

基于特征相似度降采样的交易欺诈预测模型

邹勇 林芃 马振伟

(中国银联股份有限公司 上海 200135)

摘要 在交易欺诈评估场景下,正负样本比例极其悬殊,需要对样本进行采样来解决样本不平衡。传统采样由于在采样过程中会丢失样本信息,导致模型预测的准确率不是很高。针对这类情况,提出一种基于特征相似度降采样方式的模型构建方法。该方法主要包括三个部分。(1) 依据样本数据,构建有效的与欺诈相关的特征集。(2) 通过引入样本差异度函数,在降采样时尽可能多地保留样本信息。(3) 通过多个分类器进行融合输出欺诈概率。将该方法与其他常见采样方式进行对比,实验结果表明,该方法具有更好的评估结果。

关键词 降采样 欺诈预测 集成学习

中图分类号 TP394.1

文献标志码 A

DOI:10.3969/j.issn.1000-386x.2024.04.016

TRANSACTION FRAUD PREDICTION MODEL BASED ON FEATURE SIMILARITY DOWN-SAMPLING

Zou Yong Lin Peng Ma Zhenwei

(China UnionPay Co., Ltd., Shanghai 200135, China)

Abstract In the scenario of transaction fraud evaluation, the proportion of positive and negative samples is extremely different, so it is necessary to sample the samples to solve the sample imbalance. Due to the loss of sample information in the traditional sampling process, the accuracy of model prediction is not very high. Aimed at this kind of situation, a model construction method based on feature similarity down-sampling is proposed. This method mainly included three parts. (1) According to the sample data, an effective feature set related to fraud was constructed. (2) By introducing the sample difference function, as much sample information as possible was retained when down-sampling. (3) Multiple classifiers were fused to output the fraud probability. This method was compared with other common sampling methods. Experimental results show that this method has better evaluation results.

Keywords Down-sampling Fraud prediction Integrated learning

0 引言

过去交易欺诈评估^[1]主要是基于传统专家知识和经验^[2],随着大数据和人工智能时代的到来,基于大数据和机器学习的量化评估模型^[3]已经应用到交易欺诈评估场景。在交易欺诈评估场景中,欺诈交易与正常交易的比例往往非常悬殊,是典型的非平衡数据集。直接在非平衡数据集上构建模型,一方面由于需要足够的欺诈样本数据,对应训练的正常交易数据就会非常大,需要大量的计算资源以及很长的训练时间;另一方面悬殊的样本比例会导致最后的模型出现过拟合^[4]

的情况。因此在数据建模过程中,需要对样本进行降采样处理,降低正负样本比例。传统的降采样方式会出现样本信息的丢失或者样本信息采样不够全面,导致模型学习不到有区分度的特征,进而影响模型效果。因此本文针对交易欺诈评估场景,提出一种基于特征相似度的改进采样方式,以该方法采样出的样本为训练集,运用集成学习算法^[5]构建交易欺诈评估模型。

1 欺诈评估模型构建

首先,本文着眼于机器学习的交易欺诈评估模型,基于历史的交易数据和标注的欺诈样本数据,提取有

效的特征集,筛选和欺诈相关的特征;然后通过样本差异度函数对样本进行有选择的降采样,通过分类算法结合采样后的数据进行模型训练;最后,对训练完成的交易欺诈评估模型验证,并与基于随机降采样数据训练出来的模型进行比对。

1.1 欺诈特征

特征是指区分不同交易类型的本质特点,在交易欺诈评估场景中,需要挖掘能够区分交易是否为欺诈的特征^[6],提取的特征区分度越高,模型识别欺诈交易的效果越好。

交易数据中包含的要素有:交易卡号、交易金额、交易时间等信息。本文根据交易数据,提取了众多的特征,大体的方法主要分为三类:

(1) 个性化的交易特征。每张卡片都有其习惯的交易行为特点,例如交易金额、习惯性的交易时间、经常消费的地点等,这些都属于卡片个性化的交易行为特征,本文按照卡片维度统计卡片的个性化特征,包括均值、方差、最大值、最小值等。

(2) 社会化的交易特征。不同交易的交易环境也有其不同的特点,高端商场和街边便利店交易的特点会有明显的区别,不同国家和地区风险形势也大不相同,这些都是更宏观层面行为特点,同样可以采用统计方法,构建社会化的交易特征。

(3) 依据业务知识,刻画欺诈特征。基于日常风险防控积累的经验,抽象出影响交易是否欺诈的关键特征,例如,短时异地的交易特征、连续多笔失败的场景数据特征等,同样可以依据交易数据构建这些业务经验相关的特征。

1.2 特征筛选

通过以上方法,构建了大约上千个特征。有些特征可能对于交易欺诈评估效果不明显,有些特征可能会随着时间的变化产生明显的波动,还有些特征可能相互之间关联性比较大。通过特征筛选,挑选出合适的特征集,提高模型泛化能力。特征筛选主要有以下几种方法:

(1) IV 值评估。在对变量进行分析后,需要计算变量的重要性,IV 值即 information value,中文表述为信息量或信息值,是评估变量区分度和重要性的统计量之一。IV 值越大,区分能力越强。本文设置 IV 值的阈值为 0.15,筛选掉 IV 值小于 0.15 的特征变量。

(2) 相关性过滤。相关系数用于考察两个变量或特征之间的相关程度。如果相关性过高,会导致模型重复计算。因此,需要过滤掉相关性过高的特征,本文

设定线性相关性阈值为 0.2,当两个变量相关性大于 0.2 时,保留 IV 值较大的特征变量。

(3) 稳定性过滤。在风控模型中,特征稳定性是一个非常重要的指标。如果特征分布不稳定,会导致模型效果出现较大的波动。因此需要对特征变量进行稳定性分析,过滤掉不稳定的特征变量,本文按照时间维度对特征变量稳定性进行分析,设定稳定性阈值为 0.1,筛选掉稳定性大于 0.1 的特征变量。

2 欺诈评分模型

在交易欺诈评估场景中,本文将伪冒持卡人本人或未经持卡人授权的交易定义为欺诈交易,持卡人本人或授权的交易为正常交易,为方便起见,定义欺诈交易为负样本,正常交易为正样本,在欺诈评分模型构建中,负样本占正样本的比例只有万分之一,甚至几十万分之一,是一个典型的非平衡数据集。在处理非平衡数据,最常用的方法有两种:过采样和欠采样。二者都是为了达到最终样本类别不太失衡的目的。传统的降采样方法主要有随机降采样(Random-d)和原型选择法的 NearMiss^[7]降采样(NearMiss-d)。随机降采样会导致样本信息的丢失,并且后续模型训练中无法学习到丢失的信息,影响整体模型效果;基于 NearMiss 的降采样从多数类样本中选择具有代表性的样本,简单地依靠“距离”远近,很难全面地比较样本信息。

本文提出一种方法,通过构建样本差异度函数,基于特征相似度降采样(FeatureSimilarity-d),最终达到样本类别不太失衡的同时尽量多地保留样本信息,为模型提供具有区分度的样本训练数据,以采样后的数据作为训练集,以集成学习算法^[8]作为模型训练使用的算法,最终训练出交易欺诈预测评分模型。

2.1 采样策略

每个样本都包含大量的信息,如数值数据、日期数据、类别数据和文字描述数据等。将样本信息分为两类:

(1) 数值型数据:有具体数值的,标准化处理后可量化样本间距离,也可用于采样后模型训练,如交易金额、交易笔数等数值数据。

(2) 描述型数据:难以数值化处理、具有类别区分作用的数据,如交易日期、交易类别等文字描述数据。

由于负样本比较稀缺,对正样本做降采样处理,保留全部的负样本。正样本个数记为 M ,负样本个数记为 N ,其中 M 远远大于 N 。记正样本 X_i 、负样本 Y_j 为:

$$\mathbf{X}_i = (\mathbf{X}_i^s, \mathbf{X}_i^f) \quad i=1,2,\dots,M$$

$$\mathbf{Y}_j = (\mathbf{Y}_j^s, \mathbf{Y}_j^f) \quad j=1,2,\dots,N$$

式中: \mathbf{X}_i^s 、 \mathbf{Y}_j^s 是样本数值型信息,是一个 n_s 维的数字向量,且均已做过标准化处理; \mathbf{X}_i^f 、 \mathbf{Y}_j^f 是样本的描述型信息,是一个 n_f 维的向量。

$$\mathbf{X}_i^s = (x_{i1}^s, x_{i2}^s, \dots, x_{in_s}^s), \mathbf{Y}_j^s = (y_{j1}^s, y_{j2}^s, \dots, y_{jn_s}^s)$$

$$\mathbf{X}_i^f = (x_{i1}^f, x_{i2}^f, \dots, x_{in_f}^f), \mathbf{Y}_j^f = (y_{j1}^f, y_{j2}^f, \dots, y_{jn_f}^f)$$

为了保证模型在负样本上有较好的分类效果,采样出的正样本应尽可能在负样本周围,这样模型能更好地学习负样本周围具有区分度的信息。考虑使用样本的数值型信息,计算任意两个正负样本的欧氏距离:

$$d_{ij} = D(\mathbf{X}_i, \mathbf{Y}_j) = \|\mathbf{X}_i^s - \mathbf{Y}_j^s\| \\ \left[(x_{i1}^s - y_{j1}^s)^2 + (x_{i2}^s - y_{j2}^s)^2 + \dots + (x_{in_s}^s - y_{jn_s}^s)^2 \right]^{\frac{1}{2}} \quad (1)$$

对每个负样本挑选出 m 个正样本应具有以下特点:(1)离负样本尽可能近,与负样本保持一定的相似度。(2)离负样本相同距离情况下,尽可能分散,保留更多样本信息。(3)样本信息分散与样本距离相结合,保证采样稳定。

2.2 样本差异度函数

使用样本的数值型信息和描述型信息,构造正样本差异度函数。实际上,每个样本都有 $n_s + n_f$ 个维度的信息, n_s 个维度的数值型信息, n_f 个维度的描述型信息。

考虑样本的数值型信息,针对负样本 \mathbf{Y}_j ,采样出的 m 个正样本数值型信息为:

$$\mathbf{X}_{j1}^s = (x_{j1_1}^s, x_{j1_2}^s, \dots, x_{j1_{n_s}}^s)$$

$$\mathbf{X}_{j2}^s = (x_{j2_1}^s, x_{j2_2}^s, \dots, x_{j2_{n_s}}^s)$$

⋮

$$\mathbf{X}_{jm}^s = (x_{jm_1}^s, x_{jm_2}^s, \dots, x_{jm_{n_s}}^s)$$

实际上数值型信息的每个维度均是数值,经过标准化处理后,每个维度都有取值区间。考虑所有 M 个正样本的数值型信息,针对每个维度将取值区间划分成若干小区间,观察采样出的 m 个正样本在小区间上的分布,判断样本的离散程度。

以数值型信息的第一个维度为例,所有 M 个正样本的数值型信息第一维度在取值区间上,划分成 k_1 个小区间。全部 m 个正样本第一维度的数值型信息为 $(x_{j1_1}^s, x_{j2_1}^s, \dots, x_{jm_1}^s)$,落在 k_1 个小区间上的个数分别为 s_1, s_2, \dots, s_{k_1} ,则有如下关系:

$$s_1 + s_2 + \dots + s_{k_1} = m$$

考虑 $(x_{j1_1}^s, x_{j2_1}^s, \dots, x_{jm_1}^s)$ 在 k_1 个小区间上的离散程

度,定义:

$$L_1^s = \frac{m^2(k_1 - 1)^2}{k_1 \sum_{i=1}^{k_1} (m - s_i)^2} \quad (2)$$

L_1^s 取值在 $(0, 1]$ 之间,当 s_1, s_2, \dots, s_{k_1} 均匀地分布在 k_1 个小区间上时,即:

$$s_1 = s_2 = \dots = s_{k_1} = \frac{m}{k_1}$$

此时 L_1^s 取最大值1。

针对负样本 \mathbf{Y}_j ,采样出的 m 个正样本数值型信息的离散程度定义为:

$$L^s = \frac{1}{n_s} \sum_{i=1}^{n_s} L_i^s = \frac{1}{n_s} \sum_{i=1}^{n_s} \frac{m^2(k_i - 1)^2}{k_i \sum_{i=1}^{k_i} (m - s_i^t)^2} \quad (3)$$

可以看出, L^s 取值在 $(0, 1]$ 之间,正样本数值型信息越离散, L^s 取值越接近1。

注:

k 个非负数 x_1, x_2, \dots, x_k ,其中 $x_1 + x_2 + \dots + x_k = m$,当 x_1, x_2, \dots, x_k 分别取何值时, $\sum_{i=1}^k (m - x_i)^2$ 取值最小?

可以证明,当 $x_i = \frac{m}{k} \quad i=1,2,\dots,k$ 时, $\sum_{i=1}^k (m - x_i)^2$ 有最小值 $\frac{m^2(k-1)^2}{k}$,并且此时 x_i 均匀地分布在 k 个位置,分布最为离散。

考虑样本的描述型信息,针对负样本 \mathbf{Y}_i ,采样出的 m 个正样本描述型信息为:

$$\mathbf{X}_{j1}^f = (x_{j1_1}^f, x_{j1_2}^f, \dots, x_{j1_{n_f}}^f)$$

$$\mathbf{X}_{j2}^f = (x_{j2_1}^f, x_{j2_2}^f, \dots, x_{j2_{n_f}}^f)$$

⋮

$$\mathbf{X}_{jm}^f = (x_{jm_1}^f, x_{jm_2}^f, \dots, x_{jm_{n_f}}^f)$$

实际上描述型信息的每个维度均反映了样本的属性,针对某一维度的描述型信息,样本的属性分布越离散,样本包含的信息越多。用 $U_n \{x_1, x_2, \dots, x_n\}$ 表示集合 $\{x_1, x_2, \dots, x_n\}$ 中不同元素的个数,则:

$$L_1^f = \frac{1}{m} U_n \{x_{j1_1}^f, x_{j1_2}^f, \dots, x_{j1_{n_f}}^f\} \quad (4)$$

式(4)表征了采样出的 m 个正样本第一维度的描述型信息的离散程度,且 L_1^f 的取值在 $(0, 1]$ 之间,第一维度描述型信息越离散, L_1^f 值越接近1。

针对负样本 \mathbf{Y}_j ,采样出的 m 个正样本数值型信息的离散程度定义为:

$$L^f = \frac{1}{n_f} \sum_{i=1}^{n_f} L_i^f = \frac{1}{n_f m} \sum_{i=1}^{n_f} U_n \{x_{j1_i}^f, x_{j2_i}^f, \dots, x_{jm_i}^f\} \quad (5)$$

可以看出, L^f 取值在 $(0, 1]$ 之间, 正样本数值型信息越离散, L^f 取值越接近 1。

综合样本数值型信息和描述型信息, 构造正样本离散程度函数:

$$L = \alpha_1 L^s + \alpha_2 L^f$$

式中: α_1 、 α_2 为权重参数。

考虑到从 $X_{j1}, X_{j2}, \dots, X_{jm}$ 中挑选出来的 m 个正样本与负样本 Y_j 的距离:

$$D = \beta [d(Y_j, X_{j1}) + d(Y_j, X_{j2}) + \dots + d(Y_j, X_{jm})] \quad (6)$$

即所有挑选的个正样本与负样本 Y_j 之间距离的总和, 其中 β 为权重参数。

差异度函数为:

$$S_m = \frac{L}{D} = \frac{\alpha_1 L^s + \alpha_2 L^f}{\beta [d(Y_j, X_{j1}) + d(Y_j, X_{j2}) + \dots + d(Y_j, X_{jm})]} \quad (7)$$

2.3 最优采样

针对每个负样本 Y_j , 从所有的 M 个正样本中挑选 m 个计算 S_m , 考虑全局最优, 需要计算 C_M^m 种组合, 得到差异度 S_m 最大的一组作为采样结果。

实际上, 对每个负样本采样出 m 个正样本, 最终得到 $N \times m$ 个正样本, 考虑到同一个正样本会被不同负样本采集到, 最终采样得到的正样本数量在 $[m, N \times m]$ 之间。

2.4 模型训练

基于采样的数据, 利用集成学习进行模型训练, 模型训练会输出多个成员分类器, 在模型决策时, 需要将每个成员分类器预测的结果进行融合^[9], 最后输出一个评分, 评分结果表示交易欺诈的概率:

$$S = \frac{1}{1 + e^{-\sum_{j=1}^m C_j(x)}}$$

式中: $C_j(x)$ 为第 j 个成员分类器预测的结果。

3 实验与结果分析

本文以某银行 2019 年 4 至 7 月的部分数据作为数据集, 其中: 4 至 6 月为训练集, 7 月作为测试验证集。其中确认的欺诈交易标注为负样本, 以梯度提升决策树^[10] (Gradient Boosting Decision Tree, GBDT) 集成学习算法作为模型训练阶段使用的算法。

为评估基于特征相似度的有效性, 本文选择随机降采样 (Random-d)、NearMiss 降采样 (NearMiss-d) 与

特征相似度降采样 (FeatureSimilarity-d) 进行比对, 在数据采样阶段, 三种采样方式都按照相同的比例进行采样, 在模型训练阶段, 采用同样的算法和训练参数, 最后使用同样的测试验证集进行模型验证。

3.1 采样分析

针对训练数据集, 首先对三种采样后的样本数据进行比对, 主要通过特征饱和度、正负样本相关性、特征稳定性等指标进行比对。

(1) 饱和度: 样本中, 部分特征为默认值 -1 或 0, 默认值所占比例越高, 样本包含信息越少。故而统计样本的饱和度 (非默认值特征变量个数除以总的特征变量个数), 饱和度越高, 特征变量包含信息越多, 越有用。

(2) 相关性: 正样本中特征变量与负样本特征变量之间的相关系数, 用来比较样本之间的相关性, 相关越高, 正样本越有利于识别负样本, 越有利于提高模型效果。

(3) 稳定性: Psi, 关注特征变量取值随时间推移发生波动的大小。Psi 值越小特征变量越稳定, 越有利于模型效果提升。

通过表 1 比对可以看出, 基于特征相似度的降采样 (FeatureSimilarity-d) 方法比另外两种方法在饱和度、相关性、稳定性指标上均有提高。

表 1 采样后的样本指标对比

指标	Random-d	NearMiss-d	FeatureSimilarity-d
饱和度	0.885 3	0.921 3	0.923 6
相关性	7.419 1	17.118 2	17.315 3
稳定性	6.367 0	6.135 4	6.092 4

3.2 实验结果对比

本文选择了 GBDT 集成学习算法作为模型训练使用的算法, 算法使用相同训练参数, 采样前使用相同的训练数据集, 使用相同的数据集进行验证评估, 在模型效果评估阶段, 使用召回率 (Recall) 和准确率 (Precision) 作为评价指标。

图 1 和表 2 为本文提出的基于特征相似度降采样 (FeatureSimilarity-d) 方式训练出的模型与其他采样方式训练出的模型效果对比情况。可以看出, 本文所提采样策略训练出的模型, 除了在 0.3 左右召回率的情况下, 准确率略低于基于 NearMiss 降采样训练出的模型, 在其他同样召回率的情况下, 准确率都高于其他采样算法训练出的模型。

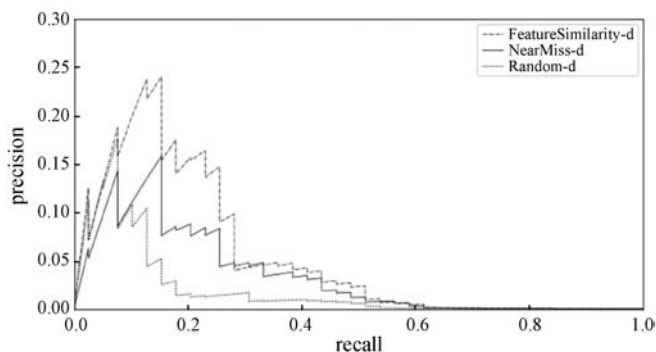


图1 模型效果对比情况

表2 模型效果对比情况(%)

召回率	准确率		
	Random-d	NearMiss-d	FeatureSimilarity-d
10	7.54	11.11	20.00
20	1.52	8.79	15.68
30	0.94	4.78	4.34
40	0.72	3.47	4.27
50	0.55	1.22	2.39

4 结 语

本文针对交易欺诈预测场景,首先提取有效的特征集,针对非平衡样本进行采样后训练建立一个用于交易欺诈预测的模型,通过与其他采样方式进行对比,验证本文所提采样方式的有效性。由于本文所提的基于特征相似度降采样(FeatureSimilarity-d)方式采样后的正样本一方面离负样本尽量近,同时正样本之间又尽量离散,保留更多样本信息与样本间组合信息,达到模型区分度高的目标。但同时本文所提的采样方法,除了需要计算每个正样本与负样本之间的距离并进行比较外,正样本之间还会进行计算和比较,所以相比随机降采样(Random-d)和原型选择法的NearMiss降采样(NearMiss-d),采样过程中计算量更大,需要的计算资源更多。

参 考 文 献

[1] 窦路路,石秀金.基于深度学习的银行卡交易反欺诈技术研究[J].智能计算机与应用,2018,8(4):85-87,91.
 [2] 李欣.大数据技术在信用卡风险管理中的应用研究[J].企业科技与发展,2019(10):104-105,108.
 [3] 曹汉平,张晓晶,祝睿杰,等.数字金融时代机器学习模型在实时反欺诈中的应用与实践[J].智能科学与技术学报,2019(4):342-351.
 [4] 刘胜兰.不平衡数据集的分类方法研究[D].北京:北京邮电大学,2019.

[5] 黄镜霖.基于集成学习的金融反欺诈模型[J].电脑知识与技术,2020,16(1):216-219.
 [6] 张燕.基于本质特征和网络特征的信用卡欺诈检测[J].微型电脑应用,2016,32(12):72-77.
 [7] 魏力,张育平.一种改进型的不平衡数据欠采样算法[J].小型微型计算机系统,2019,40(5):1094-1098.
 [8] 陈荣荣,詹国华,李志华.基于XGBoost算法模型的信用卡交易欺诈预测研究[J].计算机应用研究,2020,37(S1):111-112,115.
 [9] 钟金宏,邵晶晶,李兴国.基于组合分类策略的个人信用风险评估研究[J].合肥工业大学学报(自然科学版),2020,43(7):996-1002.
 [10] 赵金涛,邱雪涛,何东杰.基于GBDT的线上交易欺诈侦测研究[J].微型电脑应用,2017,33(10):17-18,21.

(上接第72页)

[14] 岳有军,刘英翰,赵辉,等.基于极点对称模态分解-分散熵和改进乌鸦搜索算法-核极限学习机的短期负荷区间预测[J].科学技术与工程,2020,20(22):9036-9042.
 [15] 尹新涛,杨校辉,胡道栋.居民区电动汽车充电负荷控制单元设计与实现[J].电气自动化,2020,42(6):7-9.
 [16] 黄山,程启明,张强,等.含电动汽车混合储能系统的微网多目标经济优化运行[J].高压电器,2017,53(10):142-149.

(上接第89页)

[4] 霍凯歌,张亚琦,胡志华.自动化集装箱码头多载AGV调度问题研究[J].大连理工大学学报,2016,56(3):244-251.
 [5] 朱光宇,徐文婕.考虑能耗与质量的机床构件生产线多目标柔性作业车间调度方法[J].控制与决策,2019,34(2):252-260.
 [6] Zhang L, Hu Y, Guan Y. Research on hybrid-load AGV dispatching problem for mixed-model automobile assembly line[J]. CIRP Conference on Manufacturing System, 2019, 81:1059-1064.
 [7] Mousavi M, Yap H J, Musa S N, et al. A fuzzy hybrid GA-PSO algorithm for multi-objective AGV scheduling in FMS[J]. International Journal of Simulation Modelling, 2017, 16(1):58-71.
 [8] Liu Y, Ji S, Su Z, et al. Multi-objective AGV scheduling in an automatic sorting system of an unmanned (intelligent) warehouse by using two adaptive genetic algorithms and a multi-adaptive genetic algorithm[J]. PLoS ONE, 2019, 14:e022616112.
 [9] 陶艺辉,曹小华,袁智.固定节点AGV系统调度方案分层优化算法研究[J].起重运输机械,2019(11):59-64.
 [10] 汪定伟,王俊伟,王洪峰,等.智能优化方法[M].北京:高等教育出版社,2007.