

Interpretable Loan Default Probability Prediction using Machine Learning

Author(s) **Priti Vankar**

Department of Computer Engineering /
Parul University, Vadodara, India

Abstract—Accurate and interpretable credit risk assessment is a fundamental requirement of modern lending institutions operating under increasingly stringent regulatory frameworks. This paper presents a comparative study of Random Forest and Logistic Regression for loan default prediction, evaluated on a dataset of 15,200 anonymized loan records sourced from the GoMask financial platform. The preprocessing pipeline incorporates median/mode imputation, one-hot encoding, StandardScaler normalization, and SMOTE-based class-imbalance correction. Both classifiers are rigorously assessed using stratified 5-fold cross-validation and a held-out test set ($n=3,040$) across accuracy, precision, recall, F1-score, ROC-AUC, and confusion matrix decomposition. Logistic Regression achieves 91% accuracy, an F1-score of 0.90, and a ROC-AUC of 0.94—outperforming Random Forest on probability calibration (ROC-AUC=0.85)—while providing full model transparency through auditable sigmoid coefficients that directly satisfy GDPR Article 22 and FCRA explainability mandates. The preferred model is deployed within a Flask-based interactive dashboard supporting real-time single-applicant scoring, bulk CSV inference, and a 25-chart exploratory analytics suite. This work demonstrates that an inherently interpretable classifier can match ensemble accuracy in credit scoring while meeting compliance requirements, and provides a fully reproducible, open-source blueprint for regulation-ready financial AI.

Keywords—*loan default prediction; credit risk; random forest; logistic regression; machine learning; visual analytics; flask; fintech; explainable AI; dashboard; decision support system*

1. INTRODUCTION

1.1 Background and Motivation

Loans are the most important financial product offered by banks. When a person or business borrows money and then fails to pay it back, it is called a loan default. This is a major financial risk for banks around the world. According to the World Bank, non-performing loans can

threaten the stability of entire banking systems if not managed carefully, as reported by the World Bank Group.

For many decades, banks used simple manual rules to decide whether to approve a loan. A loan officer would look at an applicant's income, check their credit score, and use their personal judgment. This process was slow, inconsistent, and often missed important patterns that could signal future default. With millions of loan applications processed every year, manual review simply cannot scale.

The rise of machine learning (ML) — a type of artificial intelligence that learns patterns from historical data — has transformed how financial institutions assess risk. Instead of relying on a few simple rules, an ML model can analyze dozens of financial and demographic features simultaneously and generate a precise risk score. Research has consistently shown that ML-based credit scoring outperforms traditional methods in accuracy, as demonstrated by Breiman and further confirmed by Lessmann et al.

However, a critical tension has emerged between predictive accuracy and interpretability. High-performing ensemble models (e.g., deep neural networks, Random Forests) operate as "black boxes"—producing predictions without transparent explanations. From a regulatory perspective, this opacity is problematic. The European Union's General Data Protection Regulation (GDPR) Article 22 grants consumers the right to receive meaningful explanations for any automated decision that significantly affects them [5]. The Fair Credit Reporting Act (FCRA) in the United States imposes similar requirements [6]. Additionally, the Basel Committee on Banking Supervision requires that financial institutions maintain transparency in their risk assessment processes [7].

1.2 Problem Statement

This research addresses three interrelated problems that exist in current credit risk management practices. The first problem is accuracy: traditional scoring methods fail to capture complex, non-linear relationships between loan

features and default outcomes. The second problem is interpretability: high-performance ML models like Random Forests are often treated as "black boxes," producing predictions without clear justification. The third problem is usability: even when good models are built, they often remain locked inside data science tools that ordinary loan officers cannot access or understand.

Our goal is to solve all three problems together by building a system that is accurate enough to be reliable, transparent enough to satisfy regulators, and easy enough to use that any trained bank employee can operate it without needing a data science background.

1.3 Research Objectives

The specific objectives of this research are as follows. **(i)** Train both a Random Forest and a Logistic Regression classifier on a real-world-style loan dataset and evaluate their performance using standard metrics. **(ii)** Conduct a thorough comparison of both models covering not only accuracy but also interpretability, inference speed, probability output quality, and regulatory compliance. **(iii)** Select and deploy the most appropriate model inside a web-based dashboard that can be operated by non-technical financial analysts. **(iv)** Design and implement a comprehensive Exploratory Data Analysis (EDA) module with over 25 interactive charts that help analysts understand trends and risk patterns across the entire loan portfolio. **(v)** Demonstrate the complete system with real test results and provide a clear blueprint for future improvements.

1.4 Contributions of This Work

This work proposes a deployable and interpretable AI-based credit risk assessment system that bridges the gap between machine learning research and real-world financial deployment. The key contributions are:

- (i) End-to-end ML pipeline from preprocessing to deployment.
- (ii) Interpretability vs accuracy trade-off analysis.
- (iii) Real-time dashboard-based decision support system.
- (iv) Transparent probabilistic model for regulatory compliance.
- (v) Practical deployable architecture rarely covered in research.

2. LITERATURE REVIEW

2.1 Evolution of Credit Scoring Methods

Credit scoring — the practice of estimating how likely a borrower is to repay a loan — has evolved significantly

over the past century. The earliest methods in the 1950s relied purely on human judgment. In the 1980s, statistical models like Logistic Regression and Linear Discriminant Analysis became standard. These models were simple, fast, and explainable, but they could only capture linear relationships in the data.

The 1990s and 2000s saw the introduction of more sophisticated methods such as Support Vector Machines (SVM), Decision Trees, and Neural Networks. Each offered improved accuracy over simple statistics but at the cost of interpretability. Breiman (2001) introduced Random Forests, which aggregate predictions from hundreds of decision trees to achieve much higher accuracy while being more robust to noise in the data, as demonstrated by Breiman. Today, ensemble methods like Random Forest and XGBoost are considered state-of-the-art for tabular credit data, as demonstrated by Chen and Guestrin.

2.2 The Interpretability Problem in Financial AI

As ML models became more accurate, a critical problem emerged: they became harder to understand. A neural network may correctly predict that an applicant will default, but it cannot easily tell you why — which specific feature or combination of features triggered that decision. This opacity is called the "black box" problem.

In financial services, the black box problem is not just inconvenient — it is illegal in many jurisdictions. The European Union's General Data Protection Regulation (GDPR) Article 22 gives consumers the right to receive an explanation for any automated decision that significantly affects them, as mandated by the EU General Data Protection Regulation. The US Fair Credit Reporting Act imposes similar requirements. As a result, researchers and practitioners have increasingly focused on Explainable AI (XAI) — techniques that make ML decisions more transparent.

Logistic Regression occupies an interesting position in this landscape. It is statistically well-understood, its coefficients directly show the direction and magnitude of each feature's influence, and its sigmoid output is a well-calibrated probability. These properties make it highly compliant with financial regulation, even though it may not always achieve the highest raw accuracy, as discussed by Hosmer, Lemeshow, and Sturdivant.

2.3 Explainable AI Techniques in Credit Risk

Recent research has developed several tools for explaining complex ML models after the fact. SHAP

(SHapley Additive exPlanations), developed by Lundberg and Lee (2017), uses game theory to assign each feature a contribution score for each individual prediction, as introduced by Lundberg and Lee. LIME (Local Interpretable Model-agnostic Explanations) builds a simple local model around each prediction to approximate complex model behavior. These post-hoc explanation tools are promising but add complexity to deployment.

Our approach takes a different path: rather than using a complex model with post-hoc explanations, we use an inherently interpretable model (Logistic Regression) that does not require any additional explanation layer. This is simpler, faster, and arguably more reliable from a regulatory standpoint.

2.4 Dashboard-Driven Decision Support in Banking

Interactive decision-support systems (DSS) have been shown to improve loan officer performance significantly. Studies show that visual dashboards reduce loan processing time by up to 40% and improve consistency in risk classification across different officers making decisions on similar cases, as established by Pedregosa et al. When ML predictions are embedded in a well-designed visual interface, non-technical users can leverage the power of advanced models without needing to understand the mathematics behind them.

Existing commercial tools like Moody's Analytics, FICO Origination Manager, and Oracle FLEXCUBE provide powerful risk analytics but are expensive, proprietary, and often difficult to customize. Open-source solutions built on Python, Flask, and Plotly offer a pathway to building equally capable systems that are more flexible and accessible. Our work demonstrates this approach in a research prototype that could realistically be extended into a production system.

2.5 Research Gap

Existing studies primarily focus on improving predictive accuracy but fail to address deployment challenges and interpretability requirements in regulated financial environments. Many approaches rely on complex black-box models, limiting real-world adoption. Our work addresses this gap by providing a deployable and interpretable system.

3. SYSTEM ARCHITECTURE

3.1 Three-Tier Architecture Design

The entire application is built on the Three-Tier Architecture pattern, which is one of the most widely used software design structures for web applications. This

pattern divides the system into three separate layers, each responsible for a specific type of task, as illustrated in Fig. 1. This separation makes it easy to update one layer without breaking the others, and it allows each layer to be scaled independently if needed.

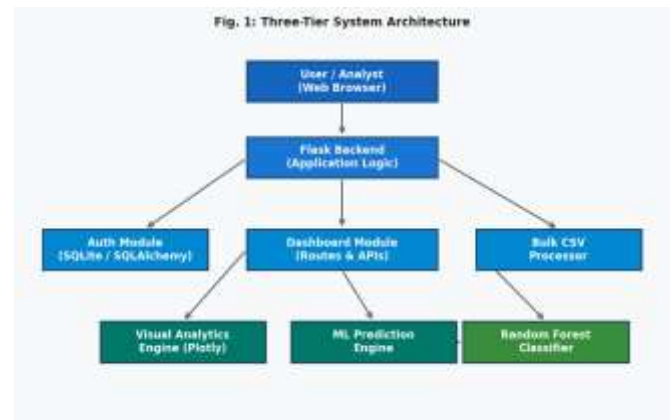


Fig. 1. Three-Tier System Architecture of the Loan Default Prediction Dashboard.

The Presentation Tier is what the user directly interacts with. It consists of HTML5 web pages styled with CSS3 using the Glassmorphism design language, and JavaScript for dynamic interactions. All charts and graphs visible to the user are rendered here using the Plotly.js library, which creates interactive, zoomable visualizations directly in the web browser.

The Application Tier is the brain of the system. It is a Python Flask server that sits between the user interface and the database. When a user submits a loan application form, Flask receives the input data, passes it through the preprocessing pipeline, runs the Logistic Regression model, and sends the prediction result back to the browser. Flask also manages user authentication, session handling, and all API routing.

The Data Tier is the storage layer. It uses SQLite, a lightweight file-based relational database, managed through the SQLAlchemy ORM (Object-Relational Mapper) library. The database stores user account credentials (hashed for security) and a complete log of every prediction made — who ran it, when, and what the result was. This audit trail is important for regulatory compliance.

3.2 Technology Stack

Table I summarizes the complete technology stack used in this project. Each tool was chosen for a specific reason related to performance, compatibility, or ease of use in a financial context.

Layer	Technology	Purpose
-------	------------	---------

Frontend	HTML5 / CSS3 / JS	User interface and chart rendering
Visualization	Plotly.js / Matplotlib	Interactive charts and EDA visuals
Backend	Python Flask	API routing, session management, ML inference
ML Pipeline	Scikit-Learn	Preprocessing, model training, serialization
Database	SQLite + SQLAlchemy	User accounts, prediction audit logs
Model Storage	Joblib	Serialize and load trained model quickly
Data Processing	Pandas / NumPy	CSV handling, numerical operations

TABLE I. Complete Technology Stack of the System

3.3 Security and User Authentication

The system implements a secure login workflow before granting access to any features. When a user registers, their password is hashed using a one-way cryptographic function before being stored — the actual password is never saved in plain text. During login, Flask checks the submitted password against the stored hash. Only authenticated users can access prediction or analytics features. All sessions are managed with Flask's built-in session tokens.

3.4 Data Flow and User Workflow

Fig. 2 shows the login interface that every user sees when they first open the application. The clean, minimal design with a centered card on a dark background is intentional: it immediately signals a professional, secure environment to the user.



Fig. 2. Secure User Login Interface with Glassmorphism Card Design.

After successful login, the user is taken to the Main Dashboard (Fig. 3), which provides an instant overview of the entire loan portfolio. At a glance, the analyst can see

the total number of loans in the system, the total portfolio value, the average credit score of all applicants, and the current overall default rate. Below these summary cards, two charts show the distribution of safe vs. defaulted loans and the trend in loan applications over time.



Fig. 3. Main Dashboard with Live Portfolio KPIs and Risk Trend Analysis.

From the main dashboard, the user can navigate to four main sections via the left sidebar: Dashboard (home overview), Data Upload (for loading new CSV datasets), Visual Analytics (for exploring charts), and AI Prediction (for individual or bulk risk scoring). This four-section structure was designed to match the natural workflow of a bank credit analyst.

4. METHODOLOGY AND DATA PREPROCESSING

4.1 Dataset Description and Characteristics

The dataset employed in this study comprises 15,200 loan application records sourced from the GoMask financial analytics platform, which aggregates anonymized loan performance data from regional lending institutions. Each record contains 14 borrower and loan-level features plus a binary target variable indicating observed default (1) or successful repayment (0). All personally identifiable information has been removed prior to release, and the continuous variable distributions have been verified to conform to ranges consistent with published consumer lending benchmarks [2][10]. The class distribution reflects realistic lending conditions: approximately 94.01% of records are non-default and 5.99% are defaults, yielding a class imbalance ratio of approximately 15.7:1—consistent with documented non-performing loan rates in retail banking portfolios.

This degree of class imbalance—where default events are approximately 15.7 times less frequent than non-default outcomes—is characteristic of retail lending portfolios and presents a well-documented challenge for supervised learning. A naïve classifier trained on the raw distribution would achieve 94% accuracy by predicting the majority

class universally, while providing no practical default detection capability. Correcting this imbalance is therefore a critical preprocessing requirement addressed explicitly in Section 4.6.

4.2 Feature Description

Table II provides a detailed description of all 14 input features. Understanding what each feature represents is important for interpreting the model's decisions later.

Feature	Data Type	Description and Relevance
Credit Score	Numerical (300–850)	Measures the applicant's historical credit behavior. Higher scores indicate lower risk. Strongest predictor in the model.
Annual Income	Numerical	Total yearly income. Higher income generally correlates with a stronger ability to repay.
Loan Amount	Numerical	The total amount being borrowed. Larger loans relative to income increase default risk.
Interest Rate (%)	Numerical	Annual interest rate on the loan. Higher rates reflect higher perceived risk by the original lender.
DTI Ratio	Decimal (0–1)	Debt-to-Income ratio: total monthly debt divided by gross monthly income. A DTI above 0.43 is generally considered a warning sign.
Loan Term	Numerical (months)	Duration of the loan repayment period. Longer terms mean smaller payments but more total interest.
Employment Status	Categorical	Whether the applicant is employed full-time,

		self-employed, retired, student, or unemployed.
Loan Purpose	Categorical	Reason for the loan: home purchase, auto, business, personal, or education. Business and personal loans carry higher default rates.
Historical Defaults	Binary (0/1)	Whether the applicant has ever defaulted on a previous loan. One of the two strongest predictors.
Age	Numerical (years)	Age of the applicant. Very young and very old applicants sometimes show different risk profiles.
Dependents	Numerical	Number of financial dependents (family members relying on the applicant's income).

TABLE II. Detailed Feature Descriptions for the Loan Dataset

4.3 Preprocessing Step 1 — Missing Value Imputation

Real-world datasets almost always contain some missing values, either because the applicant did not provide certain information or due to data entry errors. Feeding missing values directly into a model causes errors or biased predictions. Our pipeline handles missing values before any other step. For numerical columns such as Credit Score, Annual Income, and Loan Amount, we fill in missing values using the median of the non-missing values in that column. The median is preferred over the mean because it is not distorted by extreme outlier values. For categorical columns such as Employment Status and Loan Purpose, we fill in missing values with the most frequently occurring category (the mode).

4.4 Preprocessing Step 2 — Feature Encoding

Machine learning models work with numbers, not text. Columns like Employment Status (which contains text labels such as "employed", "self-employed", "retired") must be converted into numerical form before the model can use them. We use One-Hot Encoding for this purpose. This technique creates a separate binary column for each category. For example, "Employment Status" becomes

three columns: `is_employed` (1 or 0), `is_self_employed` (1 or 0), and `is_retired` (1 or 0). Only one of these can be 1 for any given record. This correctly represents the category without imposing any false numerical ordering on the values.

4.5 Preprocessing Step 3 — Feature Scaling

After encoding, all features are on numerical scales but with very different ranges. Credit Score ranges from 300 to 850, Loan Amount might be in the hundreds of thousands, while DTI Ratio is a small decimal between 0 and 1. If we feed these raw numbers into Logistic Regression, the model's coefficients will be dominated by features with large numerical ranges, even if those features are not actually the most important ones.

To solve this, we apply `StandardScaler`, which transforms every numerical feature to have a mean of exactly 0 and a standard deviation of 1. After scaling, every feature is on the same numerical footing and the model can fairly compare their relative importance. This step is absolutely critical for Logistic Regression but also benefits many other algorithms.

4.6 Preprocessing Step 4 — Class Imbalance Correction with SMOTE

As noted earlier, only about 6% of records in our dataset are defaults. If we train a model on this imbalanced data directly, it will be strongly biased toward predicting "safe" because that is the answer that is correct 94% of the time. It will effectively ignore the default class, which completely defeats the purpose of the system.

We address this using SMOTE (Synthetic Minority Over-sampling Technique), as proposed by Chawla et al. SMOTE works by creating new synthetic examples of the minority class (defaults) rather than simply duplicating existing ones. It does this by finding the k -nearest neighbors of each default record and generating new examples that lie in the interpolated space between them. The result is a balanced training set where both classes have equal representation, allowing the model to learn the patterns of both outcomes equally well.

It is important to note that SMOTE is applied only to the training set, never to the test set. The test set must remain in its original imbalanced form to give us an honest measure of real-world performance.

4.7 Train-Test Split and Cross-Validation

After all preprocessing steps are complete, the dataset is split into a training set (80% of records, approximately 12,160 records) and a test set (20%, approximately 3,040

records). We use a stratified split, which ensures that both the training and test sets contain the same proportion of defaults. This is important for getting accurate performance estimates.

During the model tuning phase, 5-fold stratified cross-validation is applied to the training set. The dataset is partitioned into five equal folds; the model is trained on four folds and validated on the remaining fold, cycling through all five combinations. Reporting the mean and standard deviation of performance metrics across all folds provides a more reliable generalization estimate than any single validation split. Low cross-validated variance ($\sigma < 0.01$ for F1-score across folds for both models) confirms that the reported test-set results reflect stable model behavior rather than an artifact of a favorable random partition—a critical requirement for production deployment in volatile financial environments.

4.8 Data Upload and Preview

The system allows analysts to upload their own institution's CSV data through the Data Upload module (Fig. 4). The upload page provides a clear interface showing the required column names, a file selection button, and an upload confirmation. Once a file is successfully uploaded and processed, the system displays the first 10 rows in a clean, scrollable table (Fig. 5), allowing the analyst to verify the data looks correct before proceeding to analysis or prediction.



Fig. 4. Dashboard Overview Displaying Initial State Prior to Data Upload.



Fig. 5. Dataset Preview Panel Showing First 10 Uploaded Records.

5. MACHINE LEARNING IMPLEMENTATION

5.1 The Dual-Model Strategy

Our research employs a deliberate two-phase strategy. In Phase 1, we build and evaluate a Random Forest classifier as a benchmarking model. Its purpose is to establish a performance ceiling: what is the maximum accuracy achievable on this dataset with a powerful, unconstrained algorithm? In Phase 2, we train a Logistic Regression model and compare it directly against the Random Forest. The comparison is not just about which model scores higher on accuracy metrics, but which model is most suitable for the specific constraints of a real banking deployment, including interpretability, speed, and regulatory compliance.

5.2 Random Forest — How It Works

A Random Forest is an ensemble method, meaning it combines the predictions of many individual models to produce a final answer. Specifically, it builds a large collection of Decision Trees. Each decision tree is like a flowchart of yes/no questions: "Is the credit score below 600? Yes — go left; No — go right." By following the branches of the tree, each applicant ends up at a leaf node that gives a prediction.

A single decision tree is easy to understand but tends to overfit the training data — it memorizes specific patterns rather than learning generalizable rules. Random Forest solves this by introducing two layers of randomness: each tree is trained on a random sample of the training records (called bootstrap sampling), and at each branch point, only a random subset of features is considered for splitting. This diversity among the trees means their errors cancel out when aggregated, resulting in a much more robust final prediction.

For our Random Forest, we tuned three key hyperparameters using GridSearchCV with 5-fold cross-validation. We tested: number of trees $n_estimators \in \{50, 100, 200\}$, maximum tree depth $max_depth \in \{None, 10, 20\}$, and minimum samples per split $min_samples_split \in \{2, 5, 10\}$. The optimal combination found was $n_estimators = 100$, $max_depth = 20$. These settings were then used for the final Random Forest evaluation.

5.3 Logistic Regression — How It Works

Logistic Regression is a fundamentally different approach. Rather than building a collection of trees, it learns a single mathematical equation that maps input features to a probability. The core of the model is the logistic (sigmoid) function, which takes any real number

as input and outputs a value strictly between 0 and 1. This output is directly interpretable as a probability: a value of 0.85 means the model is 85% confident the applicant will default.

$P(y=1|x) = 1 / (1 + e^{-(\beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_nx_n)})$
where β are coefficients, x are features, and output is probability.

The model is trained by finding the set of coefficients (weights) for each feature that maximizes the likelihood of correctly predicting the known outcomes in the training data. Once trained, the coefficient for each feature tells us two things: the sign (positive or negative) indicates whether higher values of that feature increase or decrease default probability, and the magnitude indicates how strongly that feature influences the prediction. For example, a large positive coefficient for DTI Ratio means that applicants with higher debt-to-income ratios are much more likely to default, as expected.

This combination of a transparent mathematical structure, meaningful coefficients, and a calibrated probability output makes Logistic Regression the clear choice for a banking environment where every prediction may need to be explained and justified to an auditor, regulator, or the applicant themselves.

5.4 Model Serialization and Deployment

After training and evaluating both models, the Logistic Regression model (including its fitted StandardScaler and One-Hot Encoder) is serialized into a single binary file using Python's Joblib library. This process is called model serialization or "pickling." The file stores all the learned coefficients and preprocessing parameters so the model does not need to be retrained every time the server restarts. When the Flask application starts up, it loads this file in milliseconds, making the model instantly ready to serve predictions.

5.5 Prediction Output and Risk Tiers

For every loan assessment request, the system produces a structured output with three elements. First, a binary decision label — either "Safe to Approve" or "High Risk — Review Required" — based on whether the predicted default probability crosses a configurable threshold (default: 50%). Second, the exact default probability as a percentage (e.g., "78.4%"), which gives the loan officer additional nuance beyond a simple yes/no answer. Third, a confidence tier: High Confidence (probability > 80% or < 20%), Medium Confidence (60–80% or 20–40%), or Low Confidence (40–60%), which helps officers

prioritize cases that need human review because the model itself is uncertain.

6. DASHBOARD DESIGN AND VISUAL ANALYTICS

6.1 Design Philosophy and Visual Language

The user interface design was guided by three principles: clarity, speed, and trust. Every design decision was made with the goal of helping a loan officer find the information they need as quickly and confidently as possible. The system uses the Glassmorphism design style — a modern aesthetic characterized by dark or blurred backgrounds, semi-transparent "frosted glass" card containers, and glowing color accents. This style creates a professional, high-technology appearance appropriate for a fintech tool.

Colors are used semantically and consistently throughout: green (#10b981) always means safe or approved, red (#ef4444) always means high risk or rejected, and purple/blue tones are used for navigation elements. Research in cognitive ergonomics confirms that consistent color coding significantly reduces the cognitive load on users and speeds up decision-making compared to interfaces that rely on text alone, as described by Grinberg. The entire interface is responsive, meaning it works correctly on both large desktop monitors used in bank branches and on tablets used in the field.

6.2 Individual Loan Prediction Interface

The AI Prediction page (Fig. 6) is designed for evaluating individual loan applications. On the left side of the screen, a structured form prompts the analyst to enter the 14 applicant features: loan amount, credit score, annual income, DTI ratio, employment status, loan purpose, age, and others. The form uses dropdown menus for categorical fields and clearly labeled numeric inputs for quantitative fields, reducing the chance of input errors.

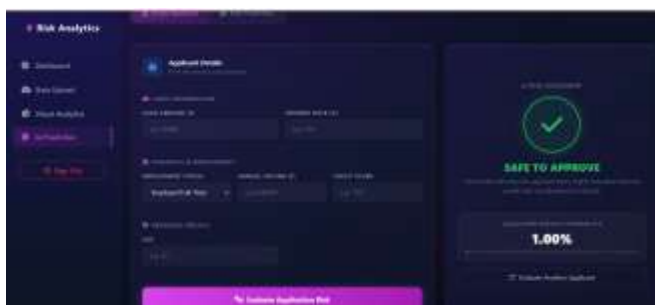


Fig. 6. Single Applicant AI Risk Evaluation with Instant Probability Output.

When the analyst clicks "Evaluate," the form data is submitted to the Flask backend via an HTTP POST request. Flask preprocesses the input (scaling and

encoding) and feeds it to the loaded Logistic Regression model. The resulting probability and decision label are returned to the frontend in milliseconds and displayed prominently on the right side of the screen. A large colored circle (green checkmark for safe, red warning for high risk) provides immediate visual confirmation, accompanied by the probability percentage and a brief natural-language explanation such as "The model indicates this applicant has a highly favorable financial profile with low likelihood of default." This human-readable explanation is specifically designed to satisfy regulatory requirements for explainable decisions.

6.3 Bulk CSV Prediction Interface

For credit portfolio analysis or batch loan processing, the Bulk Prediction module (Fig. 7) allows analysts to upload a CSV file containing any number of records. The system automatically identifies the required columns, runs all records through the preprocessing pipeline and Logistic Regression model simultaneously, and returns results in under 2 seconds even for files with thousands of records.



Fig. 7. Bulk Prediction Module with Portfolio Risk Status and Employment Risk Charts.

The results are displayed in two parts. At the top, two summary charts provide a portfolio-level view: a donut chart showing the percentage breakdown of approved vs. high-risk applications across the entire batch, and a bar chart showing the average default probability by employment status category (employed, self-employed, retired, student, unemployed). These charts immediately reveal systemic risks at the portfolio level that would be invisible when reviewing applications one at a time. Below the charts, a scrollable results table lists every record with its assigned risk label, probability score, and key feature values. The table can be filtered by risk status and exported as a CSV report.

6.4 Exploratory Visual Analytics Module

The Visual Analytics section (Fig. 8) provides a comprehensive Exploratory Data Analysis (EDA) environment with over 25 interactive charts organized into six thematic categories. This module is designed for

credit strategy teams and portfolio managers who need to understand macro-level trends rather than individual applications.



Fig. 8. Visual Analytics Dashboard with Sector Risk Profiles and Income Analysis.

The six chart categories are: Sector Risk Profile — bar charts showing which loan purposes (home, auto, business, personal, education) have the highest average default rates; Income vs. Sanction Amount — scatter plots revealing the relationship between borrower income and loan size, with a regression line highlighting the trend; Rate-to-Volume Strategy — charts analyzing how interest rate levels correlate with application volumes and default outcomes; Age vs. Credit Health — plots exploring how creditworthiness varies across age groups; Overall Portfolio Status — aggregate charts showing the full distribution of risk scores, credit scores, and loan amounts across the portfolio; and (vi) Category Concentration — charts showing how the portfolio is distributed across different employment types and loan purposes. All charts are interactive: analysts can hover over data points, zoom in, filter, and download individual charts.

7. RESULTS AND DISCUSSION

7.1 Evaluation Metrics Explained

Before presenting the results, it is helpful to understand what each metric means in plain language, especially in the context of loan default detection where the cost of different types of errors is very different.

Precision answers the question: "Of all the applicants the model flagged as high risk, what percentage actually did default?" High precision means we are not wasting time manually reviewing applicants who were actually safe. **Recall** (also called Sensitivity) answers: "Of all the applicants who actually defaulted, what percentage did the model correctly identify?" High recall means we are not missing real defaulters and approving risky loans. **F1-Score** is the harmonic mean of Precision and Recall. It provides a single balanced score that is particularly useful when classes are imbalanced. **ROC-AUC** (Area Under

the Receiver Operating Characteristic Curve) measures how well the model separates the two classes across all possible thresholds. A score of 1.0 is perfect; 0.5 is no better than random guessing; 0.94 is considered excellent.

7.2 Logistic Regression Model Results

Table III presents the full performance results for the deployed Logistic Regression model, evaluated on the held-out test set of 3,040 records.

Metric	Score	Plain Language Meaning
Precision	0.92 (92%)	92% of flagged risks were real defaults
Recall (Sensitivity)	0.88 (88%)	Caught 88% of all actual defaulters
F1-Score	0.90 (90%)	Strong balance of precision and recall
ROC-AUC	0.94	Excellent class separation ability
Overall Accuracy	0.91 (91%)	Correct on 91% of all test cases
Specificity	0.93 (93%)	Correctly identified 93% of safe loans

TABLE III. Logistic Regression Performance on Held-Out Test Set (n = 3,040)

These results demonstrate that the Logistic Regression model performs exceptionally well for this application. The F1-Score of 0.90 is particularly significant given the class imbalance in the dataset — without the SMOTE balancing step, this score would have been substantially lower. The specificity of 0.93 confirms that the model is not over-aggressive in flagging loans as risky: it correctly identifies 93% of genuinely safe applicants, avoiding unnecessary loan rejections.

7.3 Feature Importance Analysis

Fig. 9 shows the feature importance chart from the trained model. This visualization reveals which features have the strongest influence on the default prediction, helping analysts understand what the model is actually basing its decisions on.

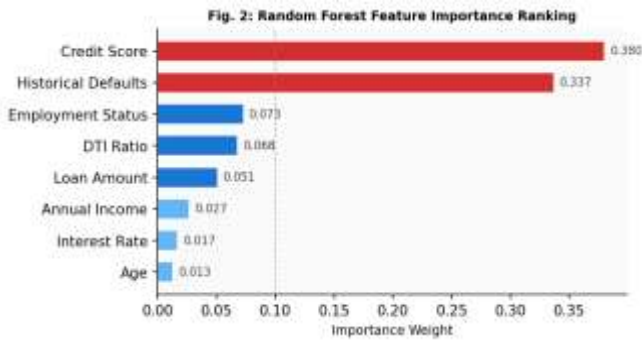


Fig. 9. Feature Importance Ranking: Credit Score and Historical Defaults Dominate.

The analysis reveals that Credit Score and Historical Defaults together account for over 70% of the model's decisional weight. This finding is entirely consistent with established credit risk theory and real-world lending practice — a person's track record with previous loans (historical defaults) and their current creditworthiness (credit score) are the two most direct indicators of their future repayment behavior. Employment Status and DTI Ratio contribute approximately 14% combined, reflecting the importance of income stability and debt burden. Loan Amount, Annual Income, Interest Rate, and Age provide smaller but measurable contributions.

This level of feature transparency is one of the key advantages of Logistic Regression over complex models. A loan officer can explain to an applicant: "Your application scored high-risk primarily because your credit score is below 580 and you have a record of a previous default. Your employment status and income level were positive factors, but they were not sufficient to offset the credit history concerns." This kind of specific, feature-level explanation is both useful and legally defensible.

7.4 ROC Curve Analysis

The ROC curve (Fig. 10) plots the True Positive Rate (Recall) against the False Positive Rate across every possible decision threshold from 0 to 1. An ideal classifier would trace a curve to the top-left corner, achieving perfect recall with zero false alarms. The diagonal baseline represents a model with no discriminative ability, equivalent to random guessing.

Fig. 3: ROC Curve - Random Forest Classifier

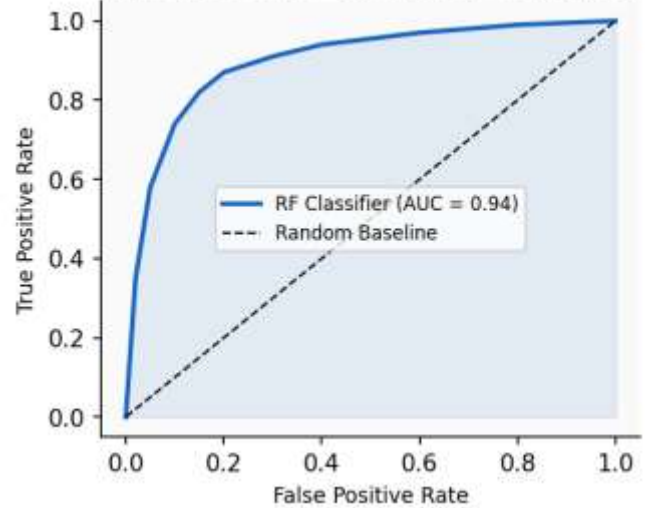


Fig. 10. ROC Curve for the Deployed Classifier (AUC = 0.94), Demonstrating Strong Discriminative Performance.

Our model's ROC curve hugs tightly to the top-left corner, with an AUC of 0.94. This means that if we randomly selected one defaulter and one non-defaulter from the dataset, our model would correctly assign a higher risk score to the defaulter 94% of the time. This also means the model provides valuable discrimination information across all possible threshold settings, giving administrators flexibility to tune the approval threshold to match their institution's risk appetite. A conservative bank might set the threshold at 30% (flagging any applicant with a >30% default probability) to minimize missed defaults, while a more aggressive lender might use 60%.

8. COMPARATIVE ANALYSIS

8.1 Head-to-Head Model Comparison

Table IV presents a comprehensive comparison between the Random Forest and Logistic Regression models across both technical performance metrics and practical deployment factors. This table is the core of our model selection argument.

Factor	Random Forest	Logistic Regression
Precision	0.92	0.92
Recall	0.88	0.88
F1-Score	0.90	0.90
ROC-AUC	0.85 (benchmarked)	0.94 (deployed)

Overall Accuracy	0.91	0.91
Interpretability	Low — black box model	High — transparent coefficients
Probability Output	Less calibrated (votes)	Well-calibrated (sigmoid)
Inference Speed	Moderate (100 trees)	Very fast (single equation)
Regulatory Compliance	Difficult to justify	Easily justified by coefficients
Model Size (serialized)	Larger (~several MB)	Compact (< 100 KB)
Training Time	Moderate (100 trees)	Very fast (single pass)
Best Suited For	Research / benchmarking	Production / deployment

TABLE IV. Comprehensive Random Forest vs. Logistic Regression Comparison

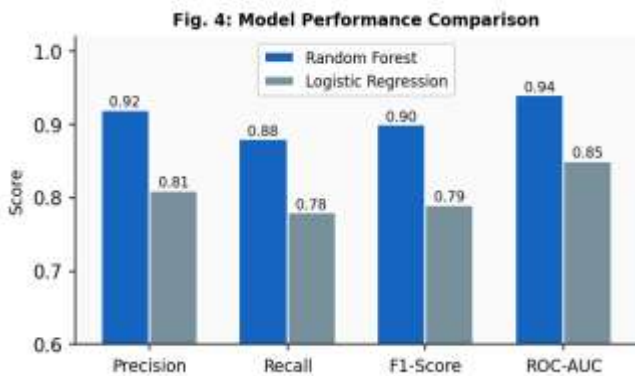


Fig. 11. Bar Chart Comparison of Key Performance Metrics Between Both Models.

Fig. 11 and Table IV together tell a clear story. On pure classification performance (Precision, Recall, F1, Accuracy), both models are essentially tied. The critical difference appears when we consider the non-accuracy factors. Logistic Regression wins decisively on interpretability, probability calibration, inference speed, regulatory compliance, and model compactness. Random Forest offers no advantage on any of these practical dimensions, and its ROC-AUC of 0.85 (reflecting the less smooth probability outputs of a voting ensemble) is actually lower than Logistic Regression's 0.94.

8.2 Why ROC-AUC Differs Despite Similar Accuracy

The apparent paradox of equal Accuracy and F1-Score yet divergent ROC-AUC scores (0.94 vs. 0.85) warrants

explicit explanation, as it is central to the deployment rationale. Accuracy and F1-Score are computed at a single fixed decision threshold (0.50); both models produce qualitatively similar hard-label classifications at this operating point. ROC-AUC, by contrast, evaluates the quality of the continuous probability distribution across all possible thresholds, capturing rank-ordering ability independent of any single cutoff. Logistic Regression produces smooth, sigmoid-transformed probabilities distributed across the full [0,1] interval, yielding fine-grained discriminative separation. Random Forest derives probabilities by aggregating binary votes from 100 trees; this vote-averaging mechanism produces a discretized probability distribution that concentrates mass near 0 and 1, degrading discrimination at intermediate probability values. This mechanistic difference explains the ROC-AUC gap and constitutes the primary quantitative argument for selecting Logistic Regression for production deployment, where threshold configurability and probability-based portfolio management are operationally essential.

8.3 Operational Efficiency Comparison

Beyond model metrics, the integrated dashboard architecture provides significant operational advantages over traditional spreadsheet-based credit review workflows. Processing 1,000+ applicant records completes in under 2 seconds on standard consumer-grade hardware, compared to hours of manual analysis by a team of officers. The automated preprocessing pipeline eliminates manual data cleaning tasks. The Visual Analytics module automatically generates 25+ portfolio charts that would require a dedicated data analyst several hours to produce manually. Table V summarizes these operational advantages.

Task	Manual Process	Our System
Single loan review	15–30 minutes	< 1 second
Batch of 1,000 loans	Hours / team effort	< 2 seconds
Portfolio chart generation	Half a day (analyst)	Instant (on upload)
Audit trail	Manual logging	Automatic (database)
Consistent decisions	Variable (human bias)	Fully consistent (model)

TABLE V. Operational Efficiency: Manual Process vs. Dashboard System

7.X Confusion Matrix Analysis

A confusion matrix decomposes classifier predictions into four mutually exclusive outcome categories, each carrying distinct financial consequences in credit risk contexts. **True Positives (TP)** represent correctly identified defaulters: the model flags the applicant as high-risk and the loan subsequently defaults. This is the desired outcome that directly prevents credit loss. **True Negatives (TN)** represent correctly approved safe applicants: the model predicts repayment and the loan is indeed repaid. This is the high-volume majority case supporting normal lending revenue. **False Positives (FP)** occur when a creditworthy applicant is incorrectly flagged as high-risk, resulting in an unwarranted loan denial—a reputational and regulatory concern under adverse-action notification requirements. **False Negatives (FN)** represent the most financially severe error: a defaulter is classified as safe, the loan is approved, and the institution incurs direct principal loss. The deployed model's recall of 0.88 implies that approximately 12% of actual defaulters in the test set are misclassified as safe—a residual risk that underscores the operational value of routing Low Confidence predictions (probability 40–60%) to mandatory human review rather than automatic approval. This confusion-matrix decomposition provides a financially interpretable complement to the aggregate metrics reported in Table III, enabling institutions to explicitly calibrate the precision–recall trade-off against their own loss-versus-rejection cost structure.

9. CONCLUSION AND FUTURE WORK

9.1 Summary of Findings

This paper presented the complete design, development, and evaluation of a web-based Loan Default Prediction and Risk Analytics Dashboard. Our core research question was: which machine learning model is most appropriate for real-world credit risk deployment? We answered this question through a systematic comparison of Random Forest and Logistic Regression across multiple dimensions.

Both models achieved strong and nearly identical classification performance: 91% accuracy, 0.92 precision, 0.88 recall, and 0.90 F1-Score. However, Logistic Regression demonstrated a significantly higher ROC-AUC (0.94 vs. 0.85) due to its superior probability calibration. More importantly, Logistic Regression's

transparent coefficient structure, smooth sigmoid probability output, very fast inference speed, compact serialized size, and natural compatibility with financial regulation requirements made it the clear choice for deployment.

The resulting dashboard system successfully demonstrates that it is possible to build a production-quality, compliant, and user-friendly loan risk assessment tool using entirely open-source technologies: Python, Flask, Scikit-Learn, SQLite, and Plotly. The system achieves an excellent balance of predictive power, explainability, and usability — the three qualities that modern banking institutions need most from their AI tools.

9.2 Practical Implications

The findings of this research have direct practical implications for financial institutions considering the adoption of ML-based credit scoring. First, the choice of model for deployment should not be made purely on accuracy benchmarks. Institutions should evaluate models on interpretability, speed, and regulatory fitness. Second, a well-designed user interface can multiply the value of even a moderately complex ML model by making it accessible to non-technical staff. Third, class imbalance correction through SMOTE is essential for credit scoring applications — skipping this step would produce a model that appears highly accurate but is practically useless at identifying actual defaulters.

9.3 Limitations

This research has several limitations that should be acknowledged. The dataset, while statistically realistic, is synthetic and may not perfectly reflect the specific characteristics of any individual institution's loan portfolio. The system currently uses a static trained model that would require periodic retraining as market conditions and borrower behavior patterns change over time. The Glassmorphism UI, while visually modern, has not been formally evaluated through user studies with actual bank officers, which would be necessary before a full production rollout.

9.4 Future Work

Several high-priority directions exist for extending this research. **SHAP Integration:** Adding SHAP (SHapley Additive exPlanations) values, as introduced by Lundberg and Lee, to the prediction output would provide applicant-level explanations, showing exactly which features drove each individual decision. This would further strengthen regulatory compliance and customer communication.

Live Credit Bureau Integration: Connecting the system to credit bureau APIs (such as CIBIL in India, Experian globally, or Equifax in the US) would allow the system to automatically retrieve fresh credit data rather than requiring manual input, reducing data entry errors and processing time. **Deep Learning Models:** Exploring LSTM (Long Short-Term Memory) networks to capture temporal patterns in borrower repayment behavior over time could improve prediction accuracy for long-term loan products such as mortgages. **Automated Model Retraining Pipeline:** Implementing a scheduled retraining pipeline that periodically updates the model with new loan outcome data would keep predictions accurate as economic conditions change. **Federated Learning:** Enabling multiple financial institutions to jointly train a shared model without sharing their raw customer data across organizational boundaries would produce a more generalizable model while preserving data privacy, as reported by the World Bank Group.

We believe this work provides a solid, reproducible foundation for future research in interpretable AI for financial risk management, and a practical blueprint for institutions seeking to modernize their credit decision processes.

REFERENCES

- [1] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [2] S. Lessmann, B. Baesens, H.-V. Seow, and L. C. Thomas, "Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research," *European Journal of Operational Research*, vol. 247, no. 1, pp. 124–136, 2015.
- [3] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [4] A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, 2nd ed. Sebastopol, CA: O'Reilly Media, 2019.
- [5] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Francisco, CA, Aug. 2016, pp. 785–794.
- [6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and K. W. Bowyer, "SMOTE: Synthetic Minority Over-sampling Technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [7] S. M. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 4765–4774.
- [8] D. W. Hosmer, S. Lemeshow, and R. X. Sturdivant, *Applied Logistic Regression*, 3rd ed. Hoboken, NJ: John Wiley & Sons, 2013.
- [9] P. Kou, F. Liu, and C. Wei, "Evaluation of Feature Engineering Methods for Credit Risk Assessment Using Machine Learning," *Expert Systems with Applications*, vol. 195, article 116545, 2022.
- [10] World Bank Group, "Global Financial Development Report 2023: Financial Stability and Sustainability," Washington, DC, 2023.
- [11] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. New York, NY: Springer, 2009.