

基于深度神经网络的庞氏骗局合约检测方法



张艳梅 楼胤成

中央财经大学信息学院 北京 100081

摘要 区块链技术的发展吸引了全球投资者的目光。目前,有数以万计的智能合约部署在以太坊上。在给金融、溯源等诸多行业带来颠覆性的创新之余,以太坊上的部分智能合约含有诸如庞氏骗局等欺诈形式,给全球投资者造成了数百万美元的损失。但是,目前针对互联网金融背景下庞氏骗局的定量识别方法较少,针对以太坊上庞氏骗局合约检测的研究较少,且检测精度有进一步提高的空间,文中提出基于深度神经网络的庞氏骗局合约检测方法。该方法提取出智能合约中有助于识别庞氏骗局的特征,如智能合约的操作码特征和账户特征,形成数据集,而后在数据集上训练模型,在测试集上检测性能。实验结果表明,基于深度神经网络的庞氏骗局合约检测方法具有 99.6% 的查准率和 96.3% 的查全率,均优于现有方法。

关键词: 区块链; 以太坊; 智能合约; 庞氏骗局; 深度神经网络

中图分类号 TP309.2

Deep Neural Network Based Ponzi Scheme Contract Detection Method

ZHANG Yan-mei and LOU Yin-cheng

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

Abstract The development of blockchain technology has attracted the attention of global investors. Currently, tens of thousands of smart contracts are deployed on Ethereum. In spite of bringing disruptive innovation to finance, traceability and many other industries, some smart contracts on Ethereum contain fraudulent forms such as Ponzi schemes, causing millions of dollars of losses to global investors. However, at present, there are few quantitative identification methods for Ponzi scheme under the background of Internet finance, few researches on detection of Ponzi scheme contract on Ethereum, and the detection accuracy needs to be improved. Therefore, a Ponzi scheme contract detection method based on deep neural network is proposed. It extracts the features of smart contract that are helpful to identify Ponzi scheme, such as operation code features and account features, to form a data set. Then, the model is trained on the dataset and performance is tested on test set. The experimental results show that the Ponzi scheme contract detection method based on deep neural network has a precision of 99.6% and a recall rate of 96.3%, which are better than that of existing methods.

Keywords Blockchain, Ethereum, Smart contract, Ponzi scheme, Deep neural network

1 引言

区块链作为脱胎于比特币的一种新兴技术^[1],在提出伊始便受到了广大投资者和普通民众的追捧^[2-3]。以太坊^[4]作为区块链 2.0 的代表,更是备受瞩目。以太坊^[5]是一个开源的公链平台,具有一套图灵完备的编程语言。以太坊将区块链技术的外延从区块链 1.0 中单纯的数字货币扩展到了金融、溯源等多种领域。智能合约^[5-6]是一套以数字形式定义的承诺。智能合约由 solidity, go 等高级语言编写完成后,需要先被编译成字节码,再通过以太坊客户端上传至以太坊公链中,才算完成合约的部署,合约方能生效。因此,以太坊公链上部署有大量的智能合约,可实现各式各样的业务^[7-9]。

区块链具有去信任化、安全性、分布式、透明性等技术优

势^[10-11],但同时也增加了区块链相关技术的复杂度,使得投资者难以解读以太坊上智能合约的具体业务逻辑,一般仅能通过智能合约上由开发者出具的少量描述性信息来了解业务的运行机制。

一些投机分子利用区块链的上述特性,将传统金融投资领域诈骗的经典形式——庞氏骗局,引入到区块链中,给投资者造成了极大的损失。智能合约中的庞氏骗局更具迷惑性^[12],因为区块链具有不可更改的特性,所以被上传至以太坊中的合约将永远不会失效。这使得很多投资者认为:一个持续运行的且能够不断获得收益的智能合约项目,不存在类似庞氏骗局的风险^[13-14]。

实则不然,庞氏骗局合约的运营者通过更加隐蔽的方式来获取利润,如:1)对每笔投资收取较高的手续费,每当投资

收稿日期:2019-11-03 返修日期:2019-12-27 本文已加入开放科学计划(OSID),请扫描上方二维码获取补充信息。

基金项目:国家自然科学基金(61602536,61773415,61672104)

This work was supported by the National Natural Science Foundation of China(61602536,61773415,61672104).

通信作者:张艳梅(jlzym0309@sina.com)

者对庞氏骗局合约进行投资时,运营者先抽取投资额的一定比例到自己的账户,再将剩余的部分转入用于发放投资收益的合约账户;2)庞氏骗局合约会按投资时间先后维护一个投资者列表,按照收益递减的方式顺序发放回报。在投入相同的投资额的情况下,在每次发放收益时,先投资的投资者比后投资的投资者收到的回报更高。因此,在一个庞氏骗局合约被发布后,骗局的运营者通常自己先进行投资,使得自身具有最高的优先级。在每次合约自动发放投资收益时,运营者均能以投资额的较大比例获得收益。

与此同时,由于一些先进入庞氏骗局合约的投资者确实会收到一定金额的收入,这使得该部分投资者信以为真,认为该合约背后有一个好的投资项目。实际上,只有极少数最早的投资者获得了超过投资金额的收入,大部分投资者均处于亏损状态。

由于区块链技术专业性高,绝大部分投资者和用户缺乏专业知识,难以对智能合约中的欺诈行为进行有效识别。并且,区块链技术具有匿名性,所有账号信息均不可追溯,这使得庞氏骗局合约的运营者可以肆无忌惮地发布类似的诈骗合约,不断作恶获利。因此,对以太坊上庞氏骗局合约的识别工作已迫在眉睫。

随着互联网金融的发展,庞氏骗局也发展出了许多新的形式^[15]。有学者估算,庞氏骗局合约大概占以太坊上所有智能合约数量的 0.03%^[16-17]。考虑到以太坊上具有多达 200 万份的智能合约,庞氏骗局合约的数目也十分庞大。目前,针对以太坊上庞氏骗局合约的研究并不多。Chen 等^[16]从机器学习的角度出发,用随机森林方法检测智能合约中是否有庞氏骗局,取得了 95% 的查准率,69% 的查全率。但是,目前的研究现状存在以下问题:第一,针对以太坊上庞氏骗局合约的研究比较少;第二,针对庞氏骗局合约的检测精度仍有一定的提升空间。

本文基于文献^[16]中的数据集进行扩展,筛选出用于识别庞氏骗局合约的相关特征,如智能合约的操作码特征和账号特征,并采用深度神经网络(DNN)的方法进行建模,实现对庞氏骗局合约的识别。实验结果显示,与现有的检测方法相比,基于深度神经网络的检测方法在真实数据集上具有更高的精度,证明了该方法的有效性。本文的主要贡献如下:

(1)创新性地引入了深度神经网络的方法对庞氏骗局合约进行检测识别。

(2)模型在测试集上的查准率为 99.6%,查全率为 96.3%,相比目前最优的随机森林方法,分别提高了 4.6% 和 27.3%,提高了识别的精度。

本文第 2 节介绍了相关工作;第 3 节详细描述了数据集、特征选取和基于深度神经网络方法的庞氏骗局合约识别模型;第 4 节阐述了数据预处理、参数选取、评估指标选取、实验结果及性能比较;最后总结全文并展望未来。

2 相关工作

总体来说,本文的相关工作大致包括 3 类:1)庞氏骗局的相关研究;2)深度学习的相关研究;3)对区块链上诈骗问题的

研究。对这三方面的相关工作进行具体介绍如下。

2.1 关于庞氏骗局的研究

庞氏骗局^[18]是一种古老但却屡试不爽的投资诈骗。庞氏骗局的组织者用新加入的投资者的资金来支付对老投资者的回报,多以低风险高回报、收益稳定等特点来对投资项目进行包装,迷惑了大量具有侥幸心理、对投资行业不甚了解的投资者。随着互联网金融的发展,打着“众筹”“P2P”等名号的庞氏骗局兴起,造成了极为恶劣的社会影响^[15]。但目前,除了从投资者角度进行宣传教育,提高他们的防范意识;还应从政府角度加强立法,完善监管。除了这两个定性的角度外^[19],并没有形成对庞氏骗局识别行之有效的定量方法。

2.2 关于深度学习的研究

随着硬件的改进和计算机算力的增强,需要进行海量计算的深度学习也取得了长足的发展^[20-21]。深度学习对特征工程的依赖较小,能处理具有高维特征的数据,适应性较强,相似的算法能应用于不同的领域^[22],在一定程度上减小了模型选择的难度。与传统的方法相比,深度学习通常具有更高的精度。针对诸如图像辨识、语音识别、自然语言处理、目标检测等诸多分类领域,目前已有多种深度学习的网络模型,如卷积神经网络、长短时记忆网络、生成式对抗网络被提出和实现,且能较好地解决相应的问题^[23]。Franscois^[24]指出,深度学习在处理具有高特征维度数据的分类问题中有着较好的应用。

从本质上来说,深度学习是一个不断分层级进行特征提取和描述的过程^[25]。目前,最为普遍且经典的深度学习框架是深度神经网络(Deep Neural Network, DNN)^[26]。DNN 在神经层的激活函数、目标函数、优化方式等多方面都采用了改进后的方案,如 ReLU, softmax 等激活函数, RMSProp, Adam 等优化方法。这使得深度神经网络的精度有了极大的提高。更重要的是,对于新输入的数据, DNN 无需将其与之前的数据整合后重新训练出一个新的模型,而是能直接对新获得数据进行训练。这使得 DNN 比一些机器学习方法更适用于连续性的在线学习^[24],无疑与我们想要实现对庞氏骗局合约的在线检测,并且通过检测数据完成对模型的学习和更新的初衷更相符。

2.3 关于区块链上诈骗问题的研究

如引言中所述,区块链在全球范围内引发了投资热潮,但也因其高技术门槛使得投资者更难对庞氏骗局等欺诈形式做出有效的识别。据统计^[27],仅 2013 年 9 月 2 日至 2014 年 9 月 9 日,比特币相关诈骗案件涉及的金额就高达 700 万美元。Vasek 等^[28]将比特币中的骗局分为庞氏骗局、挖矿骗局、钱包骗局和欺诈性交易 4 种诈骗类型。Bartoletti 等从智能合约的描述性信息、智能合约源代码以及相关交易记录 3 个方面出发,分析了庞氏智能合约的共性^[27]。Chen 等^[16]利用随机森林(RF)的方法对庞氏智能合约进行识别,取得了比其他机器学习方法更高的精度。

综上所述,目前国内外已有大量针对区块链、庞氏骗局以及深度学习的相关研究成果,但对区块链上诈骗的探索较少,而使用深度神经网络的方法对以太坊上庞氏骗局合约进行识别的研究几乎没有。

3 识别庞氏骗局合约的神经网络模型

基于神经网络识别庞氏骗局的主要思想是:首先通过 etherscan.io 接口对合约的字节码等相关数据进行爬取,得到原始数据集。在获得数据后,需要对其进行处理,提取出合约的操作码特征和账户特征作为庞氏骗局合约识别的相关特征,然后对数据集中的具体数据值进行标准化,并按照 3:1:1 的比例三分数据集,其中,训练集占 60%,验证集占

20%,测试集占 20%。而后采用 k 折验证的方法将训练集和验证集用于识别庞氏骗局合约的神经网络模型的训练。在模型训练完毕后,将测试集作为输入,获得模型判断测试集上的智能合约是否为庞氏骗局合约的识别结果。将识别结果与合约的标签值进行比较(标 0 为正常合约,标 1 为庞氏骗局合约),统计查准率、查全率等性能指标,可得神经网络方法对庞氏骗局的识别精度。图 1 给出了基于神经网络识别庞氏骗局方法的框架。

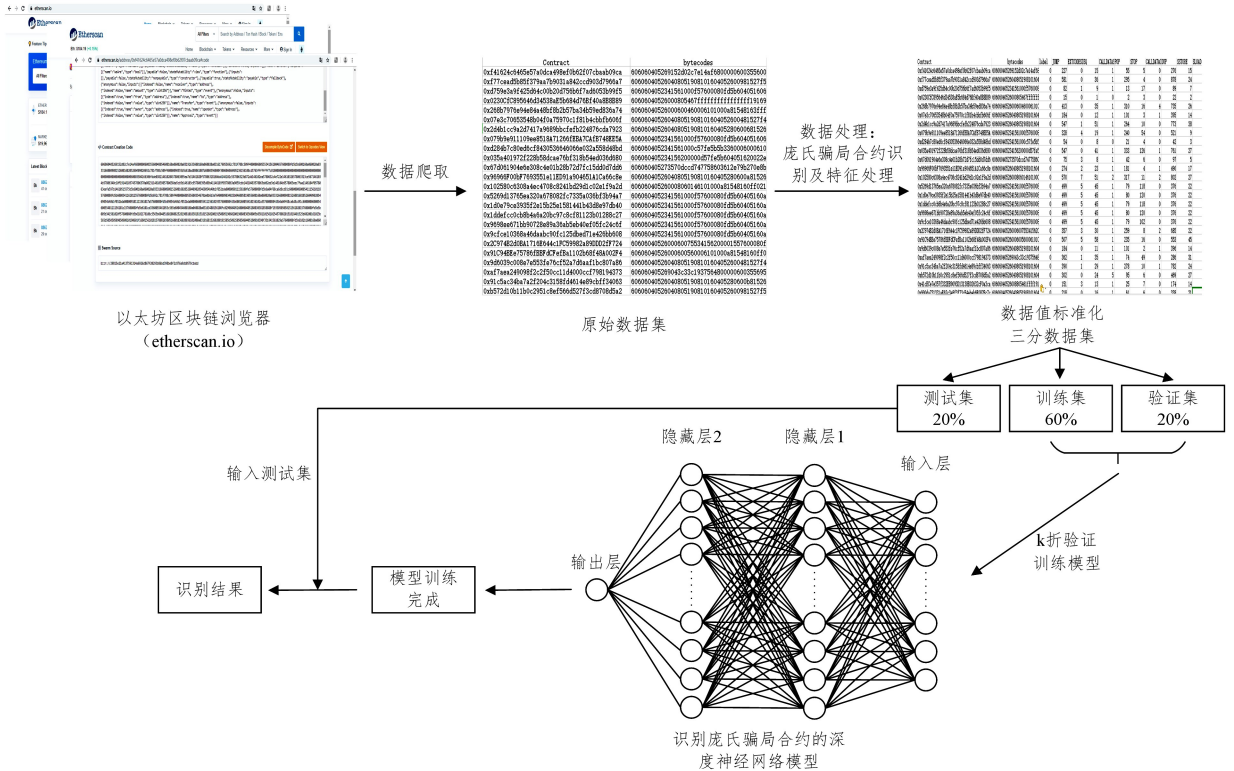


图 1 基于神经网络识别庞氏骗局的方法框架

Fig. 1 Framework of deep neural network model for identifying ponzi scheme contracts

3.1 数据集

本文采用的数据集是基于文献[17]中开源共享的数据集。其中包含有 3774 份智能合约,已经对其是否为庞氏骗局合约进行了标注。本文根据已有标签合约的地址,在 etherscan.io 上对相应合约的字节码等信息进行了爬取。通过统计可知,数据集中含有 132 份庞氏骗局合约,3 642 份正常智能合约,正反样本的比例大约为 55:2。本文在含有 3 774 份智能合约的数据集上进行模型的训练、验证和测试。

3.2 特征的选取

本文实验选取了 67 种有可能对庞氏骗局合约识别有益的特征,分别是 64 种操作码特征和 3 种账户特征,具体阐释如下。

3.2.1 操作码特征

通过对比观察庞氏骗局合约与正常智能合约的源代码可以发现两种合约的逻辑有所差别,如庞氏骗局合约有着不断向合约创建者的账户进行转账、设置有一个递减的变量,维护一个账户地址列表以进行收益的发放等诸多共性。毫无疑问,这些特征会体现在合约执行过程中所调用的操作码上。

基于这一差别,通过对 3 374 份智能合约进行统计,得到在所有智能合约中出现的操作码种类及被使用的次数。本文实验选用了 64 种不同的操作码在每份合约中被调用的次数,并将其作为识别庞氏骗局合约的依据之一。

3.2.2 账户特征

根据以太坊上庞氏骗局合约的特征可知,庞氏骗局合约与正常合约在账户的特征上也必然存在着差异,如:为了从庞氏骗局中获利,庞氏骗局合约对庞氏骗局运营者账户支付的金额会远远大于对其他正常投资者账户支付的金额。这会导致庞氏骗局合约向参与者付款的次数少于其收到投资的笔数。而且,为了使投资者相信自己能够从该合约中获利,当庞氏骗局合约的账户余额满足支付条件时,合约会立即进行回报的发放,这导致庞氏骗局合约的合约账户余额远低于正常合约的合约账户余额。

因此,本文将合约账户的余额、智能合约向所有参与者支付金额中的最大值、合约的支付率(智能合约向参与者付款的次数/参与者对智能合约投资的笔数)这 3 类账户特征作为识别庞氏骗局合约的特征。

本文使用神经网络的方法对庞氏骗局合约识别进行建模。从本质上说,庞氏骗局合约的识别是一个二分类问题。依据经验知识及已有数据集选择最优方案:选择 ReLU 函数作为中间层的激活函数,选择 sigmoid 函数作为最后一层的激活输出函数,选择二元交叉熵作为损失函数,选择 rmsprop 优化器作为模型的优化器。本实验的数据集共含有 3774 条数据,对于深度学习来说,该数据集并不算大。随着训练的层数和轮数的增大,很有可能出现过拟合。因此,在各隐藏层之间设置随机失活比率(dropout),实现正则化,以提高模型的泛化能力。所构建的模型如图 2 所示。

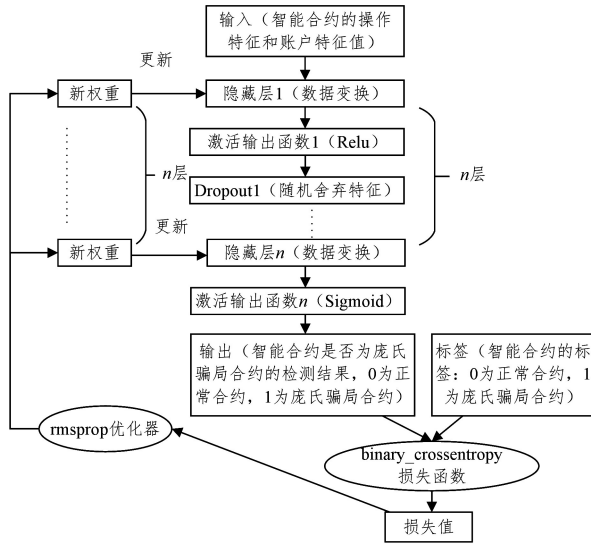


图 2 识别庞氏骗局合约的神经网络模型

Fig. 2 Deep neural network model for identifying ponzi scheme contracts

图 2 详细阐释了模型的原理:在筛选出与庞氏骗局合约相关的特征后,模型将智能合约操作特征和账户特征的具体特征值标准化,并将其作为输入。模型共含有 n 层隐藏层,输入在隐藏层 1 进行数据变换后,由激活函数 ReLU 进行非线性变化,并按照事先指定的随机失活比率(dropout)对数据特征进行舍弃,降低了模型的结构风险。将所得结果作为下一隐藏层的输入,在每一隐藏层中进行相同的操作。循环 $n-1$ 次后,进入隐藏层 n ,通过激活函数 sigmoid 得到预测值,即合约是否为庞氏骗局合约的判断。将输出与该智能合约的标签(标 0 为正常智能合约,标 1 为庞氏骗局合约)进行比较,通过损失函数衡量模型在训练集上的表现。在完成对训练集上所有数据的学习和判断后,将损失值反馈给 rmsprop 优化器,优化器对每一隐藏层中的权重进行优化更新,以进行新一轮的训练。如此循环往复,使得模型对训练集有较为精准的分类能力。最后,在测试集上对智能合约进行检测,得出模型对庞氏骗局合约的检测精度。

4 实验

4.1 数据预处理

由于智能合约的特征值具有不同的取值范围,因此在训练开始前需要对数据集进行预处理。将每个特征值进行标准化,使得每个特征的平均值为 0、标准差为 1,以消除数据取值

范围的差异,简化后续的学习过程。而后,将数据集以 3:1:1 的比例分成 3 份,分别作为训练集、验证集和测试集。训练集和验证集用于模型的训练与验证,测试集用于模型性能的测试。此外,在保证正反样例之比与整个数据集上正反样例分布相近的前提下,对每份数据集中的数据进行随机打乱,并使用 k 折验证的方式进行模型的训练,在一定程度上弥补了数据集较小的不足。

4.2 参数选取

在确定问题的本质是二分类问题时,除了能确定激活函数、损失函数和优化器这三者的选择之外,其他参数,如训练轮数、隐藏层层数、每层神经元个数、随机失活比率(dropout)和每轮训练的批尺寸(batch_size),都需要通过实验找出一个相对合适的取值。由于变量多且取值范围广,难以对每种变量组合进行遍历取最优值,因此,本次实验依据贪心算法的思想,在其他变量确定的前提下,对某个变量的取值进行选择,以期找到一个相对较优的参数组合。为了避免神经网络训练过程中的一些偶然性,使得实验结果更具一般性,在下文的参数选取实验中,数据值均为相同实验重复 10 次后去掉最大值与最小值后的平均值。

4.2.1 训练轮数

毫无疑问,随着训练轮数的增加,模型会出现过拟合的情况,因此可以设定一个较大的训练轮数,观察随着训练轮数的增加,模型验证正确率的变化情况。同时,为了更好地观察模型的走势,将每轮训练的正确率替换为前一轮正确率的指数移动平均值。设训练轮数为 100,图 3 为每轮的训练损失和验证损失,图 4 为每轮的训练正确率和验证正确率。

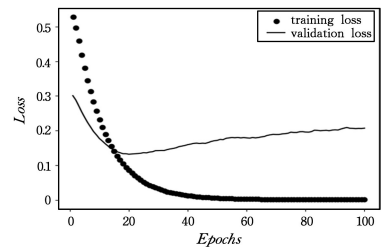


图 3 训练损失和验证损失

Fig. 3 Training loss and validation loss

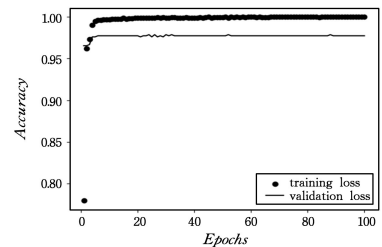


图 4 训练正确率和验证正确率

Fig. 4 Training accuracy and validation accuracy

综合图 3 和图 4 可知,在训练超过 40 轮后,训练损失有明显的下降,验证损失有明显的上升;验证正确率有明显的上升,验证正确率略有下降。这证明训练超过 40 轮后会出现过拟合的情况,因此把模型的训练轮数定为 40。

4.2.2 隐藏层层数

由于所训练的模型在本质上较为简单,因此并不需要特

别多的隐藏层,且过多的隐藏层会导致过拟合问题。本实验设隐藏层层数的取值范围为 2~10,不同层数的深度神经网络的训练正确率和验证正确率如图 5 所示。

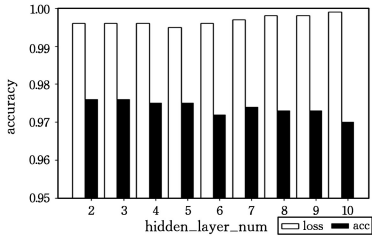


图 5 隐藏层层数不同时的训练正确率和验证正确率
Fig.5 Training accuracy and verification accuracy with different numbers of hidden layer

从图 5 可以看出,随着隐藏层层数的增加,模型的训练正确率有所上升,验证正确率有所下降。故得出结论:随着隐藏层层数的增加,模型有过拟合的趋势。由于隐藏层层数为 2 层时测试正确率已经很高,结合奥卡姆剃刀原理,把隐藏层层数定为 2 层。

4.2.3 每层神经元的个数

本实验中假设每层神经元个数的选择范围是:{8,16,32,64,128}。图 6 给出了每层神经元个数不同时模型的训练正确率和验证正确率。

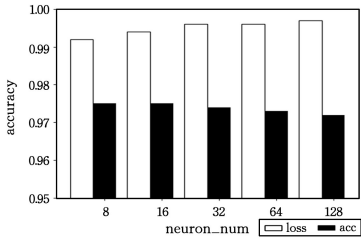


图 6 每层神经元个数不同时的训练正确率和验证正确率
Fig.6 Training accuracy with validation accuracy with different numbers of neuron in each layer

当神经元数量过少时,少量的神经元需要表示大量的特征,会导致一些特征信息丢失;当神经元过多时,深度神经网络极有可能是仅仅简单地记住了每个输入所对应的输出,容易出现过拟合的情况。根据图 6 所示的每层神经元个数不同时的训练正确率、验证正确率,并根据奥卡姆剃刀原理,将每个隐藏层的神经元个数设为 16。

4.2.4 随机失活比率(dropout)

随机失活比率(dropout)是深度神经网络训练过程中有效且常用的正则化方法。在某一隐藏层之后使用 dropout,就是在训练过程中随机将该层的一些输出特征舍弃(置 0)。随机失活比率(dropout)是指在上述过程中被舍弃的特征所占的比例,通常取值范围为 0.2~0.5。Dropout 在训练过程中引入噪声,使得训练过程中的一些偶然性结果不会被模型所记忆。但需要注意的是,随机失活比率(dropout)并不是越大越好,过大的比率会导致模型较难收敛,在训练集上便无法取得令人满意的正确率。当随机失活比率(dropout)在{0.2,0.3,0.4,0.5}中取值时,训练正确率和验证正确率如图 7 所示。

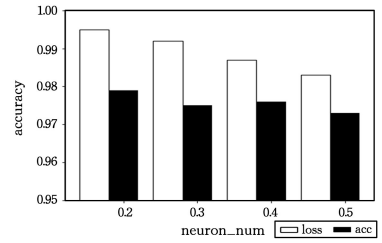


图 7 Dropout 比率不同时的训练正确率和验证正确率
Fig.7 Training accuracy and validation accuracy with different dropout ratios

由图 7 可知,随着随机失活比率(dropout)的增加,模型的训练正确率、验证正确率都有所下滑。又考虑到 dropout 值较小时,正则化项所占的比重不大,模型存在过拟合的风险,因此选取随机失活比率(dropout)为 0.3。

4.2.5 批尺寸(batch_size)

批尺寸(batch_size)与深度神经网络模型在每次学习过程中梯度下降的方向有关,太小的批尺寸会使得模型较难收敛且学习所消耗时间过长,太大的批尺寸对内存的要求较高,且会使收敛的轮数相对较长。因此,对于批尺寸大小的选取需要在训练时间和收敛所需训练轮数之间做出权衡。图 8—图 10 分别显示了批尺寸在{1,32,64,128,256,512}内取值时,模型完成训练所需的时间、模型收敛的具体轮数以及模型的训练正确率和验证正确率。

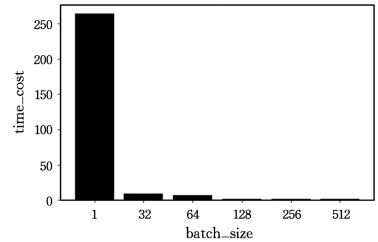


图 8 batch_size 不同时的训练时间
Fig.8 Training time with different numbers of batch_size

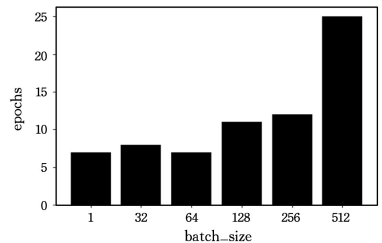


图 9 batch_size 不同时的收敛轮数
Fig.9 Convergence rounds with different numbers of batch_size

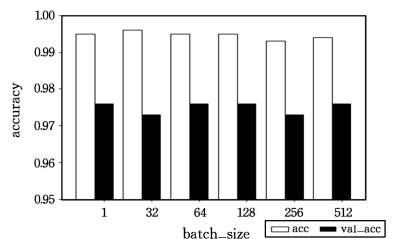


图 10 batch_size 不同时的训练正确率和验证正确率
Fig.10 Training accuracy and validation accuracy with different numbers of batch_size

通过对模型训练所需时间、模型收敛所需轮数、训练正确率和验证正确率进行权衡,本文实验选择批尺寸(batch_size)的大小为64。

4.3 指标选取

为了方便与采用其他方法对庞氏骗局合约进行检测识别的模型^[16]进行性能比较,本次实验中除了较为直观地训练正确率、验证正确率和测试正确率输出之外,还采用了查准率(Precision)、查全率(Recall)和F-score值3个指标对模型的性能进行度量。这3个指标的具体定义如下:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F\text{-score} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

其中,TP表示实际为正样本,且被判别为正样本的样本数量;FP表示实际为负样本,但被判别为正样本的样本数量;FN表示实际为正样本,但被判别为负样本的样本数量。

4.4 实验结果

本文实验采用神经网络的方法进行建模,选择ReLU函数作为中间层的激活函数,sigmoid函数作为最后一层输出的激活函数,二元交叉熵作为损失函数,rmsprop优化器作为模型的优化器,并根据第3节中所确定的参数,设定训练轮数为40,隐藏层层数为2,每层含有16个神经元,随机失活比率(dropout)为0.3,批尺寸(batch_size)为64。在对数据进行标准化的预处理后,将其拆分,60%用于训练,20%用于验证,20%用于测试。在训练集和测试集上,使用k折验证的方法对模型进行训练。训练完成后,在测试集上测试对庞氏智能合约的识别精度。实验在相同的条件下进行了10次,表1列出了每次训练后模型的性能指标值,以及10次训练后各项指标在去除了最大值和最小值后的均值。

表1 10次实验指标值及其平均值

Table 1 Index value and average of ten experiments

Experiment round	Train_accuracy	Validation_accuracy	Test_accuracy	Precision	Recall	F-score
1	0.993	0.983	0.964	1	0.96	0.982
2	0.996	0.984	0.963	0.998	0.964	0.981
3	0.995	0.986	0.962	0.997	0.964	0.980
4	0.992	0.983	0.960	0.995	0.964	0.979
5	0.983	0.994	0.963	0.998	0.963	0.981
6	0.996	0.978	0.961	0.991	0.969	0.980
7	0.996	0.985	0.963	0.998	0.963	0.981
8	0.997	0.984	0.962	0.998	0.964	0.981
9	0.994	0.983	0.962	0.998	0.964	0.981
10	0.994	0.981	0.939	0.973	0.963	0.968
Average	0.994	0.983	0.962	0.996	0.963	0.980

由表1可得,本文所构建的神经网络模型在训练集上的正确率为99.4%,在验证集上的正确率为98.3%,在测试集上的正确率为96.2%,查准率为99.6%,查全率为96.3%,F-score值为98%。可以看出,基于神经网络的模型对于庞氏骗局合约有着较为不错的识别能力。

4.5 性能比较

文献^[16]总结了采用传统机器学习方法对庞氏骗局合约的识别性能。因此,本文不仅采用了神经网络方法对庞

氏骗局合约进行检测,还复现了一些经典的机器学习分类算法,如决策树(DT)、支持向量机(SVM)、极端梯度增强(XGBoost)、一类支持向量机(OCSVM)、隔离林(IF)和随机森林(RF),以比较并衡量所采用的神经网络方法对庞氏骗局合约检测的适用性。使用查准率、查全率和F-score值对上述方法识别庞氏骗局合约的精度进行度量,结果如表2所列。

表2 不同分类方法的表现比较

Table 2 Performance comparison of different classification methods

Algorithm	Precision	Recall	F-score
DT	0.31	0.24	0.27
SVM	0.91	0.16	0.27
XGBoost	0.90	0.67	0.76
OCSVM	0.05	1.00	0.10
IF	0.02	0.05	0.04
RF	0.95	0.69	0.79
DNN	0.996	0.963	0.980

通过比较表2中各分类方法的性能指标值不难发现:在传统的机器学习方法中,随机森林(RF)方法对庞氏骗局合约的识别效果在3个指标上的表现最为稳定和优秀,其查准率为95%,查全率为69%,F-score值为79%。但与之相比,本文所采用的神经网络方法仍能使查准率略微提高4.6%,查全率提高27.3%,从而使F-score提高了19%。即,本文基于神经网络方法建模,提高了对庞氏骗局合约的识别精度。

结束语 本文将智能合约的操作码特征和账户特征作为识别庞氏智能合约的主要特征,基于神经网络的方法构建庞氏骗局合约的识别模型。依据实践经验给出模型中关键参数的取值范围,依据贪心算法的思想,通过模型在训练集、验证集和测试集上的表现选取最佳的参数,最终实现了对以太坊上庞氏骗局合约的有效识别,且相比其他现有的方法,基于神经网络的方法在检测精度和正确率上均取得了一定的提高,查准率为99.6%,查全率为96.3%。值得注意的是,即使采用了k折验证的方法来缓解小样本问题,但数据集相对较小仍不可忽视。

针对数据集较小,人工提取海量智能合约的相关特征并不现实这一问题,我们计划在接下来的工作中采用CNN直接对智能合约编译后生成的字节码进行学习。由于智能合约源代码中会存在函数和参数的调用关系,因此我们猜测,使用CNN对找出庞氏骗局合约和非庞氏骗局合约的源代码之间的区别可能会有所帮助。而且,本文使用的操作码特征也提取自合约的字节码中。考虑到每份合约字节码长度不一的事实(即数据集尺寸不一致),我们拟选择空间金字塔池化(SPP)作为研究的方法。这样只需通过爬虫在etherscan.io上获取到智能合约的字节码,即可对模型进行训练和学习,而不用人工对每个合约进行特征的提取和分类。如此能够有效地扩展数据集的大小,在现实应用中也能更方便用户对可疑智能合约进行在线检测。

将来的另一个研究方向是使用生成式对抗网络(GANs)方法。GANs可以有效地解决数据具有高维特征分布的问题,这意味着其可以用来进行异常检测,符合庞氏骗局合约识

别这一应用场景。GANs 方法的另一个优秀的特征是,它可以进行无监督学习,这意味着 GANs 能实现对正常合约的聚类,能够有效地解决在识别庞氏骗局合约过程中遇到的样本不平衡问题。

参 考 文 献

- [1] ZHENG Z B,XIE S A. Blockchain challenges and opportunities: A survey[C/OL]// International Journal of Web and Grid Services. http://https://xueshu.baidu.com/usercenter/paper/show?paperid=7e00413a964b3b16c3495eb19c64a1f4&.site=xueshu_se.
- [2] SWAN M. Blockchain:Blueprint for a New Economy[M]. Newton,MA,USA:O'Reilly Media,2015.
- [3] Bitcoin:A Peer-to-Peer Electronic Cash System. [OL]. <https://bitcoin.org/bitcoin.pdf>.
- [4] CoinDesk. Understanding Ethereum-blockchain Research Report [OL]. www.coindesk.com/research/understandingethereum-report/.
- [5] A Next-Generation Smart Contract and Decentralized Application Platform. [OL]. <https://github.com/ethereum/wiki/wiki/WhitePaper>.
- [6] SZABO N. Smart Contracts:Building Blocks for Digital Markets [OL]. http://www.fon.hum.uva.nl/rob/Courses/InformationInSpeech/CDROM/Literature/LOTwinterschool2006/szabo.best.vwh.net/smart_contracts_2.html.
- [7] BOCEK T. Digital Marketplaces Unleashed[M]. Springer-Verlag GmbH. 2017-09-15:169-184. ISBN 978-3-662-49274-1.
- [8] NORTA A. Creation of smart-contracting collaborations for decentralized autonomous organizations [OL]. https://link.springer.com/chapter/10.1007%2F978-3-319-21915-8_1.
- [9] CHRISTIDIS K,DEVETSIKIOTIS M. Blockchains and smart contracts for the internet of things[C]// IEEE Access. 2016: 2292-2303.
- [10] HE P, YU G, ZHANG Y F, et al. Survey on Blockchain Technology and Its Application Prospect [J]. Computer Science, 2017,44(4):1-7,15.
- [11] WANG Q G, HE P, NIE T Z, et al. Survey of Data Storage and Query Techniques in Blockchain Systems[J]. Computer Science, 2018,45(12):12-18.
- [12] HIGGINS S. SEC Seizes Assets from Alleged Altcoin Pyramid Scheme[OL]. <https://www.coindesk.com/sec-seizesalleged-altcoin-pyramid-scheme>.
- [13] KEIRNS G. Gemcoin Ponzi Scheme Operator Hit with \$74 Million Judgment. [OL]. <https://bitcoinwiki.co/gemcoinponzi-scheme-operator-hit-with-74-million-judgment/>.
- [14] MORRIS D Z. The Rise of Cryptocurrency Ponzi Schemes [OL]. <https://www.theatlantic.com/technology/archive/2017/05/cryptocurrency-ponzi-schemes/5286>.
- [15] ZHAO M. Identification and prevention of Ponzi scheme under the background of internet finance[J]. Zhejiang Finance, 2016(8):13-17.
- [16] CHEN W,ZHENG Z,NGAI E,et al. Exploiting Blockchain Data to Detect Smart Ponzi Schemes on Ethereum[J/OL]. IEEE Access,2019:1-1. https://www.researchgate.net/publication/331853833_Exploiting_Blockchain_Data_to_Detect_Smart_Ponzi_Schemes_on_Ethereum.
- [17] CHEN W,ZHENG Z,CUI J,et al. Detecting ponzi schemes on ethereum:Towards healthier blockchain technology[C]// Proc. World Wide Web Conf. World Wide Web,2018:1409-1418.
- [18] Wikipedia. PonziScheme[OL]. https://en.wikipedia.org/wiki/Ponzi_scheme.
- [19] YAO L,CHEN W. The Enlightenment of American P2P Supervision[J]. China Finance,2015(7):63-64.
- [20] DENG L,YU D. Deep Learning:Methods and Applications[J]. Foundations & Trends in Signal Processing,2014,7(3).
- [21] LECUN Y,BENGIO Y,HINTON G. Deep learning [OL]. <https://www.nature.com/articles/nature14539>.
- [22] SCHMIDHUBER,JÜRGEN. Deep Learning in Neural Networks:An Overview[J]. Neural Netw,2015,61:85-117.
- [23] JIAO L C,YANG S Y,LIU F,et al. Seventy Years Beyond Neural Networks:Retrospect and Prospect[J]. Chinese Journal of Computers,2016,39(8):1697-1716.
- [24] FRANSCOIS C. Deep Learning with Python[M]. Beijing:Posts and Telecommunications Press,2018.
- [25] HUANG L W,JIANG B T,LU S Y,et al. Survey on Deep Learning Based Recommender Systems[J]. Chinese Journal of Computers,2018,41(7):1619-1647.
- [26] LI C,CHAI Y M,NAN X F,et al. Research on Problem Classification Method Based on Deep Learning[J]. Computer Science, 2016,43(12):115-119.
- [27] BARTOLETTI M,CARTA S,CIMOLI T,et al. Dissecting ponzi schemes on ethereum: Identification, analysis, and impact [OL]. <https://arxiv.org/abs/1703.03779>.
- [28] VASEK M,MOORE T. There's No Free Lunch, Even Using Bitcoin: Tracking the Popularity and Profits of Virtual Currency Scams[C]// Springer Berlin Heidelberg. 2015:44-61.



ZHANG Yan-mei, born in 1976, Ph.D. professor, is a member of China Computer Federation. Her main research interests include business intelligence, service computing and blockchain.