

# Predictive Model Plan

## 1. Model Logic (Generated with GenAI)

### Model Objective:

The model aims to predict whether a customer will become delinquent (1) or not (0) based on financial behavior and account attributes. It will support Geldium's risk management by identifying high-risk customers in advance.

### Model Type:

- Primary model: **Logistic Regression**
- Backup model (optional): **Random Forest Classifier**

### Step-by-Step Predictive Pipeline Logic:

#### 1. Data Collection

Load the customer-level dataset with features like income, credit score, missed payments, and delinquency status.

#### 2. Feature Selection

Choose important features that are relevant to risk:

- Credit\_Utilization
- Missed\_Payments
- Debt\_to\_Income\_Ratio
- Account\_Tenure
- Credit\_Score

#### 3. Data Cleaning and Preprocessing

- Handle missing values: median imputation for Income, Loan\_Balance, and Credit\_Score
- Normalize numerical features
- One-hot encode categorical variables such as Employment\_Status and Credit\_Card\_Type

#### 4. Model Training

- Use **Logistic Regression** to train on the features and learn the relationship with Delinquent\_Account
- Optionally train a **Random Forest** to compare performance

#### 5. Prediction

- Model outputs a probability (e.g., 0.75) indicating the likelihood of delinquency
- Apply a threshold (e.g., 0.5) to convert probability into binary classification (0 or 1)

#### 6. Model Evaluation

- Use metrics such as Accuracy, Precision, Recall, AUC, and F1 Score
- Run fairness checks across groups (e.g., employment type)

#### Pseudo-Code :

##### # Step 1: Load and clean the data

```
data = load_data("delinquency_dataset.csv")
```

```
data = impute_missing(data)
```

```
data = encode_categorical(data)
```

##### # Step 2: Split and train

```
X_train, X_test, y_train, y_test = train_test_split(data.features, data.labels)
```

```
model = LogisticRegression()
```

```
model.fit(X_train, y_train)
```

##### # Step 3: Predict

```
y_pred = model.predict(X_test)
```

##### # Step 4: Evaluate

```
evaluate_model(y_test, y_pred)
```

## 2. Justification for Model Choice

We selected Logistic Regression as the primary model for predicting customer delinquency due to its unique balance of performance, interpretability, and industry alignment.

- **Accuracy**

- Logistic regression is well-suited for binary classification tasks like predicting Delinquent\_Account.
- It performs well on medium-sized datasets and can achieve competitive accuracy when the relationship between features and the target is linear or monotonic.

- **Transparency**

- Logistic regression offers clear, explainable outputs, showing the direction and weight of each feature's impact on the prediction.
- This is critical for Geldium's internal teams and stakeholders, especially when explaining decisions to auditors or regulators.

- **Ease of Use**

- It is easy to implement and requires relatively low computational power.
- The model is fast to train and integrates well into most existing data pipelines.

- **Relevance for Financial Prediction**

- Logistic regression is widely adopted in the financial industry due to its interpretability and regulatory friendliness.
- It aligns with common credit risk models used in banks and lending institutions.

- **Suitability for Geldium's Business Needs**

- Geldium requires a transparent, auditable, and fair model to assess customer credit risk.
- The logistic model supports probability-based risk scoring, which can help in tiering customers and prioritizing interventions.
- If needed, we can further test advanced models like Random Forests for performance improvement, while keeping Logistic Regression as the explainable baseline.

### 3. Evaluation Strategy

To ensure the model is both effective and fair, we will evaluate it using a combination of performance metrics, bias detection strategies, and ethical considerations.

#### ➤ Evaluation Metrics

We will use the following metrics to assess model performance:

- **Accuracy:** Measures overall correctness. Useful when classes are balanced.
- **Precision:** The proportion of customers predicted as delinquent who actually are delinquent — important to reduce false positives.
- **Recall:** The proportion of actual delinquents correctly identified — crucial for identifying risky customers.
- **F1 Score:** Harmonic mean of precision and recall, especially valuable for imbalanced datasets.
- **AUC (Area Under the ROC Curve):** Evaluates the model's ability to distinguish between delinquent and non-delinquent customers.

#### ➤ Interpretation of Metrics

- A high **recall** indicates the model is catching most of the risky cases.
- High **precision** shows it's not falsely labeling good customers as risky.
- **F1 score** helps ensure a good balance between precision and recall.

- A high **AUC** (closer to 1) means the model reliably separates risk vs. no-risk customers.

➤ **Bias Detection & Mitigation**

- Evaluate model performance across **different demographic or employment groups** to ensure no group is disproportionately flagged.
- Use tools like **confusion matrix disaggregation** and **demographic parity tests** to spot inconsistencies.
- If bias is detected, apply mitigation strategies such as:
  - Adjusting classification thresholds
  - Reweighting the training data
  - Removing or auditing sensitive features

➤ **Ethical Considerations**

- Predictions should be used to **support and assist** at-risk customers, not penalize them.
- The model must not reinforce existing financial inequalities or unfairly disadvantage certain customer groups.
- Transparency in model decision-making is critical to maintain trust and comply with financial regulations.