

基于堆栈降噪自编码网络的个人信用风险评估方法

杨德杰¹ 章宁¹ 袁戟² 白璐¹

(中央财经大学信息学院 北京 100081)¹ (德国慕尼黑工业大学土木-地质-环境学院 慕尼黑 80333)²

摘要 个人信用历来是银行衡量个人履约风险最重要的因素。近年来,随着我国借贷需求与日俱增,仅依据信用卡信息的传统个人信用评估方式,已不能完全满足银行业的发展需求。因此,为了构建更加丰富的用户信用画像,文中基于银行大数据提取信用风险评估特征。为了解决金融大数据带来的维度灾难和噪声问题,充分考虑了数据特征之间的相关性,对堆栈降噪自编码神经网络模型进行了改进,引入了截断的 Karhunen-Loève 展开作为噪声传入项,并在某商业银行的大数据平台上进行了一系列数据实验。实验结果显示:相比仅使用信用卡信息,利用银行大数据能使衡量正负样本分离度的指标——K-S 值提升约 11%;改进的堆栈降噪自编码神经网络方法具有更好的风险评估效果,准确率相比原模型提高了 3% 左右,验证了在银行大数据环境下进行信用风险评估的有效性。

关键词 信用风险评估,大数据,维度灾难,特征选择,堆栈降噪,深度学习

中图分类号 TP181 **文献标识码** A **DOI** 10.11896/jsjx.181102216

Individual Credit Risk Assessment Based on Stacked Denoising Autoencoder Networks

YANG De-jie¹ ZHANG Ning¹ YUAN Ji² BAI Lu¹

(School of Information, Central University of Finance and Economics, Beijing 100081, China)¹

(College of Civil, Geo and Environmental Engineering, Technical University of Munich, Munich 80333, Germany)²

Abstract Personal credit is the most important factor for banks to measure individual compliance risk. In recent years, with the increasing demand for borrowing in China, the traditional way of making credit evaluation, which is merely based on credit card transaction information, cannot fully meet the development needs of the banking industry. Therefore, this paper proposed to use the big data of personal consumption in bank as the important feature information to construct a richer user image. In order to overcome the dimensional curse and noise caused by the financial big data, a modified deep learning evaluation algorithm based on stacked denoising autoencoder neural network is proposed by considering the correlation of feature data and the truncated Karhunen-Loève expansion is applied as the noise input term, then a series of related data experiments are conducted on big data platform of a commercial bank. The experimental results show that, compared with the risk evaluation just based on credit card transaction information, the K-S value that measure the positive and negative sample resolution based on big data of bank improves 11%; the improved stack denoising autoencoder neural network method has better risk assessment results and the accuracy rate is increased by about 3% compared with the original model, thus validating the effectiveness of credit risk assessment in the big data environment of bank.

Keywords Credit risk assessment, Big data, Dimensional curse, Feature selection, Stacked denoising, Deep learning

1 引言

近 20 年来,随着经济发展、消费升级、个人消费观念的转变,个人信贷需求日益增多,相应的服务——信用卡也成为发展较快的银行业务之一。由于个人信用存在着明显差异,且借贷、还贷与银行的经济利益及业务的发展息息相关,因此在

开办信用卡业务前对个人的还贷能力进行风险评估尤为必要。显然,对在银行缺乏个人贷款记录的客户进行风险评估是困难的,因此给缺少信用数据的客户做出合理的个人借贷风险评估,成为了银行信用卡业务的痛点和难点。

人工智能的飞速发展和计算机性能的大幅提升,为算法在大数据风控中的实现提供了机会。近年来,许多学者应用

到稿日期:2018-11-29 返修日期:2019-04-15 本文受国家重点研发计划(2017YFB1400701),国家社会科学基金重点项目资助(13AXW010)资助。

杨德杰(1987—),男,博士生,高级工程师,主要研究方向为机器学习、金融风控, E-mail: yangdejiejay@163.com; 章宁(1975—),女,博士,教授,主要研究方向为金融科技、个人信息保护, E-mail: zhangning@cufe.edu.cn(通信作者); 袁戟(1985—),男,博士,助教,高级工程师,主要研究方向为贝叶斯反演分析、随机有限元方法等; 白璐(1987—),男,博士,副教授,CCF 会员,主要研究方向为机器学习、特征选择等。

基于统计学、机器学习和深度学习的方法,展开了一系列的信用风险评估模型研究^[1-2]。

机器学习方法具有可解释性强、计算效率高、鲁棒性好等优势,无论在研究中还是实际应用中都被广泛采用。Jayanthi等^[3]运用SVM模型对银行借款决策和违约率进行预测,并证实SVM预测非线性关系的效果最优。方匡南等^[4]在运用逻辑回归模型分析国内某商业银行信用卡消费违约数据时,采用了Lasso回归,其相比全变量逻辑回归模型预测,能更准确地筛选出重要变量,且效率更高。

此外,模型融合和集成学习在很大程度上能够弥补单个模型在预测精度和速度上的不足,使得预测效果超越单个模型。Lin等^[5]比较了几种不同机器学习模型在信用风险评估问题中的应用,研究结果表明集成模型具有比单一分类器更好的性能。Chen等^[6]通过集成人工蚁群算法和SVM分类器的混合模型,在美国的公司信用评分上提升了分类精度。

近年来,越来越多的深度学习算法被应用到金融领域(包括信用风险评估问题^[7])。Yu等^[8]采用Bagging方法生成多个训练子集,并在每个子集上训练单分类器,然后将输出作为深度信念网络的输入来训练深层神经网络。其在本国公开消费数据集上的测试效果优于其他机器学习方法。Sirignano等^[9]运用神经网络对美国房地产贷款逾期及提前还款风险进行预测,其效果优于逻辑回归、人工神经网络模型的效果。然而,有学者提出将深度学习应用到小数据集上未必能获得较好的效果。如Shigeyuki等^[10]在台湾的信用违约数据集上将常见的几种集成学习算法与神经网络做比较,结果验证了集成学习在此类小数据集上的效果优于深度学习模型。

随着大数据时代的到来^[11],银行数据量呈现出爆发式增长。借助于客户数据的多元化、多样化,银行可以更大程度地避免由于信息不对称导致的评估偏差。银行通过灵活运用客户信息,进行全方位的合理评估,构建了更加完整的用户画像,突破了信用特征指标单一的局限性。因此,缺少信用数据的客户的风险评估问题逐步被解决。刘新海等^[12]介绍了美国ZestFinance公司利用互联网大数据对个人客户进行信用风险评估,解决了征信缺失以及无抵押担保的小微型企业融资的需求。

Lecun等^[13]指出深度学习成功的必要条件之一是大数据,特别是对于高维数据,深度学习“端到端”学习的模式能从海量数据中学习表达对象的本质特征,从而更好地进行识别和分类。就基于大数据的信用风险评估而言,更多种类的个人数据增加了客户属性的维度,同时也带来了特征选择上的困难,即选择与信用相关或者提升信用风险评估能力的属性特征。而深度学习适合高维数据的学习能力,使得其能有效挖掘个人大数据中表征个人信用的特征组合,进而提升风险评估的效果^[14]。与此同时,大数据不可避免地带来了数据噪声问题,即数据中含有偏差、错误以及无效的信息,这些信息伴随着数据量的积累,将严重影响银行信用风险评估的精度。此外,大数据引发的维度灾难会降低计算效率,并带来特征选择上的困难^[15]。因此,在尽可能保证风险评估精度的同时提

高计算效率,成为了大数据信用风险评估极具挑战的任务。

为了更加合理地评估客户的借贷能力,本文依据客户在银行大数据平台的数据全貌,利用深度学习框架,提出改进的堆栈降噪自编码器网络算法,用于信用风险评估。与以往银行做信用评估不同,本文将银行大数据作为风险评估数据来源,丰富了信用评估特征,进而提高了评估准确率,进一步减少了信用不对称。此外,本文提出的算法在一定程度上解决了利用银行大数据带来的特征选择困难、高维数据计算复杂度以及大数据噪声多等问题。

2 算法模型

2.1 风险评估特征的选择

通常,银行对个人信用风险进行评估所使用的数据主要来源于对客户资料的收集,如央行征信数据关于个人的基本信息(性别、年龄、学历、职业等)、财务状况、信用记录、债项情况等。随着银行各项业务的拓展,尤其是对行内客户的交叉营销,客户信息不断丰富;银行客户也不再归属于某单一业务范畴,客户在银行多个业务领域留下的足迹都将成为个人信用风险评估的数据来源,这些数据也成为个人信用画像必要的补充。

本文在采用了传统评估特征即信用卡借贷和还贷数据的基础上,依托某商业银行大数据平台,加入了客户交叉业务信息,将两部分数据整合,进而构建完整的银行个人信用画像。图1给出了数据拼接整合的逻辑表述。

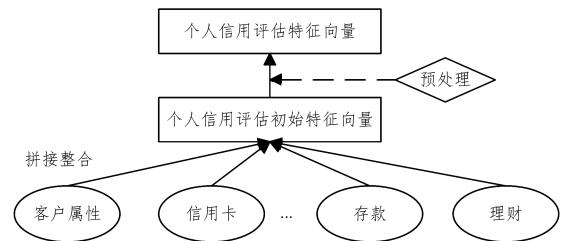


图1 某商业银行个人信用评估特征的构建

Fig. 1 Construction of personal credit evaluation features of certain commercial bank

2.2 模型和算法

本文中应用的堆栈自编码神经网络是深度学习的一种特殊结构,它主要被应用于高维数据的降维,即特征压缩^[16]。它的基本构造思路是:首先对输入项与隐含层和输出项构建浅层的神经网络,在输出项最接近输入项的条件下求解隐含层;然后将隐含层作为下一个浅层网络的输入项,重复以上过程,求解新的隐含层;最后将隐含层以堆栈形式逐层堆叠,形成深层神经网络。

考虑到大数据样本存在因数据偏差、错误以及缺失等原因带来的噪声,为增强自编码器的抗噪能力,使其更具鲁棒性和泛化性,降噪自编码器(Denoising Autoencoder, DAE)通过在原始数据上增加噪声,来生成浅层神经网络结构,压缩提取特征^[17]。该算法在多维度、高噪声数据场景下具有良好的表现^[18-19]。计算中,首先将训练样本 $x(x \in R^n)$ 进行一系列随机变换 $q_D(x'|x)$,腐蚀后得到变量 x' ,然后通过对 x' 进行自编

码器的学习训练得到重构输入 x'' 。整个学习流程示意图如图 2 所示。

在图 2 所示的神经网络中,应用梯度下降法进行反复迭代求解,使得 x'' 与原始数据 x 最为接近,此时对应的隐含层 y (或压缩特征层)为:

$$y = \arg \min \left[\sum_{i=1}^n \frac{1}{n} \mathcal{L}(x_i, x_i'') + J_w \right] \quad (1)$$

其中, $x_i \in x, x_i'' \in x'' (i=1, 2, \dots, n)$ 分别为属于原特征和重构特征的对应元素, \mathcal{L} 为均方误差或交叉熵的损失函数, J_w 是为了防止模型学习过程中产生过拟合而添加的范数约束。

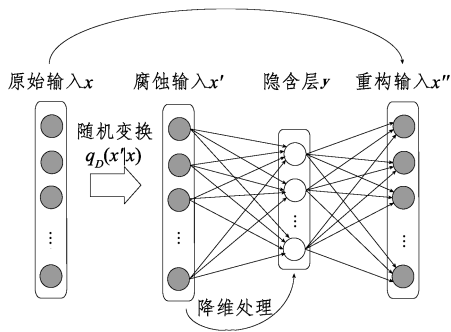


图 2 降噪自编码的计算流程图

Fig. 2 Computation flowchart of denoising autoencoder

原降噪自编码神经网络产生的随机变换只是在原数据特征上增加小幅的变动,并未考虑到输入特征之间噪声的相关性,这对于提高模型鲁棒性显然是不足的。本文在考虑样本噪声相关性的前提下,为了提高噪声生成的质量,采用 Karhunen-Loève(K-L)展开^[20]来完成随机变换 $q_D(x'|x)$,同时采用它的截断形式来提高计算效率。

首先,通过分析输入特征之间的相关性,确定其正定的相关性矩阵 Σ_{XX} 。对于任意 $n \times 1$ 的向量 e_k 和对角矩阵 $\Lambda = \text{diag}[\lambda_1, \lambda_2, \dots, \lambda_n]$,可以使得等式(2)成立:

$$\Sigma_{XX} e_k = \Lambda e_k \quad (2)$$

其中, λ_i 是相关性矩阵 Σ_{XX} 的非零特征项。Karhunen-Loève 展开可以定义为对相关矩阵的分解:

$$\Sigma_{XX} = \sum_{i=1}^n \sqrt{\lambda_i} v_i v_i^T \quad (3)$$

其中, v_i 是 λ_i 对应标准化的特征值。因此,随机变量可以表示为:

$$X = \mu_X + \sum_{i=1}^n \sqrt{\lambda_i} v_i X_i \quad (4)$$

其中, μ_X 表示随机变量 X 的均值,因在计算过程中首先会进行归一化,故可以假设 $\mu_X = 0$; X 为随机数,服从标准正态分布,即 $X \sim N(0, 1)$ 。

通常,在实际操作中为了减小计算量,不需要取得所有项 KL 谱分解的展开项,只需要取得其中的 M 项 ($M < n$),即可取得截断(truncated)近似值^[21]:

$$\hat{X} = \sum_{i=1}^M \sqrt{\lambda_i} v_i X_i \quad (5)$$

M 项需要保证得到的截断误差小于 0.05。截断误差定义为^[22]:

$$err_{r-Var} = \frac{\text{Var}[\hat{X} - X]}{\text{Var}[X]} = \frac{\sum_{i=M+1}^n \lambda_i v_i^2}{\sum_{i=1}^n \lambda_i v_i^2} \quad (6)$$

堆栈降噪自编码器(Stacked Denoising Autoencoders, SDAE)网络是由单个自编码器逐层级联形成的。具体来说,前一个自编码器训练得到的中间隐含层作为下一个自编码器的输入参与训练,这样依次堆栈,形成深层神经网络结构,如图 3 所示。

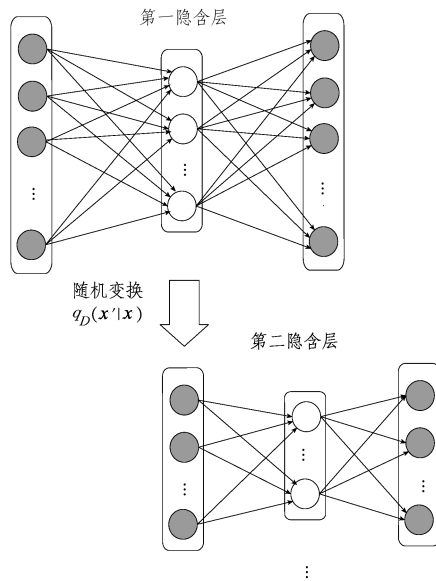


图 3 堆栈降噪自编码器网络

Fig. 3 Stacked denoising autoencoder network

基于某商业银行大数据平台,以个人信用卡数据为主数据,拼接及整合个人在大数据平台内的业务属性和交易数据,构建基于银行大数据的个人信用画像统一视图。基于银行大数据的堆栈降噪自编码神经网络的具体学习过程如图 4 所示。

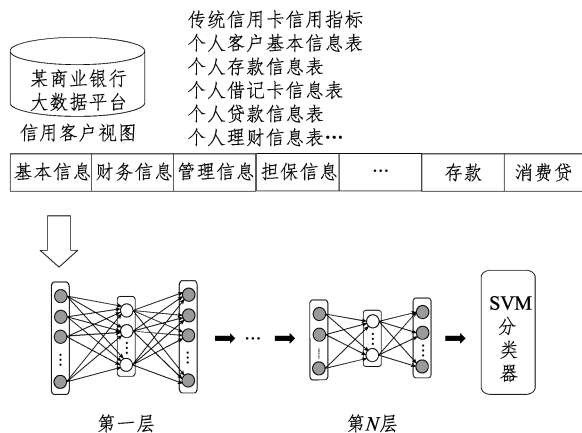


图 4 基于银行大数据的堆栈降噪自编码神经网络的学习过程

Fig. 4 Learning process of stacked denoising autoencoder neural network based on bank big data

某商业大数据平台作为基础数据集中采集、存储和处理的平台,涵盖了个人在行内的所有信息。但由于个人信息散落在各个数据源中,为了方便上层业务逻辑使用和统一管理,需要对其进行整合拼接。由于本文的目的是对个人信用风险进行评估,因此以信用维度为出发点,以个人客户为粒度,拼接和整合个人客户在银行大数据平台内的其他业务数据,如

个人财务、理财、个人贷款、担保信息等,最终形成基于银行大数据的个人信用完整视图。下面对个人信用风险评估的整体算法进行形式化表达。

(1)数据预处理。对原始数据进行抽取、整合、清洗、转换等预处理。由传统的信用评估特征形成特征向量 $[x_1, x_2, \dots, x_k]$,个人客户的行内业务数据经预处理后形成特征向量 $[x_{k+1}, x_{k+2}, \dots, x_n]$;融合两部分特征形成模型的输入特征,即以 $X=[x_1, x_2, \dots, x_n]$ 为模型训练的输入数据。

(2)模型训练。基于大数据构造输入特征集,构造改进的堆栈降噪自编码深层神经网络 KL-SDAE。

1)构造第一层网络

获取特征 X 之间的相关性矩阵 Σ_{XX} ;

通过 KL 分解进行随机扰动变换 $q_D(x' | x)$,从而得到腐蚀后的输入项 x' ;

根据降噪自编码计算流程(DAE)获得第一层隐含层 h_1 ,并将其变为首个输入层(见图2)。

2)构造第2层至第N层网络

获取隐含特征之间的相关性矩阵 Σ_{hh} ;

将上一层神经网络的隐含层作为下一个神经网络的输入项(见图3);

保留第二层以后的所有隐含层 $h_i (i=1, 2, \dots, N)$,并将其构造为深层的神经网络结构 SDAE(见图3)。

3)反向调优

使用 SVM 进行最终识别和分类,并利用 BP 对网络参数权重 w_i 和偏置 b_i 进行反向调优,即用梯度下降法调整 $[w_i, b_i](1 \leq i \leq N)$ 。

3 实验与分析

3.1 案例概述

本文数据来源于国内某商业银行大数据平台。为了得到更加全面的用户画像,建模所需数据以个人信用卡数据为中心,拼接客户在银行内的业务数据,主要包括存款、贷款、理财、担保、三方存管、基金、借记卡、电子银行交易等11个业务品种。共收集430502位客户作为建模数据,其中有421332位非违约客户(标记为0),9170位违约客户(标记为1),不良客户占比为2.13%。

3.2 数据预处理

3.2.1 数据不平衡处理

由于违约客户占比过低,导致样本分布不平衡,即违约客户数量远远小于非违约客户数量,从而影响模型对样本的学习及其对对应类别的预测能力。因此,为了不丢失违约客户的内在信息,训练样本中保留全部的违约客户,对样本占比高的非违约客户进行下采样。一般地,信用评分领域中正负样本比例最高不超过20:1,最终抽取非违约客户137550位,其与9170位违约客户共同形成146720个样本数据集。为验证模型效果,本文按7:3的比例随机将样本集划分为训练集和测试集。数据集分割情况如表1所列。

表1 模型所用数据集

	训练集	测试集
非违约客户数量	96285	41265
违约客户数量	6419	2751
占比/%	70	30

3.2.2 异常值处理

实验对原始数据进行了缺失值、异常值等的预处理,将缺失比率大于95%的数据直接删除;对连续型数值变量用0数值填充,对离散型数据变量用众数填充;采用盖帽处理极端异常值。

3.2.3 特征提取

在每类业务品种中,对关键特征(如金额、交易笔数等)进行横向衍生(如取最近3个月、最近6个月或最近一年的均值、标准差、最大值等)。通过数据拼接及特征衍生,共形成745维的特征集合。去除大部分数据缺失和常量值的特征,最终形成553维的特征集合。

3.3 实验结果与分析

本文模型基于Tensorflow实现。为保证实验结果的稳定性并防止模型过拟合,采用十折交叉验证的方法进行模型训练。因堆栈自编码器网络学习过程中涉及多个参数,本文通过网络搜索法反复实验,进而组合不同参数,最终确定了本实验数据下的最优参数,如表2所列。

表2 模型学习过程中的参数

参数	α	β	γ	numitera	batchsize	M
参数的值	0.07	0.28	0.01	200	30	233

表2中, α 表示训练模型过程中的学习率; β 表示对输入进行降噪处理的概率参数; γ 表示权重衰减项; $numitera$ 表示训练迭代次数; $batchsize$ 表示每一次迭代训练使用的批量处理样本个数; M 表示KL展开的截断项所需的特征值。

为验证基于银行大数据对个人信用风险评估结果的影响,本文分别以传统信用风险评估特征以及整合银行大数据构建的特征为模型输入,比较不同输入数据对个人信用风险评估结果的影响。按上述基础数据拼接整合及预处理流程,构建基于银行大数据的个人信用风险评估特征集合。

3.3.1 特征重要性检验

为了选择对目标变量预测能力强的特征,利用信息量(Information Value,IV)进行特征选择。IV主要用于处理目标变量是二分变量时解释变量的选择问题,且在信用风险评估领域中被广泛应用。假设某个离散自变量(包括分类变量和次序变量)有 r 个不同的取值,则IV可表示为:

$$IV = \sum_{i=1}^r (pBad_i - pGood_i) \ln \left(\frac{pBad_i}{pGood_i} \right) \quad (7)$$

其中, $pBad_i$ 表示在该离散指标取第 i 个值的样本中,违约客户数占总样本中违约客户数的比例;相同地, $pGood_i$ 表示非违约客户的占比。IV一般用于衡量解释变量对目标变量的区分能力,一般认为IV越大,该变量的区分能力越强。图5给出了基于银行大数据提取的IV排名前十的特征。

从图5的结果可以看出,排名前十的银行大数据特征依

次是:近 1 个月月末活期存款的净值总额(*cur_dpt_1m_amt*)、近半年持有产品种类(*pro_6m_amt*)、近半年个人贷款逾期次数(*loan_6m_num*)、近 3 个月线上消费次数(*onl_cons_3m_num*)、活期存款最早合约建立距今时长(*cur_dpt_fir_dura*)、近 3 个月月末三方存管余额(*thre_3m_amt_ave*)、近 3 个月理财购买金额(*fina_buy_3m_amt*)、最后动账日期距今时长(*cur_dpt_las_dura*)、柜面转账次数(*coun_trans_times*)、是否签约理财产品(*fina_sign_flg*)。其中,近 1 个月月末活期存款净值总额的 *IV* 值达到最高的 0.32,中间大部分变量的 *IV* 值在 0.2~0.3 之间,排名最后的变量即是否签约理财产品的 *IV* 值也在 0.15 以上。一般认为,*IV* 值大于 0.2 的解释变量对目标有较强的区分能力,*IV* 值在 0.1~0.2 之间的变量有一般的区分能力,*IV* 值小于 0.1 的区分能力则较弱。

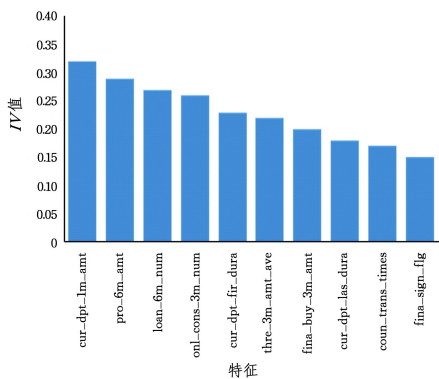


图 5 IV 值排名前十的行内交易特征

Fig. 5 Top ten IV values of bank trading features

3.3.2 单调性分析

为了更加直观地展现解释变量和目标变量之间的关系,选取 *IV* 值排名前四的变量作为候选变量,并利用信用评分中的坏样本率(Bad Rate)指标,即坏样本(违约样本)个数占总样本的比例,对解释变量与目标变量的关系进行单调性分析,结果如图 6 所示。

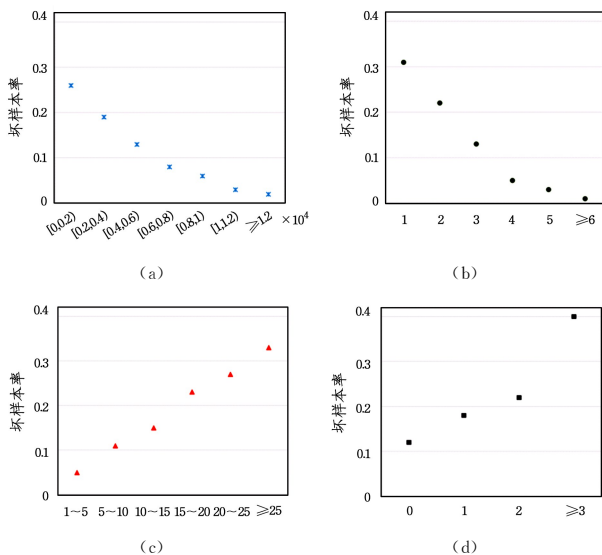


图 6 目标变量关于解释变量的单调性变化

Fig. 6 Monotonic changes of target variables due to explanatory variables

设符号 *a, b, c, d* 分别代表近 1 个月月末活期存款的净值总额、近半年持有产品种类、近 3 个月线上消费次数、近半年个人贷款逾期次数这 4 个解释变量。图 6(a)~图 6(d) 分别表示坏样本率随 *a* 变动的变化趋势,坏样本率随 *b* 变动的变化趋势,坏样本率随 *c* 变动的变化趋势,坏样本率随 *d* 变动的变化趋势。因 *a* 和 *c* 是连续型变量且取值范围较大,故对其进行分箱操作,将相邻数值区间进行合并。

由图 6 的结果可以得到,近 1 个月月末活期存款的净值总额与目标变量呈明显的负相关关系,即随着存款金额的增加,坏样本率下降,相应的违约风险逐渐降低;同样地,坏样本率随持有产品种类的增多而逐步降低。相反,坏样本率与线上消费次数、个人贷逾期次数两个变量呈现显著的正相关关系。

3.3.3 K-S 检验

采用信用评分领域中常用的非参数检验方法——Kolmogorov-Smirnov 检验,并以 K-S 值作为评价指标,来验证不同输入特征对最终个人信用风险评估结果的影响。K-S 统计量度量了两个分布之间的最大垂直距离,在评价二元分类模型的预测能力时,K-S 统计量的值越大,说明模型能够区分正负样本的程度越大。对比实验结果如图 7、图 8 所示。

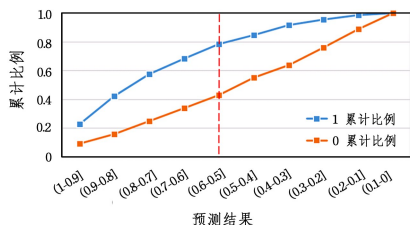


图 7 K-S 曲线——传统信用评估特征

Fig. 7 K-S curves—traditional credit evaluation features

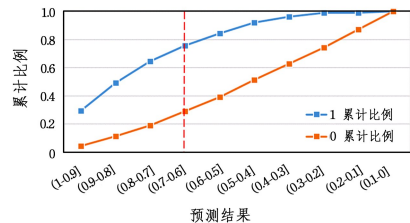


图 8 K-S 曲线——行内大数据评估特征

Fig. 8 K-S curves—bank big data evaluation features

图 7、图 8 对比展示了两种不同输入特征下对正负样本预测累计占比的结果。图 7 中采用传统信用评估特征的 K-S 最大值为 0.36;图 8 中基于银行大数据特征的 K-S 最大值为 0.47,比采用传统信用评估特征的 K-S 值高出 0.11。

综上所述可以看出,融合银行大数据特征能更好地区分违约客户和非违约客户,对个人信用评估的结果有明显的提升作用。从另一角度来看,结合银行大数据,对刻画个人信用画像起到了补充完善的作用。

为了说明本文所提算法在基于银行大数据做信用风险评估中的效果,在同时利用银行大数据的基础上对比了不同算法的效果。

3.3.4 网络层数分析

采用机器学习分类任务常用的评估指标,例如召回率

(recall, rec)、精度 (precision, pre)、正确率 (accuracy, acc)、F-Score、马修斯相关系数 mcc 和 AUC-ROC 等,来衡量和评估本文方法。

由于堆栈降噪自编码网络的 DAE 层数对模型的学习结果会产生直接影响,本文对具有不同 DAE 层数的模型的最终评估结果进行了验证,如表 3 所列。

表 3 DAE 层数的有效性验证

Table 3 Efficient validation of DAE layers

模型	rec/%	pre/%	mcc
SDAE ³	82.4	82.8	0.725
SDAE ⁴	83.2	84.6	0.748
SDAE ⁵	83.7	85.3	0.762
SDAE ⁶	83.3	84.2	0.755
SDAE ⁷	81.9	82.7	0.734

表 3 列出了 DAE 层数从 3 递增至 7 时的不同对比实验结果。由表 3 可知,随着 DAE 层数的增加,预测召回率 *rec*、精度 *pre*、马修斯相关系数 *mcc* 的值都逐渐上升;当 DAE 层数增加到 5 时,*mcc* 值达到最高的 0.762,*rec* 和 *pre* 分别达到了 83.7% 和 85.3%。这也在一定程度上说明了神经网络并非越深越好,或者网络层数越多越好,而需要根据具体业务的应用场景和数据进行动态调整 and 选择。因此,在后续的实验,采用 DAE 层数为 5 的模型进行相关实验。

3.3.5 模型对比

为了说明改进后的算法 (Karhunen-Loève Stacked Denoising Autoencoders, KL-SDAE) 在本案例研究中的优势,下面详述其与其他常见算法在个人信用风险评估中的对比结果。将本文算法与主成分分析 (Principal Component Analysis, PCA)、K 均值 (K-Means) 聚类、梯度提升决策树 (Gradient Boosting Decision Tree, GBDT) 等传统的特征选择方法,以及堆栈自编码网络 (Stacked Autoencoders, SAE), 原始堆栈降噪自编码网络 (SDAE) 进行实验比较,结果如图 9 所示。

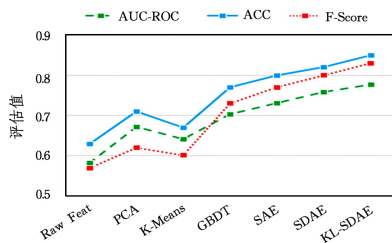


图 9 不同特征选择方法的评估结果对比

Fig. 9 Evaluation result comparison of difference features selection methods

从图 9 的结果可以看出,在同时基于银行大数据条件下,本文提出的 KL-SDAE 算法进行个人信用评估的结果好于其他传统方法,同时比原始堆栈降噪自编码网络的效果更好,准确率提升了 3% 左右。特别地,对比原始特征集 (即不经过任何特征选择, Raw Feat), *AUC-ROC*, *acc* 和 *F-Score* 分别提升了 16.5%, 19% 和 21%; GBDT 作为梯度提升树算法,常用于特征选择,其效果明显好于无监督学习方法,如 PCA 和 K-Means 等;而本文方法在各指标上较 GBDT 均略胜一筹。以上结果说明:在大数据场景下,基于深度学习框架的 KL-

SDAE 能充分提取表征个人信用风险的潜在本质特征,并且对高维稀疏特征做有效压缩和嵌入,使得在低维空间中能表达信用特征之间的关系,进而提高最终的信用评估能力。此外,对比原始自编码网络,基于降噪堆栈自编码网络在大数据环境数据质量不高的情况下提高了模型的抗噪声能力,从而取得了更好的信用评估结果。

虽然各商业银行数据存在差异,得出的结论难免有偏差,但在各商业银行争相建立大数据平台的大背景下,本文提出的方法提供了一种思路,即基于银行大数据来提高个人信用风险识别的能力。此外,本文挖掘了影响个人信用风险的若干关键因素,并将其反映和归结到个人资产状况、客户忠诚度、负债情况、消费习惯这 4 个维度。其中,个人资产状况、客户忠诚度和信用风险呈负相关关系;负债情况与信用风险呈正向关系,而通过银行大数据能有效补充负债信息,如个人贷款履约情况、其他融资类产品的购买信息等;而在消费习惯和信用风险的关系中,发现线上消费频次,如网银、手机银行的转账和消费次数等与违约风险呈正向关系。

结束语 本文针对目前个人信用风险评估问题中所选用的评估特征单一的问题,借助某商业银行大数据平台的优势,基于银行大数据构建个人信用风险评估特征,完善了个人信用风险客户画像。在大数据环境下,深度模型“端到端”学习的优势使得其能更好地挖掘大数据中表达个人信用的潜在本质特征,对个人信用刻画得更加完备,从而提升了风险评估水平。为了解决应用大数据过程中带来的高维数据计算复杂度高和噪声多等数据质量问题,本文利用深度学习算法——堆栈降噪自编码器深度网络,在原始模型基础上,充分考虑了数据特征之间的相关性,提出将截断 KL 展开作为随机噪声的输入项。最后,在某商业银行大数据集上验证了 KL-SDAE 的良好效果。

当然,未来也可引入银行外的互联网大数据 (如非结构化数据) 作为个人信用风险评估数据特征的一部分,以进一步完善客户信用画像;在算法方面,也可融合其他非神经网络模型,如 XGBoost, GBDT 等集成模型,以进一步提高分类学习的准确度。

致谢 感谢中国农业银行股份有限公司为本文研究提供了基础数据和实验环境。

参考文献

- [1] LESSMANN S, BAESENS B, SEOW H V, et al. Benchmarking State-of-the-art Classification Algorithms for Credit Scoring: An Update of Research [J]. *European Journal of Operational Research*, 2015, 247(1): 124-136.
- [2] VISHWAKARMA A C, SOLANKI R. Analysing Credit Risk using Statistical and Machine Learning Techniques [J]. *International Journal of Engineering Science and Computing*, 2018, 8(6): 18397-18404.
- [3] JAYANTHI J, JOSEPH KS, VAISHNAVI J. Bankruptcy Prediction using SVM and Hybrid SVM Survey [J]. *International Journal of Computer Application*, 2011, 33(7): 39-45.

- [4] FANG K N,ZHANG G J,ZHANG H Y. Individual Credit Risk Prediction Method: Application of a Lasso-logistic Model [J]. *The Journal of Quantitative & Technical Economics*,2014,31(2): 125-136. (in Chinese)
方匡南,章贵军,张慧颖. 基于 Lasso-logistic 模型的个人信用风险预警方法[J]. *数量经济技术经济研究*,2014,31(2):125-136.
- [5] LIN W Y,HU Y H, TSAI C F. Machine Learning in Financial Crisis Prediction:A Survey[J]. *IEEE Transactions on Systems Man & Cybernetics Part C*,2012,42(4):421-436.
- [6] CHEN M Y,CHEN C C,LIU J Y. Credit Rating Analysis with Support Vector Machines and Artificial Bee Colony Algorithm [C]//*Recent Trends in Applied Artificial Intelligence*. Amsterdam:Springer,2013:528-534.
- [7] HEATON J B,POLSON N G,WITTE J H. Deep Learning in Finance[J]. *Applied Stochastic Models in Business and Industry*,2017,33(1):561-580.
- [8] YU L,YANG Z B,TANG L. A Novel Multistage Deep Belief Network Based Extreme Learning Machine Ensemble Learning Paradigm for Credit Risk Assessment[J]. *Flexible Services & Manufacturing Journal*,2016,28(4):576-592.
- [9] SIRIGNANO J,SADHWANI A,GIESECKE K. Deep Learning for Mortgage Risk [J]. *Social Science Electronic Publishing*, 2017,22(6):134-216.
- [10] SHIGEYUKI H,MINAMI K,TAKAHIRO K, et al. Ensemble Learning or Deep Learning? Application to Default Risk Analysis[J]. *Risk and Financial Management*,2018,11(1):12-25.
- [11] MA S L,WUNIRI Q G,LI X P. Deep Learning With Big Data: State of The Art and Development [J]. *CAAI Transactions on Intelligent Systems*,2016,11(6):728-742. (in Chinese)
马世龙,乌尼日其其格,李小平. 大数据与深度学习综述[J]. *智能系统学报*,2016,11(6):728-742.
- [12] LIU X H,DING W. Big Data Credit Reporting Practices of Zest-Finance in The United States [J]. *Credit Reference*,2015, 22(8):27-32. (in Chinese)
刘新海,丁伟. 美国 ZestFinance 公司大数据征信实践 [J]. *征信*,2015,22(8):27-32.
- [13] LECUN Y,BENGIO Y,HINTON G. Deep Learning [J]. *Nature*,2015,521(7553):436-444.
- [14] CUI L X,BAI L,HANCOCK E R, et al. Identifying the most informative features using a structurally interacting elastic net [J]. *Neurocomputing*,2018,313(11):65-77.
- [15] ADDO P M,GUEGAN D,HASSANI B. Credit Risk Analysis Machine and Deep Learning Models[J]. *Risks*,2018,6(2):38-57.
- [16] HINTON G E,SALAKHUTDINOV R R. Reducing the dimensionality of data with neural networks [J]. *Science*, 2006, 313 (5786):504-507.
- [17] VINCENT P,LAROCHELLE H,LAJOIE I, et al. Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with aLocal Denoising Criterion [J]. *Journal Machine Learning Research*,2010,27(11):3371-3408.
- [18] SAGHA H,CUMMINS N,SCHULLER B. Stacked Denoising Autoencoders for Sentiment Analysis:A review[J]. *Data Mining and Knowledge Discovery*,2017,7(5):132-146.
- [19] ALHASSAN Z,MCGOUGH A,ALSHAMMARI R, et al. Stacked Denoising Autoencoders for Mortality Risk Prediction Using Imbalanced Clinical Data [C]// *IEEE International Conference on Machine Learning and Applications*. Orlando:IEEE Press,2018:396-401.
- [20] VANMARCKE E H. *Random Fields:Analysis and Synthesis* [M]. Cambridge:MIT Press,1983:92-101.
- [21] YUAN J. *Time-dependent Probabilistic Assessment of Rainfall-induced Slope Failure*[D]. Munich: Technical University of Munich,2016.
- [22] BETZ W,PAPAIIOANNOU I,STRAUB D. Numerical Methods for the Discretization of Random Fields by Means of the Karhunen-Loève Expansion[J]. *Computer Methods in Applied Mechanics and Engineering*,2014,271(0):109-129.