler
from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedShuffleSplit

!pip install scikit-learn xgboost

## Data Visualization Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.graph_objects as go
from wordcloud import WordCloud
from collections import Counter
from scipy import stats
from sklearn.tree import plot_tree
import graphviz
from IPython.display import display
from collections import Counter
```

```python
## Data Preprocessing Libraries
from sklearn.preprocessing import OrdinalEncoder
from sklearn.impute import SimpleImputer, KNNImputer
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.model_selection import train_test_split, StratifiedShuffleSplit
from sklearn.metrics import f1_score
from sklearn.tree import export_text

## Machine Learning Models and Evaluation Metrics
import xgboost as xgb
from sklearn.ensemble import RandomForestClassifier
from sklearn.utils.validation import column_or_1d
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, f1_score, precision_recall_fscore_support
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression, Lasso, Ridge
from sklearn.metrics import make_scorer
from sklearn.pipeline import make_pipeline
from sklearn.tree import export_graphviz
```

```
Collecting ipython-autotime
  Downloading ipython_autotime-0.3.2-py2.py3-none-any.whl (7.0 kB)
Requirement already satisfied: ipython in /usr/local/lib/python3.10/dist-pa
ckages (from ipython-autotime) (7.34.0)
Requirement already satisfied: setuptools>=18.5 in /usr/local/lib/python3.1
0/dist-packages (from ipython->ipython-autotime) (67.7.2)
Collecting jedi>=0.16 (from ipython->ipython-autotime)
  Downloading jedi-0.19.1-py2.py3-none-any.whl (1.6 MB)
     ──────────────────────────────────────── 1.6/1.6 MB 6.9 MB/s eta 0:00:
00
Requirement already satisfied: decorator in /usr/local/lib/python3.10/dist-
packages (from ipython->ipython-autotime) (4.4.2)
Requirement already satisfied: pickleshare in /usr/local/lib/python3.10/dis
t-packages (from ipython->ipython-autotime) (0.7.5)
Requirement already satisfied: traitlets>=4.2 in /usr/local/lib/python3.10/
dist-packages (from ipython->ipython-autotime) (5.7.1)
Requirement already satisfied: prompt-toolkit!=3.0.0,!=3.0.1,<3.1.0,>=2.0.0
in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime)
(3.0.43)
Requirement already satisfied: pygments in /usr/local/lib/python3.10/dist-p
ackages (from ipython->ipython-autotime) (2.16.1)
Requirement already satisfied: backcall in /usr/local/lib/python3.10/dist-p
ackages (from ipython->ipython-autotime) (0.2.0)
Requirement already satisfied: matplotlib-inline in /usr/local/lib/python3.
10/dist-packages (from ipython->ipython-autotime) (0.1.6)
Requirement already satisfied: pexpect>4.3 in /usr/local/lib/python3.10/dis
t-packages (from ipython->ipython-autotime) (4.9.0)
Requirement already satisfied: parso<0.9.0,>=0.8.3 in /usr/local/lib/python
3.10/dist-packages (from jedi>=0.16->ipython->ipython-autotime) (0.8.4)
Requirement already satisfied: ptyprocess>=0.5 in /usr/local/lib/python3.1
0/dist-packages (from pexpect>4.3->ipython->ipython-autotime) (0.7.0)
Requirement already satisfied: wcwidth in /usr/local/lib/python3.10/dist-pa
ckages (from prompt-toolkit!=3.0.0,!=3.0.1,<3.1.0,>=2.0.0->ipython->ipython
-autotime) (0.2.13)
Installing collected packages: jedi, ipython-autotime
Successfully installed ipython-autotime-0.3.2 jedi-0.19.1
Collecting memory-profiler
  Downloading memory_profiler-0.61.0-py3-none-any.whl (31 kB)
Requirement already satisfied: psutil in /usr/local/lib/python3.10/dist-pac
kages (from memory-profiler) (5.9.5)
Installing collected packages: memory-profiler
Successfully installed memory-profiler-0.61.0
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/di
st-packages (1.2.2)
Requirement already satisfied: xgboost in /usr/local/lib/python3.10/dist-pa
ckages (2.0.3)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.10/d
ist-packages (from scikit-learn) (1.25.2)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.10/di
st-packages (from scikit-learn) (1.11.4)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/d
ist-packages (from scikit-learn) (1.4.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/pytho
n3.10/dist-packages (from scikit-learn) (3.4.0)
time: 15.8 s (started: 2024-04-13 11:17:31 +00:00)
```

# Uploading of Dataset

```python
import pandas as pd
import gdown

# Google Drive file ID
file_id = '13fp1-YgAuSiR_bWZwetJiEcJZOeDeAtu'

# Downloading the CSV file from Google Drive
url = f'https://drive.google.com/uc?id={file_id}'
csv_file_path = 'Washington_State_HDMA_Dataset with Cluster Label'
gdown.download(url, csv_file_path, quiet=False)

# Read the CSV file into a pandas DataFrame
df = pd.read_csv(csv_file_path)

# Display the first few rows of the DataFrame to verify the data
print(df.head())
```

```
Downloading...
From: https://drive.google.com/uc?id=13fp1-YgAuSiR_bWZwetJiEcJZOeDeAtu
To: /content/Washington_State_HDMA_Dataset with Cluster Label
100%|██████████| 32.6M/32.6M [00:00<00:00, 98.5MB/s]
```

```
     S.no  tract_to_msamd_income   rate_spread   population   minority_populatio
n  \
0     1             121.690002          NaN        8381.0                23.79000
1
1     2              83.370003          NaN        4915.0                23.99000
0
2     3              91.129997          NaN        5075.0                11.82000
0
3     4             146.169998          NaN        5032.0                 8.59000
0
4     5             162.470001          NaN        5183.0                10.50000
0
```

```
     number_of_owner_occupied_units   number_of_1_to_4_family_units  \
0                            2175.0                           2660.0
1                            1268.0                           1777.0
2                            1136.0                           1838.0
3                            1525.0                           1820.0
4                            1705.0                           2104.0
```

```
     loan_amount_000s   hud_median_family_income   applicant_income_000s   ...
\
0                227                    73300.0                    116.0   ...
1                240                    57900.0                     42.0   ...
2                241                    73300.0                    117.0   ...
3                351                    73300.0                    315.0   ...
4                417                    78100.0                    114.0   ...
```

```
                              co_applicant_ethnicity_name  census_tract_number  \
0                                 Not Hispanic or Latino               413.27
1                                        No co-applicant              9208.01
2                                 Not Hispanic or Latino               414.00
3     Information not provided by applicant in mail,...               405.10
4                                 Not Hispanic or Latino               907.00
```

```
     as_of_year  application_date_indicator  applicant_sex_name  \
0        2016                            0              Female
1        2016                            0                Male
2        2016                            0                Male
3        2016                            0                Male
4        2016                            0              Female
```

```
                           applicant_ethnicity_name  \
0                             Not Hispanic or Latino
1                                 Hispanic or Latino
2                             Not Hispanic or Latino
3     Information not provided by applicant in mail,...
4                             Not Hispanic or Latino
```

```
                                         agency_name  agency_abbr  action_taken_nam
e  \
0          Consumer Financial Protection Bureau              CFPB    Loan originate
d
1     Department of Housing and Urban Development               HUD    Loan originate
d
2     Department of Housing and Urban Development               HUD    Loan originate
d
3          National Credit Union Administration              NCUA    Loan originate
d
4          Federal Deposit Insurance Corporation              FDIC    Loan originate
d
```

```
     Cluster_Label
0                 4
1                 3
2                 4
3                 4
4                 4

[5 rows x 39 columns]
time: 3.81 s (started: 2024-04-13 11:17:46 +00:00)

<ipython-input-3-f0f2a90db5ca>:13: DtypeWarning: Columns (24,25,26) have mi
xed types. Specify dtype option on import or set low_memory=False.
  df = pd.read_csv(csv_file_path)
```

In [4]:
```python
df.info()
list(df.columns)

# Assuming df is your original DataFrame
# Add your normalization or standardization code here

# Display summary statistics
df.describe()

total_records = len(df)
print(f"Total number of records: {total_records}")

# Calculate the total number of filled cells in each column
filled_cells_count = df.count()

# Sum up the counts to get the total number of filled cells in the DataFram
e
total_filled_cells = filled_cells_count.sum()

print(f"Total number of filled cells: {total_filled_cells}")

# Assuming df is your DataFrame
unique_counts = df.nunique()

# Display the number of unique values in each column
print(unique_counts)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 60000 entries, 0 to 59999
Data columns (total 39 columns):
 #   Column                        Non-Null Count   Dtype
---  ------                        --------------   -----
 0   S.no                          60000 non-null   int64
 1   tract_to_msamd_income         59878 non-null   float64
 2   rate_spread                   1866 non-null    float64
 3   population                    59878 non-null   float64
 4   minority_population           59878 non-null   float64
 5   number_of_owner_occupied_units 59876 non-null  float64
 6   number_of_1_to_4_family_units 59878 non-null   float64
 7   loan_amount_000s              60000 non-null   int64
 8   hud_median_family_income      59878 non-null   float64
 9   applicant_income_000s         53630 non-null   float64
 10  state_name                    60000 non-null   object
 11  state_abbr                    60000 non-null   object
 12  sequence_number               60000 non-null   int64
 13  respondent_id                 60000 non-null   object
 14  purchaser_type_name           60000 non-null   object
 15  property_type_name            60000 non-null   object
 16  preapproval_name              60000 non-null   object
 17  owner_occupancy_name          60000 non-null   object
 18  msamd_name                    51619 non-null   object
 19  loan_type_name                60000 non-null   object
 20  loan_purpose_name             60000 non-null   object
 21  lien_status_name              60000 non-null   object
 22  hoepa_status_name             60000 non-null   object
 23  edit_status_name              12451 non-null   object
 24  denial_reason_name_3          1 non-null       object
 25  denial_reason_name_2          43 non-null      object
 26  denial_reason_name_1          118 non-null     object
 27  county_name                   59908 non-null   object
 28  co_applicant_sex_name         60000 non-null   object
 29  co_applicant_ethnicity_name   60000 non-null   object
 30  census_tract_number           59878 non-null   float64
 31  as_of_year                    60000 non-null   int64
 32  application_date_indicator    60000 non-null   int64
 33  applicant_sex_name            60000 non-null   object
 34  applicant_ethnicity_name      60000 non-null   object
 35  agency_name                   60000 non-null   object
 36  agency_abbr                   60000 non-null   object
 37  action_taken_name             60000 non-null   object
 38  Cluster_Label                 60000 non-null   int64
dtypes: float64(9), int64(6), object(24)
memory usage: 17.9+ MB
Total number of records: 60000
Total number of filled cells: 2038780
S.no                            60000
tract_to_msamd_income           1327
rate_spread                      320
population                      1283
minority_population             1216
number_of_owner_occupied_units   996
number_of_1_to_4_family_units   1051
loan_amount_000s                1344
hud_median_family_income          15
applicant_income_000s            807
state_name                         1
state_abbr                         1
sequence_number                39340
```

```
respondent_id                      593
purchaser_type_name                 10
property_type_name                   3
preapproval_name                     3
owner_occupancy_name                 3
msamd_name                          14
loan_type_name                       4
loan_purpose_name                    3
lien_status_name                     4
hoepa_status_name                    2
edit_status_name                     1
denial_reason_name_3                 1
denial_reason_name_2                 8
denial_reason_name_1                 8
county_name                         39
co_applicant_sex_name                5
co_applicant_ethnicity_name          5
census_tract_number               1108
as_of_year                           1
application_date_indicator           2
applicant_sex_name                   4
applicant_ethnicity_name             4
agency_name                          6
agency_abbr                          6
action_taken_name                    8
Cluster_Label                        5
dtype: int64
time: 509 ms (started: 2024-04-13 11:17:50 +00:00)
```

In [5]:

```python
# Importing necessary libraries
import pandas as pd

# Assuming df is your DataFrame and 'Cluster_Label' is the column of interest

# Calculate relative frequencies
relative_frequencies = df['Cluster_Label'].value_counts(normalize=True)

# Print relative frequencies
print(relative_frequencies)

# Assuming df is your DataFrame and 'Cluster_Label' is the column of interest

# Count the occurrences of each unique value
unique_value_counts = df['Cluster_Label'].value_counts()

# Print the count of unique values
print(unique_value_counts)
```

```
Cluster_Label
3    0.236033
1    0.234150
0    0.190583
4    0.179050
2    0.160183
Name: proportion, dtype: float64
Cluster_Label
3    14162
1    14049
0    11435
4    10743
2     9611
Name: count, dtype: int64
time: 8.13 ms (started: 2024-04-13 11:17:51 +00:00)
```

```python
# Assuming df is your DataFrame
columns_list = df.columns.tolist()
columns_list

list(df.columns)
```

Out[6]:
```
['S.no',
 'tract_to_msamd_income',
 'rate_spread',
 'population',
 'minority_population',
 'number_of_owner_occupied_units',
 'number_of_1_to_4_family_units',
 'loan_amount_000s',
 'hud_median_family_income',
 'applicant_income_000s',
 'state_name',
 'state_abbr',
 'sequence_number',
 'respondent_id',
 'purchaser_type_name',
 'property_type_name',
 'preapproval_name',
 'owner_occupancy_name',
 'msamd_name',
 'loan_type_name',
 'loan_purpose_name',
 'lien_status_name',
 'hoepa_status_name',
 'edit_status_name',
 'denial_reason_name_3',
 'denial_reason_name_2',
 'denial_reason_name_1',
 'county_name',
 'co_applicant_sex_name',
 'co_applicant_ethnicity_name',
 'census_tract_number',
 'as_of_year',
 'application_date_indicator',
 'applicant_sex_name',
 'applicant_ethnicity_name',
 'agency_name',
 'agency_abbr',
 'action_taken_name',
 'Cluster_Label']
```

time: 4.18 ms (started: 2024-04-13 11:17:51 +00:00)

In [7]:
```python
# Nominal and Ordinal Columns

# Continuous and Non Continuous Columns


import pandas as pd

# Assuming df is your DataFrame
continuous_columns = df.select_dtypes(include=['float64', 'int64']).columns
non_continuous_columns = df.select_dtypes(exclude=['float64', 'int64']).col
umns

print("Continuous Columns:", list(continuous_columns))
print("Non-Continuous Columns:", list(non_continuous_columns))


# Assuming df is your DataFrame
categorical_columns = df.select_dtypes(include=['object', 'category']).colu
mns
non_categorical_columns = df.select_dtypes(exclude=['object', 'category']).
columns

print("Categorical Columns:", list(categorical_columns))
print("Non-Categorical Columns:", list(non_categorical_columns))
```

```
Continuous Columns: ['S.no', 'tract_to_msamd_income', 'rate_spread', 'popul
ation', 'minority_population', 'number_of_owner_occupied_units', 'number_of
_1_to_4_family_units', 'loan_amount_000s', 'hud_median_family_income', 'app
licant_income_000s', 'sequence_number', 'census_tract_number', 'as_of_yea
r', 'application_date_indicator', 'Cluster_Label']
Non-Continuous Columns: ['state_name', 'state_abbr', 'respondent_id', 'purc
haser_type_name', 'property_type_name', 'preapproval_name', 'owner_occupanc
y_name', 'msamd_name', 'loan_type_name', 'loan_purpose_name', 'lien_status_
name', 'hoepa_status_name', 'edit_status_name', 'denial_reason_name_3', 'de
nial_reason_name_2', 'denial_reason_name_1', 'county_name', 'co_applicant_s
ex_name', 'co_applicant_ethnicity_name', 'applicant_sex_name', 'applicant_e
thnicity_name', 'agency_name', 'agency_abbr', 'action_taken_name']
Categorical Columns: ['state_name', 'state_abbr', 'respondent_id', 'purchas
er_type_name', 'property_type_name', 'preapproval_name', 'owner_occupancy_n
ame', 'msamd_name', 'loan_type_name', 'loan_purpose_name', 'lien_status_nam
e', 'hoepa_status_name', 'edit_status_name', 'denial_reason_name_3', 'denia
l_reason_name_2', 'denial_reason_name_1', 'county_name', 'co_applicant_sex_
name', 'co_applicant_ethnicity_name', 'applicant_sex_name', 'applicant_ethn
icity_name', 'agency_name', 'agency_abbr', 'action_taken_name']
Non-Categorical Columns: ['S.no', 'tract_to_msamd_income', 'rate_spread',
'population', 'minority_population', 'number_of_owner_occupied_units', 'num
ber_of_1_to_4_family_units', 'loan_amount_000s', 'hud_median_family_incom
e', 'applicant_income_000s', 'sequence_number', 'census_tract_number', 'as_
of_year', 'application_date_indicator', 'Cluster_Label']
time: 37 ms (started: 2024-04-13 11:17:51 +00:00)
```

In [8]:
```python
###  Missing Data Statistics and Treatment
###  Missing Data Statistics: Records

# Assuming df is your DataFrame

# Count the missing values in each column
missing_data = df.isnull().sum()

# Create a DataFrame to display missing data statistics
missing_data_stats = pd.DataFrame({
    'Column': missing_data.index,
    'Missing Records': missing_data.values,
    'Percentage Missing': (missing_data / len(df)) * 100
})

# Sort the DataFrame by the percentage of missing values in descending orde
r
missing_data_stats = missing_data_stats.sort_values(by='Percentage Missin
g', ascending=False)

# Print the missing data statistics
print(missing_data_stats)
```

```
                                                            Column  \
denial_reason_name_3                        denial_reason_name_3
denial_reason_name_2                        denial_reason_name_2
denial_reason_name_1                        denial_reason_name_1
rate_spread                                          rate_spread
edit_status_name                                edit_status_name
msamd_name                                            msamd_name
applicant_income_000s                      applicant_income_000s
number_of_owner_occupied_units    number_of_owner_occupied_units
tract_to_msamd_income                      tract_to_msamd_income
hud_median_family_income                hud_median_family_income
number_of_1_to_4_family_units      number_of_1_to_4_family_units
minority_population                        minority_population
population                                            population
census_tract_number                        census_tract_number
county_name                                          county_name
co_applicant_ethnicity_name        co_applicant_ethnicity_name
co_applicant_sex_name                      co_applicant_sex_name
S.no                                                        S.no
as_of_year                                            as_of_year
application_date_indicator        application_date_indicator
applicant_sex_name                        applicant_sex_name
agency_name                                          agency_name
agency_abbr                                          agency_abbr
action_taken_name                            action_taken_name
applicant_ethnicity_name            applicant_ethnicity_name
loan_type_name                                    loan_type_name
hoepa_status_name                            hoepa_status_name
lien_status_name                              lien_status_name
loan_purpose_name                            loan_purpose_name
owner_occupancy_name                        owner_occupancy_name
preapproval_name                              preapproval_name
property_type_name                          property_type_name
purchaser_type_name                        purchaser_type_name
respondent_id                                      respondent_id
sequence_number                                  sequence_number
state_abbr                                            state_abbr
state_name                                            state_name
loan_amount_000s                              loan_amount_000s
Cluster_Label                                      Cluster_Label

                                Missing Records   Percentage Missing
denial_reason_name_3                       59999            99.998333
denial_reason_name_2                       59957            99.928333
denial_reason_name_1                       59882            99.803333
rate_spread                                58134            96.890000
edit_status_name                           47549            79.248333
msamd_name                                  8381            13.968333
applicant_income_000s                       6370            10.616667
number_of_owner_occupied_units               124             0.206667
tract_to_msamd_income                        122             0.203333
hud_median_family_income                     122             0.203333
number_of_1_to_4_family_units                122             0.203333
minority_population                          122             0.203333
population                                   122             0.203333
census_tract_number                          122             0.203333
county_name                                   92             0.153333
co_applicant_ethnicity_name                    0             0.000000
co_applicant_sex_name                          0             0.000000
S.no                                           0             0.000000
as_of_year                                     0             0.000000
```

```
application_date_indicator              0              0.000000
applicant_sex_name                      0              0.000000
agency_name                             0              0.000000
agency_abbr                             0              0.000000
action_taken_name                       0              0.000000
applicant_ethnicity_name                0              0.000000
loan_type_name                          0              0.000000
hoepa_status_name                       0              0.000000
lien_status_name                        0              0.000000
loan_purpose_name                       0              0.000000
owner_occupancy_name                    0              0.000000
preapproval_name                        0              0.000000
property_type_name                      0              0.000000
purchaser_type_name                     0              0.000000
respondent_id                           0              0.000000
sequence_number                         0              0.000000
state_abbr                              0              0.000000
state_name                              0              0.000000
loan_amount_000s                        0              0.000000
Cluster_Label                           0              0.000000
time: 102 ms (started: 2024-04-13 11:17:51 +00:00)
```

In [9]:
```python
# List of columns to drop
columns_to_drop = ['state_name', 'state_abbr','denial_reason_name_3', 'deni
al_reason_name_2', 'denial_reason_name_1', 'rate_spread', 'edit_status_nam
e', 'as_of_year']

# Drop columns with more than 50% missing values
df_cleaned = df.drop(columns=columns_to_drop)

# Print the cleaned DataFrame
df1 = df_cleaned

# Count the missing values in each column
missing_data = df1.isnull().sum()

# Create a DataFrame to display missing data statistics
missing_data_stats = pd.DataFrame({
    'Column': missing_data.index,
    'Missing Records': missing_data.values,
    'Percentage Missing': (missing_data / len(df)) * 100
})
# Sort the DataFrame by the percentage of missing values in descending orde
r
missing_data_stats = missing_data_stats.sort_values(by='Percentage Missin
g', ascending=False)

# Print the missing data statistics
print(missing_data_stats)
```

|  | Column \ |
|---|---|
| msamd_name | msamd_name |
| applicant_income_000s | applicant_income_000s |
| number_of_owner_occupied_units | number_of_owner_occupied_units |
| census_tract_number | census_tract_number |
| population | population |
| minority_population | minority_population |
| number_of_1_to_4_family_units | number_of_1_to_4_family_units |
| tract_to_msamd_income | tract_to_msamd_income |
| hud_median_family_income | hud_median_family_income |
| county_name | county_name |
| agency_name | agency_name |
| applicant_ethnicity_name | applicant_ethnicity_name |
| applicant_sex_name | applicant_sex_name |
| application_date_indicator | application_date_indicator |
| hoepa_status_name | hoepa_status_name |
| action_taken_name | action_taken_name |
| co_applicant_ethnicity_name | co_applicant_ethnicity_name |
| co_applicant_sex_name | co_applicant_sex_name |
| agency_abbr | agency_abbr |
| S.no | S.no |
| lien_status_name | lien_status_name |
| loan_purpose_name | loan_purpose_name |
| loan_type_name | loan_type_name |
| owner_occupancy_name | owner_occupancy_name |
| preapproval_name | preapproval_name |
| property_type_name | property_type_name |
| purchaser_type_name | purchaser_type_name |
| respondent_id | respondent_id |
| sequence_number | sequence_number |
| loan_amount_000s | loan_amount_000s |
| Cluster_Label | Cluster_Label |

|  | Missing Records | Percentage Missing |
|---|---|---|
| msamd_name | 8381 | 13.968333 |
| applicant_income_000s | 6370 | 10.616667 |
| number_of_owner_occupied_units | 124 | 0.206667 |
| census_tract_number | 122 | 0.203333 |
| population | 122 | 0.203333 |
| minority_population | 122 | 0.203333 |
| number_of_1_to_4_family_units | 122 | 0.203333 |
| tract_to_msamd_income | 122 | 0.203333 |
| hud_median_family_income | 122 | 0.203333 |
| county_name | 92 | 0.153333 |
| agency_name | 0 | 0.000000 |
| applicant_ethnicity_name | 0 | 0.000000 |
| applicant_sex_name | 0 | 0.000000 |
| application_date_indicator | 0 | 0.000000 |
| hoepa_status_name | 0 | 0.000000 |
| action_taken_name | 0 | 0.000000 |
| co_applicant_ethnicity_name | 0 | 0.000000 |
| co_applicant_sex_name | 0 | 0.000000 |
| agency_abbr | 0 | 0.000000 |
| S.no | 0 | 0.000000 |
| lien_status_name | 0 | 0.000000 |
| loan_purpose_name | 0 | 0.000000 |
| loan_type_name | 0 | 0.000000 |
| owner_occupancy_name | 0 | 0.000000 |
| preapproval_name | 0 | 0.000000 |
| property_type_name | 0 | 0.000000 |
| purchaser_type_name | 0 | 0.000000 |

```
respondent_id                          0           0.000000
sequence_number                        0           0.000000
loan_amount_000s                       0           0.000000
Cluster_Label                          0           0.000000
time: 71 ms (started: 2024-04-13 11:17:51 +00:00)
```

In [10]:
```python
### Missing Records (ROWS)

# Count the missing values in each row
missing_rows = df1.isnull().sum(axis=1)

# Count the number of rows with at least one missing value
num_rows_with_missing = len(missing_rows[missing_rows > 0])

# Print the number of rows with missing data
print("Number of rows with missing data:", num_rows_with_missing)


# Calculate the percentage of missing values in each row
missing_percentage_rows = (df1.isnull().sum(axis=1) / len(df1.columns)) * 1
00

# Count the number of rows with more than 50% missing data
num_rows_more_than_50_percent_missing = len(missing_percentage_rows[missing
_percentage_rows > 50])

# Print the number of rows with more than 50% missing data
print("Number of rows with more than 50% missing data:", num_rows_more_than
_50_percent_missing)
```

```
Number of rows with missing data: 13629
Number of rows with more than 50% missing data: 0
time: 131 ms (started: 2024-04-13 11:17:51 +00:00)
```

```
In [11]: # DIVIDING DF1 into Cat and Non Cat

         # Assuming df1 is your DataFrame
         cat_columns = df1.select_dtypes(include=['object']).columns
         noncat_columns = df1.select_dtypes(exclude=['object']).columns

         # Creating categorical and non-categorical DataFrames
         catdf1 = df1[cat_columns]
         noncatdf1 = df1[noncat_columns]

         #print(list(catdf.columns))
         #print(list(noncatdf.columns))
         print(list(catdf1.columns))
         print(list(noncatdf1.columns))

         #20
         #list(noncatdf.columns)
```

```
['respondent_id', 'purchaser_type_name', 'property_type_name', 'preapproval
_name', 'owner_occupancy_name', 'msamd_name', 'loan_type_name', 'loan_purpo
se_name', 'lien_status_name', 'hoepa_status_name', 'county_name', 'co_appli
cant_sex_name', 'co_applicant_ethnicity_name', 'applicant_sex_name', 'appli
cant_ethnicity_name', 'agency_name', 'agency_abbr', 'action_taken_name']
['S.no', 'tract_to_msamd_income', 'population', 'minority_population', 'num
ber_of_owner_occupied_units', 'number_of_1_to_4_family_units', 'loan_amount
_000s', 'hud_median_family_income', 'applicant_income_000s', 'sequence_numb
er', 'census_tract_number', 'application_date_indicator', 'Cluster_Label']
time: 28.7 ms (started: 2024-04-13 11:17:51 +00:00)
```

# PreProcessing of Data

```
In [12]: # Data Bifurcation
         catdf = df[['S.no', 'respondent_id', 'purchaser_type_name', 'property_type_
         name', 'preapproval_name', 'owner_occupancy_name', 'msamd_name', 'loan_type
         _name', 'loan_purpose_name',
                     'lien_status_name', 'hoepa_status_name', 'county_name', 'co_app
         licant_sex_name', 'co_applicant_ethnicity_name', 'applicant_sex_name', 'app
         licant_ethnicity_name', 'agency_name',
                     'agency_abbr', 'action_taken_name']] # Categorical Data [Nomina
         l | Ordinal]

         noncatdf = df[['S.no', 'tract_to_msamd_income', 'population', 'minority_pop
         ulation', 'number_of_owner_occupied_units', 'number_of_1_to_4_family_unit
         s', 'loan_amount_000s',
                        'hud_median_family_income', 'applicant_income_000s', 'sequen
         ce_number', 'census_tract_number', 'application_date_indicator', 'Cluster_L
         abel']] # Non-Categorical Data
```

```
time: 25.6 ms (started: 2024-04-13 11:17:51 +00:00)
```

In [13]:
```python
####  STATISTICS OF CAT DATASET

# Count and frequency statistics for each column in catdf
catdf_stats = pd.DataFrame()

for column in catdf.columns:
    col_count = catdf[column].value_counts().reset_index()
    col_count.columns = [column, 'Frequency']
    catdf_stats = pd.concat([catdf_stats, col_count], axis=1)

# Display the count and frequency statistics
#print(catdf_stats)

# Summary for each column in catdf
catdf_summary = catdf.describe(include='all').transpose()

# Display the summary
print(catdf_summary)


# Calculate the proportion (relative frequency) for each categorical column
#proportion_stats = catdf.apply(lambda x: x.value_counts(normalize=True).id
xmax() + ': ' + "{:.2%}".format(x.value_counts(normalize=True).max()))

# Display the proportion statistics
#print(proportion_stats)
```

```
                               count  unique  \
S.no                          60000.0    NaN
respondent_id                  60000    593
purchaser_type_name            60000     10
property_type_name             60000      3
preapproval_name               60000      3
owner_occupancy_name           60000      3
msamd_name                     51619     14
loan_type_name                 60000      4
loan_purpose_name              60000      3
lien_status_name               60000      4
hoepa_status_name              60000      2
county_name                    59908     39
co_applicant_sex_name          60000      5
co_applicant_ethnicity_name    60000      5
applicant_sex_name             60000      4
applicant_ethnicity_name       60000      4
agency_name                    60000      6
agency_abbr                    60000      6
action_taken_name              60000      8
```

```
top  \
S.no
NaN
respondent_id                                                      32
489
purchaser_type_name          Loan was not originated or was not sold in ca
l...
property_type_name           One-to-four family dwelling (other than manuf
a...
preapproval_name                                           Not applica
ble
owner_occupancy_name                       Owner-occupied as a principal dwell
ing
msamd_name                                     Seattle, Bellevue, Everett -
WA
loan_type_name                                                 Conventio
nal
loan_purpose_name                                              Refinanc
ing
lien_status_name                                     Secured by a first l
ien
hoepa_status_name                                          Not a HOEPA l
oan
county_name                                                    King Cou
nty
co_applicant_sex_name                                      No co-applic
ant
co_applicant_ethnicity_name                                No co-applic
ant
applicant_sex_name                                                    M
ale
applicant_ethnicity_name                            Not Hispanic or Lat
ino
agency_name                  Department of Housing and Urban Developm
ent
agency_abbr
HUD
action_taken_name                                          Loan origina
ted
```

|  | freq | mean | std | min | 25% \ |
| --- | --- | --- | --- | --- | --- |
| S.no | NaN | 30000.5 | 17320.652413 | 1.0 | 15000.75 |
| respondent_id | 5006 | NaN | NaN | NaN | NaN |
| purchaser_type_name | 16112 | NaN | NaN | NaN | NaN |
| property_type_name | 57630 | NaN | NaN | NaN | NaN |
| preapproval_name | 47832 | NaN | NaN | NaN | NaN |
| owner_occupancy_name | 53940 | NaN | NaN | NaN | NaN |
| msamd_name | 17965 | NaN | NaN | NaN | NaN |
| loan_type_name | 42917 | NaN | NaN | NaN | NaN |
| loan_purpose_name | 28576 | NaN | NaN | NaN | NaN |
| lien_status_name | 57046 | NaN | NaN | NaN | NaN |
| hoepa_status_name | 59996 | NaN | NaN | NaN | NaN |
| county_name | 12915 | NaN | NaN | NaN | NaN |
| co_applicant_sex_name | 26987 | NaN | NaN | NaN | NaN |
| co_applicant_ethnicity_name | 26987 | NaN | NaN | NaN | NaN |
| applicant_sex_name | 37070 | NaN | NaN | NaN | NaN |
| applicant_ethnicity_name | 44014 | NaN | NaN | NaN | NaN |
| agency_name | 27514 | NaN | NaN | NaN | NaN |
| agency_abbr | 27514 | NaN | NaN | NaN | NaN |
| action_taken_name | 55815 | NaN | NaN | NaN | NaN |

|  | 50% | 75% | max |
| --- | --- | --- | --- |
| S.no | 30000.5 | 45000.25 | 60000.0 |
| respondent_id | NaN | NaN | NaN |
| purchaser_type_name | NaN | NaN | NaN |
| property_type_name | NaN | NaN | NaN |
| preapproval_name | NaN | NaN | NaN |
| owner_occupancy_name | NaN | NaN | NaN |
| msamd_name | NaN | NaN | NaN |
| loan_type_name | NaN | NaN | NaN |
| loan_purpose_name | NaN | NaN | NaN |
| lien_status_name | NaN | NaN | NaN |
| hoepa_status_name | NaN | NaN | NaN |
| county_name | NaN | NaN | NaN |
| co_applicant_sex_name | NaN | NaN | NaN |
| co_applicant_ethnicity_name | NaN | NaN | NaN |
| applicant_sex_name | NaN | NaN | NaN |
| applicant_ethnicity_name | NaN | NaN | NaN |
| agency_name | NaN | NaN | NaN |
| agency_abbr | NaN | NaN | NaN |
| action_taken_name | NaN | NaN | NaN |

time: 460 ms (started: 2024-04-13 11:17:51 +00:00)

```
In [14]: #### STATISTICS OF NONCAT DATASET

# Display descriptive statistics for non-categorical variables
noncatdf_descriptive_stats = noncatdf.describe()

# Print the descriptive statistics
print(noncatdf_descriptive_stats)
```

|       | S.no | tract_to_msamd_income | population | minority_populati on \ |
|-------|------|----------------------|------------|------------------------|
| count | 60000.000000 | 59878.000000 | 59878.000000 | 59878.0000 00 |
| mean  | 30000.500000 | 107.617351 | 5278.782157 | 23.2444 42 |
| std   | 17320.652413 | 28.233471 | 1716.101490 | 14.4162 09 |
| min   | 1.000000 | 14.050000 | 98.000000 | 2.0400 00 |
| 25%   | 15000.750000 | 88.970001 | 4070.000000 | 12.9500 00 |
| 50%   | 30000.500000 | 105.550003 | 5145.000000 | 19.4200 00 |
| 75%   | 45000.250000 | 123.330002 | 6382.000000 | 29.6800 00 |
| max   | 60000.000000 | 257.140015 | 13025.000000 | 94.7900 01 |

|       | number_of_owner_occupied_units | number_of_1_to_4_family_units \ |
|-------|-------------------------------|--------------------------------|
| count | 59876.000000 | 59878.000000 |
| mean  | 1399.044375 | 1873.281456 |
| std   | 518.330561 | 738.505184 |
| min   | 15.000000 | 27.000000 |
| 25%   | 1034.000000 | 1414.000000 |
| 50%   | 1359.000000 | 1770.000000 |
| 75%   | 1722.000000 | 2249.000000 |
| max   | 2997.000000 | 5893.000000 |

|       | loan_amount_000s | hud_median_family_income | applicant_income_000s \ |
|-------|------------------|--------------------------|-------------------------|
| count | 60000.000000 | 59878.000000 | 53630.000000 |
| mean  | 291.358717 | 73869.411136 | 112.822301 |
| std   | 604.958183 | 12811.243390 | 122.862496 |
| min   | 1.000000 | 48700.000000 | 1.000000 |
| 25%   | 170.000000 | 63100.000000 | 61.000000 |
| 50%   | 242.000000 | 73300.000000 | 89.000000 |
| 75%   | 337.000000 | 90300.000000 | 132.000000 |
| max   | 55000.000000 | 90300.000000 | 6161.000000 |

|       | sequence_number | census_tract_number | application_date_indicator \ |
|-------|-----------------|---------------------|------------------------------|
| count | 6.000000e+04 | 59878.000000 | 60000.000000 |
| mean  | 7.752647e+04 | 1750.597252 | 0.025867 |
| std   | 1.505157e+05 | 3359.676740 | 0.225976 |
| min   | 1.000000e+00 | 1.000000 | 0.000000 |
| 25%   | 3.231500e+03 | 114.020000 | 0.000000 |
| 50%   | 1.648100e+04 | 403.020000 | 0.000000 |
| 75%   | 7.276225e+04 | 713.100000 | 0.000000 |
| max   | 1.241590e+06 | 9757.000000 | 2.000000 |

|       | Cluster_Label |
|-------|---------------|
| count | 60000.000000 |
| mean  | 1.978817 |
| std   | 1.395815 |
| min   | 0.000000 |
| 25%   | 1.000000 |
| 50%   | 2.000000 |
| 75%   | 3.000000 |
| max   | 4.000000 |

time: 70.6 ms (started: 2024-04-13 11:17:52 +00:00)

In [15]:
```python
# Missing Data Statistics: Non-Categorical Variables or Features

# Calculate missing data statistics for non-categorical columns
missing_data_non_categorical = noncatdf.isnull().sum().reset_index()
missing_data_non_categorical.columns = ['Feature', 'Missing_Records']

# Display the missing data statistics
print(missing_data_non_categorical)
```

```
                           Feature  Missing_Records
0                             S.no                0
1              tract_to_msamd_income              122
2                         population              122
3                  minority_population              122
4         number_of_owner_occupied_units          124
5          number_of_1_to_4_family_units          122
6                      loan_amount_000s            0
7               hud_median_family_income          122
8                  applicant_income_000s          6370
9                      sequence_number            0
10                 census_tract_number           122
11             application_date_indicator           0
12                       Cluster_Label            0
time: 8.33 ms (started: 2024-04-13 11:17:52 +00:00)
```

In [16]:
```python
#  Missing Data Treatment: Non-Categorical Variables or Features

# Dataset Used : df_noncat

si_noncat = SimpleImputer(missing_values=np.nan, strategy='mean') # Other S
trategy : mean | median | most_frequent | constant
si_noncat_fit = si_noncat.fit_transform(noncatdf)
imputed_data_non_categorical = pd.DataFrame(si_noncat_fit, columns=noncatd
f.columns); # Missing Non-Categorical Data Imputed Subset using Simple Impu
ter
imputed_data_non_categorical.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 60000 entries, 0 to 59999
Data columns (total 13 columns):
 #   Column                          Non-Null Count  Dtype
---  ------                          --------------  -----
 0   S.no                            60000 non-null  float64
 1   tract_to_msamd_income           60000 non-null  float64
 2   population                      60000 non-null  float64
 3   minority_population             60000 non-null  float64
 4   number_of_owner_occupied_units  60000 non-null  float64
 5   number_of_1_to_4_family_units   60000 non-null  float64
 6   loan_amount_000s                60000 non-null  float64
 7   hud_median_family_income        60000 non-null  float64
 8   applicant_income_000s           60000 non-null  float64
 9   sequence_number                 60000 non-null  float64
 10  census_tract_number             60000 non-null  float64
 11  application_date_indicator      60000 non-null  float64
 12  Cluster_Label                   60000 non-null  float64
dtypes: float64(13)
memory usage: 6.0 MB
time: 36 ms (started: 2024-04-13 11:17:52 +00:00)
```

```python
In [17]:  # Calculate standard deviation for non-categorical columns
          std_deviation_non_categorical = imputed_data_non_categorical.std()

          # Creating a DataFrame to display the results
          dispersion_non_categorical_df = pd.DataFrame({
              'Variable': imputed_data_non_categorical.columns,
              'Standard Deviation': std_deviation_non_categorical.values
          })

          print(dispersion_non_categorical_df)
```

```
                          Variable  Standard Deviation
0                             S.no        17320.652413
1              tract_to_msamd_income          28.204752
2                       population         1714.355870
3                minority_population           14.401545
4      number_of_owner_occupied_units          517.794667
5        number_of_1_to_4_family_units          737.753976
6                   loan_amount_000s          604.958183
7             hud_median_family_income        12798.211781
8                applicant_income_000s          116.157479
9                    sequence_number       150515.678152
10               census_tract_number         3356.259273
11           application_date_indicator            0.225976
12                      Cluster_Label            1.395815
time: 17.6 ms (started: 2024-04-13 11:17:52 +00:00)
```

In [18]:
```python
# Dataset Used : df_cat

si_cat = SimpleImputer(missing_values=np.nan, strategy='most_frequent') # S
trategy = median [When Odd Number of Categories Exists]
si_cat_fit = si_cat.fit_transform(catdf)
imputed_data_categorical = pd.DataFrame(si_cat_fit, columns=catdf.columns);
# Missing Categorical Data Imputed Subset
imputed_data_categorical.info()
imputed_data_categorical.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 60000 entries, 0 to 59999
Data columns (total 19 columns):
 #   Column                     Non-Null Count   Dtype
---  ------                     --------------   -----
 0   S.no                       60000 non-null   object
 1   respondent_id              60000 non-null   object
 2   purchaser_type_name        60000 non-null   object
 3   property_type_name         60000 non-null   object
 4   preapproval_name           60000 non-null   object
 5   owner_occupancy_name       60000 non-null   object
 6   msamd_name                 60000 non-null   object
 7   loan_type_name             60000 non-null   object
 8   loan_purpose_name          60000 non-null   object
 9   lien_status_name           60000 non-null   object
 10  hoepa_status_name          60000 non-null   object
 11  county_name                60000 non-null   object
 12  co_applicant_sex_name      60000 non-null   object
 13  co_applicant_ethnicity_name 60000 non-null  object
 14  applicant_sex_name         60000 non-null   object
 15  applicant_ethnicity_name   60000 non-null   object
 16  agency_name                60000 non-null   object
 17  agency_abbr                60000 non-null   object
 18  action_taken_name          60000 non-null   object
dtypes: object(19)
memory usage: 8.7+ MB
```

Out[18]:

| | S.no | respondent_id | purchaser_type_name | property_type_name | preapproval_name | owner_ |
|---|---|---|---|---|---|---|
| **0** | 1 | 480228 | Freddie Mac (FHLMC) | One-to-four family dwelling (other than manufa... | Not applicable | O |
| **1** | 2 | 7257500009 | Life insurance company, credit union, mortgage... | One-to-four family dwelling (other than manufa... | Not applicable | O |
| **2** | 3 | 72-1545376 | Loan was not originated or was not sold in cal... | One-to-four family dwelling (other than manufa... | Not applicable | O |
| **3** | 4 | 4878 | Loan was not originated or was not sold in cal... | One-to-four family dwelling (other than manufa... | Not applicable | O |
| **4** | 5 | 32489 | Freddie Mac (FHLMC) | One-to-four family dwelling (other than manufa... | Not applicable | O |

```
time: 981 ms (started: 2024-04-13 11:17:52 +00:00)
```

In [19]:
```python
# ENCODING
# Converting Categorical Variable into Numeric

# Calculate the number of unique values in each column
unique_values_categorical = imputed_data_categorical.nunique().reset_index()
unique_values_categorical.columns = ['Feature', 'Number_of_Unique_Values']

# Display the number of unique values
print(unique_values_categorical)

# Initialize LabelEncoder
label_encoder = LabelEncoder()

# Create a copy of the imputed_data_categorical dataframe to avoid modifying the original
encoded_data_categorical = imputed_data_categorical.copy()

# Columns to exclude from encoding
exclude_columns = ["S.no", "Cluster_Label"]

# Iterate through each column in the dataframe
mapping = {}  # To store the mapping of variable names to numeric representation

for column in encoded_data_categorical.columns:
    if column not in exclude_columns:
        # Perform numerical encoding
        encoded_data_categorical[column] = label_encoder.fit_transform(encoded_data_categorical[column])

        # Store the mapping information
        mapping[column] = dict(zip(label_encoder.classes_, label_encoder.transform(label_encoder.classes_)))

# Display the mapping
for variable, variable_mapping in mapping.items():
    print(f"\nMapping for {variable}:")
    print(variable_mapping)

# Display the encoded data
print(encoded_data_categorical)
```

```
                        Feature    Number_of_Unique_Values
0                         S.no                        60000
1                respondent_id                          593
2          purchaser_type_name                           10
3           property_type_name                            3
4             preapproval_name                            3
5          owner_occupancy_name                           3
6                   msamd_name                           14
7               loan_type_name                            4
8            loan_purpose_name                            3
9             lien_status_name                            4
10           hoepa_status_name                            2
11                  county_name                           39
12          co_applicant_sex_name                          5
13    co_applicant_ethnicity_name                          5
14             applicant_sex_name                          4
15       applicant_ethnicity_name                          4
16                   agency_name                           6
17                   agency_abbr                           6
18              action_taken_name                          8
```

Mapping for respondent_id:
{'01-0681100': 0, '01-0726495': 1, '02-0793125': 2, '03-0488052': 3, '04-32
12636': 4, '04-3568208': 5, '04-3660901': 6, '04-7534967': 7, '05-0402708':
8, '06-1016329': 9, '1015560': 10, '10257': 11, '1047': 12, '1097500000': 1
3, '1099800006': 14, '11-3399725': 15, '11-3412303': 16, '11-3714032': 17,
'11162': 18, '112837': 19, '11443': 20, '1146500007': 21, '11729': 22, '117
34': 23, '12135': 24, '1216826': 25, '12219': 26, '1227300009': 27, '1231
1': 28, '1265': 29, '1281': 30, '13-3222578': 31, '13-3602661': 32, '13-375
3941': 33, '13-4225190': 34, '13-4362989': 35, '13-6131491': 36, '13232': 3
7, '13303': 38, '13392': 39, '13778': 40, '14-1841762': 41, '14206': 42, '1
4252': 43, '143662': 44, '1461700004': 45, '14662': 46, '146672': 47, '1474
0': 48, '14843': 49, '151': 50, '15219': 51, '15732': 52, '16-1686740': 53,
'16243': 54, '1635900004': 55, '16450': 56, '16814': 57, '169653': 58, '175
87': 59, '17672': 60, '177957': 61, '17874': 62, '17884': 63, '1842065': 6
4, '19307': 65, '19628': 66, '197478': 67, '19899': 68, '19976': 69, '20-01
42846': 70, '20-0192872': 71, '20-0304793': 72, '20-0640473': 73, '20-07401
51': 74, '20-1255434': 75, '20-1832276': 76, '20-2053401': 77, '20-235529
6': 78, '20-2470783': 79, '20-2471369': 80, '20-2485875': 81, '20-2693054':
82, '20-2718340': 83, '20-2752826': 84, '20-2928975': 85, '20-3702275': 86,
'20-3828708': 87, '20-4136310': 88, '20-4224234': 89, '20-4255880': 90, '20
-4866754': 91, '20-5238443': 92, '20-5239910': 93, '20-5741925': 94, '20-80
06279': 95, '20-8083209': 96, '20-8544905': 97, '20-8745846': 98, '20-88034
49': 99, '20-8921389': 100, '2003500009': 101, '20061': 102, '20068': 103,
'20214': 104, '20516': 105, '20624': 106, '20774': 107, '210434': 108, '211
22': 109, '212465': 110, '2137100009': 111, '2149009991': 112, '21717': 11
3, '2191': 114, '2193616': 115, '22-3039688': 116, '22-3470404': 117, '22-3
554558': 118, '22-3626426': 119, '22-3747694': 120, '22-3887207': 121, '221
34': 122, '22157': 123, '22407': 124, '22444': 125, '22637': 126, '2285': 1
27, '22939': 128, '23-2470039': 129, '23-2769131': 130, '23041': 131, '2317
700005': 132, '23216': 133, '23416': 134, '23521': 135, '23850': 136, '2392
2': 137, '23957': 138, '24077': 139, '24080': 140, '24107': 141, '24169': 1
42, '24224': 143, '24235': 144, '24326': 145, '24382': 146, '24671': 147,
'24708': 148, '24713': 149, '24719': 150, '24753': 151, '24760': 152, '2483
1': 153, '24849': 154, '2489805': 155, '25080': 156, '25093': 157, '25103':
158, '2562164': 159, '2590037': 160, '26-0012825': 161, '26-0021318': 162,
'26-0335190': 163, '26-0360466': 164, '26-0362771': 165, '26-0423240': 166,
'26-0455770': 167, '26-0508430': 168, '26-0595342': 169, '26-0707492': 170,
'26-1242154': 171, '26-1334020': 172, '26-1589507': 173, '26-1773722': 174,
'26-2049351': 175, '26-2261031': 176, '26-2593704': 177, '26-2689428': 178,
'26-2916887': 179, '26-3264687': 180, '26-3416474': 181, '26-3780954': 182,

```
'26-4193875': 183, '26-4461592': 184, '26-4558390': 185, '26-4599244': 186,
'2618780': 187, '26610': 188, '267100004': 189, '27-0222046': 190, '27-0267
182': 191, '27-0684906': 192, '27-0812934': 193, '27-1190043': 194, '27-143
8405': 195, '27-2389039': 196, '27-3459350': 197, '27-4023565': 198, '27-46
26537': 199, '27280': 200, '2734': 201, '2735146': 202, '27601': 203, '2765
79': 204, '280110': 205, '28116': 206, '28151': 207, '28316': 208, '28405':
209, '28453': 210, '28454': 211, '28489': 212, '28599': 213, '2888798': 21
4, '29012': 215, '29058': 216, '29209': 217, '2945': 218, '3027509990': 21
9, '3076248': 220, '30788': 221, '30810': 222, '31-1197926': 223, '31-16900
08': 224, '31-1712553': 225, '311845': 226, '3121': 227, '31286': 228, '315
0447': 229, '319': 230, '32-0016270': 231, '3212149': 232, '32178': 233, '3
2489': 234, '3284070': 235, '33-0231744': 236, '33-0397503': 237, '33-04199
92': 238, '33-0594693': 239, '33-0750812': 240, '33-0816610': 241, '33-0828
099': 242, '33-0941669': 243, '33-0962918': 244, '33-0975529': 245, '3318
3': 246, '33508': 247, '3374412': 248, '33806': 249, '33826': 250, '33985
8': 251, '34-1194858': 252, '34-1633105': 253, '34-1716542': 254, '34-17196
15': 255, '34-2000096': 256, '34106': 257, '34214': 258, '34585': 259, '346
07': 260, '34627': 261, '34953': 262, '35-2486440': 263, '35013': 264, '350
14': 265, '35139': 266, '35261': 267, '35355': 268, '3537897': 269, '3540
6': 270, '36-4327855': 271, '365325': 272, '37-1493496': 273, '37-1542226':
274, '38-2434249': 275, '38-2749215': 276, '38-2750395': 277, '38-2799035':
278, '38-3564305': 279, '39-2001010': 280, '3918898': 281, '3938186': 282,
'3956': 283, '4004574': 284, '41-1795868': 285, '41-1914032': 286, '41-2277
737': 287, '4114567': 288, '413208': 289, '4142': 290, '42-1554181': 291,
'42-1739728': 292, '4239': 293, '43-0951349': 294, '43-1965151': 295, '43
5': 296, '45-3143795': 297, '45-3508295': 298, '451965': 299, '46-0530020':
300, '46-1728831': 301, '46-1836968': 302, '46-2477192': 303, '46-2888046':
304, '46-3435079': 305, '46-3528270': 306, '46-3676810': 307, '46-5671661':
308, '47-0873092': 309, '47-0912342': 310, '47-0933090': 311, '47-3632618':
312, '471809999': 313, '476810': 314, '48-1148159': 315, '48-1236121': 316,
'480228': 317, '4878': 318, '491224': 319, '501105': 320, '504713': 321, '5
1-0488301': 322, '51-0517525': 323, '52-2091594': 324, '52-2276553': 325,
'52-2321476': 326, '52-2323186': 327, '5380': 328, '54-0259290': 329, '54-1
094297': 330, '54-1994393': 331, '54-2070914': 332, '542409990': 333, '5426
49': 334, '546571': 335, '55130': 336, '5556209999': 337, '5582': 338, '558
8': 339, '56-2103469': 340, '56-2237729': 341, '56-2471041': 342, '566': 34
3, '57-1175755': 344, '57033': 345, '57071': 346, '57074': 347, '57167': 34
8, '57451': 349, '57542': 350, '57607': 351, '57614': 352, '57633': 353, '5
7776': 354, '57777': 355, '57949': 356, '57955': 357, '57978': 358, '58-186
5166': 359, '58305': 360, '58322': 361, '58341': 362, '58380': 363, '5877
8': 364, '59-3378746': 365, '5912': 366, '5913': 367, '592448': 368, '59527
0': 369, '595869': 370, '60042': 371, '60059': 372, '60079': 373, '601050':
374, '60143': 375, '60438': 376, '606046': 377, '60613': 378, '6158': 379,
'6161': 380, '617677': 381, '619877': 382, '62-0997810': 383, '62-1494087':
384, '62-1532940': 385, '624': 386, '6248': 387, '62659': 388, '62665': 38
9, '62745': 390, '6288': 391, '63-1052225': 392, '63069': 393, '63194': 39
4, '63196': 395, '6328': 396, '63315': 397, '63440': 398, '63799': 399, '64
103': 400, '644': 401, '6443809990': 402, '64482': 403, '64546': 404, '6455
2': 405, '65595': 406, '656377': 407, '65644': 408, '656733': 409, '66157':
410, '66328': 411, '66331': 412, '66337': 413, '66349': 414, '66373': 415,
'66399': 416, '665': 417, '66734': 418, '66751': 419, '66841': 420, '6720
1': 421, '67262': 422, '67264': 423, '67389': 424, '675332': 425, '67911':
426, '67955': 427, '68-0151632': 428, '68-0295876': 429, '68-0309242': 430,
'68061': 431, '68095': 432, '68186': 433, '68187': 434, '68196': 435, '6820
3': 436, '68205': 437, '68222': 438, '68223': 439, '68237': 440, '68239': 4
41, '68253': 442, '68255': 443, '68271': 444, '68278': 445, '68284': 446,
'68293': 447, '68298': 448, '68304': 449, '68315': 450, '68362': 451, '6837
5': 452, '68394': 453, '68423': 454, '68457': 455, '68465': 456, '68517': 4
57, '68530': 458, '68546': 459, '685676': 460, '68576': 461, '68598': 462,
'68601': 463, '68613': 464, '68622': 465, '694904': 466, '697633': 467, '70
2825': 468, '702889': 469, '703927': 470, '704347': 471, '7056000000': 472,
```

'706707': 473, '706809': 474, '707674': 475, '708146': 476, '708186': 477, '7101': 478, '712504': 479, '713570': 480, '713964': 481, '714970': 482, '715018': 483, '716195': 484, '7162800002': 485, '716456': 486, '717936': 487, '7197000003': 488, '72-1545376': 489, '7257500009': 490, '7272800006': 491, '7275700004': 492, '73-1374559': 493, '73-1545233': 494, '73-1577221': 495, '74-2508160': 496, '75-2585326': 497, '75-2695327': 498, '75-2921540': 499, '75-3170028': 500, '7505400005': 501, '7516800003': 502, '75633': 503, '76-0236067': 504, '76-0503625': 505, '76-0561995': 506, '76-0629353': 507, '7635500004': 508, '7638200000': 509, '7667200009': 510, '7674000006': 511, '77-0158990': 512, '77-0605392': 513, '77-0672274': 514, '77-0717225': 515, '7748': 516, '7756300009': 517, '7810600004': 518, '7811300008': 519, '7927200007': 520, '7983500003': 521, '7992700007': 522, '80-0233937': 523, '80-0312140': 524, '80-0860209': 525, '80122': 526, '804963': 527, '8100': 528, '817824': 529, '83-0171636': 530, '83-0368862': 531, '84-0927358': 532, '84-1040263': 533, '84-1412422': 534, '84-1496821': 535, '84-1564935': 536, '84-1594306': 537, '85-0260899': 538, '852218': 539, '852320': 540, '857': 541, '86-0415227': 542, '86-0431588': 543, '86-0634557': 544, '86-0860478': 545, '8663': 546, '87-0623581': 547, '87-0675992': 548, '87-0682600': 549, '87-0691650': 550, '8796': 551, '8797': 552, '88-0209429': 553, '88-0508228': 554, '8854': 555, '90-0790926': 556, '91-1374387': 557, '91-1395192': 558, '91-1441009': 559, '91-1465333': 560, '91-1529683': 561, '91-1569077': 562, '91-1780488': 563, '91-1841798': 564, '91-1913382': 565, '91-2006136': 566, '915878': 567, '9289': 568, '93-1231049': 569, '93-1248952': 570, '93-1296762': 571, '93-1301081': 572, '934329': 573, '9366': 574, '936855': 575, '9373': 576, '94-3195577': 577, '9483': 578, '9486': 579, '95-3821253': 580, '95-3990375': 581, '95-4196389': 582, '95-4234730': 583, '95-4267987': 584, '95-4462959': 585, '95-4482547': 586, '95-4523866': 587, '95-4623407': 588, '95-4762204': 589, '95-4769926': 590, '95-4866828': 591, '972590': 592}

Mapping for purchaser_type_name:
{'Affiliate institution': 0, 'Commercial bank, savings bank or savings association': 1, 'Fannie Mae (FNMA)': 2, 'Farmer Mac (FAMC)': 3, 'Freddie Mac (FHLMC)': 4, 'Ginnie Mae (GNMA)': 5, 'Life insurance company, credit union, mortgage bank, or finance company': 6, 'Loan was not originated or was not sold in calendar year covered by register': 7, 'Other type of purchaser': 8, 'Private securitization': 9}

Mapping for property_type_name:
{'Manufactured housing': 0, 'Multifamily dwelling': 1, 'One-to-four family dwelling (other than manufactured housing)': 2}

Mapping for preapproval_name:
{'Not applicable': 0, 'Preapproval was not requested': 1, 'Preapproval was requested': 2}

Mapping for owner_occupancy_name:
{'Not applicable': 0, 'Not owner-occupied as a principal dwelling': 1, 'Owner-occupied as a principal dwelling': 2}

Mapping for msamd_name:
{'Bellingham - WA': 0, 'Bremerton, Silverdale - WA': 1, 'Kennewick, Richland - WA': 2, 'Lewiston - ID, WA': 3, 'Longview - WA': 4, 'Mount Vernon, Anacortes - WA': 5, 'Olympia, Tumwater - WA': 6, 'Portland, Vancouver, Hillsboro - OR, WA': 7, 'Seattle, Bellevue, Everett - WA': 8, 'Spokane, Spokane Valley - WA': 9, 'Tacoma, Lakewood - WA': 10, 'Walla Walla - WA': 11, 'Wenatchee - WA': 12, 'Yakima - WA': 13}

Mapping for loan_type_name:
{'Conventional': 0, 'FHA-insured': 1, 'FSA/RHS-guaranteed': 2, 'VA-guaranteed': 3}

```
Mapping for loan_purpose_name:
{'Home improvement': 0, 'Home purchase': 1, 'Refinancing': 2}

Mapping for lien_status_name:
{'Not applicable': 0, 'Not secured by a lien': 1, 'Secured by a first lie
n': 2, 'Secured by a subordinate lien': 3}

Mapping for hoepa_status_name:
{'HOEPA loan': 0, 'Not a HOEPA loan': 1}

Mapping for county_name:
{'Adams County': 0, 'Asotin County': 1, 'Benton County': 2, 'Chelan Count
y': 3, 'Clallam County': 4, 'Clark County': 5, 'Columbia County': 6, 'Cowli
tz County': 7, 'Douglas County': 8, 'Ferry County': 9, 'Franklin County': 1
0, 'Garfield County': 11, 'Grant County': 12, 'Grays Harbor County': 13, 'I
sland County': 14, 'Jefferson County': 15, 'King County': 16, 'Kitsap Count
y': 17, 'Kittitas County': 18, 'Klickitat County': 19, 'Lewis County': 20,
'Lincoln County': 21, 'Mason County': 22, 'Okanogan County': 23, 'Pacific C
ounty': 24, 'Pend Oreille County': 25, 'Pierce County': 26, 'San Juan Count
y': 27, 'Skagit County': 28, 'Skamania County': 29, 'Snohomish County': 30,
'Spokane County': 31, 'Stevens County': 32, 'Thurston County': 33, 'Wahkiak
um County': 34, 'Walla Walla County': 35, 'Whatcom County': 36, 'Whitman Co
unty': 37, 'Yakima County': 38}

Mapping for co_applicant_sex_name:
{'Female': 0, 'Information not provided by applicant in mail, Internet, or
telephone application': 1, 'Male': 2, 'No co-applicant': 3, 'Not applicabl
e': 4}

Mapping for co_applicant_ethnicity_name:
{'Hispanic or Latino': 0, 'Information not provided by applicant in mail, I
nternet, or telephone application': 1, 'No co-applicant': 2, 'Not Hispanic
or Latino': 3, 'Not applicable': 4}

Mapping for applicant_sex_name:
{'Female': 0, 'Information not provided by applicant in mail, Internet, or
telephone application': 1, 'Male': 2, 'Not applicable': 3}

Mapping for applicant_ethnicity_name:
{'Hispanic or Latino': 0, 'Information not provided by applicant in mail, I
nternet, or telephone application': 1, 'Not Hispanic or Latino': 2, 'Not ap
plicable': 3}

Mapping for agency_name:
{'Consumer Financial Protection Bureau': 0, 'Department of Housing and Urba
n Development': 1, 'Federal Deposit Insurance Corporation': 2, 'Federal Res
erve System': 3, 'National Credit Union Administration': 4, 'Office of the
Comptroller of the Currency': 5}

Mapping for agency_abbr:
{'CFPB': 0, 'FDIC': 1, 'FRS': 2, 'HUD': 3, 'NCUA': 4, 'OCC': 5}

Mapping for action_taken_name:
{'Application approved but not accepted': 0, 'Application denied by financi
al institution': 1, 'Application withdrawn by applicant': 2, 'File closed f
or incompleteness': 3, 'Loan originated': 4, 'Loan purchased by the institu
tion': 5, 'Preapproval request approved but not accepted': 6, 'Preapproval
request denied by financial institution': 7}
        S.no  respondent_id  purchaser_type_name  property_type_name  \
0          1            317                    4                   2
```

| | | | | |
|---|---|---|---|---|
| 1 | 2 | 490 | 6 | 2 |
| 2 | 3 | 489 | 7 | 2 |
| 3 | 4 | 318 | 7 | 2 |
| 4 | 5 | 234 | 4 | 2 |
| ... | ... | ... | ... | ... |
| 59995 | 59996 | 472 | 8 | 2 |
| 59996 | 59997 | 488 | 4 | 2 |
| 59997 | 59998 | 55 | 5 | 2 |
| 59998 | 59999 | 116 | 5 | 2 |
| 59999 | 60000 | 47 | 2 | 2 |

| | preapproval_name | owner_occupancy_name | msamd_name | loan_type_name |
|---|---|---|---|---|
| 0 | 0 | 2 | 7 | 0 |
| 1 | 0 | 2 | 11 | 1 |
| 2 | 0 | 2 | 7 | 0 |
| 3 | 0 | 2 | 7 | 0 |
| 4 | 0 | 2 | 1 | 0 |
| ... | ... | ... | ... | ... |
| 59995 | 0 | 2 | 7 | 0 |
| 59996 | 0 | 1 | 5 | 0 |
| 59997 | 0 | 2 | 0 | 1 |
| 59998 | 0 | 2 | 1 | 3 |
| 59999 | 0 | 2 | 7 | 0 |

| | loan_purpose_name | lien_status_name | hoepa_status_name | county_name |
|---|---|---|---|---|
| 0 | 2 | 2 | 1 | 5 |
| 1 | 1 | 2 | 1 | 35 |
| 2 | 2 | 2 | 1 | 5 |
| 3 | 2 | 2 | 1 | 5 |
| 4 | 0 | 2 | 1 | 17 |
| ... | ... | ... | ... | ... |
| 59995 | 2 | 2 | 1 | 5 |
| 59996 | 2 | 2 | 1 | 28 |
| 59997 | 0 | 2 | 1 | 36 |
| 59998 | 2 | 2 | 1 | 17 |
| 59999 | 2 | 2 | 1 | 5 |

| | co_applicant_sex_name | co_applicant_ethnicity_name | applicant_sex_name |
|---|---|---|---|
| 0 | 2 | 3 | 0 |
| 1 | 3 | 2 | 2 |
| 2 | 0 | 3 | 2 |
| 3 | 0 | 1 | 2 |
| 4 | 2 | 3 | 0 |
| ... | ... | ... | ... |
| 59995 | 0 | 0 | 2 |
| 59996 | 3 | 2 | 0 |
| 59997 | 0 | 3 | 2 |
| 59998 | 0 | 3 | 2 |

```
59999                        0                        3
2
```

| | applicant_ethnicity_name | agency_name | agency_abbr | action_taken_name |
|---|---|---|---|---|
| 0 4 | 2 | 0 | 0 | |
| 1 4 | 0 | 1 | 3 | |
| 2 4 | 2 | 1 | 3 | |
| 3 4 | 1 | 4 | 4 | |
| 4 4 | 2 | 2 | 1 | |
| ... ... | ... | ... | ... | |
| 59995 4 | 2 | 1 | 3 | |
| 59996 4 | 1 | 1 | 3 | |
| 59997 4 | 2 | 1 | 3 | |
| 59998 4 | 2 | 0 | 0 | |
| 59999 4 | 2 | 0 | 0 | |

```
[60000 rows x 19 columns]
time: 396 ms (started: 2024-04-13 11:17:53 +00:00)
```

In [20]:
```python
print(imputed_data_non_categorical.columns)
```

```
Index(['S.no', 'tract_to_msamd_income', 'population', 'minority_populatio
n',
       'number_of_owner_occupied_units', 'number_of_1_to_4_family_units',
       'loan_amount_000s', 'hud_median_family_income', 'applicant_income_00
0s',
       'sequence_number', 'census_tract_number', 'application_date_indicato
r',
       'Cluster_Label'],
      dtype='object')
time: 1.02 ms (started: 2024-04-13 11:17:53 +00:00)
```

In [21]:
```python
def identify_outliers(column):
    Q1 = np.percentile(column, 25)
    Q3 = np.percentile(column, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = (column < lower_bound) | (column > upper_bound)
    return outliers

# Apply the function to each column to get a DataFrame of True/False values
outliers = imputed_data_non_categorical.apply(identify_outliers)

# Display the number of outliers for each column
outlier_counts = outliers.sum()
print(outlier_counts)
```

```
S.no                               0
tract_to_msamd_income           1309
population                       553
minority_population             2641
number_of_owner_occupied_units   478
number_of_1_to_4_family_units   2150
loan_amount_000s                2467
hud_median_family_income           0
applicant_income_000s           3765
sequence_number                 7898
census_tract_number             9391
application_date_indicator       776
Cluster_Label                      0
dtype: int64
time: 59.1 ms (started: 2024-04-13 11:17:53 +00:00)
```

In [22]:
```python
# Iterate through each column and print count of unique values
for column in imputed_data_non_categorical.columns:
    unique_count = imputed_data_non_categorical[column].nunique()
    print(f"Count of unique values in {column} column: {unique_count}")
```

```
Count of unique values in S.no column: 60000
Count of unique values in tract_to_msamd_income column: 1328
Count of unique values in population column: 1284
Count of unique values in minority_population column: 1217
Count of unique values in number_of_owner_occupied_units column: 997
Count of unique values in number_of_1_to_4_family_units column: 1052
Count of unique values in loan_amount_000s column: 1344
Count of unique values in hud_median_family_income column: 16
Count of unique values in applicant_income_000s column: 808
Count of unique values in sequence_number column: 39340
Count of unique values in census_tract_number column: 1109
Count of unique values in application_date_indicator column: 2
Count of unique values in Cluster_Label column: 5
time: 34.1 ms (started: 2024-04-13 11:17:53 +00:00)
```

In [23]:
```python
# Initialize the StandardScaler
scaler = StandardScaler()

# Apply Standard Scaling to your dataset
scaled_data = scaler.fit_transform(imputed_data_non_categorical)

def identify_outliers(column):
    Q1 = np.percentile(column, 25)
    Q3 = np.percentile(column, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = (column < lower_bound) | (column > upper_bound)
    return outliers

# Apply the function to each column in the scaled dataset
outliers_scaled = pd.DataFrame(scaled_data, columns=imputed_data_non_categorical.columns).apply(identify_outliers)

# Display the number of outliers for each column in the scaled dataset
outlier_counts_scaled = outliers_scaled.sum()
print(outlier_counts_scaled)
```

```
S.no                                0
tract_to_msamd_income            1309
population                        553
minority_population              2641
number_of_owner_occupied_units    478
number_of_1_to_4_family_units    2150
loan_amount_000s                 2467
hud_median_family_income            0
applicant_income_000s            3765
sequence_number                  7898
census_tract_number              9391
application_date_indicator        776
Cluster_Label                       0
dtype: int64
time: 59.9 ms (started: 2024-04-13 11:17:54 +00:00)
```

In [24]:
```python
# Initialize the RobustScaler
scaler = RobustScaler()

# Apply Robust Scaling to your dataset
scaled_data_robust = scaler.fit_transform(imputed_data_non_categorical)

# Check for outliers in the scaled dataset
outliers_robust = pd.DataFrame(scaled_data_robust, columns=imputed_data_non
_categorical.columns).apply(identify_outliers)

# Display the number of outliers for each column in the scaled dataset
outlier_counts_robust = outliers_robust.sum()
print(outlier_counts_robust)
```

```
S.no                               0
tract_to_msamd_income           1309
population                       553
minority_population             2641
number_of_owner_occupied_units   478
number_of_1_to_4_family_units   2150
loan_amount_000s                2467
hud_median_family_income           0
applicant_income_000s           3765
sequence_number                 7898
census_tract_number             9391
application_date_indicator       776
Cluster_Label                      0
dtype: int64
time: 76.6 ms (started: 2024-04-13 11:17:54 +00:00)
```

In [25]:
```python
# Define columns to exclude from normalization
columns_to_exclude = ['S.no','hud_median_family_income', 'application_date_
indicator', 'Cluster_Label']

# Create a copy of the DataFrame with excluded columns
data_to_scale = imputed_data_non_categorical.drop(columns=columns_to_exclud
e)

# Initialize the MinMaxScaler
scaler = MinMaxScaler()

# Apply Min-Max Scaling to the selected columns
scaled_data = scaler.fit_transform(data_to_scale)

# Create a DataFrame with scaled data and original column names
scaled_df = pd.DataFrame(scaled_data, columns=data_to_scale.columns)

# Add back the excluded columns to the scaled DataFrame
scaled_df[columns_to_exclude] = imputed_data_non_categorical[columns_to_exc
lude]

# Display the scaled DataFrame
print(scaled_df)
```

```
        tract_to_msamd_income   population   minority_population   \
0                    0.442799     0.640752              0.234501
1                    0.285162     0.372631              0.236658
2                    0.317084     0.385008              0.105445
3                    0.543502     0.381682              0.070620
4                    0.610556     0.393363              0.091213
...                       ...          ...                   ...
59995                0.394915     0.299373              0.149434
59996                0.445679     0.369923              0.060162
59997                0.418734     0.609964              0.105553
59998                0.336419     0.179392              0.322803
59999                0.442799     0.640752              0.234501

        number_of_owner_occupied_units   number_of_1_to_4_family_units   \
0                             0.724346                        0.448858
1                             0.420188                        0.298329
2                             0.375922                        0.308728
3                             0.506372                        0.305660
4                             0.566734                        0.354074
...                                ...                             ...
59995                         0.370557                        0.240709
59996                         0.556673                        0.383396
59997                         0.783032                        0.558814
59998                         0.195171                        0.116263
59999                         0.724346                        0.448858

        loan_amount_000s   applicant_income_000s   sequence_number   \
0               0.004109                0.018669          0.096625
1               0.004346                0.006656          0.042368
2               0.004364                0.018831          0.005001
3               0.006364                0.050974          0.000158
4               0.007564                0.018344          0.026241
...                  ...                     ...               ...
59995           0.002927                0.012013          0.024452
59996           0.002564                0.007955          0.268165
59997           0.004546                0.014123          0.020504
59998           0.004655                0.018153          0.289037
59999           0.005309                0.013149          0.032247

        census_tract_number      S.no   hud_median_family_income   \
0                  0.042258       1.0                     73300.0
1                  0.943728       2.0                     57900.0
2                  0.042333       3.0                     73300.0
3                  0.041421       4.0                     73300.0
4                  0.092866       5.0                     78100.0
...                     ...       ...                         ...
59995              0.041926   59996.0                     73300.0
59996              0.963715   59997.0                     61400.0
59997              0.000724   59998.0                     69900.0
59998              0.082002   59999.0                     78100.0
59999              0.042258   60000.0                     73300.0

        application_date_indicator   Cluster_Label
0                              0.0             4.0
1                              0.0             3.0
2                              0.0             4.0
3                              0.0             4.0
4                              0.0             4.0
...                            ...             ...
59995                          0.0             1.0
59996                          0.0             1.0
```

```
59997                          0.0        1.0
59998                          0.0        1.0
59999                          0.0        1.0

[60000 rows x 13 columns]
time: 24.8 ms (started: 2024-04-13 11:17:54 +00:00)
```

In [26]:
```python
def identify_outliers(column):
    Q1 = np.percentile(column, 25)
    Q3 = np.percentile(column, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = (column < lower_bound) | (column > upper_bound)
    return outliers

# Apply the function to each column in the scaled dataset
scaled_outliers = scaled_df.apply(identify_outliers)

# Display the number of outliers for each column in the scaled dataset
scaled_outlier_counts = scaled_outliers.sum()
print(scaled_outlier_counts)
```

```
tract_to_msamd_income          1309
population                      553
minority_population            2641
number_of_owner_occupied_units  478
number_of_1_to_4_family_units  2150
loan_amount_000s               2467
applicant_income_000s          3765
sequence_number                7898
census_tract_number            9391
S.no                              0
hud_median_family_income          0
application_date_indicator      776
Cluster_Label                     0
dtype: int64
time: 45.4 ms (started: 2024-04-13 11:17:54 +00:00)
```

In [27]:
```python
# Calculate standard deviation for non-categorical columns
std_deviation_non_categorical1 = scaled_df.std()

# Creating a DataFrame to display the results
dispersion_non_categorical_df1 = pd.DataFrame({
    'Variable': scaled_df.columns,
    'Standard Deviation': std_deviation_non_categorical1.values
})

print(dispersion_non_categorical_df1)
```

```
                          Variable  Standard Deviation
0           tract_to_msamd_income            0.116026
1                      population            0.132618
2              minority_population            0.155273
3     number_of_owner_occupied_units          0.173640
4       number_of_1_to_4_family_units         0.125768
5                  loan_amount_000s            0.010999
6              applicant_income_000s           0.018857
7                  sequence_number            0.121228
8              census_tract_number            0.344020
9                             S.no       17320.652413
10           hud_median_family_income     12798.211781
11         application_date_indicator          0.225976
12                    Cluster_Label            1.395815
time: 16.5 ms (started: 2024-04-13 11:17:54 +00:00)
```

In [28]:
```python
# Pre-Processed Dataset
combined_data = pd.merge(encoded_data_categorical, scaled_df, on='S.no')

# Display the Pre-Processed Dataset
%memit
combined_data
```

```
peak memory: 420.05 MiB, increment: 0.09 MiB
```

Out[28]:

|       | S.no  | respondent_id | purchaser_type_name | property_type_name | preapproval_name | c |
|-------|-------|---------------|---------------------|--------------------|------------------|---|
| 0     | 1     | 317           | 4                   | 2                  | 0                |   |
| 1     | 2     | 490           | 6                   | 2                  | 0                |   |
| 2     | 3     | 489           | 7                   | 2                  | 0                |   |
| 3     | 4     | 318           | 7                   | 2                  | 0                |   |
| 4     | 5     | 234           | 4                   | 2                  | 0                |   |
| ...   | ...   | ...           | ...                 | ...                | ...              |   |
| 59995 | 59996 | 472           | 8                   | 2                  | 0                |   |
| 59996 | 59997 | 488           | 4                   | 2                  | 0                |   |
| 59997 | 59998 | 55            | 5                   | 2                  | 0                |   |
| 59998 | 59999 | 116           | 5                   | 2                  | 0                |   |
| 59999 | 60000 | 47            | 2                   | 2                  | 0                |   |

60000 rows × 31 columns

◄ ▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮                                                                                          ▶

```
time: 346 ms (started: 2024-04-13 11:17:54 +00:00)
```

```
In [29]:  # Get the index of the 'Cluster_Label' column
          cluster_label_index = combined_data.columns.get_loc('Cluster_Label')

          # Reorder the columns to move 'Cluster_Label' to the extreme right
          combined_data = combined_data[[col for col in combined_data if col != 'Clus
          ter_Label'] + ['Cluster_Label']]

          # Display the updated dataset
          print(combined_data.head())
```

```
# Get the index of the 'Cluster_Label' column
cluster_label_index = combined_data.columns.get_loc('Cluster_Label')

# Reorder the columns to move 'Cluster_Label' to the extreme right
combined_data = combined_data[[col for col in combined_data if col != 'Clus
ter_Label'] + ['Cluster_Label']]

# Display the updated dataset
print(combined_data.head())
```

```
     S.no   respondent_id   purchaser_type_name   property_type_name   \
0     1             317                       4                     2
1     2             490                       6                     2
2     3             489                       7                     2
3     4             318                       7                     2
4     5             234                       4                     2

     preapproval_name   owner_occupancy_name   msamd_name   loan_type_name   \
0                   0                      2            7                  0
1                   0                      2           11                  1
2                   0                      2            7                  0
3                   0                      2            7                  0
4                   0                      2            1                  0

     loan_purpose_name   lien_status_name   ...   minority_population   \
0                     2                  2   ...              0.234501
1                     1                  2   ...              0.236658
2                     2                  2   ...              0.105445
3                     2                  2   ...              0.070620
4                     0                  2   ...              0.091213

     number_of_owner_occupied_units   number_of_1_to_4_family_units   \
0                          0.724346                        0.448858
1                          0.420188                        0.298329
2                          0.375922                        0.308728
3                          0.506372                        0.305660
4                          0.566734                        0.354074

     loan_amount_000s   applicant_income_000s   sequence_number   \
0            0.004109                0.018669          0.096625
1            0.004346                0.006656          0.042368
2            0.004364                0.018831          0.005001
3            0.006364                0.050974          0.000158
4            0.007564                0.018344          0.026241

     census_tract_number   hud_median_family_income   application_date_indicato
r   \
0               0.042258                    73300.0                            0.
0
1               0.943728                    57900.0                            0.
0
2               0.042333                    73300.0                            0.
0
3               0.041421                    73300.0                            0.
0
4               0.092866                    78100.0                            0.
0

     Cluster_Label
0              4.0
1              3.0
2              4.0
3              4.0
4              4.0

[5 rows x 31 columns]
time: 22.9 ms (started: 2024-04-13 11:17:54 +00:00)
```

```
In [30]:  df_ppd_subset = combined_data.copy()
```

time: 16.8 ms (started: 2024-04-13 11:17:54 +00:00)

```
In [31]:  list(combined_data.columns)
```

Out[31]:  ['S.no',
           'respondent_id',
           'purchaser_type_name',
           'property_type_name',
           'preapproval_name',
           'owner_occupancy_name',
           'msamd_name',
           'loan_type_name',
           'loan_purpose_name',
           'lien_status_name',
           'hoepa_status_name',
           'county_name',
           'co_applicant_sex_name',
           'co_applicant_ethnicity_name',
           'applicant_sex_name',
           'applicant_ethnicity_name',
           'agency_name',
           'agency_abbr',
           'action_taken_name',
           'tract_to_msamd_income',
           'population',
           'minority_population',
           'number_of_owner_occupied_units',
           'number_of_1_to_4_family_units',
           'loan_amount_000s',
           'applicant_income_000s',
           'sequence_number',
           'census_tract_number',
           'hud_median_family_income',
           'application_date_indicator',
           'Cluster_Label']

time: 4.03 ms (started: 2024-04-13 11:17:54 +00:00)

# DT

```
In [32]:  #####  DT
```

time: 354 µs (started: 2024-04-13 11:17:54 +00:00)

```
In [33]:  df1 = df_ppd_subset.copy()
```

time: 12.1 ms (started: 2024-04-13 11:17:54 +00:00)

In [34]: `df1.columns`

Out[34]:
```
Index(['S.no', 'respondent_id', 'purchaser_type_name', 'property_type_nam
e',
       'preapproval_name', 'owner_occupancy_name', 'msamd_name',
       'loan_type_name', 'loan_purpose_name', 'lien_status_name',
       'hoepa_status_name', 'county_name', 'co_applicant_sex_name',
       'co_applicant_ethnicity_name', 'applicant_sex_name',
       'applicant_ethnicity_name', 'agency_name', 'agency_abbr',
       'action_taken_name', 'tract_to_msamd_income', 'population',
       'minority_population', 'number_of_owner_occupied_units',
       'number_of_1_to_4_family_units', 'loan_amount_000s',
       'applicant_income_000s', 'sequence_number', 'census_tract_number',
       'hud_median_family_income', 'application_date_indicator',
       'Cluster_Label'],
      dtype='object')
```

time: 3.87 ms (started: 2024-04-13 11:17:54 +00:00)

In [35]: `df1.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 60000 entries, 0 to 59999
Data columns (total 31 columns):
 #   Column                          Non-Null Count  Dtype
---  ------                          --------------  -----
 0   S.no                            60000 non-null  object
 1   respondent_id                   60000 non-null  int64
 2   purchaser_type_name             60000 non-null  int64
 3   property_type_name              60000 non-null  int64
 4   preapproval_name                60000 non-null  int64
 5   owner_occupancy_name            60000 non-null  int64
 6   msamd_name                      60000 non-null  int64
 7   loan_type_name                  60000 non-null  int64
 8   loan_purpose_name               60000 non-null  int64
 9   lien_status_name                60000 non-null  int64
 10  hoepa_status_name               60000 non-null  int64
 11  county_name                     60000 non-null  int64
 12  co_applicant_sex_name           60000 non-null  int64
 13  co_applicant_ethnicity_name     60000 non-null  int64
 14  applicant_sex_name              60000 non-null  int64
 15  applicant_ethnicity_name        60000 non-null  int64
 16  agency_name                     60000 non-null  int64
 17  agency_abbr                     60000 non-null  int64
 18  action_taken_name               60000 non-null  int64
 19  tract_to_msamd_income           60000 non-null  float64
 20  population                      60000 non-null  float64
 21  minority_population             60000 non-null  float64
 22  number_of_owner_occupied_units  60000 non-null  float64
 23  number_of_1_to_4_family_units   60000 non-null  float64
 24  loan_amount_000s                60000 non-null  float64
 25  applicant_income_000s           60000 non-null  float64
 26  sequence_number                 60000 non-null  float64
 27  census_tract_number             60000 non-null  float64
 28  hud_median_family_income        60000 non-null  float64
 29  application_date_indicator      60000 non-null  float64
 30  Cluster_Label                   60000 non-null  float64
dtypes: float64(12), int64(18), object(1)
memory usage: 14.2+ MB
time: 22.4 ms (started: 2024-04-13 11:17:54 +00:00)
```

```
In [36]: df1_inputs_all = df1[['respondent_id', 'purchaser_type_name', 'property_typ
         e_name',
                 'preapproval_name', 'owner_occupancy_name', 'msamd_name',
                 'loan_type_name', 'loan_purpose_name', 'lien_status_name',
                 'hoepa_status_name', 'county_name', 'co_applicant_sex_name',
                 'co_applicant_ethnicity_name', 'applicant_sex_name',
                 'applicant_ethnicity_name', 'agency_name', 'agency_abbr',
                 'tract_to_msamd_income', 'population',
                 'minority_population', 'number_of_owner_occupied_units',
                 'number_of_1_to_4_family_units', 'loan_amount_000s',
                 'applicant_income_000s', 'sequence_number', 'census_tract_number',
                 'hud_median_family_income', 'application_date_indicator',
                 'Cluster_Label']]
         df1_output = df1[['action_taken_name']]

         df1_inputs_all_names = df1_inputs_all.columns
         df1_output_labels = df1_output['action_taken_name'].unique().astype(str)
```

time: 12.8 ms (started: 2024-04-13 11:17:54 +00:00)

```
In [37]: # Initialize StratifiedShuffleSplit with desired test size and random state
         stratified_split = StratifiedShuffleSplit(n_splits=1, test_size=0.25, rando
         m_state=45005)

         # Perform the stratified split to get training and testing indices
         for train_index, test_index in stratified_split.split(df1_inputs_all, df1_o
         utput):
             df1_inputs_all_train = df1_inputs_all.iloc[train_index]
             df1_inputs_all_test = df1_inputs_all.iloc[test_index]
             df1_output_train = df1_output.iloc[train_index]
             df1_output_test = df1_output.iloc[test_index]
```

time: 290 ms (started: 2024-04-13 11:17:54 +00:00)

```
In [38]: # Decision Tree : Model (Training Subset)
         dtc_all = DecisionTreeClassifier(criterion='gini', random_state=45005,max_d
         epth=3) # Other Criteria : Entropy,  Log Loss
         dtc_model_all = dtc_all.fit(df1_inputs_all_train, df1_output_train); dtc_mo
         del_all
```

```
Out[38]:               ▼         DecisionTreeClassifier

         DecisionTreeClassifier(max_depth=3, random_state=45005)
```

time: 135 ms (started: 2024-04-13 11:17:55 +00:00)

In [39]:
```python
# Decision Tree : Feature Importance
dtc_all_imp_features = pd.DataFrame({'feature': df1_inputs_all_names, 'impo
rtance': np.round(dtc_model_all.feature_importances_, 3)})
dtc_all_imp_features.sort_values('importance', ascending=False, inplace=Tru
e); dtc_all_imp_features
```

Out[39]:

| | feature | importance |
|---|---|---|
| 27 | application_date_indicator | 0.479 |
| 5 | msamd_name | 0.292 |
| 1 | purchaser_type_name | 0.229 |
| 0 | respondent_id | 0.000 |
| 16 | agency_abbr | 0.000 |
| 26 | hud_median_family_income | 0.000 |
| 25 | census_tract_number | 0.000 |
| 24 | sequence_number | 0.000 |
| 23 | applicant_income_000s | 0.000 |
| 22 | loan_amount_000s | 0.000 |
| 21 | number_of_1_to_4_family_units | 0.000 |
| 20 | number_of_owner_occupied_units | 0.000 |
| 19 | minority_population | 0.000 |
| 18 | population | 0.000 |
| 17 | tract_to_msamd_income | 0.000 |
| 14 | applicant_ethnicity_name | 0.000 |
| 15 | agency_name | 0.000 |
| 13 | applicant_sex_name | 0.000 |
| 12 | co_applicant_ethnicity_name | 0.000 |
| 11 | co_applicant_sex_name | 0.000 |
| 10 | county_name | 0.000 |
| 9 | hoepa_status_name | 0.000 |
| 8 | lien_status_name | 0.000 |
| 7 | loan_purpose_name | 0.000 |
| 6 | loan_type_name | 0.000 |
| 4 | owner_occupancy_name | 0.000 |
| 3 | preapproval_name | 0.000 |
| 2 | property_type_name | 0.000 |
| 28 | Cluster_Label | 0.000 |

```
time: 18.4 ms (started: 2024-04-13 11:17:55 +00:00)
```

```python
# Subset df1 based on Inputs as {mpg, hp, cyl, vs} & Output as {am}
df1_inputs = df1[['application_date_indicator', 'msamd_name', 'purchaser_ty
pe_name', 'loan_type_name', 'loan_purpose_name', 'hud_median_family_incom
e', 'loan_amount_000s']];
df1_inputs
df1_output = df1[['action_taken_name']]; df1_output

df1_inputs_names = df1_inputs.columns; df1_inputs_names
df1_output_labels = df1_output['action_taken_name'].unique().astype(str); d
f1_output_labels
```

Out[40]:  `array(['4', '0', '1', '2', '3', '5', '6', '7'], dtype='<U21')`

time: 9.78 ms (started: 2024-04-13 11:17:55 +00:00)

```python
# Initialize StratifiedShuffleSplit with desired test size and random state
stratified_split = StratifiedShuffleSplit(n_splits=1, test_size=0.25, rando
m_state=45005)

# Perform the stratified split to get training and testing indices
for train_index, test_index in stratified_split.split(df1_inputs, df1_outpu
t):
    df1_inputs_train = df1_inputs.iloc[train_index]
    df1_inputs_test = df1_inputs.iloc[test_index]
    df1_output_train = df1_output.iloc[train_index]
    df1_output_test = df1_output.iloc[test_index]
```

time: 262 ms (started: 2024-04-13 11:17:55 +00:00)

In [42]:
```python
from sklearn.linear_model import LogisticRegression
from sklearn.feature_selection import SelectFromModel
import numpy as np

# Initialize Logistic Regression model with L1 regularization
logreg_l1 = LogisticRegression(penalty='l1', solver='liblinear', random_state=45005)

# Fit the model on the training data
logreg_l1.fit(df1_inputs_train, df1_output_train.values.ravel())

# Get feature importances from the fitted model
feature_importances = np.abs(logreg_l1.coef_).flatten()

# Calculate the threshold as 20% of the maximum feature importance
threshold = 0.2 * np.max(feature_importances)

# Create a selector object to select features based on non-zero coefficients
selector = SelectFromModel(logreg_l1, threshold=threshold)

# Transform the training and testing input data to select features
df1_inputs_train_selected = selector.transform(df1_inputs_train)
df1_inputs_test_selected = selector.transform(df1_inputs_test)

# Get the selected features
selected_features = df1_inputs_names[selector.get_support()]

# Print the selected features and the calculated threshold
print("Selected Features:", selected_features)
print("Threshold:", threshold)
```

```
Selected Features: Index(['application_date_indicator', 'purchaser_type_name',
       'loan_purpose_name'],
      dtype='object')
Threshold: 1.8129252376051945
time: 2.81 s (started: 2024-04-13 11:17:55 +00:00)

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:432: UserWarning: X
has feature names, but SelectFromModel was fitted without feature names
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:432: UserWarning: X
has feature names, but SelectFromModel was fitted without feature names
  warnings.warn(
```

In [43]:
```python
from sklearn.linear_model import LogisticRegression
from sklearn.feature_selection import SelectFromModel
import numpy as np

# Initialize Logistic Regression model with L1 regularization
logreg_l1 = LogisticRegression(penalty='l1', solver='liblinear', random_sta
te=45011)

# Fit the model on the training data
logreg_l1.fit(df1_inputs_train, df1_output_train.values.ravel())

# Get feature importances from the fitted model
feature_importances = np.abs(logreg_l1.coef_).flatten()

# Calculate the threshold as 20% of the maximum feature importance
threshold = 0.2 * np.max(feature_importances)

# Create a selector object to select features based on non-zero coefficient
s
selector = SelectFromModel(logreg_l1, threshold=threshold)

# Transform the training and testing input data to select features
df1_inputs_train_selected = selector.transform(df1_inputs_train)
df1_inputs_test_selected = selector.transform(df1_inputs_test)

# Get the selected features
selected_features = df1_inputs_names[selector.get_support()]

# Print the selected features and the calculated threshold
print("Selected Features:", selected_features)
print("Threshold:", threshold)
```

```
Selected Features: Index(['application_date_indicator', 'purchaser_type_nam
e',
        'loan_purpose_name'],
      dtype='object')
Threshold: 1.7010443233194732
time: 3.24 s (started: 2024-04-13 11:17:58 +00:00)

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:432: UserWarning: X
has feature names, but SelectFromModel was fitted without feature names
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:432: UserWarning: X
has feature names, but SelectFromModel was fitted without feature names
  warnings.warn(
```

In [44]:
```python
# Decision Tree : Model (Training Subset)
dtc = DecisionTreeClassifier(criterion='gini', random_state=45005,max_depth
=3) # Other Criteria : Entropy,  Log Loss
dtc_model = dtc.fit(df1_inputs_train, df1_output_train); dtc_model
```

Out[44]:
```
▼                   DecisionTreeClassifier

DecisionTreeClassifier(max_depth=3, random_state=45005)
```

```
time: 34.9 ms (started: 2024-04-13 11:18:01 +00:00)
```

In [45]:
```python
# Decision Tree : Model Rules
dtc_model_rules = export_text(dtc_model, feature_names = list(df1_inputs_na
mes)); print(dtc_model_rules)
```

```
|--- application_date_indicator <= 1.00
|    |--- purchaser_type_name <= 6.50
|    |    |--- class: 4
|    |--- purchaser_type_name >  6.50
|    |    |--- msamd_name <= 9.50
|    |    |    |--- class: 4
|    |    |--- msamd_name >  9.50
|    |    |    |--- class: 4
|--- application_date_indicator >  1.00
|    |--- class: 5

time: 1.97 ms (started: 2024-04-13 11:18:01 +00:00)
```

In [46]:
```python
# Decision Tree : Feature Importance
dtc_imp_features = pd.DataFrame({'feature': df1_inputs_names, 'importance':
np.round(dtc_model.feature_importances_, 3)})
dtc_imp_features.sort_values('importance', ascending=False, inplace=True);
dtc_imp_features
```

Out[46]:

|   | feature | importance |
|---|---|---|
| 0 | application_date_indicator | 0.479 |
| 1 | msamd_name | 0.292 |
| 2 | purchaser_type_name | 0.229 |
| 3 | loan_type_name | 0.000 |
| 4 | loan_purpose_name | 0.000 |
| 5 | hud_median_family_income | 0.000 |
| 6 | loan_amount_000s | 0.000 |

```
time: 16.6 ms (started: 2024-04-13 11:18:01 +00:00)
```

In [47]:
```python
from sklearn.tree import DecisionTreeClassifier, export_text
import numpy as np
import pandas as pd

# Initialize the Decision Tree classifier with Gini coefficient criterion
dtc_gini = DecisionTreeClassifier(criterion='gini', random_state=45005, max
_depth=3)

# Train the Decision Tree model using the training subset for Gini coeffici
ent
dtc_model_gini = dtc_gini.fit(df1_inputs_train, df1_output_train)

# Print the trained Decision Tree model for Gini coefficient (optional)
print(dtc_model_gini)

# Get the rules of the trained Decision Tree model for Gini coefficient
dtc_model_rules_gini = export_text(dtc_model_gini, feature_names=list(df1_i
nputs_names))
print(dtc_model_rules_gini)

# Calculate feature importance based on Gini coefficient
dtc_imp_features_gini = pd.DataFrame({'feature': df1_inputs_names, 'importa
nce': np.round(dtc_model_gini.feature_importances_, 3)})
dtc_imp_features_gini.sort_values('importance', ascending=False, inplace=Tr
ue)
print(dtc_imp_features_gini)

# Initialize the Decision Tree classifier with entropy criterion
dtc_entropy = DecisionTreeClassifier(criterion='entropy', random_state=4500
5, max_depth=3)

# Train the Decision Tree model using the training subset for entropy
dtc_model_entropy = dtc_entropy.fit(df1_inputs_train, df1_output_train)

# Print the trained Decision Tree model for entropy (optional)
print(dtc_model_entropy)

# Get the rules of the trained Decision Tree model for entropy
dtc_model_rules_entropy = export_text(dtc_model_entropy, feature_names=list
(df1_inputs_names))
print(dtc_model_rules_entropy)

# Calculate feature importance based on entropy
dtc_imp_features_entropy = pd.DataFrame({'feature': df1_inputs_names, 'impo
rtance': np.round(dtc_model_entropy.feature_importances_, 3)})
dtc_imp_features_entropy.sort_values('importance', ascending=False, inplace
=True)
print(dtc_imp_features_entropy)
```

```
DecisionTreeClassifier(max_depth=3, random_state=45005)
|--- application_date_indicator <= 1.00
|    |--- purchaser_type_name <= 6.50
|    |    |--- class: 4
|    |--- purchaser_type_name >  6.50
|    |    |--- msamd_name <= 9.50
|    |    |    |--- class: 4
|    |    |--- msamd_name >  9.50
|    |    |    |--- class: 4
|--- application_date_indicator >  1.00
|    |--- class: 5
```

```
                      feature   importance
0   application_date_indicator      0.479
1                  msamd_name      0.292
2          purchaser_type_name      0.229
3              loan_type_name      0.000
4           loan_purpose_name      0.000
5      hud_median_family_income      0.000
6             loan_amount_000s      0.000
```

```
DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=4500
5)
|--- purchaser_type_name <= 6.50
|    |--- application_date_indicator <= 1.00
|    |    |--- class: 4
|    |--- application_date_indicator >  1.00
|    |    |--- class: 5
|--- purchaser_type_name >  6.50
|    |--- hud_median_family_income <= 72800.00
|    |    |--- msamd_name <= 6.50
|    |    |    |--- class: 4
|    |    |--- msamd_name >  6.50
|    |    |    |--- class: 4
|    |--- hud_median_family_income >  72800.00
|    |    |--- msamd_name <= 7.50
|    |    |    |--- class: 4
|    |    |--- msamd_name >  7.50
|    |    |    |--- class: 4
```

```
                      feature   importance
2          purchaser_type_name      0.389
0   application_date_indicator      0.304
5      hud_median_family_income      0.198
1                  msamd_name      0.109
3              loan_type_name      0.000
4           loan_purpose_name      0.000
6             loan_amount_000s      0.000
time: 68.9 ms (started: 2024-04-13 11:18:01 +00:00)
```

```python
from sklearn.metrics import log_loss
from sklearn.tree import DecisionTreeClassifier

# Assuming you have already trained the DecisionTreeClassifier model
# dtc_model_entropy = DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=45005)
# dtc_model_gini = DecisionTreeClassifier(criterion='gini', max_depth=3, random_state=45005)

# Calculate entropy
y_pred_proba_entropy = dtc_model_entropy.predict_proba(df1_inputs_test)
entropy = log_loss(df1_output_test, y_pred_proba_entropy)

# Calculate Gini impurity
y_pred_proba_gini = dtc_model_gini.predict_proba(df1_inputs_test)
gini_impurity = 1 - (y_pred_proba_gini ** 2).sum(axis=1).mean()

print("Entropy:", entropy)
print("Gini Impurity:", gini_impurity)
```

```
Entropy: 0.15333397314417907
Gini Impurity: 0.08053758117824039
time: 30.1 ms (started: 2024-04-13 11:18:01 +00:00)
```

In [49]:
```python
# Decision Tree : Model Prediction (Training Subset)
dtc_model_predict = dtc_model.predict(df1_inputs_train); dtc_model_predict
```

Out[49]:
```
array([4, 4, 4, ..., 4, 4, 4])

time: 12.8 ms (started: 2024-04-13 11:18:01 +00:00)
```

In [50]:
```python
# Decision Tree : Prediction (Testing Subset)
dtc_predict = dtc_model.predict(df1_inputs_test); dtc_predict
```

Out[50]:
```
array([4, 4, 4, ..., 4, 4, 4])

time: 8.09 ms (started: 2024-04-13 11:18:01 +00:00)
```

In [51]:
```python
# Decision Tree : Model Evaluation (Training Subset)
dtc_model_conf_mat = pd.DataFrame(confusion_matrix(df1_output_train, dtc_mo
del_predict)); dtc_model_conf_mat
dtc_model_perf = classification_report(df1_output_train, dtc_model_predic
t); print(dtc_model_perf)
```

```
              precision    recall  f1-score   support

           0       0.00      0.00      0.00        63
           1       0.00      0.00      0.00       121
           2       0.00      0.00      0.00      2000
           3       0.00      0.00      0.00       334
           4       0.94      1.00      0.97     41861
           5       1.00      1.00      1.00       582
           6       0.00      0.00      0.00        13
           7       0.00      0.00      0.00        26

    accuracy                           0.94     45000
   macro avg       0.24      0.25      0.25     45000
weighted avg       0.89      0.94      0.92     45000


time: 79.5 ms (started: 2024-04-13 11:18:01 +00:00)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:
1344: UndefinedMetricWarning: Precision and F-score are ill-defined and bei
ng set to 0.0 in labels with no predicted samples. Use `zero_division` para
meter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:
1344: UndefinedMetricWarning: Precision and F-score are ill-defined and bei
ng set to 0.0 in labels with no predicted samples. Use `zero_division` para
meter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:
1344: UndefinedMetricWarning: Precision and F-score are ill-defined and bei
ng set to 0.0 in labels with no predicted samples. Use `zero_division` para
meter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

In [52]:
```python
import matplotlib.pyplot as plt
from sklearn.tree import plot_tree

# Set a larger figure size for better clarity
plt.figure(figsize=(10, 10))

# Plot the decision tree
train_subset_dtc_plot = plot_tree(dtc_model, feature_names=df1_inputs_names, class_names=df1_output_labels, rounded=True, filled=True, fontsize=20)

# Show the plot
plt.show()
```

```
application_date_indicator <= 1.0
gini = 0.132
samples = 45000
value = [63, 121, 2000, 334, 41861, 582, 13, 26]
class = 3
```

```
purchaser_type_name
gini = 0.11
samples = 444
value = [63, 121, 2000, 334,
class = 3
```

```
gini = 0.0
samples = 582
value = [0, 0, 0, 0, 0, 582, 0, 0]
class = 5
```

```
gin
sample
value = [0, 0, 0
cla
```

```
msamd_name <= 9.5
gini = 0.31
samples = 14024
value = [63, 121, 2000, 334, 11467, 0, 13, 26]
class = 3
```

```
gini = 0.
samples =
value = [63, 109, 845, 2
class =
```

```
gini = 0.543
samples = 2433
value = [0, 12, 1155, 99, 1166, 0, 1, 0]
class = 3
```

In [53]:
```python
# Set up the plot
ax = plt.axes()

# Plot the confusion matrix with annotations in integer format
sns.heatmap(dtc_model_conf_mat, annot=True, fmt='d', cmap='Paired')

# Set labels and title
ax.set_xlabel('Predicted Label')
ax.set_ylabel('True Label')
ax.set_title('Decision Tree : Confusion Matrix')

# Show the plot
plt.show()
```

time: 514 ms (started: 2024-04-13 11:18:02 +00:00)



Decision Tree : Confusion Matrix

time: 669 ms (started: 2024-04-13 11:18:02 +00:00)

In [54]:
```python
# Cross Validation
from sklearn.model_selection import cross_val_score

# Define your decision tree classifier with desired parameters
dtc_cv = DecisionTreeClassifier(criterion='gini', random_state=45005)

# Perform 5-fold cross-validation
cv_scores = cross_val_score(dtc_cv, df1_inputs, df1_output.values.ravel(),
cv=20)
print("Cross-Validation Scores:", cv_scores)
print("Average Cross-Validation Score:", np.mean(cv_scores))
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:7
00: UserWarning: The least populated class in y has only 17 members, which
is less than n_splits=20.
  warnings.warn(

Cross-Validation Scores: [0.92333333 0.931       0.94633333 0.90266667 0.905
0.919
 0.88566667 0.916       0.932       0.94466667 0.96        0.96566667
 0.96866667 0.96766667 0.96733333 0.97233333 0.97233333 0.97166667
 0.95766667 0.95533333]
Average Cross-Validation Score: 0.9432166666666667
time: 1.35 s (started: 2024-04-13 11:18:03 +00:00)
```

In [55]:
```python
from sklearn.metrics import f1_score

# Compute F1 score
f1 = f1_score(df1_output_test, dtc_predict, average='macro')  # or 'weighte
d' for weighted F1 score
print("F1 Score:", f1)

# Weighted F1 score
weighted_f1 = f1_score(df1_output_test, dtc_predict, average='weighted')
print("Weighted F1 Score:", weighted_f1)
```

```
F1 Score: 0.24629694019471488
Weighted F1 Score: 0.9156413351877607
time: 28.5 ms (started: 2024-04-13 11:18:04 +00:00)
```

In [56]:

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, StratifiedShuffleSpli
t, cross_val_score
from sklearn.tree import DecisionTreeClassifier, export_text, plot_tree
from sklearn.metrics import accuracy_score, classification_report, confusio
n_matrix, f1_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.feature_selection import SelectFromModel
from sklearn.svm import SVC
import numpy as np
import time
import psutil

# Function to measure memory usage
def memory_usage():
    process = psutil.Process()
    return process.memory_info().rss / 1024 ** 2  # Memory usage in MB

# Start time
start_time = time.time()

# Data preprocessing and splitting
# Assuming you have your data loaded into cars_inputs and cars_output
df1_inputs_train, df1_inputs_test, df1_output_train, df1_output_test = trai
n_test_split(df1_inputs, df1_output, test_size=0.2, random_state=42)

# End time
end_time = time.time()

# Time taken for data preprocessing and splitting
data_preprocessing_time = end_time - start_time

# Memory usage after data preprocessing
data_preprocessing_memory = memory_usage()

# Decision Tree
dt_start_time = time.time()
dt_model = DecisionTreeClassifier(criterion='gini', random_state=45005, max
_depth=3)
dt_model.fit(df1_inputs_train, df1_output_train)
dt_training_time = time.time() - dt_start_time
dt_memory_used = memory_usage()
dt_pred = dt_model.predict(df1_inputs_test)
dt_accuracy = accuracy_score(df1_output_test, dt_pred)

# Cross-validation for Decision Tree
dtc_cv_start_time = time.time()
dtc_cv = DecisionTreeClassifier(criterion='gini', random_state=45007)
cv_scores_dtc = cross_val_score(dtc_cv, df1_inputs, df1_output.values.ravel
(), cv=20)
dtc_cv_time = time.time() - dtc_cv_start_time
dtc_cv_accuracy = np.mean(cv_scores_dtc)

print("Decision Tree:")
print(f"  - Training Time (s): {dt_training_time}")
print(f"  - Memory Used (MB): {dt_memory_used}")
print(f"  - Single Split Accuracy: {dt_accuracy}")
```

```
print(f"  - Cross Validation Accuracy: {dtc_cv_accuracy}")
print()
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:7
00: UserWarning: The least populated class in y has only 17 members, which
is less than n_splits=20.
  warnings.warn(

Decision Tree:
   - Training Time (s): 0.026858806610107422
   - Memory Used (MB): 506.37109375
   - Single Split Accuracy: 0.9454166666666667
   - Cross Validation Accuracy: 0.9432499999999999

time: 1.36 s (started: 2024-04-13 11:18:04 +00:00)
```

# Random Forest

In [57]:
```python
## Data Visualization Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.graph_objects as go
from wordcloud import WordCloud
from collections import Counter
from scipy import stats
from sklearn.tree import plot_tree
import graphviz
from IPython.display import display
from collections import Counter

## Machine Learning Models and Evaluation Metrics
from sklearn.ensemble import RandomForestClassifier
from sklearn.utils.validation import column_or_1d
from sklearn.metrics import accuracy_score, classification_report, confusio
n_matrix, f1_score, precision_recall_fscore_support
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression, Lasso, Ridge
from sklearn.metrics import make_scorer
from sklearn.pipeline import make_pipeline
from sklearn.tree import export_graphviz
```

```
time: 1.69 ms (started: 2024-04-13 11:18:05 +00:00)
```

In [58]:
```python
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=4500
5)
```

```
time: 617 µs (started: 2024-04-13 11:18:06 +00:00)
```

In [59]:
```python
rf_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])
```

Out[59]:
```
  ▼           RandomForestClassifier

RandomForestClassifier(random_state=45005)
```

```
time: 2.75 s (started: 2024-04-13 11:18:06 +00:00)
```

In [60]:
```python
y_train_pred_rf = rf_classifier.predict(df1_inputs_train)
y_test_pred_rf = rf_classifier.predict(df1_inputs_test)
```

time: 1.37 s (started: 2024-04-13 11:18:08 +00:00)

In [61]:
```python
from sklearn.metrics import log_loss
from sklearn.ensemble import RandomForestClassifier

# Assuming you have already trained the RandomForestClassifier model
# rf_classifier = RandomForestClassifier(n_estimators=100, random_state=450
05)
# rf_classifier.fit(df1_inputs_train, df1_output_train['action_taken_nam
e'])

# Calculate entropy
y_pred_proba_rf = rf_classifier.predict_proba(df1_inputs_test)
entropy_rf = log_loss(df1_output_test, y_pred_proba_rf)

# Calculate Gini impurity
gini_impurity_rf = 1 - (y_pred_proba_rf ** 2).sum(axis=1).mean()

print("Entropy for Random Forest:", entropy_rf)
print("Gini Impurity for Random Forest:", gini_impurity_rf)
```

Entropy for Random Forest: 0.4448700986425233
Gini Impurity for Random Forest: 0.03587907150255809
time: 483 ms (started: 2024-04-13 11:18:10 +00:00)

In [62]:
```python
# Train the Random Forest classifier
rf_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])

# Print feature importances
feature_importances = rf_classifier.feature_importances_
feature_importance_df = pd.DataFrame({'Feature': df1_inputs_train.columns,
'Importance': feature_importances})
sorted_feature_importance_df = feature_importance_df.sort_values(by='Import
ance', ascending=False)
print("Feature Importances:")
print(sorted_feature_importance_df)
```

```
Feature Importances:
                    Feature  Importance
6           loan_amount_000s    0.333727
2        purchaser_type_name    0.201634
0  application_date_indicator    0.195447
5    hud_median_family_income    0.120389
1                 msamd_name    0.096998
3             loan_type_name    0.027990
4          loan_purpose_name    0.023816
time: 6.03 s (started: 2024-04-13 11:18:10 +00:00)
```

In [63]:
```python
# For training set
print("Training Set Confusion Matrix:")
print(confusion_matrix(df1_output_train['action_taken_name'], y_train_pred_
rf))

print("\nTraining Set Classification Report:")
print(classification_report(df1_output_train['action_taken_name'], y_train_
pred_rf))
```

```
Training Set Confusion Matrix:
[[   31     0    13     1    21     0     0     0]
 [    1    69    21     2    43     0     0     0]
 [    1     8  1825    26   289     0     0     1]
 [    0     3    51   229    75     0     0     0]
 [    2    10   260    27 44307     0     0     0]
 [    0     0     0     0     0   642     0     0]
 [    0     0     1     0     3     0    11     0]
 [    0     0     1     0     2     0     0    24]]

Training Set Classification Report:
              precision    recall  f1-score   support

           0       0.89      0.47      0.61        66
           1       0.77      0.51      0.61       136
           2       0.84      0.85      0.84      2150
           3       0.80      0.64      0.71       358
           4       0.99      0.99      0.99     44606
           5       1.00      1.00      1.00       642
           6       1.00      0.73      0.85        15
           7       0.96      0.89      0.92        27

    accuracy                           0.98     48000
   macro avg       0.91      0.76      0.82     48000
weighted avg       0.98      0.98      0.98     48000

time: 165 ms (started: 2024-04-13 11:18:16 +00:00)
```

In [64]:
```python
# For testing set
print("\nTesting Set Confusion Matrix:")
print(confusion_matrix(df1_output_test['action_taken_name'], y_test_pred_r
f))

print("\nTesting Set Classification Report:")
print(classification_report(df1_output_test['action_taken_name'], y_test_pr
ed_rf))
```

```
Testing Set Confusion Matrix:
[[    0     0     5     3    10     0     0     0]
 [    0     1     7     1    16     0     0     0]
 [    6    10   272    33   194     0     0     1]
 [    0     5    36    16    31     0     0     0]
 [    2     7   196    17 10984     0     3     0]
 [    0     0     0     0     0   134     0     0]
 [    0     0     1     0     1     0     0     0]
 [    0     0     2     0     2     0     1     3]]

Testing Set Classification Report:
              precision    recall  f1-score   support

           0       0.00      0.00      0.00        18
           1       0.04      0.04      0.04        25
           2       0.52      0.53      0.53       516
           3       0.23      0.18      0.20        88
           4       0.98      0.98      0.98     11209
           5       1.00      1.00      1.00       134
           6       0.00      0.00      0.00         2
           7       0.75      0.38      0.50         8

    accuracy                           0.95     12000
   macro avg       0.44      0.39      0.41     12000
weighted avg       0.95      0.95      0.95     12000

time: 88 ms (started: 2024-04-13 11:18:16 +00:00)
```

In [65]:
```python
# Assuming rf is your trained Random Forest classifier
for i in range(3):
    tree = rf_classifier.estimators_[i]  # Assuming rf_classifier is your t
rained Random Forest model
    dot_data = export_graphviz(tree,
                                  feature_names=df1_inputs_train.columns,
                                  filled=True,
                                  max_depth=2,
                                  impurity=False,
                                  proportion=True)
    graph = graphviz.Source(dot_data)
    display(graph)
```



time: 381 ms (started: 2024-04-13 11:18:17 +00:00)

```
In [66]:  # Export and visualize the first three decision trees
          for i in range(3):
              tree = rf_classifier.estimators_[i]
              dot_data = export_graphviz(tree,
                                         feature_names=df1_inputs_train.columns,
                                         filled=True,
                                         max_depth=2,
                                         impurity=False,
                                         proportion=True)
              graph = graphviz.Source(dot_data)
              display(graph)
```







time: 154 ms (started: 2024-04-13 11:18:17 +00:00)

In [67]:
```python
# Initialize a dictionary to store tree frequencies and rules
tree_frequency = Counter()
tree_rules = {}

# Loop through the trees and count their frequency while storing rules
for i in range(len(rf_classifier.estimators_)):
    tree = rf_classifier.estimators_[i]
    tree_str = export_text(tree, feature_names=list(df1_inputs_train.column
s))
    tree_frequency[tree_str] += 1
    tree_rules[tree_str] = tree

# Get the three most frequent trees
top_trees = tree_frequency.most_common(3)

# Print the rules for the top three trees
for tree_str, frequency in top_trees:
    print(f"Tree Frequency: {frequency}")
    print(tree_str)
    print("\n")
```

```
Tree Frequency: 1
|--- application_date_indicator <= 1.00
|   |--- hud_median_family_income <= 72800.00
|   |   |--- purchaser_type_name <= 6.50
|   |   |   |--- class: 4.0
|   |   |--- purchaser_type_name >  6.50
|   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |--- hud_median_family_income <= 71100.00
|   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |--- hud_median_family_income <= 67850.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- msamd_name <= 4.50
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  2.50
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- msamd_name >  4.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- hud_median_family_income >  67850.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 59850.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 52150.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
17
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  52150.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
29
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
28
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
13
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
13
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
17
|   |   |   |   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |--- hud_median_family_income >  59850.00
|   |   |   |   |   |   |--- hud_median_family_income <= 62450.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |--- hud_median_family_income >  62450.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
18
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |--- hud_median_family_income >  71100.00
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.35
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.35
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_type_name >  1.50
```

```
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |--- class: 4.0
|   |--- hud_median_family_income >  72800.00
|   |   |--- msamd_name <= 7.50
|   |   |   |--- msamd_name <= 6.50
|   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 75
850.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  75
850.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
```

```
|   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |--- msamd_name >  7.50
|   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |---  class: 4.0
|   |   |   |--- purchaser_type_name >  6.50
```

```
|   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 0.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
```

```
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.
71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 6.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.
71
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.02
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.02
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 6.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 6.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 7.0
|   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
```

```
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|--- application_date_indicator >  1.00
|   |--- class: 5.0


Tree Frequency: 1
|--- hud_median_family_income <= 58650.00
|   |--- loan_purpose_name <= 1.50
|   |   |--- purchaser_type_name <= 6.50
|   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |--- class: 5.0
|   |   |   |--- loan_type_name >  0.50
|   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |--- class: 5.0
|   |   |--- purchaser_type_name >  6.50
|   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |--- msamd_name <= 12.00
|   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
20
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
16
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
```

```
11
|   |   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |--- msamd_name >  12.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
```

```
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 52150.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  52150.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |--- class: 4.0
|   |--- loan_purpose_name >  1.50
|   |   |--- application_date_indicator <= 1.00
|   |   |   |--- msamd_name <= 12.00
|   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
24
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
17
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.02
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.02
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.02
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.02
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- purchaser_type_name <= 5.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- purchaser_type_name >  5.50
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.00
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |--- msamd_name >  12.00
|   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |   |---- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |--- application_date_indicator >  1.00
|   |   |   |--- class: 5.0
|--- hud_median_family_income >  58650.00
|   |--- purchaser_type_name <= 6.50
|   |   |--- application_date_indicator <= 1.00
|   |   |   |--- class: 4.0
|   |   |--- application_date_indicator >  1.00
|   |   |   |--- class: 5.0
|   |--- purchaser_type_name >  6.50
|   |   |--- purchaser_type_name <= 7.50
|   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |--- hud_median_family_income <= 72800.00
|   |   |   |   |   |--- hud_median_family_income <= 71100.00
|   |   |   |   |   |   |--- msamd_name <= 10.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 63800.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 61
600.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  61
600.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
```

```
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  63800.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- msamd_name >  1.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  10.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |--- hud_median_family_income >  71100.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator <=
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator >
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
16
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
24
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
14
|   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |--- hud_median_family_income >  72800.00
|   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |--- hud_median_family_income <= 73450.00
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |--- hud_median_family_income >  73450.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator <=
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator >
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 75
850.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  75
850.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
                                                     |
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- application_date_indicator <=
1.00
|   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |--- application_date_indicator >
1.00
|   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |--- loan_purpose_name <= 1.00
|   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |--- loan_purpose_name >  1.00
|   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
```

```
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
18
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |--- hud_median_family_income <= 72800.00
|   |   |   |   |   |--- loan_amount_000s <= 0.40
|   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |--- hud_median_family_income <= 71100.00
|   |   |   |   |   |--- msamd_name <= 10.50
|   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- msamd_name >  10.50
|   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |--- hud_median_family_income >  71100.00
|   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |--- hud_median_family_income <= 71100.
00
|   |   |   |   |   |--- hud_median_family_income <= 62
450.00
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- hud_median_family_income >  62
450.00
|   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |--- hud_median_family_income >  71100.
00
|   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |--- truncated branch of depth
```

```
10
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.40
|   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |--- hud_median_family_income >  72800.00
|   |   |   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |--- msamd_name <= 9.00
|   |   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82
084.71
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82
084.71
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  9.00
|   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |--- purchaser_type_name >  7.50
|   |   |   |--- class: 4.0


Tree Frequency: 1
|--- purchaser_type_name <= 6.50
|   |--- application_date_indicator <= 1.00
|   |   |--- class: 4.0
|   |--- application_date_indicator >  1.00
|   |   |--- class: 5.0
|--- purchaser_type_name >  6.50
|   |--- application_date_indicator <= 1.00
|   |   |--- purchaser_type_name <= 7.50
|   |   |   |--- hud_median_family_income <= 72800.00
|   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- msamd_name <= 8.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |--- msamd_name >  8.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 67
700.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  67
700.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 60
500.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  60
500.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 52150.
00
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  52150.
00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
```

```
|   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 63950.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  63950.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 4.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  4.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 61350.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  61350.00
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  1.00
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name >  1.00
```

```
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 58700.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
21
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
24
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |--- hud_median_family_income >  58700.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 58700.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
16
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  58700.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |--- msamd_name <= 11.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 10.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |   |   |--- msamd_name >  10.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
25
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |--- msamd_name >  11.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 55900.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
17
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
14
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  55900.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |--- hud_median_family_income >  72800.00
|   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |--- hud_median_family_income <= 75850.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 73450.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  73450.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 73450.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  73450.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- hud_median_family_income >  75850.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 7.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82
084.71
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82
084.71
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
10
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.
71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 0.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.
71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82
084.71
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82
084.71
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |--- purchaser_type_name >  7.50
|   |   |   |--- class: 4.0
|   |--- application_date_indicator >  1.00
|   |   |--- class: 5.0
```

time: 2.24 s (started: 2024-04-13 11:18:17 +00:00)

In [68]:
```python
# Initialize a dictionary to store tree frequencies
tree_frequency = Counter()

# Loop through the trees and count their frequency
for i in range(len(rf_classifier.estimators_)):
    tree = rf_classifier.estimators_[i]
    tree_str = export_graphviz(tree, feature_names=df1_inputs_train.columns,
                                filled=True, max_depth=2, impurity=False, proportion=True)
    tree_frequency[tree_str] += 1

# Get the three most frequent trees
top_trees = tree_frequency.most_common(3)

# Plot the top three trees
for tree_str, frequency in top_trees:
    graph = graphviz.Source(tree_str)
    display(graph)
```







```
time: 356 ms (started: 2024-04-13 11:18:19 +00:00)
```

In [69]:
```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.graph_objects as go
from wordcloud import WordCloud
from collections import Counter
from scipy import stats
from sklearn.tree import plot_tree, export_graphviz
import graphviz
from IPython.display import display
from collections import Counter
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import time
import psutil

# Function to measure memory usage
def memory_usage():
    process = psutil.Process()
    return process.memory_info().rss / 1024 ** 2  # Memory usage in MB

# Start time
start_time = time.time()

# Initialize Random Forest classifier
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=45005)

# Train the Random Forest classifier
rf_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])

# Predictions
y_train_pred_rf = rf_classifier.predict(df1_inputs_train)
y_test_pred_rf = rf_classifier.predict(df1_inputs_test)

# End time
end_time = time.time()

# Time taken
execution_time = end_time - start_time

# Memory usage
memory_used = memory_usage()

# Accuracy
accuracy_train = accuracy_score(df1_output_train['action_taken_name'], y_train_pred_rf)
accuracy_test = accuracy_score(df1_output_test['action_taken_name'], y_test_pred_rf)

# Print time, memory usage, and accuracy
print("Time taken (seconds):", execution_time)
print("Memory used (MB):", memory_used)
print("Training Set Accuracy:", accuracy_train)
print("Testing Set Accuracy:", accuracy_test)

# Print feature importances
feature_importances = rf_classifier.feature_importances_
```

```python
feature_importance_df = pd.DataFrame({'Feature': df1_inputs_train.columns,
'Importance': feature_importances})
sorted_feature_importance_df = feature_importance_df.sort_values(by='Import
ance', ascending=False)
print("\nFeature Importances:")
print(sorted_feature_importance_df)

# Print confusion matrix and classification report for training set
print("\nTraining Set Confusion Matrix:")
print(confusion_matrix(df1_output_train['action_taken_name'], y_train_pred_
rf))
print("\nTraining Set Classification Report:")
print(classification_report(df1_output_train['action_taken_name'], y_train_
pred_rf))

# Print confusion matrix and classification report for testing set
print("\nTesting Set Confusion Matrix:")
print(confusion_matrix(df1_output_test['action_taken_name'], y_test_pred_r
f))
print("\nTesting Set Classification Report:")
print(classification_report(df1_output_test['action_taken_name'], y_test_pr
ed_rf))

# Export and visualize the first three decision trees
for i in range(3):
    tree = rf_classifier.estimators_[i]
    dot_data = export_graphviz(tree,
                               feature_names=df1_inputs_train.columns,
                               filled=True,
                               max_depth=2,
                               impurity=False,
                               proportion=True)
    graph = graphviz.Source(dot_data)
    display(graph)

# Initialize a dictionary to store tree frequencies and rules
tree_frequency = Counter()
tree_rules = {}

# Loop through the trees and count their frequency while storing rules
for i in range(len(rf_classifier.estimators_)):
    tree = rf_classifier.estimators_[i]
    tree_str = export_text(tree, feature_names=list(df1_inputs_train.column
s))
    tree_frequency[tree_str] += 1
    tree_rules[tree_str] = tree

# Get the three most frequent trees
top_trees = tree_frequency.most_common(3)

# Print the rules for the top three trees
for tree_str, frequency in top_trees:
    print(f"Tree Frequency: {frequency}")
    print(tree_str)
    print("\n")

# Initialize a dictionary to store tree frequencies
tree_frequency = Counter()

# Loop through the trees and count their frequency
for i in range(len(rf_classifier.estimators_)):
```

```python
    tree = rf_classifier.estimators_[i]
    tree_str = export_graphviz(tree, feature_names=df1_inputs_train.column
s,
                               filled=True, max_depth=2, impurity=False, pr
oportion=True)
    tree_frequency[tree_str] += 1

# Get the three most frequent trees
top_trees = tree_frequency.most_common(3)

# Plot the top three trees
for tree_str, frequency in top_trees:
    graph = graphviz.Source(tree_str)
    display(graph)
```

```
Time taken (seconds): 7.921905517578125
Memory used (MB): 591.46484375
Training Set Accuracy: 0.9820416666666667
Testing Set Accuracy: 0.9508333333333333
```

Feature Importances:

| | Feature | Importance |
|---|---|---|
| 6 | loan_amount_000s | 0.333727 |
| 2 | purchaser_type_name | 0.201634 |
| 0 | application_date_indicator | 0.195447 |
| 5 | hud_median_family_income | 0.120389 |
| 1 | msamd_name | 0.096998 |
| 3 | loan_type_name | 0.027990 |
| 4 | loan_purpose_name | 0.023816 |

Training Set Confusion Matrix:
```
[[   31     0    13     1    21     0     0     0]
 [    1    69    21     2    43     0     0     0]
 [    1     8  1825    26   289     0     0     1]
 [    0     3    51   229    75     0     0     0]
 [    2    10   260    27 44307     0     0     0]
 [    0     0     0     0     0   642     0     0]
 [    0     0     1     0     3     0    11     0]
 [    0     0     1     0     2     0     0    24]]
```

Training Set Classification Report:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.89 | 0.47 | 0.61 | 66 |
| 1 | 0.77 | 0.51 | 0.61 | 136 |
| 2 | 0.84 | 0.85 | 0.84 | 2150 |
| 3 | 0.80 | 0.64 | 0.71 | 358 |
| 4 | 0.99 | 0.99 | 0.99 | 44606 |
| 5 | 1.00 | 1.00 | 1.00 | 642 |
| 6 | 1.00 | 0.73 | 0.85 | 15 |
| 7 | 0.96 | 0.89 | 0.92 | 27 |
| | | | | |
| accuracy | | | 0.98 | 48000 |
| macro avg | 0.91 | 0.76 | 0.82 | 48000 |
| weighted avg | 0.98 | 0.98 | 0.98 | 48000 |

Testing Set Confusion Matrix:
```
[[    0     0     5     3    10     0     0     0]
 [    0     1     7     1    16     0     0     0]
 [    6    10   272    33   194     0     0     1]
 [    0     5    36    16    31     0     0     0]
 [    2     7   196    17 10984     0     3     0]
 [    0     0     0     0     0   134     0     0]
 [    0     0     1     0     1     0     0     0]
 [    0     0     2     0     2     0     1     3]]
```

Testing Set Classification Report:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.00 | 0.00 | 0.00 | 18 |
| 1 | 0.04 | 0.04 | 0.04 | 25 |
| 2 | 0.52 | 0.53 | 0.53 | 516 |
| 3 | 0.23 | 0.18 | 0.20 | 88 |
| 4 | 0.98 | 0.98 | 0.98 | 11209 |
| 5 | 1.00 | 1.00 | 1.00 | 134 |

| | | | | |
|---|---|---|---|---|
| 6 | 0.00 | 0.00 | 0.00 | 2 |
| 7 | 0.75 | 0.38 | 0.50 | 8 |
| | | | | |
| accuracy | | | 0.95 | 12000 |
| macro avg | 0.44 | 0.39 | 0.41 | 12000 |
| weighted avg | 0.95 | 0.95 | 0.95 | 12000 |

```
Tree Frequency: 1
|--- application_date_indicator <= 1.00
|   |--- hud_median_family_income <= 72800.00
|   |   |--- purchaser_type_name <= 6.50
|   |   |   |--- class: 4.0
|   |   |--- purchaser_type_name >  6.50
|   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |--- hud_median_family_income <= 71100.00
|   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |--- hud_median_family_income <= 67850.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- msamd_name <= 4.50
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  2.50
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- msamd_name >  4.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- hud_median_family_income >  67850.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 59850.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 52150.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
17
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  52150.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
29
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
28
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
13
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
13
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
17
|   |   |   |   |   |   |   |   |   |---  msamd_name >  9.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |--- hud_median_family_income >  59850.00
|   |   |   |   |   |   |   |--- hud_median_family_income <= 62450.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |--- hud_median_family_income >  62450.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
18
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |--- hud_median_family_income >  71100.00
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.35
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.35
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_type_name >  1.50
```

```
|   |   |   |   |   |   |   |   |---  loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |    |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |    |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |    |    |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |    |    |    |--- class: 2.0
|   |   |   |   |   |   |   |   |    |    |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |    |    |    |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |    |    |    |    |--- class: 1.0
|   |   |   |   |   |   |   |   |    |    |    |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |    |    |    |    |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |    |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |    |    |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |    |    |    |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |    |    |    |    |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |    |    |    |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |    |    |    |    |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |    |    |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |    |    |    |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |    |    |    |    |--- class: 4.0
|   |   |   |   |   |   |   |   |    |    |    |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |    |    |    |    |--- truncated branch of depth
10
|   |   |   |--- purchaser_type_name >  7.50
|   |   |   |    |--- class: 4.0
|   |--- hud_median_family_income >  72800.00
|   |   |--- msamd_name <= 7.50
|   |   |   |--- msamd_name <= 6.50
|   |   |   |    |--- purchaser_type_name <= 6.50
|   |   |   |    |    |--- class: 4.0
|   |   |   |    |--- purchaser_type_name >  6.50
|   |   |   |    |    |--- purchaser_type_name <= 7.50
|   |   |   |    |    |    |--- loan_purpose_name <= 1.50
|   |   |   |    |    |    |    |--- loan_amount_000s <= 0.00
|   |   |   |    |    |    |    |    |--- loan_purpose_name <= 0.50
|   |   |   |    |    |    |    |    |    |--- class: 4.0
|   |   |   |    |    |    |    |--- loan_purpose_name >  0.50
|   |   |   |    |    |    |    |    |--- class: 1.0
|   |   |   |    |    |    |--- loan_amount_000s >  0.00
|   |   |   |    |    |    |    |--- msamd_name <= 3.50
|   |   |   |    |    |    |    |    |--- loan_amount_000s <= 0.00
|   |   |   |    |    |    |    |    |    |--- loan_type_name <= 0.50
|   |   |   |    |    |    |    |    |    |    |--- truncated branch of depth
4
|   |   |   |    |    |    |    |    |    |--- loan_type_name >  0.50
|   |   |   |    |    |    |    |    |    |    |--- truncated branch of depth
3
|   |   |   |    |    |    |    |    |--- loan_amount_000s >  0.00
|   |   |   |    |    |    |    |    |    |--- class: 4.0
|   |   |   |    |    |    |    |--- msamd_name >  3.50
|   |   |   |    |    |    |    |    |--- loan_amount_000s <= 0.00
|   |   |   |    |    |    |    |    |    |--- loan_purpose_name <= 0.50
|   |   |   |    |    |    |    |    |    |    |--- class: 4.0
|   |   |   |    |    |    |    |    |    |--- loan_purpose_name >  0.50
|   |   |   |    |    |    |    |    |    |    |--- truncated branch of depth
2
|   |   |   |    |    |    |    |    |--- loan_amount_000s >  0.00
|   |   |   |    |    |    |    |    |    |--- class: 4.0
```

```
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 75
850.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  75
850.00
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
```

```
|   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |--- msamd_name >  7.50
|   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |--- class: 4.0
|   |   |   |--- purchaser_type_name >  6.50
```

```
|   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 0.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
```

```
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.
71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 6.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.
71
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.02
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.02
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 6.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 6.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 7.0
|   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
```

```
|   |   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|--- application_date_indicator >  1.00
|   |--- class: 5.0


Tree Frequency: 1
|--- hud_median_family_income <= 58650.00
|   |--- loan_purpose_name <= 1.50
|   |   |--- purchaser_type_name <= 6.50
|   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |--- class: 5.0
|   |   |   |--- loan_type_name >  0.50
|   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |--- class: 5.0
|   |   |--- purchaser_type_name >  6.50
|   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |--- msamd_name <= 12.00
|   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
20
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
16
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |    |--- truncated branch of depth
```

```
11
|   |   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 56
750.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  56
750.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |--- msamd_name >  12.00
|   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
```

```
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 52150.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  52150.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |--- class: 4.0
|   |--- loan_purpose_name >  1.50
|   |   |--- application_date_indicator <= 1.00
|   |   |   |--- msamd_name <= 12.00
|   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
24
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
17
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.02
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.02
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.02
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.02
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |--- purchaser_type_name <= 5.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- purchaser_type_name >  5.50
|   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.00
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |--- msamd_name >  12.00
|   |   |   |   |--- purchaser_type_name <= 6.50
|   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- purchaser_type_name >  6.50
|   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |   |---  loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |--- purchaser_type_name <= 7.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |--- purchaser_type_name >  7.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |--- application_date_indicator >  1.00
|   |   |   |--- class: 5.0
|--- hud_median_family_income >  58650.00
|   |--- purchaser_type_name <= 6.50
|   |   |--- application_date_indicator <= 1.00
|   |   |   |--- class: 4.0
|   |   |--- application_date_indicator >  1.00
|   |   |   |--- class: 5.0
|   |--- purchaser_type_name >  6.50
|   |   |--- purchaser_type_name <= 7.50
|   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |--- hud_median_family_income <= 72800.00
|   |   |   |   |   |--- hud_median_family_income <= 71100.00
|   |   |   |   |   |   |--- msamd_name <= 10.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 63800.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 61
600.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  61
600.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
```

```
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  63800.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- msamd_name >  1.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  10.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |--- hud_median_family_income >  71100.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator <=
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator >
1.00
|   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
16
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
24
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
14
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |--- hud_median_family_income >  72800.00
|   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |--- hud_median_family_income <= 73450.00
|   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |--- hud_median_family_income >  73450.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator <=
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator >
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 75
850.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  75
850.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
2
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator <=
1.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- application_date_indicator >
1.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
18
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.50
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.50
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |--- hud_median_family_income <= 72800.00
|   |   |   |   |   |--- loan_amount_000s <= 0.40
|   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 71100.00
|   |   |   |   |   |   |   |   |--- msamd_name <= 10.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- msamd_name >  10.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  71100.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 71100.
00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 62
450.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  62
450.00
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  71100.
00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
10
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |   |   |--- loan_amount_000s >  0.40
|   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |--- hud_median_family_income >  72800.00
|   |   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |   |--- application_date_indicator <= 1.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |--- application_date_indicator >  1.00
|   |   |   |   |   |   |   |   |--- class: 5.0
|   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |--- msamd_name <= 9.00
|   |   |   |   |   |   |--- msamd_name <= 7.50
|   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82
084.71
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82
084.71
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  9.00
|   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |--- purchaser_type_name >  7.50
|   |   |   |--- class: 4.0


Tree Frequency: 1
|--- purchaser_type_name <= 6.50
|   |--- application_date_indicator <= 1.00
|   |   |--- class: 4.0
|   |--- application_date_indicator >  1.00
|   |   |--- class: 5.0
|--- purchaser_type_name >  6.50
|   |--- application_date_indicator <= 1.00
|   |   |--- purchaser_type_name <= 7.50
|   |   |   |--- hud_median_family_income <= 72800.00
|   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- msamd_name <= 8.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |---  class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |--- msamd_name >  8.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 67
700.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  67
700.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
6
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 60
500.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  60
500.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 52150.
00
|   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  52150.
00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
```

```
|   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 63950.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  63950.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 4.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  4.50
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 61350.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  61350.00
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  1.00
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- msamd_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 1.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |   |--- msamd_name >  1.00
```

```
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- msamd_name <= 3.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- msamd_name >  3.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |--- msamd_name <= 9.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- hud_median_family_income <= 58700.00
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
21
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
24
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |--- hud_median_family_income >  58700.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 58700.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
16
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  58700.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
5
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- msamd_name >  9.50
|   |   |   |   |   |   |   |--- msamd_name <= 11.50
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |--- msamd_name <= 10.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
19
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |   |   |--- msamd_name >  10.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
10
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
25
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
15
|   |   |   |   |   |   |   |--- msamd_name >  11.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 55900.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
17
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
14
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  55900.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
9
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
11
|   |   |   |       |--- hud_median_family_income >  72800.00
|   |   |   |       |   |--- msamd_name <= 7.50
|   |   |   |       |   |   |--- hud_median_family_income <= 75850.00
|   |   |   |       |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |       |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |       |   |   |   |   |   |--- hud_median_family_income <= 73450.00
|   |   |   |       |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |       |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |       |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |       |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |       |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |       |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |       |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  73450.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 73450.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  73450.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |--- msamd_name <= 6.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 3.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- msamd_name >  6.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
```

```
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- hud_median_family_income >  75850.00
|   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |--- msamd_name >  7.50
|   |   |   |   |   |--- loan_purpose_name <= 0.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |--- class: 4.0
```

```
|   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
4
|   |   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |--- loan_purpose_name >  0.50
|   |   |   |   |   |   |--- loan_type_name <= 0.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 7.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82
084.71
|   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82
084.71
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
```

```
10
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.
71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 0.0
|   .|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.
71
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
12
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |--- loan_type_name >  0.50
|   |   |   |   |   |   |   |--- loan_purpose_name <= 1.50
|   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 1.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name <= 1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_type_name >  1.50
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
3
|   |   |   |   |   |   |   |   |--- hud_median_family_income >  82084.71
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
8
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |--- loan_purpose_name >  1.50
|   |   |   |   |   |   |   |   |--- loan_type_name <= 2.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.01
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
7
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.01
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |--- loan_type_name >  2.00
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s <= 0.00
|   |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
2
|   |   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |   |--- class: 4.0
|   |   |   |   |   |   |   |   |   |--- loan_amount_000s >  0.00
|   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income <= 82
084.71
```

```
|  |   |   |   |   |   |   |   |   |   |   |--- class: 2.0
|  |   |   |   |   |   |   |   |   |   |   |--- hud_median_family_income >  82
084.71
|  |   |   |   |   |   |   |   |   |   |   |--- truncated branch of depth
5
|  |   |--- purchaser_type_name >  7.50
|  |   |   |--- class: 4.0
|  |--- application_date_indicator >  1.00
|  |   |--- class: 5.0
```



time: 11.9 s (started: 2024-04-13 11:18:20 +00:00)

```python
In [70]: # Random Forest
rf_start_time = time.time()
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=4500
5)
rf_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])
rf_training_time = time.time() - rf_start_time
rf_memory_used = memory_usage()
y_train_pred_rf = rf_classifier.predict(df1_inputs_train)
y_test_pred_rf = rf_classifier.predict(df1_inputs_test)
rf_accuracy = accuracy_score(df1_output_test, y_test_pred_rf)

# Cross-validation for Random Forest
rf_cv_start_time = time.time()
rf_cv = RandomForestClassifier(n_estimators=100, random_state=45005)
cv_scores_rf = cross_val_score(rf_cv, df1_inputs, df1_output.values.ravel
(), cv=20)
rf_cv_time = time.time() - rf_cv_start_time
rf_cv_accuracy = np.mean(cv_scores_rf)

print("Random Forest:")
print(f"  - Training Time (s): {rf_training_time}")
print(f"  - Memory Used (MB): {rf_memory_used}")
print(f"  - Single Split Accuracy: {rf_accuracy}")
print(f"  - Cross Validation Accuracy: {rf_cv_accuracy}")
print()
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:7
00: UserWarning: The least populated class in y has only 17 members, which
is less than n_splits=20.
  warnings.warn(

Random Forest:
  - Training Time (s): 5.206611394882202
  - Memory Used (MB): 628.19140625
  - Single Split Accuracy: 0.9508333333333333
  - Cross Validation Accuracy: 0.94575

time: 1min 12s (started: 2024-04-13 11:18:32 +00:00)
```

```python
import pandas as pd
import numpy as np
import time
import psutil
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Function to measure memory usage
def memory_usage():
    process = psutil.Process()
    return process.memory_info().rss / 1024 ** 2  # Memory usage in MB

# Data preprocessing and splitting
# Assuming you have your data loaded into cars_inputs and cars_output
df1_inputs_train, df1_inputs_test, df1_output_train, df1_output_test = trai
n_test_split(df1_inputs, df1_output, test_size=0.2, random_state=42)

# Initialize lists to store results
models = []
training_times = []
memory_used = []
single_split_accuracies = []
cross_validation_accuracies = []

# Decision Tree
dt_start_time = time.time()
dt_model = DecisionTreeClassifier(criterion='gini', random_state=45011, max
_depth=3)
dt_model.fit(df1_inputs_train, df1_output_train)
dt_training_time = time.time() - dt_start_time
dt_memory_used = memory_usage()
dt_pred = dt_model.predict(df1_inputs_test)
dt_accuracy = accuracy_score(df1_output_test, dt_pred)

# Cross-validation for Decision Tree
dtc_cv_start_time = time.time()
dtc_cv = DecisionTreeClassifier(criterion='gini', random_state=45011)
cv_scores_dtc = cross_val_score(dtc_cv, df1_inputs, df1_output.values.ravel
(), cv=20)
dtc_cv_time = time.time() - dtc_cv_start_time
dtc_cv_accuracy = np.mean(cv_scores_dtc)

# Append Decision Tree results to lists
models.append('Decision Tree')
training_times.append(dt_training_time)
memory_used.append(dt_memory_used)
single_split_accuracies.append(dt_accuracy)
cross_validation_accuracies.append(dtc_cv_accuracy)

# Random Forest
rf_start_time = time.time()
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=4500
5)
rf_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])
rf_training_time = time.time() - rf_start_time
rf_memory_used = memory_usage()
y_train_pred_rf = rf_classifier.predict(df1_inputs_train)
y_test_pred_rf = rf_classifier.predict(df1_inputs_test)
```

```python
rf_accuracy_train = accuracy_score(df1_output_train['action_taken_name'], y
_train_pred_rf)
rf_accuracy_test = accuracy_score(df1_output_test['action_taken_name'], y_t
est_pred_rf)

# Append Random Forest results to lists
models.append('Random Forest')
training_times.append(rf_training_time)
memory_used.append(rf_memory_used)
single_split_accuracies.append(rf_accuracy_test)  # Using test accuracy as
we already calculated it
cross_validation_accuracies.append(np.nan)  # Cross-validation accuracy not
calculated here

# Create DataFrame
results_df = pd.DataFrame({
    'Model': models,
    'Training Time (s)': training_times,
    'Memory Used (MB)': memory_used,
    'Single Split Accuracy': single_split_accuracies,
    'Cross Validation Accuracy': cross_validation_accuracies
})

# Print DataFrame
print(results_df)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:7
00: UserWarning: The least populated class in y has only 17 members, which
is less than n_splits=20.
  warnings.warn(

          Model  Training Time (s)  Memory Used (MB)  Single Split Accurac
y  \
0  Decision Tree           0.047592           683.375                0.94541
7
1  Random Forest           2.410989           683.375                0.95083
3

   Cross Validation Accuracy
0                     0.9431
1                        NaN
time: 5.11 s (started: 2024-04-13 11:19:44 +00:00)
```

# XGBoost

```python
In [72]: dtrain = xgb.DMatrix(df1_inputs_train, label=df1_output_train['action_taken
_name'])
dtest = xgb.DMatrix(df1_inputs_test, label=df1_output_test['action_taken_na
me'])
```

```
time: 32.2 ms (started: 2024-04-13 11:19:50 +00:00)
```

In [73]:
```python
# Define XGBoost parameters
params = {
    'objective': 'multi:softmax',  # For multi-class classification
    'num_class': len(df1_output_train['action_taken_name'].unique()),  # Nu
mber of unique classes in the output
    'max_depth': 5,
    'learning_rate': 0.1,
    'n_estimators': 1000,
    'eval_metric': 'merror'  # Use 'merror' for multiclass classification e
rror
}

# Initialize the XGBoost classifier
xgb_classifier = xgb.XGBClassifier(**params)

# Train the classifier
xgb_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])
```

Out[73]:
```
                              XGBClassifier

XGBClassifier(base_score=None, booster=None, callbacks=None,
              colsample_bylevel=None, colsample_bynode=None,
              colsample_bytree=None, device=None, early_stopping_rounds
=None,
              enable_categorical=False, eval_metric='merror',
              feature_types=None, gamma=None, grow_policy=None,
              importance_type=None, interaction_constraints=None,
              learning_rate=0.1, max_bin=None, max_cat_threshold=None,
              max_cat_to_onehot=None, max_delta_step=None, max_depth=5,
              max_leaves=None, min_child_weight=None, missing=nan,
```

```
time: 32.2 s (started: 2024-04-13 11:19:50 +00:00)
```

In [74]:
```python
import xgboost as xgb
from sklearn.metrics import log_loss

# Assuming you have already trained the XGBoost classifier
# xgb_classifier = xgb.XGBClassifier(**params)
# xgb_classifier.fit(df1_inputs_train, df1_output_train['action_taken_nam
e'])

# Predict probabilities for each class
y_pred_proba_xgb = xgb_classifier.predict_proba(df1_inputs_test)

# Calculate entropy
entropy_xgb = log_loss(df1_output_test, y_pred_proba_xgb)

# Calculate Gini impurity
gini_impurity_xgb = 1 - (y_pred_proba_xgb ** 2).sum(axis=1).mean()

print("Entropy for XGBoost:", entropy_xgb)
print("Gini Impurity for XGBoost:", gini_impurity_xgb)
```

```
Entropy for XGBoost: 0.11704169720141437
Gini Impurity for XGBoost: 0.04694253206253052
time: 2.33 s (started: 2024-04-13 11:20:22 +00:00)
```

```python
In [75]:  # Print feature importances
          feature_importances = xgb_classifier.feature_importances_
          feature_importance_df = pd.DataFrame({'Feature': df1_inputs_train.columns,
          'Importance': feature_importances})
          sorted_feature_importance_df = feature_importance_df.sort_values(by='Import
          ance', ascending=False)
          print(sorted_feature_importance_df)
```

```
                        Feature  Importance
0  application_date_indicator    0.765484
2          purchaser_type_name    0.199578
5     hud_median_family_income    0.014517
1                   msamd_name    0.012562
3               loan_type_name    0.003441
4            loan_purpose_name    0.003015
6             loan_amount_000s    0.001402
time: 13.1 ms (started: 2024-04-13 11:20:24 +00:00)
```

```python
In [76]:  # Make predictions on the training set
          y_train_pred = xgb_classifier.predict(df1_inputs_train)

          # Print classification report for training set
          print("Training Set Classification Report:")
          print(classification_report(df1_output_train['action_taken_name'], y_train_
          pred))
```

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

           0       0.88      0.21      0.34        66
           1       0.84      0.20      0.32       136
           2       0.73      0.74      0.74      2150
           3       0.75      0.38      0.50       358
           4       0.98      0.99      0.99     44606
           5       1.00      1.00      1.00       642
           6       0.88      0.47      0.61        15
           7       0.96      0.89      0.92        27

    accuracy                           0.97     48000
   macro avg       0.88      0.61      0.68     48000
weighted avg       0.97      0.97      0.97     48000

time: 6.57 s (started: 2024-04-13 11:20:24 +00:00)
```

```python
In [77]:  # Print confusion matrix for training set
          print("Training Set Confusion Matrix:")
          print(confusion_matrix(df1_output_train['action_taken_name'], y_train_pre
          d))
```

```
Training Set Confusion Matrix:
[[   14     1    16     2    33     0     0     0]
 [    0    27    33     1    74     0     1     0]
 [    0     2  1599    28   520     0     0     1]
 [    0     0   107   136   115     0     0     0]
 [    2     2   425    15 44162     0     0     0]
 [    0     0     0     0     0   642     0     0]
 [    0     0     2     0     6     0     7     0]
 [    0     0     1     0     2     0     0    24]]
time: 15.1 ms (started: 2024-04-13 11:20:31 +00:00)
```

In [78]:
```python
# Make predictions on the test set
y_pred = xgb_classifier.predict(df1_inputs_test)

# Evaluate the model
print(classification_report(df1_output_test['action_taken_name'], y_pred))
```

```
              precision    recall  f1-score   support

           0       0.00      0.00      0.00        18
           1       0.12      0.04      0.06        25
           2       0.55      0.59      0.57       516
           3       0.24      0.11      0.15        88
           4       0.98      0.98      0.98     11209
           5       1.00      1.00      1.00       134
           6       0.00      0.00      0.00         2
           7       0.50      0.12      0.20         8

    accuracy                           0.95     12000
   macro avg       0.42      0.36      0.37     12000
weighted avg       0.95      0.95      0.95     12000

time: 1.42 s (started: 2024-04-13 11:20:31 +00:00)

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:
1344: UndefinedMetricWarning: Precision and F-score are ill-defined and bei
ng set to 0.0 in labels with no predicted samples. Use `zero_division` para
meter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:
1344: UndefinedMetricWarning: Precision and F-score are ill-defined and bei
ng set to 0.0 in labels with no predicted samples. Use `zero_division` para
meter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:
1344: UndefinedMetricWarning: Precision and F-score are ill-defined and bei
ng set to 0.0 in labels with no predicted samples. Use `zero_division` para
meter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

In [79]:
```python
# Make predictions on the test set
y_test_pred = xgb_classifier.predict(df1_inputs_test)

# Print confusion matrix for testing set
print("\nTesting Set Confusion Matrix:")
print(confusion_matrix(df1_output_test['action_taken_name'], y_test_pred))
```

```
Testing Set Confusion Matrix:
[[    0     0     8     2     8     0     0     0]
 [    0     1     7     0    17     0     0     0]
 [    4     3   305    20   183     0     0     1]
 [    0     2    41    10    35     0     0     0]
 [    1     2   188    10 11008     0     0     0]
 [    0     0     0     0     0   134     0     0]
 [    0     0     1     0     1     0     0     0]
 [    0     0     5     0     2     0     0     1]]
time: 1.43 s (started: 2024-04-13 11:20:32 +00:00)
```

In [80]:
```python
# XGBoost
xgb_start_time = time.time()

# Define XGBoost parameters
params = {
    'objective': 'multi:softmax',  # For multi-class classification
    'num_class': len(df1_output_train['action_taken_name'].unique()),  # Nu
mber of unique classes in the output
    'max_depth': 5,
    'learning_rate': 0.1,
    'n_estimators': 1000,
    'eval_metric': 'merror'  # Use 'merror' for multiclass classification e
rror
}

# Initialize the XGBoost classifier
xgb_classifier = xgb.XGBClassifier(**params)

# Train the classifier
xgb_classifier.fit(df1_inputs_train, df1_output_train['action_taken_name'])

xgb_training_time = time.time() - xgb_start_time
xgb_memory_used = memory_usage()
y_train_pred_xgb = xgb_classifier.predict(df1_inputs_train)
y_test_pred_xgb = xgb_classifier.predict(df1_inputs_test)
xgb_accuracy = accuracy_score(df1_output_test, y_test_pred_xgb)

# Cross-validation for XGBoost
xgb_cv_start_time = time.time()
xgb_cv = xgb.XGBClassifier(**params)
cv_scores_xgb = cross_val_score(xgb_cv, df1_inputs, df1_output.values.ravel
(), cv=20)
xgb_cv_time = time.time() - xgb_cv_start_time
xgb_cv_accuracy = np.mean(cv_scores_xgb)

print("XGBoost:")
print(f"  - Training Time (s): {xgb_training_time}")
print(f"  - Memory Used (MB): {xgb_memory_used}")
print(f"  - Single Split Accuracy: {xgb_accuracy}")
print(f"  - Cross Validation Accuracy: {xgb_cv_accuracy}")
print()
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_split.py:7
00: UserWarning: The least populated class in y has only 17 members, which
is less than n_splits=20.
  warnings.warn(

XGBoost:
  - Training Time (s): 34.26365828514099
  - Memory Used (MB): 748.02734375
  - Single Split Accuracy: 0.9549166666666666
  - Cross Validation Accuracy: 0.9495833333333336

time: 14min 3s (started: 2024-04-13 11:20:34 +00:00)
```

In [81]:
```python
# Plot the first tree in the XGBoost model
xgb.plot_tree(xgb_classifier, num_trees=0, rankdir='LR')  # num_trees=0 plo
ts the first tree
plt.rcParams['figure.figsize'] = [30, 30]  # Adjust the figure size if need
ed
plt.show()
```



time: 1.72 s (started: 2024-04-13 11:34:37 +00:00)

In [81]:

time: 1.73 s (started: 2024-04-13 11:34:37 +00:00)