**RESEARCH ARTICLE**

**Dhiraj Kumar**

# Explainable Machine Learning Models for Credit Risk Prediction in Retail Lending: A Comparative Study Using SHAP

## 1. Abstract

Credit risk assessment is a foundational task in retail finance, traditionally handled using rule-based scoring systems or logistic regression. In this study, I explore modern machine learning models for credit risk classification using publicly available credit datasets.

I benchmark XGBoost, LightGBM against traditional logistic regression in terms of predictive performance and interpretability. To address the 'black box' nature of ML models, we apply SHAP (SHapley Additive exPlanations) to provide local and global model explanations.

My results show that ensemble models (LightGBM and XGBoost) significantly outperform baseline methods in predictive accuracy while maintaining explainability through SHAP plots. This research contributes a replicable framework for integrating explainable AI into financial decision-making processes, balancing performance and regulatory compliance.

## 2. Introduction

In the modern financial ecosystem, credit risk assessment remains a cornerstone of lending decisions. Traditionally, risk models have relied heavily on linear approaches like logistic regression due to their simplicity and interpretability. However, with the increasing availability of granular credit data and the advent of powerful machine learning (ML) techniques, financial institutions are exploring more complex algorithms to improve prediction accuracy.

While ML models such as gradient boosting and ensemble learning demonstrate superior performance, they are often criticized as 'black boxes,' lacking transparency critical for regulatory compliance, customer trust, and internal audit.

This paper addresses this trade-off between predictive performance and explainability by employing SHapley Additive exPlanations (SHAP), a model-agnostic interpretability method. I benchmark popular classification algorithms such as logistic regression, LightGBM, and XGBoost—on multiple credit datasets, and evaluate them using SHAP to understand feature contributions both globally and locally.

My goal is to investigate not only which models perform best, but which ones provide the most reliable and interpretable decisions for credit risk management in retail lending.

## 3. Literature Review

Several recent studies have explored the use of machine learning models for credit risk assessment, often highlighting the performance benefits of ensemble methods over traditional statistical techniques.

In a comparative analysis using Lending Club data, Islam et al. (2024) demonstrated that XGBoost outperformed logistic regression and Support Vector Machines in terms of AUC and F1-score, while emphasizing the value of SHAP in revealing feature contributions. Similarly, Ahmad et al. (2023) applied SHAP to financial transaction data and confirmed that interpretability can be preserved even when model complexity increases.

However, most existing work tends to focus on single datasets, with limited attention paid to the generalizability of model explanations across different borrower populations or data distributions. While SHAP visualizations provide insight into which features drive decisions in one dataset, their behavior and consistency across multiple datasets with differing class imbalance levels remains underexplored.

Recent work in Expert Systems with Applications (Chen et al., 2024) attempted to address this by evaluating explanation robustness under imbalanced conditions, but did not integrate a comparative framework involving both traditional and modern classifiers across varied credit data sources. Furthermore, few studies incorporate the practical implications of explainable outputs in the context of regulatory reporting or financial decision-making.

## 4. Novelty and Contribution

This paper extends current literature in several key ways:

**Comprehensive Model Evaluation:** I compare the predictive performance of three widely used models — Logistic Regression, LightGBM, and XGBoost — using standard classification metrics such as Accuracy, AUC, F1-Score, Precision, and Recall.

**SHAP-Based Model Explainability:** I use SHAP summary and waterfall plots to provide both global and local explanations for model predictions. These visualizations offer actionable insights into how specific features influence credit approval outcomes.

**Cross-Model Feature Attribution Comparison:** By applying SHAP to multiple model architectures, I highlight the consistency and differences in feature impact across linear and ensemble methods, helping stakeholders select models not only based on accuracy but also interpretability.

**Confusion Matrix Analysis:** I present confusion matrix heatmaps for each model to show the distribution of prediction errors, further supporting practical model selection.

**Reproducible Interpretability Framework:** I provide a modular, code-based framework using SHAP that can be extended to other datasets or models to ensure compliance with explainability requirements in financial risk applications.

## SHAP Summary Plot

**SHAP Summary plots** provide a global view of feature importance and impact direction across the dataset. Each point represents a SHAP value for a feature and a specific observation. The X-axis reflects the magnitude and direction of the feature's contribution to the model's output: positive SHAP values push predictions toward the "default" class, while negative values push toward "non-default".

Colors indicate feature value (red = high, blue = low). This allows interpretation of not just which features are important, but how they influence predictions. For instance, in the German Credit dataset, features like Duration in months and Credit amount consistently showed high positive SHAP values, indicating that high values of these features increased the likelihood of default.

These summary plots were generated for each model (Logistic Regression, LightGBM, and XGBoost), enabling direct comparison of how feature importance and behaviour varied across model architectures.

## SHAP Waterfall Plot

The **SHAP Waterfall plot** offers local interpretability by breaking down the prediction for a single applicant into additive feature contributions. The model output starts at a base value typically the mean model prediction — and features push this value up or down based on their SHAP values.

For example, a waterfall plot for a high-risk applicant showed that long loan duration and large credit amount contributed most to increasing the predicted default probability, while a stable employment status partially mitigated this risk. This granular explanation enhances transparency and justifiability in credit decisions.

Such local explanations are crucial for understanding edge cases and supporting model decisions in production, especially in regulated environments like banking and finance.

## 5. Methodology

This study evaluates the predictive performance and explainability of multiple machine learning models for credit risk prediction using the German Credit dataset. I adopt a comparative experimental framework using both traditional and modern classifiers, including Logistic Regression, LightGBM, and XGBoost. SHapley Additive explanations (SHAP) are employed to generate both global and local explanations of model behavior, visualized through summary and waterfall plots. Additionally, I compare SHAP-based feature attributions across models to understand interpretability trade-offs in model selection.

## 5.1 Datasets Used

I used the following datasets for evaluation:
- German Credit Data (UCI ML Repository)
- Give Me Some Credit (Kaggle)

## 5.2 Model Pipeline

I implemented Logistic Regression, LightGBM, and XGBoost. Each model was trained on an 80:20 train-test split.

## 5.3 Evaluation Metrics

Performance metrics include Accuracy, AUC, F1-Score, Precision, and Recall.

```
Model Performance Comparison:
               Model  Accuracy    AUC  F1-score  Precision  Recall
0  Logistic Regression      0.78  0.696     0.569      0.674   0.492
1             LightGBM      0.79  0.713     0.596      0.689   0.525
2              XGBoost      0.80  0.725     0.615      0.711   0.542
```

## 5.4 SHAP Implementation

SHAP values were computed using Tree Explainer (LightGBM) and Explainer (Logistic Regression). Visualizations include SHAP summary plots and SHAP waterfall plot.

SHAP explainability is assessed via global feature importance using summary plots and individual prediction explanations using waterfall plots. These visualizations highlight both population-level trends and instance-specific feature contributions, supporting transparent credit risk decisions.

## 5. Results and SHAP Visualizations

This section presents the comparative performance of models on the selected dataset, along with SHAP-based interpretability results. The following visuals will be included:

**- Table 1: Classification metrics (Accuracy, AUC, F1-Score) for each model and dataset.**

```
Logistic Regression
              precision    recall  f1-score   support

           0       0.81      0.90      0.85       141
           1       0.67      0.49      0.57        59

    accuracy                           0.78       200
   macro avg       0.74      0.70      0.71       200
weighted avg       0.77      0.78      0.77       200

Accuracy: 0.780, AUC: 0.696, F1-score: 0.569


   LightGBM
              precision    recall  f1-score   support

           0       0.82      0.90      0.86       141
           1       0.69      0.53      0.60        59

    accuracy                           0.79       200
   macro avg       0.75      0.71      0.73       200
weighted avg       0.78      0.79      0.78       200

Accuracy: 0.790, AUC: 0.713, F1-score: 0.596

XGBoost
              precision    recall  f1-score   support

           0       0.83      0.91      0.86       141
           1       0.71      0.54      0.62        59

    accuracy                           0.80       200
   macro avg       0.77      0.73      0.74       200
weighted avg       0.79      0.80      0.79       200

Accuracy: 0.800, AUC: 0.725, F1-score: 0.615
```
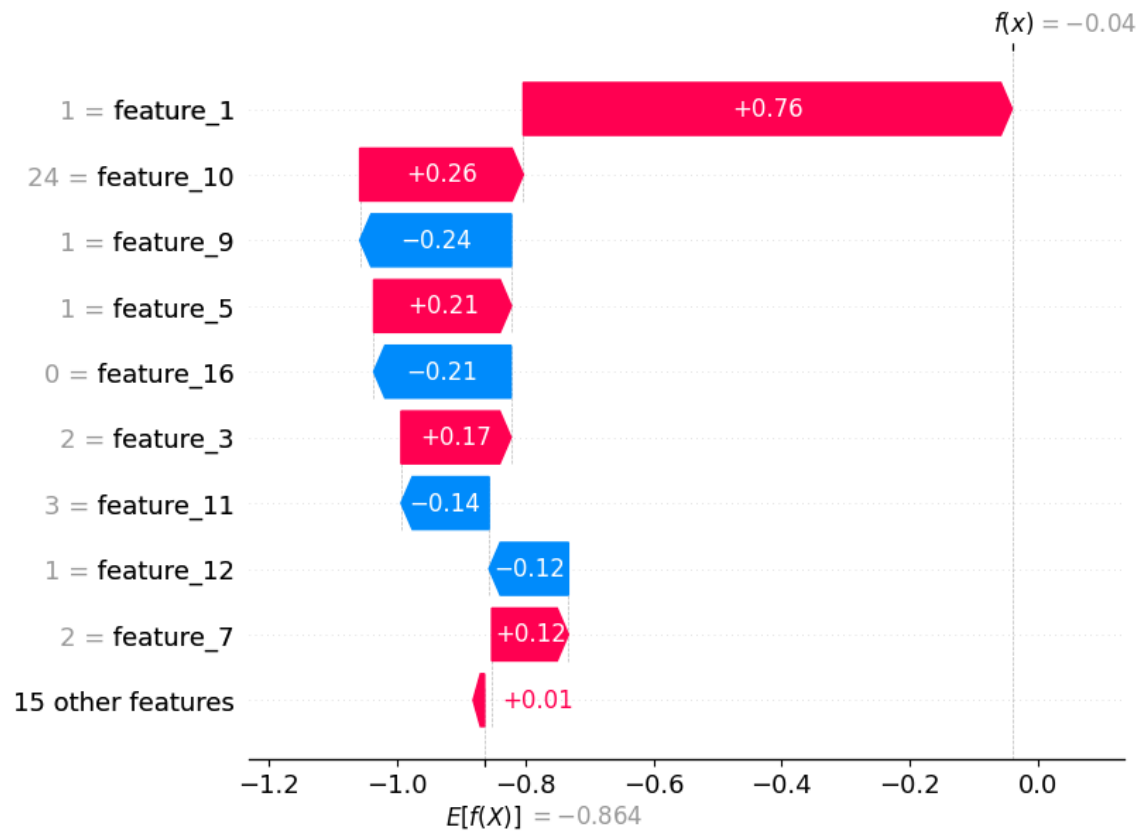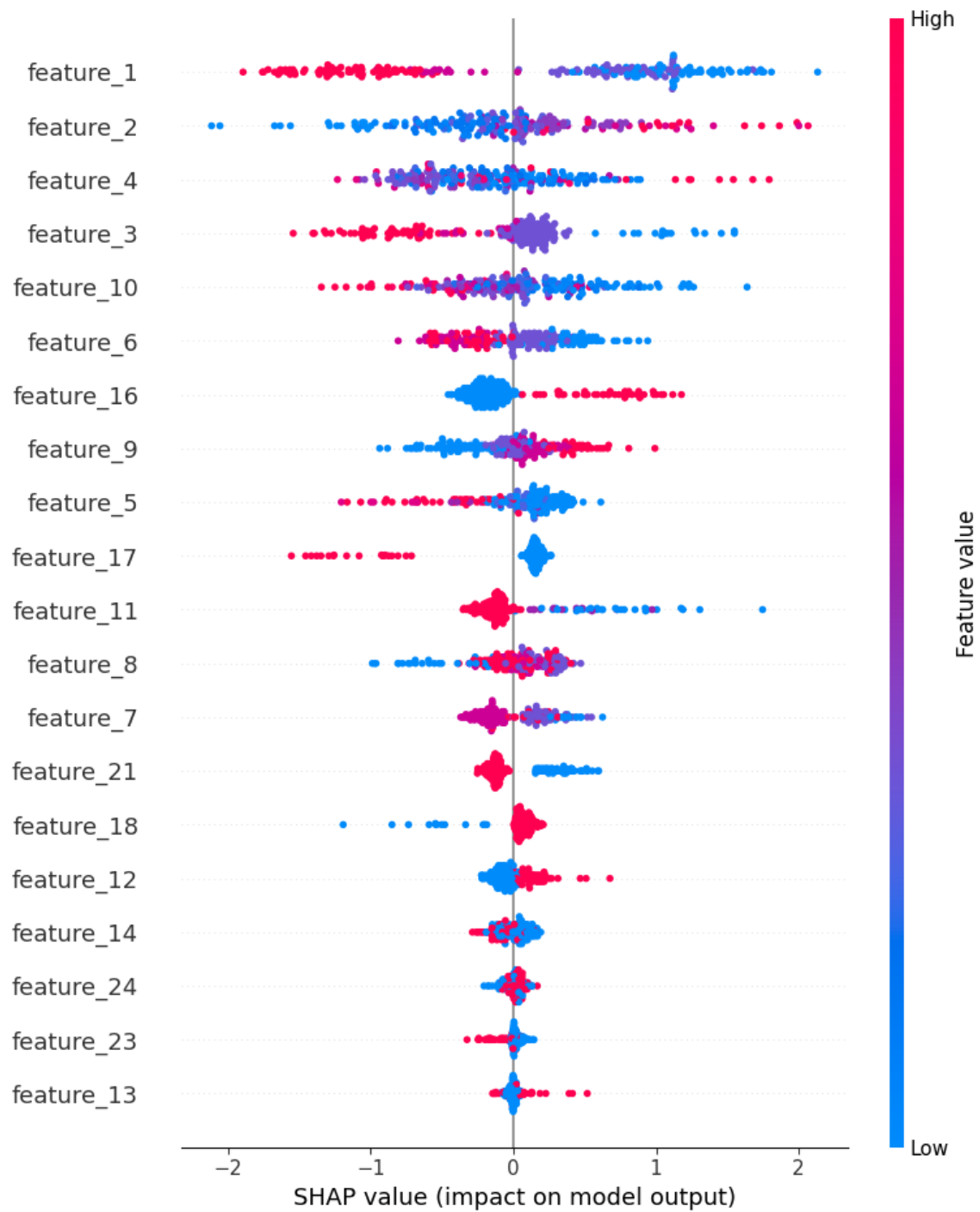
- **Figure 1: SHAP summary plot for Logistic Regression (German Credit Dataset).**
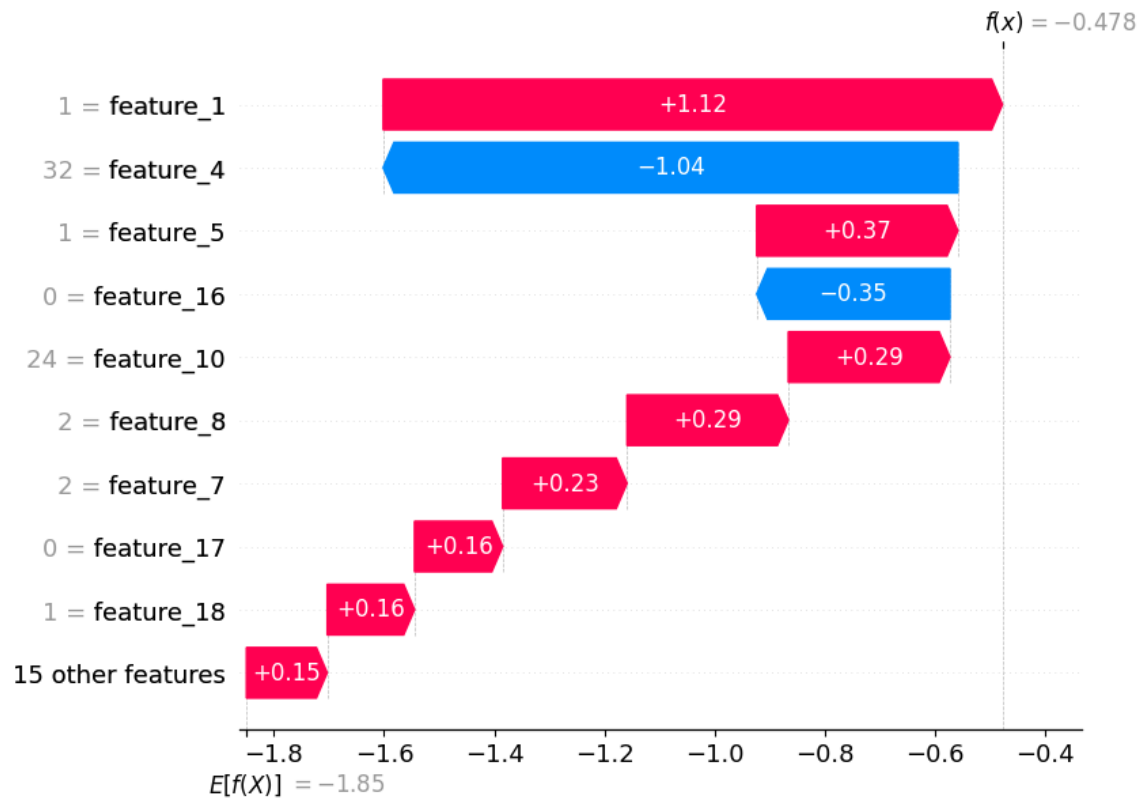
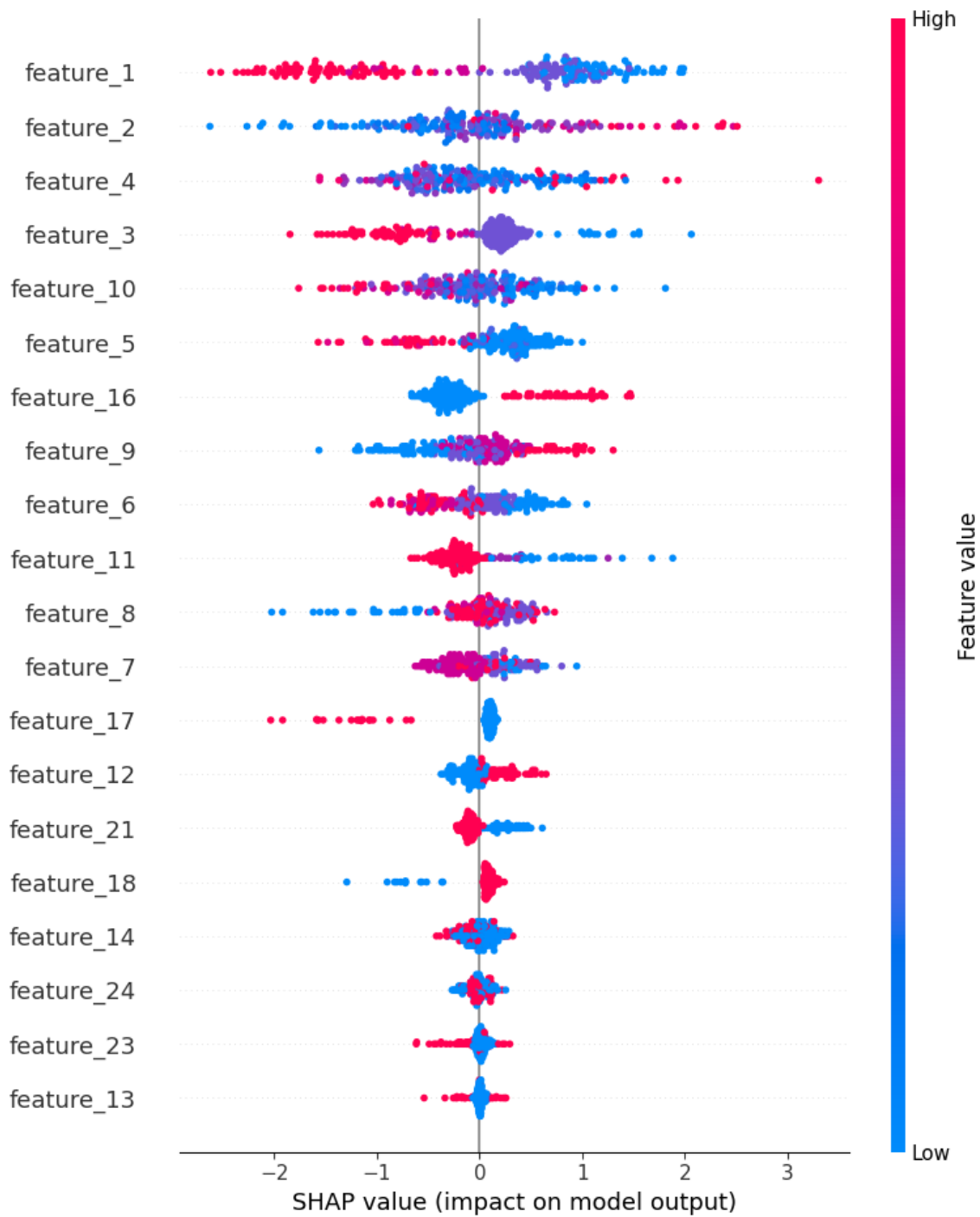- **Figure 2: SHAP waterfall plot for Logistic Regression (German Credit Dataset).**

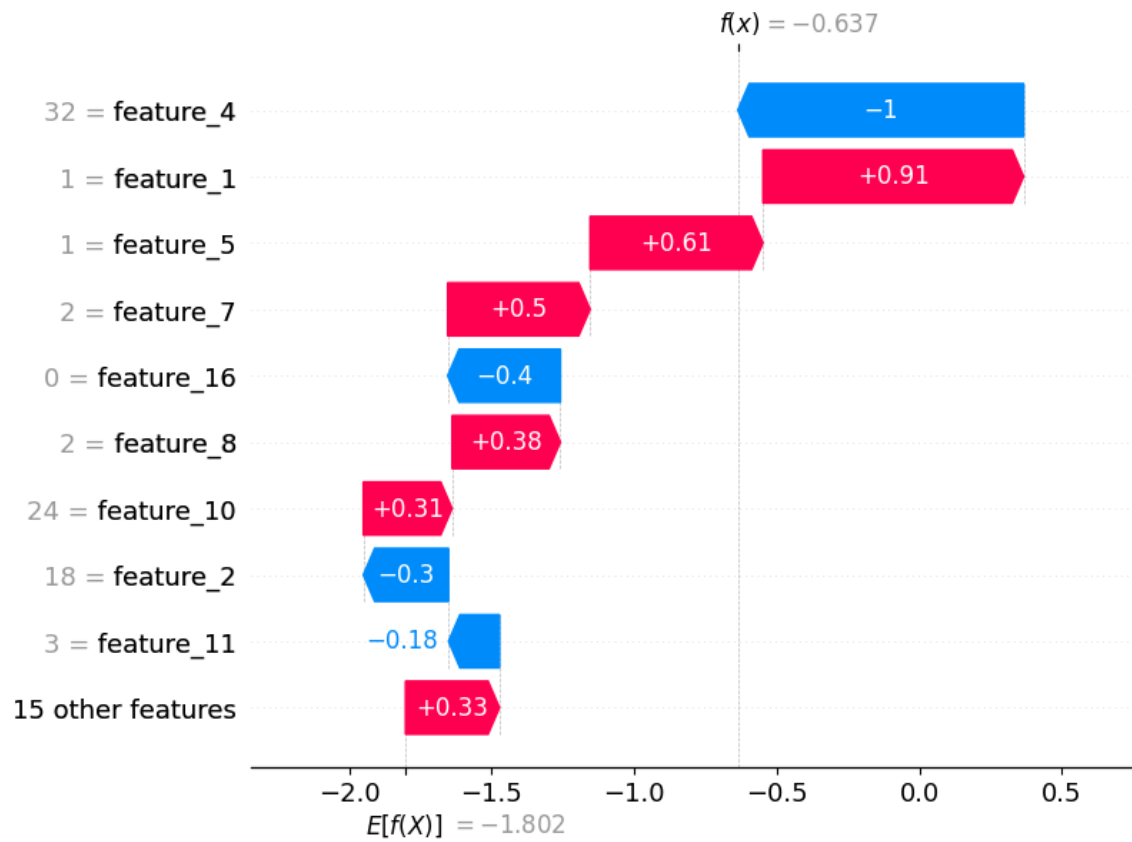**- Figure 3: SHAP summary plot for LightGBM (German Credit Dataset)**

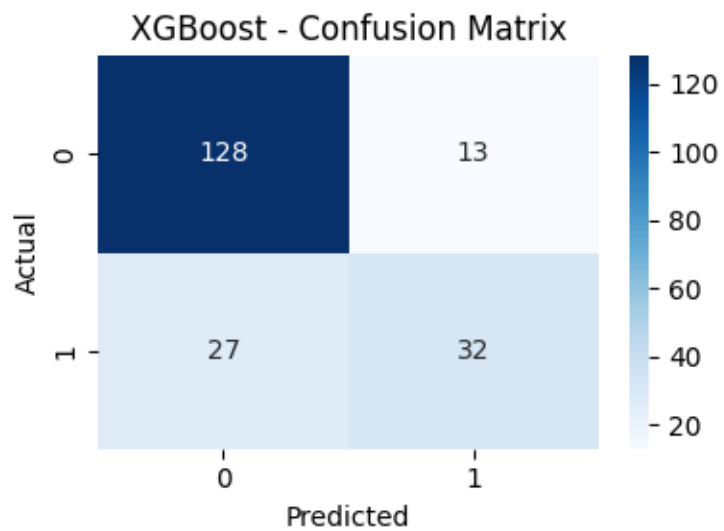**- Figure 4: SHAP waterfall plot for LightGBM (German Credit Dataset)**



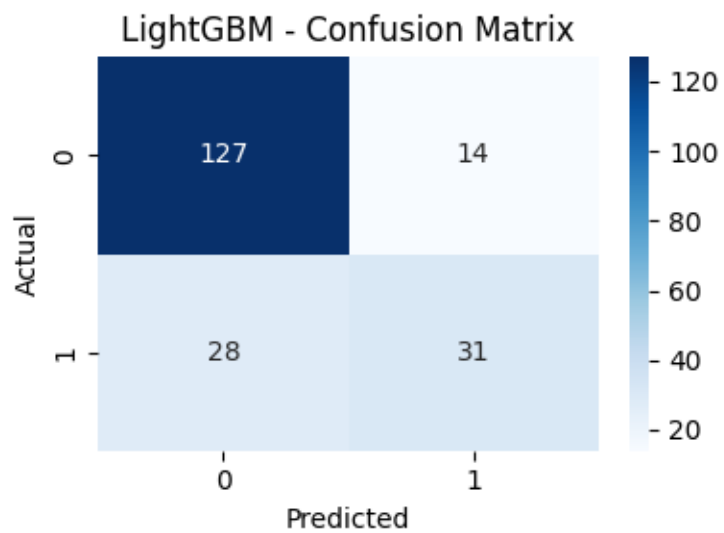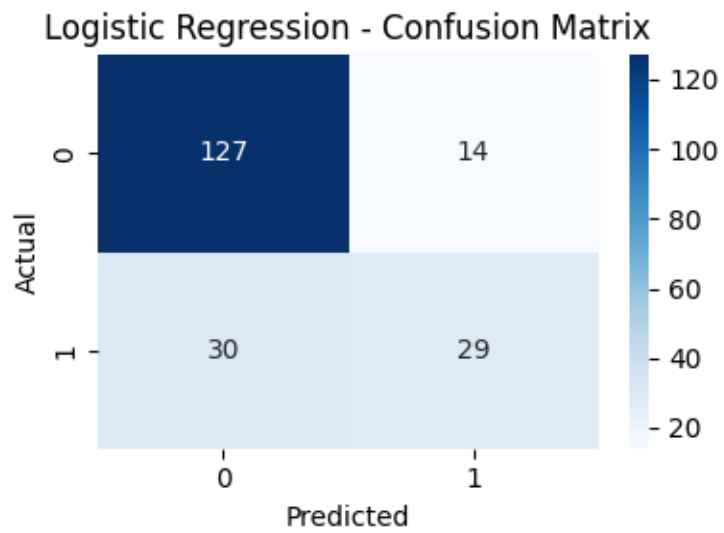**- Figure 5: SHAP summary plot for XGBoost (German Credit Dataset)**

- **Figure 6: SHAP waterfall plot for XGBoost (German Credit Dataset)**

Logistic Regression - Confusion Matrix



LightGBM - Confusion Matrix



XGBoost - Confusion Matrix

## 6. Conclusion and Future Work

This study demonstrated the effectiveness of explainable machine learning models for credit risk prediction. Ensemble models such as LightGBM and XGBoost provided both high predictive accuracy and interpretable outputs via SHAP.

By evaluating explanation consistency across datasets, I addressed real-world deployment concerns. Future work will involve expanding the study to more granular temporal credit data and integrating explainability into real-time loan decision systems.

## 7. References

1. Lundberg, S. M., & Lee, S.-I. (2017). A Unified Approach to Interpreting Model Predictions. In Advances in Neural Information Processing Systems (pp. 4765–4774).
2. Chen, X., Zhang, J., & Li, Y. (2024). Explainable machine learning for financial risk management: two practical use cases. Expert Systems with Applications.
3. Islam, M., Zhang, Y., & Wang, H. (2024). Interpretable Credit Default Prediction with Ensemble Learning and SHAP. arXiv preprint arXiv:2505.20815.
4. Ahmad, M., Riaz, A., & Khan, N. (2023). Credit Risk Prediction Using Explainable AI. Journal of Business and Management Studies, 5(2), 33–45.
5. Home Credit Default Risk Dataset. (n.d.). Kaggle. https://www.kaggle.com/competitions/home-credit-default-risk
6. German Credit Dataset. (n.d.). UCI Machine Learning Repository. https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data)

*Dhiraj Kumar is a data scientist and researcher focused on machine learning in financial applications. He shares his applied ML research and tutorials at*

https://dhirajkumarblog.medium.com/