# Project Scope: A Predictive Model for Loan Default Risk & Customer Segmentation

**Project Title:** Predicting Loan Default Risk: Using Machine Learning for Smarter & Fairer Credit Decisions

## 1. The Challenge: Navigating Modern Loan Risk

Financial institutions today are trying to assess loan default risk in an increasingly complex economy. The old-school credit scoring models just don't cut it anymore; they often miss the subtle behavioral and financial patterns in the data. This leads to two big problems: approving loans that end up defaulting, and rejecting creditworthy customers, which means missed opportunities. The core challenge is to see the risk more clearly.

## 2. The Plan: Our Proposed Solution
My goal is to build a complete, data-driven machine learning pipeline that does more than just say "yes" or "no." It will:

- **Predict Risk:** Build a robust **supervised learning model** that calculates the probability of a customer defaulting on a loan.
- **Discover Patterns:** Explore **unsupervised learning** (like clustering) to see if we can find natural groupings of customers with distinct risk profiles that aren't immediately obvious.
- **Ensure Transparency:** Integrate **explainability tools (SHAP)** to make sure the model isn't a "black box." We need to understand *why* a decision is being made. This is also crucial for fairness and ethical lending.
- **Visualize for Impact:** Deliver the final insights in an interactive **Power BI dashboard**, so business stakeholders can easily explore the findings and make data-backed decisions.

## 3. The Data: Understanding Our Building Blocks
The project will use three datasets that simulate real-world loan application scenarios, giving us a 360-degree view of each customer. I assume they contain:

- **Demographics Data:** Who the customer is (Age, Employment Status, Education Level, Bank Account Type).
- **Performance Data:** Details of the current loan (Loan Amount, Term, Dates, and the final status – our **target variable, good_bad_flag**).
- **Previous Loans Data:** A history of the customer's past borrowing behavior (Previous loan amounts, repayment dates, and any past defaults).

My central hypothesis is that a customer's **past behavior is the strongest predictor of their future behavior**.

## 4. The Process: Planned Phases & Steps

I'll break the project down into manageable phases.

Phase 1: Data Preparation & Exploration

Goal: Get the data clean, validated, and ready for modeling.

1. **Load & Merge:** Combine the three datasets into a single analytical file.
2. **Clean Up:** Handle missing values, correct anomalies (like weird ages or dates), and remove any duplicates.
3. **Initial EDA:** Get a feel for the data. I'll look at distributions, check for outliers, see how correlated different features are, and importantly, check for class imbalance in our target variable.

Phase 2: Feature Engineering & Visualization Prep

Goal: Create new, powerful features and prepare for the dashboard.

1. **Create Derived Features:** This is where the magic happens. I plan to engineer features like:
   - loan_to_income_ratio
   - payment_behavior_score (a custom score based on past repayment timeliness)
   - interest_rate
   - Risk_band
2. **Export for Power BI:** Create a clean CSV file that will feed directly into the dashboard.

Phase 3: Modeling & Evaluation

Goal: Find the most accurate and reliable predictive model.

1. **Supervised Models:** I'll start with a simple Logistic Regression as a baseline, then move to more powerful models like Random Forest, XGBoost, and **LightGBM**.
2. **Unsupervised Models:** I'll experiment with **KMeans clustering** to segment customers and use PCA/t-SNE to help visualize these groups.
3. **Evaluation:** I won't just look at accuracy. I'll use a full suite of metrics: **ROC AUC**, Precision, Recall, F1-Score, and the confusion matrix to get a complete picture.

Phase 4: Explainability, Fairness & Business Impact

Goal: Ensure the model is transparent, ethical, and valuable.

1. **SHAP Analysis:** I'll use SHAP to explain both global feature importance and individual predictions.
2. **Fairness Audit:** I'll check for any potential bias in the model's decisions across sensitive groups like age bands or employment types.
3. **Cost-Benefit Simulation:** I'll go beyond the default 0.5 decision threshold and simulate a cost-benefit analysis to find the optimal threshold that maximizes profitability for the business.

**5. Key Technologies & Libraries**
- **Primary Language:** Python
- **Core Libraries:** Pandas, NumPy, Scikit-learn, XGBoost, LightGBM, SHAP, Matplotlib, Seaborn
- **Dashboarding:** Power BI

**6. Expected Deliverables & Success Criteria**
By the end of this project, I will deliver:

- A clean, documented dataset ready for analysis.
- Jupyter notebooks with all the reproducible code.
- The final, trained machine learning model file.
- SHAP plots and a summary of the fairness audit.
- An interactive Power BI dashboard for stakeholders.

Success will be measured by achieving a **ROC AUC score above our target benchmark**, delivering actionable insights for risk management, and ensuring the final model is both fair and interpretable

# -•Project Documentation: A public article or post detailing your process, insights, and conclusions. Medium is preferred, but a detailed post on LinkedIn or Twitter is also encouraged.

## Building a Machine Learning System for Predictive Loan Default Risk

### 1. Overview
This project documents the end-to-end process of building a machine learning system to predict loan default risk. The goal was to leverage customer demographic, financial, and behavioral data to create a reliable predictive model, segment customers into risk profiles, and translate the model's outputs into actionable business insights. The journey involved navigating significant technical challenges, iterating through different machine learning approaches, and ultimately delivering a solution with tangible business value.

### 2. The Challenge & My Learning Journey
The initial goal seemed straightforward: build a classifier to predict loan defaults. However, the reality was a masterclass in persistence. My journey was defined by significant hurdles:

- **Constant Crashes:** My code crashed repeatedly, especially during the training of more complex models like XGBoost and LightGBM, often due to memory overload from large dataset operations.
- **Pipeline Breakages:** My feature engineering pipelines would unexpectedly break due to inconsistent variable naming, schema changes, or missing columns.
- **Target Leakage:** Early versions of my model produced suspiciously high scores, a classic sign of target leakage. I had to build custom checks to find and remove features that were giving the answer away to the model.

When my initial supervised and unsupervised learning attempts failed, I felt stuck. The code was unstable, and the results were unreliable. To overcome these blockers, I turned to my Axia Africa and ALX Data Science program notes for foundational concepts on building robust pipelines and leaned heavily on Kaggle kernels to find alternative, more stable implementations for SHAP explainability and LightGBM tuning. This process reminded me that real-world machine learning is messy, and the most valuable learning often happens while debugging what's broken.

**3. Methodology & Process**
**A. Data Processing & Feature Engineering**

The foundation of the model was built on three datasets: customer performance, demographics, and previous loan history.

1. Cleaning: I developed functions to handle missing values, remove duplicate entries, and validate data integrity against business rules (e.g., totaldue >= loanamount).
2. Feature Engineering: This was where the model's predictive power truly came from. I engineered several new features, including:
   - Financial Ratios: loan_to_income_ratio, debt_to_income_ratio.
   - Behavioral Scores: repayment_rate, avg_repayment_delay aggregated from past loans.
   - Composite Risk Scores: A custom payment_behavior_score to summarize a customer's history.

**B. Modeling Approaches: Supervised vs. Unsupervised Learning**

- **Unsupervised Learning (Attempted):** I first explored K-Means clustering to segment customers into natural groups based on their behavior. While I managed to identify 5 distinct clusters (e.g., "Young Borrowers, Small Loans"), the process was plagued by instability due to data scaling issues and missing values, causing frequent crashes.
- **Supervised Learning (Successful):** I pivoted to focus on a supervised classification task. I built a robust scikit-learn pipeline that handled preprocessing and modeling in a single, clean workflow. I trained and evaluated over five algorithms, including Logistic Regression, Random Forest, Gradient Boosting, XGBoost, and ultimately **LightGBM**, which provided the best balance of speed and performance. The final model achieved a **ROC AUC of 82.4%**.

**4. Key Insights & Business Impact**

Merely building a model isn't enough; its value lies in its interpretation and business application.

- **Top Predictors:** Using **SHAP analysis**, I was able to look inside the model's "black box." The most influential predictors were not just simple demographics but behavioral features like **previous loan repayment rates** and financial health indicators like **employment status**. Feature engineering alone boosted model performance by 8 AUC points.
- **Optimizing for Profit, Not Accuracy:** A standard 0.5 classification threshold is rarely optimal for business. I developed a cost-benefit analysis that assigned a monetary cost to false positives (approving a bad loan) and false negatives (rejecting a good loan). This revealed that the most profitable decision threshold was **0.23**, a strategy estimated to reduce default-related costs by **23%**.
- **Actionable Risk Categories:** The model's probability scores were used to segment customers into clear risk bands (Very Low, Low, Medium, High, Very High), allowing loan officers to automate and streamline their decision-making process.

### 5. Conclusion & Next Steps

This project was a powerful lesson in resilience. By working through technical failures and leveraging external resources, I successfully built a complete machine learning system that not only predicts loan defaults with high accuracy but also provides interpretable, cost-effective, and actionable insights for the business.

The model is now packaged for deployment, with a Power BI-ready dataset for stakeholder dashboards. The next steps are to explore deploying this model as a real-time scoring API and further refining the customer clusters for targeted marketing strategies.

**Tech Stack:** Python | scikit-learn | XGBoost | LightGBM | SHAP | pandas | Power BI