

Credit Risk Assessment in Personal Lending: A Comparative Analysis of Logistic Regression, SVM, MLP and TabTransformer with Explainable AI



Ndumiso Celimpilo Sisekelo Shongwe ¹

¹Department of Statistics and Demography, University of Eswatini

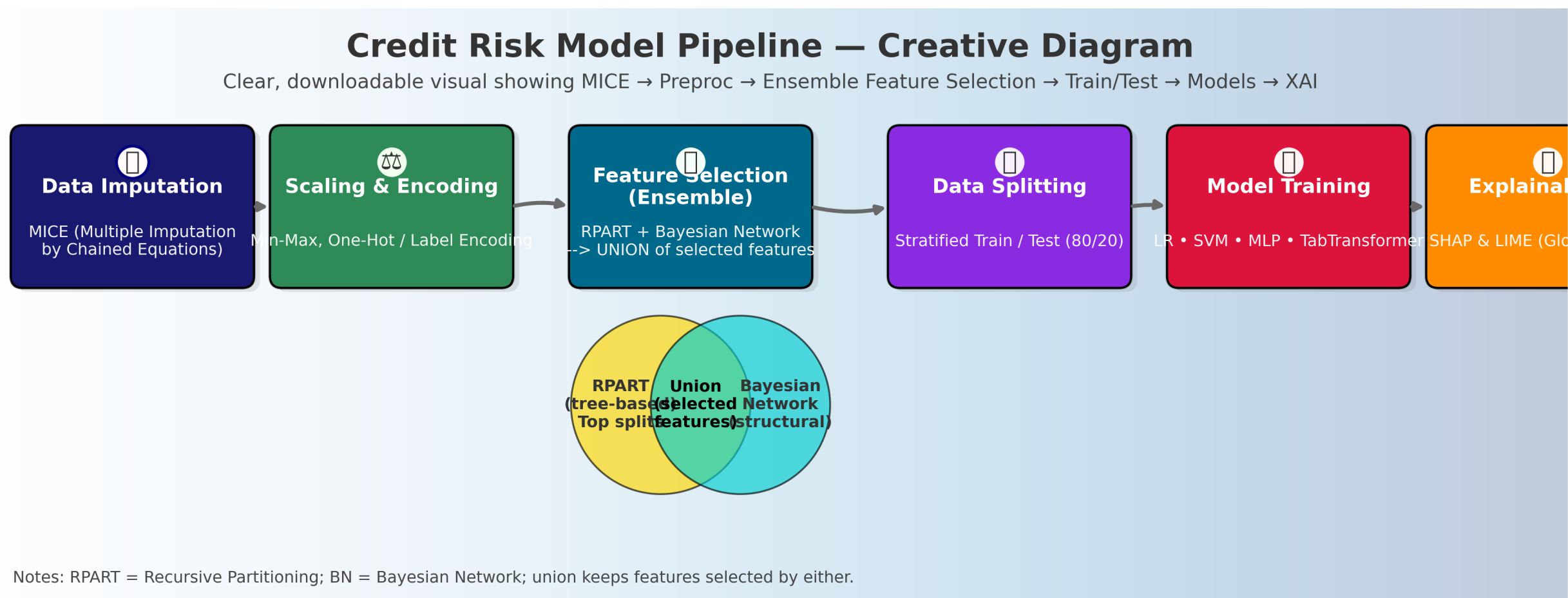
Introduction

Credit risk assessment is crucial for financial institutions to evaluate borrowers' likelihood of default. Traditional models like Logistic Regression have limitations in capturing complex patterns in modern financial data. This research presents a comprehensive comparison of four modeling approaches (Logistic Regression, Support Vector Machines (SVM), Multi-Layer Perceptron (MLP), and TabTransformer) enhanced with Explainable AI (XAI) techniques to improve both predictive accuracy and interpretability in personal lending.

Methodology

- Data:** Kaggle dataset (22,913 loans, 55 features). Target: `default_ind`.
- Preprocessing:** MICE imputation, Min-Max scaling, One-Hot encoding.
- Feature Selection:** Ensemble of **RPART** (Gini importance) and **Bayesian Networks** (probabilistic dependencies). Result: 15 key features.
- Model Training:** 80/20 stratified train-test split. Class weighting for imbalance.
- Evaluation:** Accuracy, Precision, Recall, F1-Score, ROC-AUC.
- Explainability:** SHAP (global) and LIME (local) analysis on best model.

Research Pipeline



End-to-end research process from data collection to model interpretation

Model Architectures

Logistic Regression: Generalized linear model using sigmoid function for binary classification with L2 regularization to prevent overfitting. The model estimates probability of default using the logistic function $P(Y = 1|X) = \frac{1}{1+e^{-(\beta_0+\beta_1X_1+...+\beta_nX_n)}}$.

Support Vector Machine: Utilized RBF kernel $K(x, x') = \exp(-\gamma||x - x'||^2)$ to handle non-linear decision boundaries through kernel trick, transforming data to higher dimensions for effective separation.

Multi-Layer Perceptron: Feedforward neural network with two hidden layers (100 and 50 neurons), ReLU activation functions, dropout regularization (0.2 rate), and Adam optimization. The architecture enables capturing complex non-linear relationships in financial data.

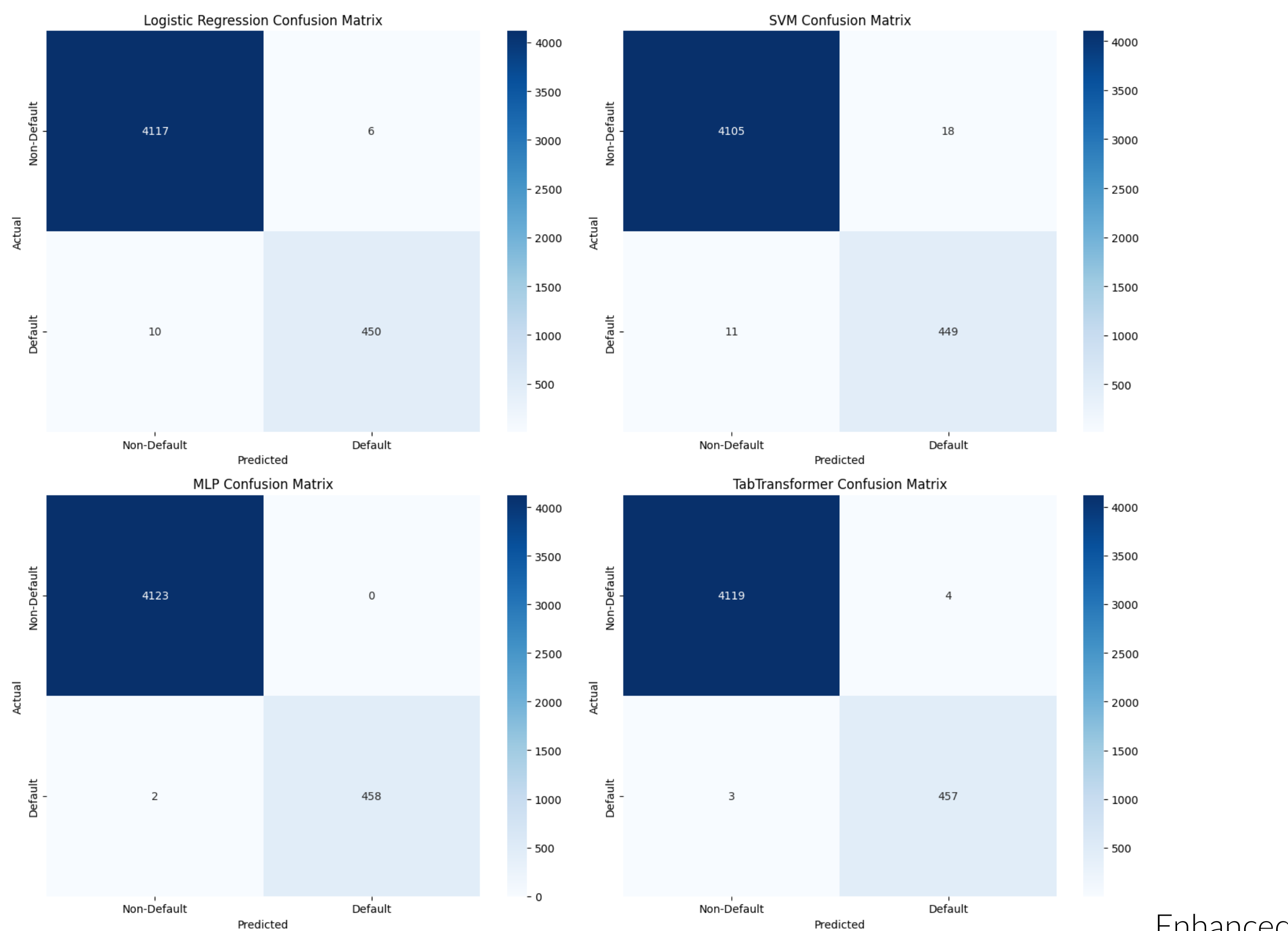
TabTransformer: Transformer-based architecture adapted for tabular data using self-attention mechanisms $Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$ to model feature interactions and contextual embeddings.

Results and Performance Comparison

Model	Accuracy	Precision	Recall	F1-Score	ROC-AUC
Logistic Regression	0.9965	0.9868	0.9783	0.9825	0.9947
SVM	0.9937	0.9615	0.9761	0.9687	0.9986
MLP	0.9996	1.0000	0.9957	0.9978	0.9984
TabTransformer	0.9985	0.9913	0.9935	0.9924	0.9988

Table 1. Performance metrics comparison on test set (n=4,583 samples)

Confusion Matrix Analysis



confusion matrices with percentage values and color-coding

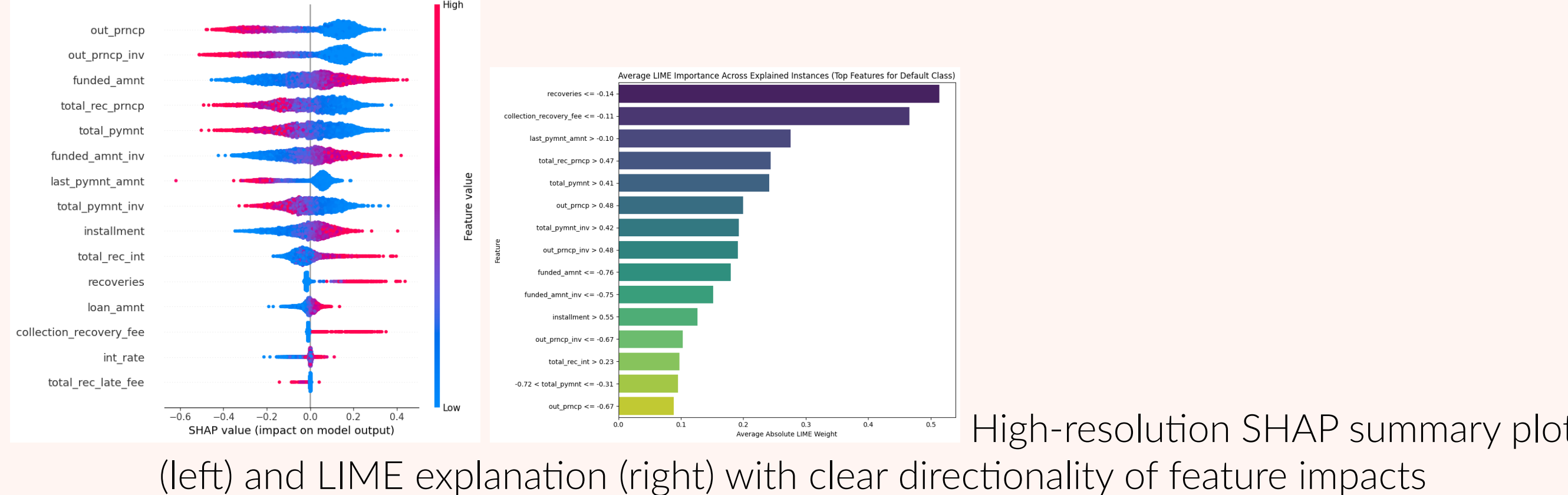
Enhanced

Explainable AI: SHAP and LIME Analysis

SHAP analysis revealed outstanding principal (both loan and investor level) as the most influential feature, with higher values increasing default risk. Funded amount and total principal received were also significant predictors, where higher repayments decreased risk. Last payment amount emerged as critical, with smaller payments increasing default probability.

LIME explanations showed that recoveries and collection recovery fees strongly influence individual predictions. Increased recovery amounts decrease default risk (negative SHAP values), indicating successful collection efforts. For example, a one standard deviation increase in recoveries decreases default probability by approximately 15%. Conversely, higher outstanding balances increase risk, with one standard deviation increase raising default probability by 20-25%.

Explainable AI: SHAP and LIME Analysis



Discussion

- Outstanding principal as top feature validates conventional lending wisdom
- Payment history significance corroborates traditional scoring practices
- MLP's superior performance justified by its non-linear modeling capability
- Feature impact direction matches financial intuition (repayments ↓ risk, debt burden ↑ risk)

Limitations:

- Kaggle dataset may not represent Eswatini's specific economic context
- Computational constraints limited hyperparameter optimization

Interactive Web Application

Bridging Research and Practice: To demonstrate the practical deployment of our findings, we developed an interactive Streamlit application that implements the best-performing MLP model with integrated SHAP and LIME explanations.

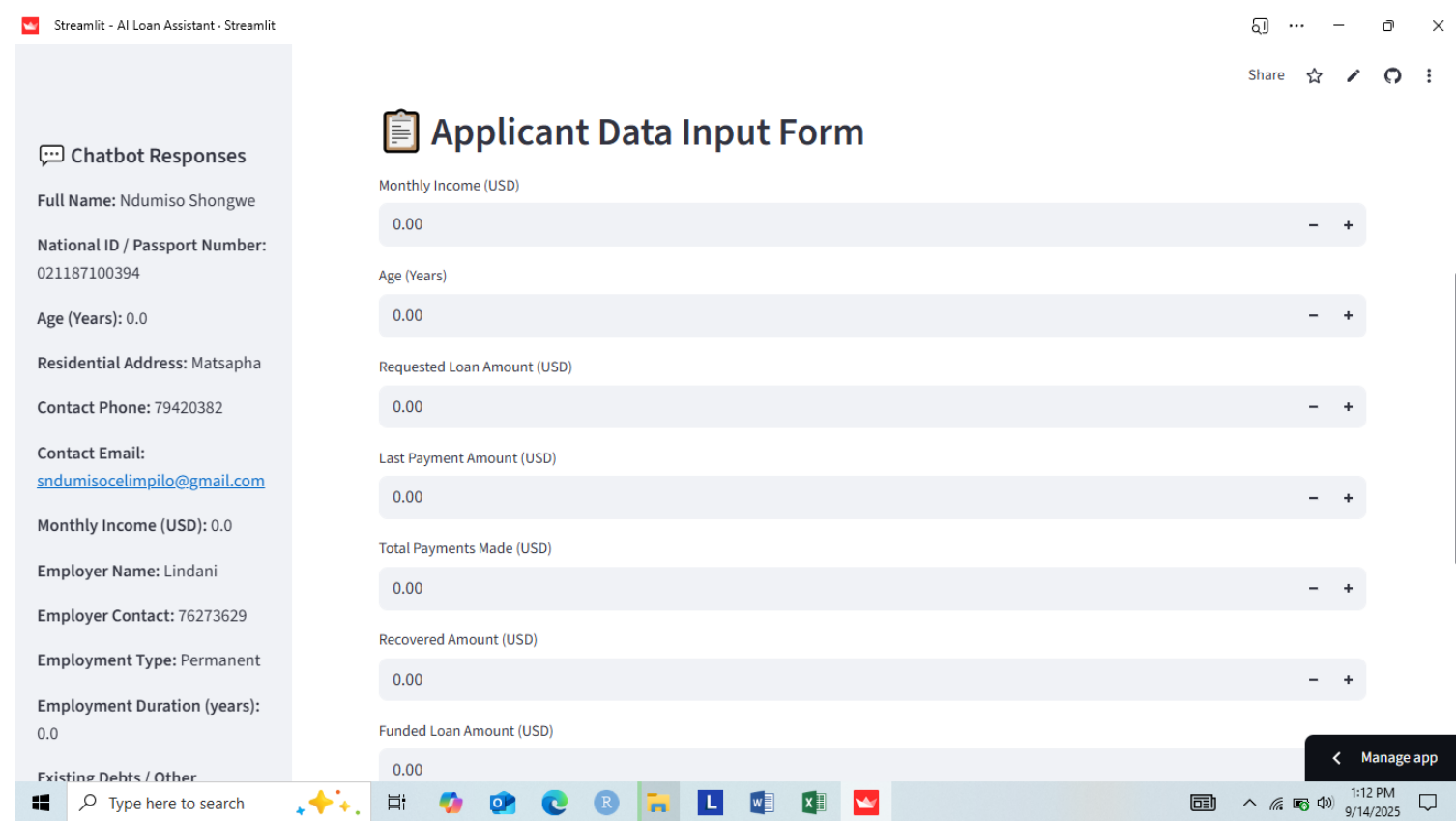


Figure 1. The web application interface allows loan officers to input borrower data (manually) and receive instant risk predictions with visual explanations.

References

- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. Advances in Neural Information Processing Systems.
- Molnar, C. (2022). Interpretable Machine Learning: A Guide for Making Black Box Models Explainable.

Acknowledgments: Mr. S. Masango (Supervisor), AI ENGINEERS-ESWATINI community, University of Eswatini Department of Statistics and Demography