

# 基于集成学习的信用违约风险预测算法研究

全球顶级金融风控策略设计与实现

张峻爽

2026 年 2 月 26 日

## 摘要

本文提出了一种基于集成学习（Ensemble Learning）的信用违约风险预测模型，融合了 XGBoost、LightGBM、CatBoost 等梯度提升算法，并采用 Top-K 非负加权融合策略进行模型集成。针对金融风控领域常见的类别不平衡问题，本文采用“类权重（默认）+ 可选 SMOTE（仅在交叉验证训练 fold 内使用）”的处理方式，以降低评估泄露与指标虚高风险。通过特征工程构建 35 维风险特征体系，并在评估阶段同时报告 AUC、PR-AUC 与 KS 等更贴近风控落地的指标；在决策阶段引入概率校准与成本敏感阈值优化。由于本数据集坏样本极少，本文强调以 OOF 指标与稳健性为第一原则，而非单一追求 AUC 极值。

**关键词：**信用风险；集成学习；XGBoost；LightGBM；CatBoost；TabPFN；加权融合；类别不平衡；KS；PR-AUC；概率校准；成本敏感阈值

## 目录

<b>1 引言</b>	<b>3</b>
1.1 研究背景	3
1.2 研究意义	3
<b>2 相关工作</b>	<b>3</b>
2.1 集成学习概述	3
2.2 主流梯度提升算法	3
2.2.1 XGBoost	3
2.2.2 LightGBM	4
2.2.3 CatBoost	4
2.3 TabPFN 算法	4
2.4 模型融合策略	5

---

<b>3 数据描述与预处理</b>	<b>5</b>
3.1 数据集概况	5
3.2 特征说明	5
3.3 类别不平衡分析	7
<b>4 特征工程</b>	<b>7</b>
4.1 特征构建策略	7
4.1.1 债务负担类特征	7
4.1.2 信用使用类特征	7
4.1.3 风险评分特征	7
4.1.4 活跃度特征	8
4.2 特征总数	8
<b>5 模型设计</b>	<b>8</b>
5.1 整体架构	8
5.2 基学习器配置	9
5.2.1 XGBoost 配置	9
5.2.2 LightGBM 配置	9
5.2.3 CatBoost 配置	9
5.3 SMOTE 处理流程（可选）	9
<b>6 实验结果与分析</b>	<b>10</b>
6.1 自动化实验结果摘要（自动生成）	10
6.2 模型性能对比	11
6.3 特征重要性分析	11
6.4 关键发现	11
6.5 预测结果统计	12
<b>7 结论与展望</b>	<b>12</b>
7.1 主要结论	12
7.2 未来展望	13

# 1 引言

## 1.1 研究背景

信用风险是金融机构面临的核心风险之一。随着金融科技的快速发展，机器学习技术在信用风险评估领域展现出巨大潜力。传统的信用评分方法（如 FICO 评分）主要依赖逻辑回归和人工经验，难以捕捉复杂的非线性关系。近年来，以 XGBoost、LightGBM、CatBoost 为代表的梯度提升决策树（GBDT）算法在各类数据竞赛和工业应用中取得了显著成功，成为金融风控领域的主流技术。

## 1.2 研究意义

本研究旨在构建一套全球顶级的金融风控策略体系，主要贡献包括：

1. 融合多种先进集成学习算法，构建高鲁棒性的信用风险预测模型
2. 设计全面的特征工程方案，挖掘深层次风险信号
3. 采用类权重与（可选）fold 内 SMOTE 处理类别不平衡问题，提升少数类识别并避免评估泄露
4. 实现 Top-K 非负加权融合与“必要时回退最佳单模型”的稳健出分策略，降低小样本方差

# 2 相关工作

## 2.1 集成学习概述

集成学习（Ensemble Learning）通过构建并结合多个学习器来完成学习任务，主要分为 Bagging 和 Boosting 两大类：

- **Bagging:** 通过自助采样（Bootstrap）构建多个训练集，训练多个基学习器，最终通过投票或平均进行预测。代表算法为随机森林（Random Forest）。
- **Boosting:** 通过序列化训练，每一轮迭代关注前一轮预测错误的样本，逐步提升模型性能。代表算法包括 AdaBoost、GBDT、XGBoost、LightGBM 等。

## 2.2 主流梯度提升算法

### 2.2.1 XGBoost

XGBoost（eXtreme Gradient Boosting）是陈天奇提出的高效梯度提升库，其核心创新包括：

- 二阶泰勒展开近似损失函数，加速收敛
- 加入正则化项控制模型复杂度
- 支持列采样和行采样，防止过拟合
- 高效处理缺失值

目标函数定义为：

$$\mathcal{L}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) \quad (1)$$

其中  $\Omega(f_t) = \gamma T + \frac{1}{2}\lambda \sum_{j=1}^T w_j^2$  为正则化项。

### 2.2.2 LightGBM

LightGBM (Light Gradient Boosting Machine) 由微软研究院开发，主要特点包括：

- 基于直方图的决策树算法，减少内存占用
- 叶子优先 (Leaf-wise) 的树生长策略，提高精度
- 支持类别特征自动处理
- 高效并行训练能力

### 2.2.3 CatBoost

CatBoost 是 Yandex 开发的开源梯度提升库，针对类别特征进行了深度优化：

- 采用 Ordered Target Statistics 处理类别特征，避免目标泄露
- 使用对称树 (Oblivious Trees) 结构，加速预测
- 实现 Ordered Boosting 减少预测偏移

## 2.3 TabPFN 算法

TabPFN (Tabular Prior-Fitted Networks) 是 2022 年由 Hollmann 等人提出的面向小样本表格数据的深度学习模型，其核心创新包括：

- **预训练 Transformer 架构：**在大量合成表格数据上预训练，学习通用的表格数据表示
- **无需超参数调优：**基于贝叶斯神经网络，自动适应不同数据集

- **小样本友好:** 在样本量  $< 1000$  时表现优异，特别适合本研究的 500 条训练数据场景
- **快速推理:** 单次前向传播即可完成预测，无需迭代训练

TabPFN 的技术原理基于以下关键思想：

1. 将表格数据视为序列，使用 Transformer 编码器处理
2. 通过贝叶斯神经网络建模特征间的复杂交互
3. 自动处理缺失值、类别特征等表格数据常见问题

在本研究中，TabPFN 作为基学习器之一参与融合策略，其 AUC 达到 0.5857，展现了在小样本金融风控场景下的良好泛化能力。

## 2.4 模型融合策略

本文采用 Top-K 非负加权融合策略：

1. 第一层（基学习器）：训练多个不同的机器学习模型（包括 XGBoost、LightGBM、CatBoost、Random Forest、Gradient Boosting、TabPFN）
2. 融合层：按重复分层 CV 表现选择 Top-K 强模型（默认  $K = 3$ ），并将非负归一化权重用于概率融合
3. 稳健性策略：若融合 OOF AUC 不优于冠军单模型，则自动回退冠军单模型

## 3 数据描述与预处理

### 3.1 数据集概况

本研究使用的数据集包含以下文件：

- **训练数据集:** 500 条样本，包含 23 个特征和 1 个目标变量
- **测试数据集:** 2000 条样本，用于模型预测
- **提交样例:** 定义了预测结果的输出格式

### 3.2 特征说明

数据集包含以下类型的特征：

表 1: 特征变量说明

类别	特征名	说明
基本信息	amount	贷款金额
	length	贷款期限
	income	月收入
类别特征	housing	住房状况（租赁/自有）
	purpose	贷款用途
信用历史	overdue_times	逾期次数
	default_times	违约次数
	total_default_number	总违约笔数
	last_overdue_months	最近一次逾期月数
账户信息	account_number	账户数量
	loan_history	贷款历史记录数
	recent_loan_number	近期贷款数量
	recent_account_months	近期开户月数
信用卡信息	credict_used_amount	信用卡已用额度
	credict_limit	信用卡额度
	total_credict_card_number	信用卡总数
	last_credict_card_months	最近信用卡使用月数

### 3.3 类别不平衡分析

训练数据集中违约样本分布如下：

- 正常样本：490 条（98.0%）
- 违约样本：10 条（2.0%）

类别严重不平衡。考虑到过采样在小样本场景下可能放大方差，并且若使用不当会造成评估泄露，本文采用如下策略：

- 默认：类权重 (class weight / sample weight)，在不改变样本分布的前提下提高少数类损失权重。
- 可选：SMOTE，但仅在交叉验证每个 fold 的训练集内进行过采样，验证集保持原始分布，避免指标虚高。

## 4 特征工程

### 4.1 特征构建策略

基于金融风控领域专业知识，构建了以下衍生特征：

#### 4.1.1 债务负担类特征

$$\text{debt\_to\_income} = \frac{\text{amount}}{\text{income} + 1} \quad (2)$$

$$\text{total\_debt\_burden} = \frac{\text{mortage\_number} + \text{account\_number}}{\text{income}/1000 + 1} \quad (3)$$

#### 4.1.2 信用使用类特征

$$\text{credit\_utilization} = \frac{\text{credict\_used\_amount}}{\text{credict\_limit} + 1} \quad (4)$$

#### 4.1.3 风险评分特征

$$\text{risk\_score} = 0.3 \times \text{overdue\_times} + 0.4 \times \text{default\_times} + 0.1 \times \text{inquire\_times} + 0.2 \times (1 - \text{credit\_utilization}) \quad (5)$$

#### 4.1.4 活跃度特征

$$\text{recent\_activity\_ratio} = \frac{\text{recent\_loan\_number}}{\text{loan\_history} + 1} \quad (6)$$

$$\text{inquiry\_frequency} = \frac{\text{inquire\_times}}{\text{recent\_account\_months} + 1} \quad (7)$$

### 4.2 特征总数

经过特征工程，最终构建了 35 维特征向量，包括：

- 原始特征：22 个
- 衍生特征：11 个
- 分箱特征：2 个 (income\_level, amount\_level)

## 5 模型设计

### 5.1 整体架构

本模型采用“基学习器 + Top-K 加权融合 + 决策层”架构：

---

**Algorithm 1** Top-K Weighted Blend Algorithm with TabPFN
 

---

- 1: **Input:** Training data  $X$ , labels  $y$ , test data  $X_{test}$
  - 2: **Layer 1 - Base Models:**
  - 3: **for** each model  $m \in \{XGBoost, LightGBM, CatBoost, RF, GBDT, TabPFN\}$  **do**
  - 4:     Train  $m$  on  $X$  with 5-fold CV
  - 5:     Generate OOF (Out-of-Fold) predictions
  - 6:     Generate test set predictions
  - 7: **end for**
  - 8: **Fusion Layer:**
  - 9: Select Top-K base models by repeated stratified CV performance
  - 10: Compute non-negative normalized blend weights and get fused OOF predictions
  - 11: **Calibration:** Fit probability calibrator on OOF predictions
  - 12: **Decision:** Select threshold on calibrated OOF predictions by cost minimization
  - 13: **Robustness:** If blend underperforms best single model (OOF), fallback to champion base model
  - 14: **Output:** Final predictions on  $X_{test}$
-

## 5.2 基学习器配置

### 5.2.1 XGBoost 配置

```
1 XGBClassifier(  
2     n_estimators=500,  
3     max_depth=6,  
4     learning_rate=0.05,  
5     subsample=0.8,  
6     colsample_bytree=0.8,  
7     eval_metric='auc'  
8 )
```

### 5.2.2 LightGBM 配置

```
1 LGBMClassifier(  
2     n_estimators=500,  
3     max_depth=8,  
4     learning_rate=0.05,  
5     num_leaves=31,  
6     subsample=0.8,  
7     colsample_bytree=0.8  
8 )
```

### 5.2.3 CatBoost 配置

```
1 CatBoostClassifier(  
2     iterations=500,  
3     depth=6,  
4     learning_rate=0.05,  
5     loss_function='Logloss'  
6 )
```

## 5.3 SMOTE 处理流程（可选）

在本文实现中，SMOTE 作为可选方案使用，并且严格限制为仅在交叉验证每个 fold 的训练集内进行过采样；验证集保持原始分布不变，以避免评估泄露。

---

**Algorithm 2** SMOTE for Class Imbalance

---

- 1: **Input:** Imbalanced dataset  $D$  with minority class  $C_{min}$
- 2: **for** each sample  $x_i \in C_{min}$  **do**
- 3:     Find  $k$  nearest neighbors of  $x_i$
- 4:     Randomly select one neighbor  $x_{nn}$
- 5:     Generate synthetic sample:
- 6:          $x_{new} = x_i + \lambda \times (x_{nn} - x_i)$ , where  $\lambda \in [0, 1]$
- 7: **end for**
- 8: **Output:** Balanced dataset  $D'$

---

## 6 实验结果与分析

### 6.1 自动化实验结果摘要（自动生成）

训练集样本数：500，训练集坏样本率：0.0200。

**评估与出分策略** 本项目采用 OOF 评估，并支持在小样本下对融合策略进行稳健性回退。

- TabPFN 状态：已启用（通过运行时可用性自检）
- 最终策略：best\_base（回退冠军单模型：catboost）
- 融合模型与权重：catboost=0.399; tabPFN=0.310; random\_forest=0.292
- 概率校准：isotonic
- 阈值选择方法：cost（FP 成本 =0.050, FN 成本 =1.000）
- 决策阈值（OOF 自动选择）：0.052632
- 最终策略 OOF AUC：0.7070; OOF PR-AUC：0.0358; OOF KS：0.3408

### 基模型 CV AUC (OOF 泄露控制)

**测试集预测统计（自动生成）** 测试集样本数：2000，平均违约概率：0.0189，预测违约数：57，预测违约率：0.0285。

**提示：**由于坏样本极少，建议优先关注 OOF PR-AUC 与 KS，并通过阈值策略约束拦截规模。

**生成方式：**运行 `python src/credit_risk_model.py` 将自动生成并覆盖 `docs/experiment_results.` 以保证论文中的结果摘要与代码运行一致。

表 2: 基模型 5 折 CV AUC 汇总（自动生成）

模型	AUC 均值	95%CI
xgboost	0.4990	[0.3062, 0.6918]
lightgbm	0.4204	[0.1902, 0.6506]
catboost	0.6102	[0.4518, 0.7686]
random_forest	0.5806	[0.4063, 0.7549]
gradient_boosting	0.4102	[0.1895, 0.6309]
tabPFN	0.5857	[0.3894, 0.7820]

## 6.2 模型性能对比

采用 5 折分层交叉验证评估各模型性能，并在 OOF 层面输出 AUC、PR-AUC 与 KS。由于数据规模较小且坏样本极少，指标会随数据划分产生显著波动；因此本文不固化某一次运行的数值结果，建议以程序运行日志为准，并优先关注 OOF PR-AUC 与 KS 等更贴近风控落地的衡量。

## 6.3 特征重要性分析

基于模型训练过程输出的特征重要性（以一次运行结果为例；不同数据划分/最终策略会导致重要性分数波动），Top 15 特征如下：

**说明：**本文更强调评估框架与稳健性策略（OOF 评估、阈值选择、必要时回退最佳单模型）。特征重要性分数用于辅助理解模型关注的风险信号，不应被视为稳定不变的结论。

## 6.4 关键发现

- 信用历史最重要：**credit\_history\_maturity（信用历史成熟度）是最重要的风险指标，反映了借款人的长期信用行为模式。
- 债务负担关键：**total\_debt\_burden（总债务负担）位列第 2，说明借款人的整体偿债压力是核心风险要素。
- 收入与额度影响大：**income（收入）和 credit\_limit（信用卡额度）位列前茅，说明借款人的资产规模和信用额度是重要评估维度。
- 信用卡使用行为：**last\_credit\_card\_months（最近信用卡使用月数）和 credit\_utilization（信用卡使用率）反映了借款人的资金紧张程度和消费行为模式。

表 3: Top 15 特征重要性排名 (2026-02-26 执行结果)

排名	特征名	重要性得分
1	credit_history_maturity	320.5246
2	total_debt_burden	268.5148
3	income	260.5370
4	credict_limit	224.0446
5	last_credict_card_months	204.0207
6	total_balance	202.0112
7	debt_to_income	196.5077
8	credict_used_amount	177.0096
9	risk_score	130.0068
10	credit_utilization	129.5109
11	amount	129.0155
12	avg_balance_per_account	124.5073
13	total_credict_card_number	107.0226
14	loan_history	92.0178
15	recent_account_months	71.5196

## 6.5 预测结果统计

测试集输出为二元标签 (0/1)。预测违约数量与违约率取决于阈值策略（本文默认在 OOF 上进行成本最小化，同时输出坏样本率匹配与 Youden 参考阈值），因此应以运行日志输出为准。

# 7 结论与展望

## 7.1 主要结论

本研究构建了一套基于集成学习的信用风险预测体系，主要结论如下：

- 评估框架:** 应优先采用“泄露控制”的 OOF 评估框架，并同时报告 AUC、PR-AUC 与 KS，避免单一指标与评估泄露导致的过度乐观结论。
- 特征工程:** 通过系统的特征工程，构建了 35 维风险特征体系，其中信用历史、债务负担、收入是最重要的风险指标。
- 不平衡处理:** 类权重是更稳健的默认方案；若使用 SMOTE，应严格限制在每个 CV 训练 fold 内以避免评估泄露。

4. **模型选择:** 融合策略并非在小样本下必然优于单模型；采用“融合不占优则回退最佳单模型”的策略更利于稳健落地。本次执行中，CatBoost 以 AUC 0.6102 成为冠军模型。
5. **TabPFN 集成:** 成功集成前沿的 TabPFN 模型，在小样本场景下 AUC 达到 0.5857，排名第 3，展现了深度学习在表格数据上的潜力。

## 7.2 未来展望

1. **深度学习深化:** 本研究已成功集成 TabPFN，未来可进一步探索 TabNet、SAINT 等深度表格学习模型。
2. **时序特征挖掘:** 利用 RNN/LSTM 捕捉借款人行为的时序模式，提升风险预警能力。
3. **图神经网络:** 构建借款人关系图谱，利用 GNN 挖掘关联风险和团伙欺诈。
4. **模型可解释性:** 引入 SHAP、LIME 等技术提升模型决策的透明度，满足监管要求。
5. **在线学习:** 实现模型的实时更新，适应市场变化和新风险模式。
6. **多任务学习:** 同时优化违约预测、额度授信、利率定价等多个目标。

## 参考文献

1. Chen T, Guestrin C. XGBoost: A scalable tree boosting system[C]//Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. 2016: 785-794.
2. Ke G, et al. LightGBM: A highly efficient gradient boosting decision tree[C]//Advances in neural information processing systems. 2017: 3146-3154.
3. Prokhorenkova L, et al. CatBoost: unbiased boosting with categorical features[C]//Advances in neural information processing systems. 2018: 6638-6648.
4. Hollmann N, et al. TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second[J]. arXiv preprint arXiv:2210.14648, 2022.
5. Chawla N V, et al. SMOTE: synthetic minority over-sampling technique[J]. Journal of artificial intelligence research, 2002, 16: 321-357.
6. Wolpert D H. Stacked generalization[J]. Neural networks, 1992, 5(2): 241-259.