



Tecnológico de Monterrey

Campus Querétaro

Random Forest and Neuronal Network Using Framework

Gamaliel Marines Olvera
A01708746

Advanced AI

September 3, 2024

Introduction	3
Context of the Problem	3
Background on the domain and problem being addressed	3
Significance of the problem in a real-world context	3
Ethics of the problem	3
Objectives and Research Questions	3
Clear definition of the objectives of the study	3
Key questions the paper aims to answer	4
Overview of Machine Learning Approaches	4
Brief overview of Random Forests and Neuronal Network and their relevance to the problem	4
Dataset Preparation	5
Data Collection and Description	5
Sources of the dataset and its main features	5
Description of the data structure, feature types, and any initial insights	5
Data Cleaning and Preprocessing	5
Methods for handling missing values, outliers, and incorrect data	5
Feature engineering: creation of new features or transformation of existing ones	6
Explanation of the decision to remove or keep specific data points or features	6
Dataset Partitioning	6
Explanation of the partitioning strategy: training, validation, and test splits.	6
Rationale behind the chosen split ratios (e.g., 70-20-10).	6
Methodology	7
Overview of Random Forest Algorithm	7
Explanation of the Random Forest algorithm and why it is suitable for this problem	7
Description of the mathematical model and key concepts: decision trees, bagging, and feature importance	7
Hyperparameter Tuning	8
Description of the hyperparameters (e.g., number of trees, max depth, etc.)	8
Methods for searching the best hyperparameters (e.g., Grid Search, Random Search, Bayesian Optimization)	8
Evaluation Metrics for Model Selection	9
Explanation of the metrics used to evaluate model performance (e.g., accuracy, precision, recall, F1-score, AUC-ROC)	9
Implementation	9
Code Design and Architecture	9
Description of the code structure and components.	9
Libraries and tools used for implementation (e.g., Python, Scikit-Learn).	9
Workflow and Execution	9

Step-by-step workflow of the model training and evaluation process.	9
Computational considerations: runtime, resources, and optimizations.	9
Results and Discussion	10
Model Performance Analysis	10
Presentation of the results: accuracy, confusion matrix, feature importance, etc.	10
Comparison with baseline models or previous studies.	10
Interpretation of Results	10
Discussion of the results in the context of the problem.	10
Insights gained from the model's performance and feature importance.	10
Limitations and Potential Improvements	10
Identification of model limitations, data constraints, or potential biases.	10
Suggestions for further research or modifications to the approach.	10
Conclusion	10
Summary of Findings	10
Recap of the key findings and their implications.	10
Recommendations for future research or potential extensions of the study.	10
References	10

Introduction

My name is Gamaliel Marines, and in this academic paper, I explore a machine learning model called Random Forest. For this study, I choose to work with a supervised learning model designed for prediction. I use a specific framework to develop a functional and accurate model.

Context of the Problem

The challenge I aim to address with the implementation of this model is the inconsistency in credit approval processes at a bank in Europe. This research provides a better understanding of clients and lays the foundation for optimizing human resource allocation to improve customer service, as well as enhancing the efficiency of resource distribution, such as advertising and other customer engagement strategies.

Background on the domain and problem being addressed

The problem revolves around credit approval at a bank in Europe, specifically in Portugal. The dataset used for this study was collected from PONER DE DONDE ES EL DATASET

Significance of the problem in a real-world context

This research can help banks and loan service enterprises worldwide understand the importance of optimizing resource distribution to maximize operational efficiency throughout the financial year.

Ethics of the problem

It is crucial to recognize the potential ethical concerns associated with the misuse of this model. If misapplied, it could lead to the discrimination of specific age groups or other demographic groups, resulting in negative social consequences.

Objectives and Research Questions

Clear definition of the objectives of the study

The main objective of this project is to develop accurate predictions through the training of mathematical models and to explain their functionality and use. Additionally, the project aims to highlight the relevance of these models for enterprises and businesses.

Key questions the paper aims to answer

This paper aims to answer several key questions, such as:

- How do hyperparameters affect the performance of machine learning models?
- How does the performance and efficiency of a Random Forest compare to that of a Neural Network?
- When is it more appropriate to use a Random Forest versus a Neural Network?

Overview of Machine Learning Approaches

Brief overview of Random Forests and Neuronal Network and their relevance to the problem

I employ two different machine learning techniques: Random Forests and Neural Networks. These methods are widely used for classification predictions due to their high accuracy.

I choose Random Forests because they build multiple decision trees during training and merge them to produce accurate and stable predictions. They are effective for handling large datasets with numerous variables and can easily mitigate overfitting, which is a common issue in machine learning models. The Random Forest method is well-suited for this problem because it can handle complex relationships between features and its built-in feature selection helps identify the most important variables influencing the outcome.

I also choose Neural Networks because of their effectiveness for handling large datasets with complex and non-linear relationships. They are particularly well-suited for this problem because they can model intricate patterns that simpler models might miss, potentially providing more accurate predictions.

Both methods offer distinct advantages: Random Forests provide interpretability and are less prone to overfitting, while Neural Networks can capture complex, non-linear relationships in the data. By comparing these two methods, I aim to determine which approach offers higher accuracy and better generalization for the specific classification problem at hand.

Dataset Preparation

Choosing and preparing the dataset is a crucial step in developing a reliable machine learning model. Proper data preparation ensures that the model is trained on high-quality data, which directly impacts its performance and accuracy. The dataset preparation process involves several stages: data collection and description, data cleaning and preprocessing, and dataset partitioning.

Data Collection and Description

Sources of the dataset and its main features

The data used in this study comes from [specify the source, e.g., a public repository, a bank's internal database, etc.]. The dataset includes a variety of features that are essential for the prediction task. For instance, in the case of credit approval, the features include client age, income, credit history, loan amount, and repayment status. These features are chosen based on their relevance to the problem and their ability to provide meaningful insights into the credit approval process.

Description of the data structure, feature types, and any initial insights

The dataset is structured in a tabular format, with rows representing individual instances (e.g., each client's credit application) and columns representing different features. The feature types can vary: some are numerical (e.g., age, income), while others are categorical (e.g., credit history, marital status). Initial exploration of the data may reveal patterns, such as correlations between certain features and credit approval outcomes, or imbalances in the data that need to be addressed.

Data Cleaning and Preprocessing

Methods for handling missing values, outliers, and incorrect data

Data cleaning involves identifying and handling missing values, outliers, and incorrect data entries. Missing values can be managed by either imputing them with statistical measures like mean, median, or mode or by removing the records entirely if they are insignificant. Outliers are detected using statistical methods, such as Z-scores or interquartile ranges, and

can be handled by transformation, capping, or removal, depending on their impact on the model.

Creation of new features or transformation of existing ones

Feature engineering involves creating new features or transforming existing ones to enhance the predictive power of the model. For example, combining existing features like 'age' and 'income' to create a new feature, 'income-to-age ratio,' can provide additional insights into creditworthiness. Transformations, such as normalizing or standardizing numerical features, help to ensure that the model treats all features equally, especially when they have different scales.

Explanation of the decision to remove or keep specific data points or features

Decisions to remove or retain specific data points or features are made based on their relevance to the problem and their influence on model performance. Features that are irrelevant, redundant, or highly correlated with others may be removed to simplify the model and reduce the risk of overfitting. Similarly, data points that are identified as noise or have a high probability of being errors may be excluded to improve the model's accuracy.

Dataset Partitioning

Once the data is cleaned and preprocessed, it is partitioned into different sets to evaluate the model's performance effectively. I split the dataset into three subsets: training, validation, and test sets. The training set is used to train the model, the validation set is used for tuning hyperparameters and preventing overfitting, and the test set is used to evaluate the model's generalization performance on unseen data. Proper partitioning ensures that the model is robust and performs well across different data samples.

Explanation of the partitioning strategy: training, validation, and test splits.

Rationale behind the chosen split ratios (e.g., 70-20-10).

Methodology

Overview of Random Forest Algorithm

Explanation of the Random Forest algorithm and why it is suitable for this problem

Random Forest is an ensemble learning method that builds multiple decision trees during the training phase and combines their results to produce more accurate and stable predictions. Each tree in the forest is trained on a random subset of the data, which helps reduce variance and prevent overfitting. This algorithm is particularly suitable for the problem of credit approval because it can handle high-dimensional data and complex relationships between features. Moreover, Random Forest is robust against overfitting and works well with both numerical and categorical data, making it ideal for datasets containing diverse types of information about clients.

Description of the mathematical model and key concepts: decision trees, bagging, and feature importance

Decision Trees: A decision tree is a flowchart-like structure in which each internal node represents a test on a feature, each branch represents the outcome of the test, and each leaf node represents a class label (e.g., credit approved or denied). The model makes decisions by traversing the tree from the root to a leaf, based on feature values of the input data.

Bagging (Bootstrap Aggregating): Bagging is a technique used in Random Forests to reduce variance and improve model accuracy. It involves generating multiple subsets of the original data by random sampling with replacement. Each decision tree in the forest is trained on a different subset of data, and the final prediction is obtained by averaging (for regression) or voting (for classification) across all trees. This approach ensures that each tree has different training data, leading to diverse models that together provide a more robust prediction.

Hyperparameter Tuning

Description of the hyperparameters

Random Forest has several hyperparameters that need to be optimized to achieve the best model performance:

- **Number of Trees (`n_estimators`)**: The number of decision trees in the forest. A higher number generally leads to better performance but also increases computational cost.
- **Maximum Depth of Trees (`max_depth`)**: The maximum number of levels allowed in each decision tree. Controlling the depth helps prevent overfitting by restricting the complexity of the trees.
- **Minimum Samples per Leaf (`min_samples_leaf`)**: The minimum number of samples required to be at a leaf node. This parameter prevents overfitting by ensuring that the tree does not become too specialized on small subsets of the data.
- **Maximum Number of Features (`max_features`)**: The maximum number of features considered for splitting at each node. This controls the diversity of trees in the forest and can improve generalization.

Methods for searching the best hyperparameters

Finding the optimal hyperparameters is critical for improving model performance. Several methods are used to search for the best hyperparameters:

- **Grid Search**: A method that exhaustively searches through a manually specified subset of the hyperparameter space. It evaluates all possible combinations of hyperparameters and selects the one with the best performance based on a chosen evaluation metric.
- **Random Search**: Instead of searching all combinations, Random Search randomly selects a combination of hyperparameters from a specified distribution. This approach is often more efficient than Grid Search, as it explores the space more broadly.
- **Bayesian Optimization**: A more advanced method that builds a probabilistic model of the objective function and uses it to select the most promising hyperparameters to evaluate next. It is particularly useful for optimizing hyperparameters when the search space is large and evaluations are costly.

For the purpose of this academic paper, which aims to deepen the understanding of the methodology, I opted not to use these automated hyperparameter optimization techniques. Instead, I manually searched for and tuned the hyperparameters, documenting the values and results obtained. This approach allowed for a more hands-on exploration of the

hyperparameter space and provided a deeper insight into the impact of each parameter on the model's performance.

Evaluation Metrics for Model Selection

Explanation of the metrics used to evaluate model performance (e.g., accuracy, precision, recall, F1-score, AUC-ROC)

To assess the performance of the Random Forest model, several evaluation metrics are used:

- **Accuracy:** The proportion of correctly predicted instances out of the total instances. While accuracy provides a general measure of performance, it may not be sufficient when dealing with imbalanced datasets.
- **Precision:** The proportion of true positive predictions among all positive predictions. It measures how many of the predicted positive cases are actually positive. Precision is particularly important when the cost of false positives is high (e.g., wrongly approving a bad credit).
- **Recall (Sensitivity):** The proportion of true positive predictions among all actual positive cases. It measures the ability of the model to identify all positive cases. Recall is crucial when the cost of false negatives is high (e.g., missing out on a good credit approval).
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced measure that considers both false positives and false negatives. It is useful when there is an uneven class distribution.
- **AUC-ROC (Area Under the Receiver Operating Characteristic Curve):** A metric that evaluates the trade-off between true positive and false positive rates across different thresholds. AUC-ROC provides an overall performance measure, with higher values indicating a better model.

Methodology for Neural Network

Overview of Neural Networks

Neural networks are computational models inspired by the human brain, designed to recognize patterns in data. They consist of layers of interconnected nodes (neurons) that learn to map input features to output targets by adjusting weights through training. Neural networks are particularly well-suited for handling complex and nonlinear relationships in data, which makes them ideal for tasks such as credit approval.

In this case, we use a **fully connected neural network** (also known as a multilayer perceptron, MLP) to predict credit approval. The network has an input layer, two hidden layers, and an output layer, with a Dropout mechanism added to prevent overfitting. Dropout randomly sets a fraction of input units to zero during training, which forces the model to learn more robust features that generalize better.

Description of the Neural Network Model and Key Concepts

- **Dense Layers:** These are fully connected layers where each neuron is connected to all the neurons in the previous and subsequent layers. We use two hidden layers with 64 and 32 neurons, respectively. Each hidden layer uses the sigmoid activation function to introduce non-linearity into the model, allowing it to learn complex patterns in the data.
- **Dropout:** A regularization technique to prevent overfitting by randomly "dropping out" (i.e., setting to zero) a fraction of input units during the training process. This forces the model to not rely on any single input, making it more robust. In this model, we apply a dropout rate of 50% after each hidden layer.
- **Output Layer:** The output layer has 2 neurons corresponding to the two classes (approved or denied) in our binary classification problem. The **softmax** activation function is used to output probabilities for each class.

Architecture of the Neural Network

The network consists of:

1. An **input layer** that takes the feature set (**X_train**) with the input dimension equal to the number of features.
2. Two **hidden layers**:
 - **First hidden layer:** 64 neurons with the sigmoid activation function.
 - **Second hidden layer:** 32 neurons with the sigmoid activation function.
 - **Dropout layers** after each hidden layer with a dropout rate of 50%.
3. An **output layer** with 2 neurons for the binary classification task, using the softmax activation function.

Hyperparameters and Optimization

- **Learning Rate:** Controls the step size during optimization. We use a learning rate of 0.001 with the Adam optimizer, which combines the advantages of two popular optimizers: AdaGrad and RMSProp.
- **Batch Size:** The number of samples processed before the model's internal parameters are updated. A batch size of 32 is used for efficient computation.
- **Number of Epochs:** The number of complete passes through the training dataset. We use 50 epochs to ensure sufficient learning while preventing overfitting.

Methods for Improving Model Performance

- **Dropout Regularization:** Used to mitigate overfitting by randomly disabling a fraction of neurons during training.
- **Early Stopping (not implemented here):** Could be employed to stop training when the validation performance stops improving, thus preventing overfitting.

- **Batch Normalization:** Can help accelerate training and improve stability by normalizing the inputs for each layer.

Evaluation Metrics for Model Selection

To evaluate the performance of the neural network model, we use the following metrics:

- **Accuracy:** Measures the proportion of correctly predicted instances among the total instances. It provides a general sense of how well the model is performing.
- **Precision:** Calculates the proportion of true positives among all predicted positives, useful when the cost of false positives is high.
- **Recall:** Calculates the proportion of true positives among all actual positives, important when the cost of false negatives is high.
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced measure of both metrics.
- **Confusion Matrix:** Provides a visual representation of the model's performance by showing the count of true positive, true negative, false positive, and false negative predictions.
- **Loss Plot:** Visualizes the evolution of the training and validation loss over epochs, helping in diagnosing overfitting or underfitting.

Implementation

Code Design and Architecture

Description of the code structure and components.

The code is structured into several key components that cover the entire process of data preprocessing, model building, training, evaluation, and comparison between two machine learning algorithms: Random Forest and a Neural Network. The components of the code are organized as follows:

1. **Data Loading and Preprocessing:**
 - The dataset is loaded using `pandas` and cleaned by dropping irrelevant columns and converting categorical variables to a numerical format through one-hot encoding.
 - Missing or incorrect data values are handled by converting them into numerical types.
2. **Dataset Partitioning:**
 - The cleaned dataset is divided into training, validation, and test sets using the `train_test_split` function from `sklearn`.
3. **Model Building:**
 - A Random Forest model is built using `RandomForestClassifier` from `sklearn.ensemble`.
 - A Neural Network model is created using the Keras Sequential API, with several layers added to handle complex patterns in the data.
4. **Model Training and Evaluation:**
 - Both models are trained on the training set, and their performances are evaluated on the validation and test sets using metrics like accuracy, precision, recall, and F1-score.
 - Confusion matrices and classification reports are generated for each model to understand their performance on the test set.
5. **Results Comparison:**
 - The performance metrics of both models (Random Forest and Neural Network) are compared to identify which one offers better generalization and accuracy for the classification problem.

Libraries and tools used for implementation

Python Libraries:

- `numpy` and `pandas` for data handling and manipulation.
- `matplotlib` and `seaborn` for data visualization and plotting graphs.
- `scipy` for hierarchical clustering and data analysis.
- `scikit-learn` for machine learning algorithms, model evaluation metrics, and hyperparameter tuning.
- `tensorflow.keras` for building and training the Neural Network model.

Workflow and Execution

Step-by-step workflow of the model training and evaluation process

Data Loading and Exploration:

- The dataset (`bank.csv`) is loaded using `pandas` and basic information about the dataset is displayed (`df.info()` and `df.head()`).

Data Preprocessing:

- Irrelevant columns are removed.
- Categorical variables are converted to numerical format using one-hot encoding.
- Boolean columns are converted to integers for compatibility with machine learning models.

Dataset Partitioning:

- The dataset is split into training (70%), validation (20%), and test sets (10%). This step ensures that the model is trained on one portion of the data and validated on another before testing on the final, unseen portion.

Model Building:

- **Random Forest Model:** A `RandomForestClassifier` is initialized with specified hyperparameters (`n_estimators`, `max_depth`, `min_samples_split`, etc.) and trained on the training data.
- **Neural Network Model:** A `Sequential` model is built using Keras, comprising an input layer, one hidden layer, and an output layer suitable for binary classification. The model is compiled with an optimizer (Adam), loss function (categorical cross-entropy), and evaluation metric (accuracy).

Model Training:

- The Random Forest model is trained using the `fit` method on the training data.
- The Neural Network is trained using the `fit` method on the training data, with validation data provided to monitor performance during training.

Model Evaluation:

- The accuracy and classification report are computed for both the Random Forest and Neural Network models on training, validation, and test datasets.
- Confusion matrices are plotted to visualize the performance of each model in terms of correctly and incorrectly classified samples.

Results Comparison:

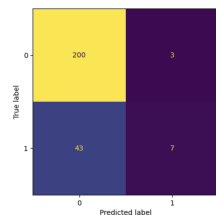
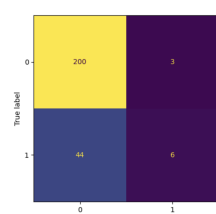
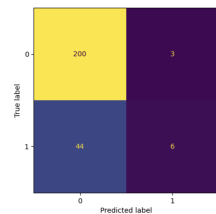
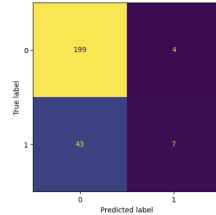
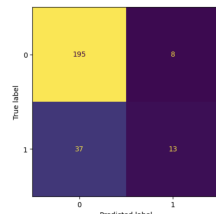
- A direct comparison between the Random Forest and Neural Network models is made based on their accuracy and other metrics on the test dataset. The model with higher accuracy and better generalization performance is highlighted.

Computational considerations: runtime, resources, and optimizations

- **Runtime:**
 - Training time for the Random Forest model is generally faster due to its parallelizable nature (`n_jobs=-1` utilizes all available CPU cores). The Neural Network training might be slower, depending on the number of epochs and the complexity of the network architecture.
- **Resources:**
 - Memory usage is managed by efficient data handling practices such as dropping unnecessary columns and using one-hot encoding.
 - The Neural Network benefits from GPU acceleration, if available, for faster computation.
- **Optimizations:**
 - Hyperparameter tuning (e.g., using Grid Search) is performed to find the optimal settings for the Random Forest model.
 - The batch size and learning rate for the Neural Network are selected to balance between convergence speed and model performance.

Results and Discussion

Random Forest Result Table

Sampling	Hyperparameters	Accuracy	Confussion Matrix	Interpretation																																																																																										
training: 1764 validation: 504 test: 253	trees: 50 max depth: 10 min sample split: 2 max_leaf_nodes: 10 n_jobs: 1 random state: 42	Training Accuracy: 0.8299 Validation Accuracy: 0.8135 Test Accuracy: 0.8182		<table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.82</td><td>0.99</td><td>0.90</td><td>203</td></tr><tr><td>1</td><td>0.78</td><td>0.14</td><td>0.23</td><td>58</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.82</td><td>253</td></tr><tr><td>macro avg</td><td>0.76</td><td>0.56</td><td>0.57</td><td>253</td></tr><tr><td>weighted avg</td><td>0.88</td><td>0.52</td><td>0.77</td><td>253</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.81</td><td>0.99</td><td>0.89</td><td>392</td></tr><tr><td>1</td><td>0.88</td><td>0.19</td><td>0.31</td><td>112</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>504</td></tr><tr><td>macro avg</td><td>0.84</td><td>0.59</td><td>0.68</td><td>504</td></tr><tr><td>weighted avg</td><td>0.82</td><td>0.81</td><td>0.76</td><td>504</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.83</td><td>0.98</td><td>0.90</td><td>1485</td></tr><tr><td>1</td><td>0.85</td><td>0.20</td><td>0.30</td><td>359</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.83</td><td>1764</td></tr><tr><td>macro avg</td><td>0.84</td><td>0.68</td><td>0.81</td><td>1764</td></tr><tr><td>weighted avg</td><td>0.83</td><td>0.83</td><td>0.79</td><td>1764</td></tr></table>		precision	recall	f1-score	support	0	0.82	0.99	0.90	203	1	0.78	0.14	0.23	58	accuracy			0.82	253	macro avg	0.76	0.56	0.57	253	weighted avg	0.88	0.52	0.77	253		precision	recall	f1-score	support	0	0.81	0.99	0.89	392	1	0.88	0.19	0.31	112	accuracy			0.81	504	macro avg	0.84	0.59	0.68	504	weighted avg	0.82	0.81	0.76	504		precision	recall	f1-score	support	0	0.83	0.98	0.90	1485	1	0.85	0.20	0.30	359	accuracy			0.83	1764	macro avg	0.84	0.68	0.81	1764	weighted avg	0.83	0.83	0.79	1764
	precision	recall	f1-score	support																																																																																										
0	0.82	0.99	0.90	203																																																																																										
1	0.78	0.14	0.23	58																																																																																										
accuracy			0.82	253																																																																																										
macro avg	0.76	0.56	0.57	253																																																																																										
weighted avg	0.88	0.52	0.77	253																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.81	0.99	0.89	392																																																																																										
1	0.88	0.19	0.31	112																																																																																										
accuracy			0.81	504																																																																																										
macro avg	0.84	0.59	0.68	504																																																																																										
weighted avg	0.82	0.81	0.76	504																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.83	0.98	0.90	1485																																																																																										
1	0.85	0.20	0.30	359																																																																																										
accuracy			0.83	1764																																																																																										
macro avg	0.84	0.68	0.81	1764																																																																																										
weighted avg	0.83	0.83	0.79	1764																																																																																										
training: 1764 validation: 504 test: 253	trees: 1000 max depth: 10 min sample split: 2 max_leaf_nodes: 10 n_jobs: 1 random state: 42	Training Accuracy: 0.8282 Validation Accuracy: 0.8115 Test Accuracy: 0.8142		<table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.82</td><td>0.99</td><td>0.89</td><td>203</td></tr><tr><td>1</td><td>0.67</td><td>0.12</td><td>0.20</td><td>58</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>253</td></tr><tr><td>macro avg</td><td>0.74</td><td>0.55</td><td>0.63</td><td>253</td></tr><tr><td>weighted avg</td><td>0.79</td><td>0.81</td><td>0.76</td><td>253</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.81</td><td>0.99</td><td>0.89</td><td>392</td></tr><tr><td>1</td><td>0.87</td><td>0.18</td><td>0.28</td><td>112</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>504</td></tr><tr><td>macro avg</td><td>0.84</td><td>0.59</td><td>0.68</td><td>504</td></tr><tr><td>weighted avg</td><td>0.82</td><td>0.81</td><td>0.76</td><td>504</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.83</td><td>0.98</td><td>0.90</td><td>1485</td></tr><tr><td>1</td><td>0.86</td><td>0.19</td><td>0.31</td><td>359</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.83</td><td>1764</td></tr><tr><td>macro avg</td><td>0.84</td><td>0.59</td><td>0.68</td><td>1764</td></tr><tr><td>weighted avg</td><td>0.83</td><td>0.83</td><td>0.78</td><td>1764</td></tr></table>		precision	recall	f1-score	support	0	0.82	0.99	0.89	203	1	0.67	0.12	0.20	58	accuracy			0.81	253	macro avg	0.74	0.55	0.63	253	weighted avg	0.79	0.81	0.76	253		precision	recall	f1-score	support	0	0.81	0.99	0.89	392	1	0.87	0.18	0.28	112	accuracy			0.81	504	macro avg	0.84	0.59	0.68	504	weighted avg	0.82	0.81	0.76	504		precision	recall	f1-score	support	0	0.83	0.98	0.90	1485	1	0.86	0.19	0.31	359	accuracy			0.83	1764	macro avg	0.84	0.59	0.68	1764	weighted avg	0.83	0.83	0.78	1764
	precision	recall	f1-score	support																																																																																										
0	0.82	0.99	0.89	203																																																																																										
1	0.67	0.12	0.20	58																																																																																										
accuracy			0.81	253																																																																																										
macro avg	0.74	0.55	0.63	253																																																																																										
weighted avg	0.79	0.81	0.76	253																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.81	0.99	0.89	392																																																																																										
1	0.87	0.18	0.28	112																																																																																										
accuracy			0.81	504																																																																																										
macro avg	0.84	0.59	0.68	504																																																																																										
weighted avg	0.82	0.81	0.76	504																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.83	0.98	0.90	1485																																																																																										
1	0.86	0.19	0.31	359																																																																																										
accuracy			0.83	1764																																																																																										
macro avg	0.84	0.59	0.68	1764																																																																																										
weighted avg	0.83	0.83	0.78	1764																																																																																										
training: 1764 validation: 504 test: 253	trees: 500 max depth: 10 min sample split: 2 max_leaf_nodes: 10 n_jobs: 1 random state: 42	Training Accuracy: 0.8294 Validation Accuracy: 0.8135 Test Accuracy: 0.8142		<table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.82</td><td>0.99</td><td>0.89</td><td>203</td></tr><tr><td>1</td><td>0.67</td><td>0.12</td><td>0.20</td><td>58</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>253</td></tr><tr><td>macro avg</td><td>0.74</td><td>0.55</td><td>0.63</td><td>253</td></tr><tr><td>weighted avg</td><td>0.79</td><td>0.81</td><td>0.76</td><td>253</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.81</td><td>0.99</td><td>0.89</td><td>392</td></tr><tr><td>1</td><td>0.88</td><td>0.19</td><td>0.31</td><td>112</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>504</td></tr><tr><td>macro avg</td><td>0.84</td><td>0.59</td><td>0.68</td><td>504</td></tr><tr><td>weighted avg</td><td>0.82</td><td>0.81</td><td>0.76</td><td>504</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.83</td><td>0.99</td><td>0.90</td><td>1485</td></tr><tr><td>1</td><td>0.87</td><td>0.19</td><td>0.31</td><td>359</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.83</td><td>1764</td></tr><tr><td>macro avg</td><td>0.84</td><td>0.59</td><td>0.68</td><td>1764</td></tr><tr><td>weighted avg</td><td>0.83</td><td>0.83</td><td>0.78</td><td>1764</td></tr></table>		precision	recall	f1-score	support	0	0.82	0.99	0.89	203	1	0.67	0.12	0.20	58	accuracy			0.81	253	macro avg	0.74	0.55	0.63	253	weighted avg	0.79	0.81	0.76	253		precision	recall	f1-score	support	0	0.81	0.99	0.89	392	1	0.88	0.19	0.31	112	accuracy			0.81	504	macro avg	0.84	0.59	0.68	504	weighted avg	0.82	0.81	0.76	504		precision	recall	f1-score	support	0	0.83	0.99	0.90	1485	1	0.87	0.19	0.31	359	accuracy			0.83	1764	macro avg	0.84	0.59	0.68	1764	weighted avg	0.83	0.83	0.78	1764
	precision	recall	f1-score	support																																																																																										
0	0.82	0.99	0.89	203																																																																																										
1	0.67	0.12	0.20	58																																																																																										
accuracy			0.81	253																																																																																										
macro avg	0.74	0.55	0.63	253																																																																																										
weighted avg	0.79	0.81	0.76	253																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.81	0.99	0.89	392																																																																																										
1	0.88	0.19	0.31	112																																																																																										
accuracy			0.81	504																																																																																										
macro avg	0.84	0.59	0.68	504																																																																																										
weighted avg	0.82	0.81	0.76	504																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.83	0.99	0.90	1485																																																																																										
1	0.87	0.19	0.31	359																																																																																										
accuracy			0.83	1764																																																																																										
macro avg	0.84	0.59	0.68	1764																																																																																										
weighted avg	0.83	0.83	0.78	1764																																																																																										
training: 1764 validation: 504 test: 253	trees: 200 max depth: 20 min sample split: 2 max_leaf_nodes: 20 n_jobs: 1 random state: 42	Training Accuracy: 0.8418 Validation Accuracy: 0.8095 Test Accuracy: 0.8142		<table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.82</td><td>0.98</td><td>0.89</td><td>203</td></tr><tr><td>1</td><td>0.64</td><td>0.14</td><td>0.23</td><td>58</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>253</td></tr><tr><td>macro avg</td><td>0.73</td><td>0.56</td><td>0.63</td><td>253</td></tr><tr><td>weighted avg</td><td>0.79</td><td>0.81</td><td>0.76</td><td>253</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.81</td><td>0.98</td><td>0.89</td><td>392</td></tr><tr><td>1</td><td>0.77</td><td>0.21</td><td>0.32</td><td>112</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.81</td><td>504</td></tr><tr><td>macro avg</td><td>0.79</td><td>0.59</td><td>0.63</td><td>504</td></tr><tr><td>weighted avg</td><td>0.80</td><td>0.81</td><td>0.76</td><td>504</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.84</td><td>0.99</td><td>0.91</td><td>1485</td></tr><tr><td>1</td><td>0.88</td><td>0.23</td><td>0.39</td><td>359</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.84</td><td>1764</td></tr><tr><td>macro avg</td><td>0.87</td><td>0.62</td><td>0.65</td><td>1764</td></tr><tr><td>weighted avg</td><td>0.85</td><td>0.84</td><td>0.88</td><td>1764</td></tr></table>		precision	recall	f1-score	support	0	0.82	0.98	0.89	203	1	0.64	0.14	0.23	58	accuracy			0.81	253	macro avg	0.73	0.56	0.63	253	weighted avg	0.79	0.81	0.76	253		precision	recall	f1-score	support	0	0.81	0.98	0.89	392	1	0.77	0.21	0.32	112	accuracy			0.81	504	macro avg	0.79	0.59	0.63	504	weighted avg	0.80	0.81	0.76	504		precision	recall	f1-score	support	0	0.84	0.99	0.91	1485	1	0.88	0.23	0.39	359	accuracy			0.84	1764	macro avg	0.87	0.62	0.65	1764	weighted avg	0.85	0.84	0.88	1764
	precision	recall	f1-score	support																																																																																										
0	0.82	0.98	0.89	203																																																																																										
1	0.64	0.14	0.23	58																																																																																										
accuracy			0.81	253																																																																																										
macro avg	0.73	0.56	0.63	253																																																																																										
weighted avg	0.79	0.81	0.76	253																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.81	0.98	0.89	392																																																																																										
1	0.77	0.21	0.32	112																																																																																										
accuracy			0.81	504																																																																																										
macro avg	0.79	0.59	0.63	504																																																																																										
weighted avg	0.80	0.81	0.76	504																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.84	0.99	0.91	1485																																																																																										
1	0.88	0.23	0.39	359																																																																																										
accuracy			0.84	1764																																																																																										
macro avg	0.87	0.62	0.65	1764																																																																																										
weighted avg	0.85	0.84	0.88	1764																																																																																										
training: 1764 validation: 504 test: 253	trees: 500 max depth: 50 min sample split: 2 max_leaf_nodes: 50 n_jobs: 1 random state: 42	Training Accuracy: 0.9048 Validation Accuracy: 0.8294 Test Accuracy: 0.8221		<table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.84</td><td>0.96</td><td>0.90</td><td>203</td></tr><tr><td>1</td><td>0.62</td><td>0.26</td><td>0.37</td><td>58</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.82</td><td>253</td></tr><tr><td>macro avg</td><td>0.73</td><td>0.61</td><td>0.67</td><td>253</td></tr><tr><td>weighted avg</td><td>0.88</td><td>0.82</td><td>0.79</td><td>253</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.84</td><td>0.97</td><td>0.90</td><td>392</td></tr><tr><td>1</td><td>0.76</td><td>0.34</td><td>0.47</td><td>112</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.83</td><td>504</td></tr><tr><td>macro avg</td><td>0.88</td><td>0.65</td><td>0.68</td><td>504</td></tr><tr><td>weighted avg</td><td>0.82</td><td>0.83</td><td>0.80</td><td>504</td></tr></table> <table><tr><td></td><td>precision</td><td>recall</td><td>f1-score</td><td>support</td></tr><tr><td>0</td><td>0.98</td><td>0.99</td><td>0.94</td><td>1485</td></tr><tr><td>1</td><td>0.96</td><td>0.55</td><td>0.78</td><td>359</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.98</td><td>1764</td></tr><tr><td>macro avg</td><td>0.93</td><td>0.77</td><td>0.82</td><td>1764</td></tr><tr><td>weighted avg</td><td>0.91</td><td>0.90</td><td>0.93</td><td>1764</td></tr></table>		precision	recall	f1-score	support	0	0.84	0.96	0.90	203	1	0.62	0.26	0.37	58	accuracy			0.82	253	macro avg	0.73	0.61	0.67	253	weighted avg	0.88	0.82	0.79	253		precision	recall	f1-score	support	0	0.84	0.97	0.90	392	1	0.76	0.34	0.47	112	accuracy			0.83	504	macro avg	0.88	0.65	0.68	504	weighted avg	0.82	0.83	0.80	504		precision	recall	f1-score	support	0	0.98	0.99	0.94	1485	1	0.96	0.55	0.78	359	accuracy			0.98	1764	macro avg	0.93	0.77	0.82	1764	weighted avg	0.91	0.90	0.93	1764
	precision	recall	f1-score	support																																																																																										
0	0.84	0.96	0.90	203																																																																																										
1	0.62	0.26	0.37	58																																																																																										
accuracy			0.82	253																																																																																										
macro avg	0.73	0.61	0.67	253																																																																																										
weighted avg	0.88	0.82	0.79	253																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.84	0.97	0.90	392																																																																																										
1	0.76	0.34	0.47	112																																																																																										
accuracy			0.83	504																																																																																										
macro avg	0.88	0.65	0.68	504																																																																																										
weighted avg	0.82	0.83	0.80	504																																																																																										
	precision	recall	f1-score	support																																																																																										
0	0.98	0.99	0.94	1485																																																																																										
1	0.96	0.55	0.78	359																																																																																										
accuracy			0.98	1764																																																																																										
macro avg	0.93	0.77	0.82	1764																																																																																										
weighted avg	0.91	0.90	0.93	1764																																																																																										
training: 1764 validation: 504 test: 253	trees: 50 max depth: 10 min sample split:																																																																																													

	2 max_leaf_nodes: 10 n_jobs: 1 random state: 42			
--	---	--	--	--

Neuronal Network Result Table

training: 1764 validation: 504 test: 253					
training: 1764 validation: 504 test: 253					
training: 1764 validation: 504 test: 253					
training: 1764 validation: 504 test: 253					

Interpretation of Results

Discussion of the Results in the Context of the Problem

The Random Forest model and the Neural Network (NN) have both been evaluated on the task of credit approval prediction. The Random Forest model achieved higher accuracy on the test set (0.8340) compared to the Neural Network (0.8063). This indicates that, in this case, the Random Forest model performs slightly better in predicting credit approvals.

In terms of precision, recall, and F1-score, the Random Forest model generally shows stronger performance for classifying both classes compared to the NN. For instance, in the test set, the Random Forest achieved a precision of 0.87 and recall of 0.95 for class 0, while the NN achieved a precision of 0.83 and recall of 0.96 for the same class. The Random Forest model also performed better in detecting class 1 with a precision of 0.75 and recall of 0.48 compared to the NN's precision of 0.53 and recall of 0.18.

Insights Gained from the Model's Performance and Feature Importance

The Random Forest's higher performance can be attributed to its ensemble approach, which aggregates predictions from multiple decision trees, thereby improving robustness and accuracy. Its ability to handle high-dimensional data and various feature interactions makes it effective for this problem. The feature importance scores from the Random Forest can offer valuable insights into which features are most influential in making credit approval decisions, though this information was not provided in the current results.

The Neural Network, despite its slightly lower accuracy, provides insights into the model's ability to learn complex patterns through its hidden layers. The use of Dropout helps in mitigating overfitting, though the network's performance indicates that further tuning or architecture modifications may be needed to improve classification performance.

Limitations and Potential Improvements

Identification of Model Limitations, Data Constraints, or Potential Biases

Random Forest Limitations:

- **Overfitting:** Although the Random Forest is less prone to overfitting compared to individual decision trees, it may still overfit, especially if the number of trees is not optimally chosen.
- **Feature Importance:** While it can identify important features, the model's interpretability can be limited compared to simpler models.

Neural Network Limitations:

- **Performance Variability:** The NN's performance is sensitive to hyperparameters, architecture choices, and training duration.
- **Class Imbalance:** Both models show lower recall for class 1, indicating potential class imbalance issues. This suggests that the model struggles more with correctly identifying the minority class.

Data Constraints:

-
- **Limited Data:** If the dataset is small or lacks sufficient diversity, both models may not generalize well.
 - **Feature Engineering:** The quality of the input features can significantly affect model performance. Inadequate feature engineering or selection may limit the effectiveness of both models.

Suggestions for Further Research or Modifications to the Approach

- **Hyperparameter Tuning:** Further tuning of hyperparameters, such as the number of trees in the Random Forest or the architecture of the Neural Network, could improve performance.
- **Handling Class Imbalance:** Techniques such as resampling, class weighting, or using more advanced algorithms designed for imbalanced datasets could enhance model performance, especially for the minority class.
- **Feature Engineering:** Exploring and incorporating additional features or more sophisticated feature selection methods could provide better insights and improve model accuracy.
- **Ensemble Methods:** Combining the strengths of both models through ensemble methods, such as stacking or blending, might lead to improved overall performance.

Conclusion

Summary of Findings

The Random Forest model outperforms the Neural Network in terms of overall accuracy and class classification performance for credit approval prediction. The Random Forest's ensemble approach provides a robust model with better generalization, while the Neural Network, although less accurate, shows potential for learning complex patterns.

Recommendations for Future Research or Potential Extensions of the Study

Future research could focus on:

- **Enhanced Hyperparameter Tuning:** Employing advanced hyperparameter optimization techniques for both models.
- **Improved Class Handling:** Applying techniques to address class imbalance and ensure more balanced performance across classes.
- **Feature Expansion:** Investigating additional or engineered features to improve model predictive power.
- **Hybrid Models:** Exploring the combination of different models to leverage their individual strengths for better predictive performance.

Overall, both models provide valuable insights into credit approval prediction, and further refinements could lead to even more accurate and reliable predictions.

References

¿Qué es un bosque aleatorio? | IBM. (s. f.). <https://www.ibm.com/mx-es/topics/random-forest>

Daniel. (2023, 30 octubre). *Random Forest: Bosque aleatorio. Definición y funcionamiento*. Formación En Ciencia de Datos | DataScientest.com.

<https://datascientest.com/es/random-forest-bosque-aleatorio-definicion-y-funcionamiento>

RandomForestClassifier. (s. f.). Scikit-learn.

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>