# Explainable AI Module for Financial Machine Learning Models

Your Name

July 2, 2025

## 1. Motivation: Interpretability in Finance

Financial institutions must not only build models with high predictive accuracy, but also ensure they are interpretable and *explainable* for:

- **Regulatory compliance:** Models (e.g., for credit approval) must satisfy regulators that decisions are fair and non-discriminatory.

- **Stakeholder trust:** Loan officers, auditors, and customers need human-readable reasons behind each prediction.

- **Risk management:** Understanding drivers of risk scores helps in debugging and improving models.

## 2. SHAP and LIME: Local Feature Attributions

**SHAP (SHapley Additive exPlanations):** Based on cooperative game theory. Computes each feature's contribution to the prediction by averaging over all coalitions of features.
*Advantages:* Global consistency, exact for tree models (TreeSHAP), additive explanation.

**LIME (Local Interpretable Model-agnostic Explanations):** Fits a sparse, interpretable surrogate model (e.g., linear) around the neighborhood of a single instance.
*Advantages:* Model-agnostic, intuitive for small perturbations.

## 3. Algorithm Outline

1. **Data preparation:** Standardize or encode financial features (age, income, credit_score, debt, . . . ).

2. **Model training:** Fit a black-box classifier (e.g., XGBoost, RandomForest, LogisticRegression).

3. **Choose explainer:**
   - **SHAP:** Use `shap.TreeExplainer` for tree-based, or `shap.KernelExplainer` otherwise.
   - **LIME:** Use `lime.lime_tabular.LimeTabularExplainer`.

4. **Compute explanations:**
   - *SHAP:*

   $$\hat{f}(x) = \phi_0 + \sum_{i=1}^{d} \phi_i, \quad \phi_i = \sum_{S \subseteq \{1,\ldots,d\} \setminus \{i\}} \frac{|S|! \, (d - |S| - 1)!}{d!} \big[ f(x_{S \cup \{i\}}) - f(x_S) \big].$$

   - *LIME:*

   $$\underset{g \in G}{\arg\min} \, \mathcal{L}\big(f, \, g, \, \pi_x\big) + \Omega(g),$$

   where $\pi_x$ is a proximity kernel around $x$.

5. **Visualize results:** Force plots, summary beeswarm (SHAP), or bar charts (LIME).

6. **Integrate into pipeline:** Generate explanations for new loan applications automatically.

# 4. Example Integration in Python

```python
import numpy as np
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
import shap
import matplotlib.pyplot as plt

# 1) Synthetic financial dataset
np.random.seed(0)
n = 500
X = pd.DataFrame({
    'Age': np.random.randint(21, 70, n),
    'Income': np.random.normal(60000, 15000, n),
    'Credit_Score': np.random.randint(300, 850, n),
    'Debt': np.random.normal(15000, 5000, n),
    'Years_Employed': np.random.randint(0, 40, n)
})
# Binary target: risk (1=high risk, 0=low risk)
y = (0.3*(X['Debt']/X['Income']) +
     0.2*(800 - X['Credit_Score'])/500 +
     0.1*(40 - X['Years_Employed'])/40 +
     np.random.randn(n)*0.05 > 0.3).astype(int)

# 2) Train/test split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42
)

# 3) Fit black-box model
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# 4) Explain predictions with SHAP
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)

# 5) Global summary plot
shap.summary_plot(shap_values[1], X_test)

# 6) Force plot for a single sample
idx = 0
shap.initjs()
force_plot = shap.force_plot(
    explainer.expected_value[1],
    shap_values[1][idx],
    X_test.iloc[idx]
)
```