

# 面向征信数据安全共享的 SVM 训练机制

沈蒙<sup>1)</sup> 张杰<sup>1)</sup> 祝烈煌<sup>1)</sup> 徐恪<sup>2),3)</sup> 张开翔<sup>4)</sup> 李辉忠<sup>4)</sup> 唐湘云<sup>1)</sup>

<sup>1)</sup>(北京理工大学计算机学院 北京 100081)

<sup>2)</sup>(清华大学计算机科学与技术系 北京 100084)

<sup>3)</sup>(北京信息科学与技术国家研究中心 北京 100084)

<sup>4)</sup>(深圳前海微众银行股份有限公司 深圳 518052)

**摘要** 在征信行业中,征信数据的丰富性和多样性对信用评价极为重要.然而,征信机构尤其是小型征信机构拥有的征信数据存在内容不完整、种类不全、数量不充足等问题.同时由于征信数据价值高、隐私性强、易被非授权复制,征信机构之间难以直接共享数据.针对这一问题,本文提出了面向征信数据安全共享的 SVM 训练机制.首先共享数据经同态加密后存储在区块链上,保证数据不可篡改以及隐私安全.其次使用基于安全多方计算的支持向量机(SVM)在共享的加密数据上进行运算,保证在不泄露原始数据的条件下,训练信用评价模型.最后,通过真实数据集上的实验对本文所提出机制的可用性和性能进行验证.实验结果显示,相比于基于明文数据集训练出的模型,本文提出的机制在可接受时间内训练出的模型无准确率损失.同时,与其他同类隐私训练方案相比,本机制在实验数据集上的计算耗时小于对比实验的 5%,且无需可信第三方协助计算.

**关键词** 联盟链;征信数据;支持向量机;隐私保护;同态加密

**中图法分类号** TP309 **DOI号** 10.11897/SP.J.1016.2021.00696

## SVM Training Mechanism for Secure Sharing of Credit Data

SHEN Meng<sup>1)</sup> ZHANG Jie<sup>1)</sup> ZHU Lie-Huang<sup>1)</sup> XU Ke<sup>2),3)</sup> ZHANG Kai-Xiang<sup>4)</sup>

LI Hui-Zhong<sup>4)</sup> TANG Xiang-Yun<sup>1)</sup>

<sup>1)</sup>(School of Computer Science, Beijing Institute of Technology, Beijing 100081)

<sup>2)</sup>(Department of Computer Science and Technology, Tsinghua University, Beijing 100084)

<sup>3)</sup>(Beijing National Research Center for Information Science and Technology, Beijing 100084)

<sup>4)</sup>(Shenzhen Qianhai Micro Public Bank Co., Ltd., Shenzhen 518052)

**Abstract** In the credit reporting industry, the richness and diversity of credit reporting data is extremely important for the development of credit evaluation. However, credit data owned by credit reporting agencies, especially small credit reporting agencies, has issues like incomplete content, incomplete types, and insufficient instance numbers. Therefore, data sharing among credit reporting agencies is very important. In practical application scenarios, credit data has the characteristics of high value, strong privacy, and easy to be copied without authorization. These characteristics will cause great security challenges when sharing credit data. To solve this problem, this paper proposes a SVM training mechanism for secure sharing of credit data. Meanwhile we design a system prototype based on this mechanism, as showed in Figure 3 in the manuscript. This mechanism is based on the consortium blockchain and the addition homomorphic encryption

收稿日期:2019-09-24;在线发布日期:2020-02-20. 本课题得到国家重点研发计划(2018YFB0803405)、国家自然科学基金(61902039, 61872041, 61932016)、北京市自然科学基金(4192050)和 CCF-腾讯犀牛鸟基金微众银行专项基金资助. 沈蒙, 博士, 副教授, 中国计算机学会(CCF)会员, 主要研究方向为云计算隐私保护、区块链技术与应用. E-mail: shenmeng@bit.edu.cn. 张杰, 硕士研究生, 主要研究方向为区块链、隐私保护. 祝烈煌(通信作者), 博士, 教授, 国家重点研发计划首席科学家, 中国计算机学会(CCF)会员, 主要研究领域为密码学、网络与信息安全. E-mail: liehuangz@bit.edu.cn. 徐恪, 博士, 教授, 国家杰出青年科学基金入选者, 中国计算机学会(CCF)会员, 主要研究领域为网络安全可信、区块链应用. 张开翔, 本科, 微众银行区块链首席架构师, 主要研究方向为区块链. 李辉忠, 硕士, 微众银行区块链高级架构师, 中国计算机学会(CCF)会员, 主要研究方向为区块链底层技术. 唐湘云, 博士研究生, 主要研究方向为应用密码学.

scheme Paillier. With the decentralization of blockchain technology, this mechanism does not need to rely on any trusted third party during model training. At the same time, through secure collaborative computing between credit reporting agencies, the mechanism can meet the credit evaluation needs of the model trainer without revealing data privacy. Firstly, the shared data is stored on the blockchain and is encrypted to ensure that the data is secure and cannot be tampered. This process is completed through smart contracts, without the need for a third party as a data sharing platform. Secondly, based on the addition homomorphic encryption algorithm Paillier, this paper implements various secure operations in the SVM training process based on the stochastic gradient descent algorithm, and designs a secure SVM training algorithm according to the training process. The algorithm flow is shown in Algorithm 2. Based on this algorithm, the credit reporting agencies participating to the calculation can perform operations on the shared encrypted data, ensuring that the model trainer can train the credit evaluation model without leaking the original data. During the training process, only the data provider and a model trainer participate in the calculation. The calculation based on the encrypted data does not require the assistance of a third party, which avoids the risk of privacy leakage caused by the introduction of a third party. The mechanism proposed in this paper is verified by security analysis. In the threat model, neither the model parameters of the model trainer nor the original data of the data provider will have the problem of privacy leakage. At the same time, this paper verifies the usability and performance of the proposed mechanism through experiments on real-world datasets. The experimental results show that compared with the model trained on the plaintext data set under normal conditions, the model trained by the proposed mechanism has no loss of accuracy and the training time is acceptable. In order to further evaluate the advantages of the scheme in this paper, a comparative experiment with other similar privacy training schemes is carried out. The experimental results show that the computation time of this mechanism on the experimental dataset is less than 5% of the comparison mechanism. At the same time, relying on the characteristics of decentralized training, the scheme in this paper has prospects in practical application scenarios.

**Keywords** consortium blockchain; credit data; support vector machine; privacy preserving; homomorphic encryption

## 1 引 言

征信是指采集、保存、整理和使用企业和个人的信用信息,并对其资信状况进行评价的专业活动<sup>[1-2]</sup>. 征信是国家经济健康运转和发展的保障,也是金融市场规避交易风险的重要手段. 于企业征信机构而言,客户信用状况的评估结果决定着商业决策的制定和商业活动的开展. 因而评价原始征信数据的信用状况是征信活动过程的关键步骤,建立一个准确率高、效果好的信用评价模型就尤为重要. 目前主流的征信机构一般使用两类信用评价模型训练方法: 统计方法和机器学习方法<sup>[3]</sup>. 统计方法包括逻辑回归和线性判别分析等,机器学习方法包括 SVM、神经网络等,其中, SVM 具有广泛的应用<sup>[4-8]</sup>. 通过

SVM 建立信用评价模型能够解决征信过程中部分环节效率不高、效果不好的问题,实现更加准确高效的信用评价方案.

训练信用评价模型所使用的训练集与其准确率紧密相关. 一方面,由于小型征信机构面临征信数据量少、质量不高等问题,需要得到其他机构的数据支持. 另一方面,大型征信机构之间存在用户群体差异,因而需要相互共享数据形成数据集上的互补,一定程度上完善和更正征信数据. 由此可见,理想的信用评价模型的训练离不开征信机构间的数据共享. 近年来网络技术尤其是移动网络技术的发展与完善为数据的采集、传输和共享提供了关键的技术支持<sup>[9-10]</sup>. 然而在征信领域,征信数据的高隐私、易非授权复制、高价值等特性,令彼此独立的征信机构之间形成的数据孤岛问题仍旧无法解决,数据的直接

共享面临许多挑战。目前,国内成立的百行征信试图连接这些孤岛,但是成员机构间数据共享的发展并不顺利,包括数据安全、机构内部的利益公平分配、征信数据格式等在内的诸多问题依然缺乏令人信服解决方案<sup>[11-12]</sup>。

本文提出面向征信数据安全共享的 SVM 训练机制,旨在解决数据共享及数据协同计算过程中的隐私保护问题,即在数据安全的条件下训练出机器学习模型。该机制引入区块链技术和安全多方计算技术,避免原始征信数据的直接共享,联合多方征信机构基于同态加密后的征信数据协同训练一个基于 SVM 的信用评价模型。

首先,区块链技术凭借其不可篡改、去中心化等优势技术特性,逐渐被应用到不同领域的数据安全共享方案中,例如智慧医疗、智能电网、智慧城市等。其中多数的方案采用复制型数据共享方式,即数据共享的结果是数据请求方得到数据提供方的原始数据。在征信行业也不乏区块链结合数据共享的研究,然而对于征信机构而言,在数据共享过程中,高隐私、易被非授权复制、高价值的原始征信数据不希望被其他机构直接获取。为了解决数据传输和存储过程中面临的隐私泄露问题,该机制要求征信数据在共享之前利用同态加密算法完成加密,且在存储和计算期间一直保持加密状态。除此之外,我们基于联盟链建立了征信数据共享平台,相关征信机构以该平台为核心完成数据共享和计算。相比于公有链,联盟链更小的开放程度和规模为方案提供了更多的隐私保护,区块链共享账本的不可篡改特性保证任何征信机构都无法对账本中的征信数据进行修改,进一步保护了在存储状态下征信数据<sup>[13]</sup>。

同态加密和差分隐私是两种主流的隐私保护算法,在应用到机器学习数据隐私保护的过程中在安全性或效率上面临挑战。差分隐私算法的隐私保护程度不高,且会给训练结果引入偏差。同态加密算法具备高隐私性和高可靠性,然而在运行效率方面存在明显劣势,同时多数基于同态加密的隐私保护方案需要引入可信第三方<sup>[14]</sup>。为了训练密态数据上的机器学习模型,本文基于部分同态加密算法的同态性质,实现了机器学习模型训练过程中必要的计算组件。不同于其他基于同态加密的机器学习隐私保护方案,在本文的方案中,训练过程无需第三方的参与,且训练时间开销低、模型准确率无损失。在具体使用过程中,当有机构向其他征信机构发起数据共享请求以扩充自己的数据集规模时,数据请求方会

协调多个数据提供方在加密的数据上进行基于同态加密的安全多方计算,完成基于密态征信数据的信用评价模型训练。在训练过程中,本文提出的机制能够保证不泄露各个参与机构原始数据隐私以及中间计算结果。也就意味着,在机制运行的过程中,模型训练者和其他参与计算的征信机构无法得知某一征信机构的原始征信数据<sup>[15]</sup>。

本文的主要贡献如下:

(1) 基于联盟链设计并实现安全可靠的数据共享平台,解决了征信机构面对数据共享时互不信任的问题。密态征信数据以及数据计算过程中的中间值记录到区块链上,既保护链上数据隐私安全,不可篡改,又维护了数据计算过程的透明性。

(2) 设计并实现一种隐私保护的 SVM 模型训练算法,借助同态加密技术,在无须第三方参与的情况下,在密态征信数据上训练 SVM 模型。

(3) 基于本文提出的机制,面向征信行业需求设计出征信数据信用评价模型训练系统。同时使用真实的数据集开展实验,训练并测试信用评价模型,验证机制的可用性和计算性能。结果显示在保证征信数据隐私的同时,相比于无保护的训练方案,系统在可接受时间内训练出的信用评价模型没有准确率损失。

本文第 2 节介绍国内外相关研究工作;第 3 节介绍背景知识,涉及到区块链和智能合约以及密码学基础知识;第 4 节对问题进行了定义;第 5 节是本文提出该机制下信用评价模型训练系统的详细介绍;第 6 节针对信用评价模型开展验证实验以及与同类方法的对比分析;第 7 节是结论。

## 2 相关工作

国内外许多领域将区块链应用到数据共享方案中,比如医疗行业<sup>[16-19]</sup>、能源行业<sup>[20-21]</sup>等。相比于各自领域传统的数据共享方案,区块链技术的应用一定程度上改善了数据共享过程中的信任问题和安全问题。

在征信领域,也不乏“区块链+征信”的研究<sup>[22-23]</sup>,塔琳等人建立了一个基于区块链的跨平台征信数据共享模型<sup>[23]</sup>。在此方案中,征信数据虽然在存储过程中经过加密,但是在共享给其他机构时,原始的征信数据还会直接暴露给对方,考虑到征信数据的易复制性,该方案对数据隐私的保护有待进一步加强。但是仅将加密的数据共享出来,数据的可

用性将无法得到保障。

为了同时满足数据隐私安全性和可用性, 相关研究也有开展. 董祥千等人<sup>[24]</sup>借助安全多方计算和差分隐私技术保障数据所有者计算和输出隐私. 通过安全多方计算, 数据在受到保护的条件下依旧能够计算出有效的输出结果. Yue 等人<sup>[25]</sup>使用区块链用于存储个人的医疗数据信息, 同时借助安全多方计算在没有可信第三方的情况下完成对加密医疗数据的计算工作, 这些研究证明了区块链与安全多方计算结合的合理性和有效性.

机器学习模型训练过程中的隐私保护问题一直是研究重点. 差分隐私和同态加密是最常用的两种隐私保护方案. 差分隐私技术在原始数据上添加精心计算的扰动, 由此来保证公开数据的隐私. Abadi 等人<sup>[26]</sup>将差分隐私用于深度学习中对数据集中敏感信息的保护. 虽然差分隐私是一种高效的隐私保护方法, 但是扰动的引入将会给最终模型准确率带来损失. 由于全同态加密技术计算复杂, 部分同态加密技术的应用更为广泛. 然而部分同态加密技术支持的计算操作有限, 无法保证模型训练过程中出现的各种计算, 因此在多数使用部分同态加密的隐私保护方案中, 需可信第三方的协助. Gonzalez-Serrano 等人<sup>[27]</sup>基于部分同态加密技术, 解决多个数据提供方训练 SVM 模型时的隐私保护问题. 在该方案中, 除了数据提供方和模型训练方(应用方), 又引入一个认证方协助进行密钥分发, 同时完成一些部分同态加密无法实现的操作. 第三方的引入带来的安全问题以及部分同态加密在计算上的局限性给基于同态加密的隐私保护方案带来极大的挑战. 在数据的安全分类领域, 也有相关的研究. De Cock 等人<sup>[28]</sup>和 Bost 等人<sup>[29]</sup>提出安全的算法用于在隐私保护的场景中完成机器学习分类模型的测试.

## 3 背景知识

### 3.1 区块链及智能合约

#### 3.1.1 区块链

Peer-to-Peer(P2P)技术凭借其去中心化、可扩展性高等优势被广泛应用于文件共享等领域<sup>[30]</sup>. 区块链本质上是一个运行在点对点网络上的公开账本<sup>[31]</sup>, 具备去中心化、安全性、不可篡改的特点, 常用于解决互不信任实体之间的信任问题<sup>[32-33]</sup>.

按照区块链的去中心化程度从高到低或是参与方从多到少, 区块链可以被分为三类: 公有链、联盟

链和私有链<sup>[34-35]</sup>.

(1) 公有链是完全公开的区块链, 用户可以随时加入区块链网络, 查看区块链上的数据. 公有链最典型的应用是比特币和以太坊.

(2) 联盟链的开放程度小于公有链, 只有经过认证的联盟成员才能够加入区块链网络, 并且设有访问控制来限制成员访问区块的权限. 超级账本以及 FISCO BCOS 平台都属于联盟链.

(3) 私有链一般应用在单个企业、机构、组织内部, 去中心化程度最低, 具备很高级别的访问控制和权限设置能力.

从区块链的分类情况不难发现, 联盟链的特征更适用于征信机构之间的关系, 更能够满足去中心化和数据安全上的需求. 因而, 本文提出的征信数据安全协同计算机制运行在联盟链上.

#### 3.1.2 智能合约

智能合约是运行在区块链上的一组代码, 由区块链上的节点调用, 调用后按照编写之初的规则自动执行, 实现特定功能. 智能合约的出现令区块链的应用场景从数字货币领域扩展到其他领域. 通常, 智能合约被认为能够一定程度上替代第三方, 从而简化中心化系统在运行过程中的繁琐环节, 提高系统运行效率, 减少运行成本. 最重要的是, 智能合约一经发布, 就不被人所干预, 提高了安全性.

### 3.2 同态加密

加密系统包括三个算法: 密钥生成算法( $Gen$ )、加密算法( $Enc$ )和解密算法( $Dec$ ). 在公钥加密系统中使用一对密钥对( $PK, SK$ ), 其中公钥  $PK$  用于加密, 私钥  $SK$  用于解密. 某些加密系统拥有同态的性质被称为同态加密. 同态性质是指加密算法可以将密文上的操作映射到相应的明文而不需要对应的密文信息. 正式的同态加密的定义在定义 1 中给出.

**定义 1**(同态)<sup>[36]</sup>. 一个公钥加密算法( $Gen, Enc, Dec$ )是同态加密当且仅当对于所有的密钥对( $PK, SK$ ), 该加密算法依赖  $PK$  定义了两个群  $M, C$  有如下操作:

(1) 明文空间是  $M$ , 所有的密文通过  $Enc_{pk}$  产生, 并且在群  $C$  中.

(2) 对于任意信息  $m_1, m_2 \in M, c_1$  和  $c_2$  是算法  $Enc_{pk}$  输出的相应的密文有如下操作:  $Dec_{sk}(o(c_1, c_2)) = \sigma(m_1, m_1)$ .

Paillier: 在我们的方案中, 我们应用加法同态算法 Paillier<sup>[37]</sup>. Paillier 是一种具备同态性质的公钥加密算法, 它基于复合剩余类的困难问题的公钥

加密系统. 经过 Paillier 加密后的数据成为密文, 除了用对应的私钥解密, 密文被认为是安全的. 假设  $p$  和  $q$  是两个  $n$  比特长的大素数,  $N = pq$ . 则在 Paillier 中公钥是  $N$ , 私钥是  $(N, \varnothing(N))$ . 加密算法表达为  $c := [(1+N)^m r^N \bmod N^2]$ , 其中  $m \in \mathbb{Z}_N$ . 解密算法表达为

$$m := \left\lceil \frac{[c^{\varnothing(N)} \bmod N^2 - 1]}{N} \cdot \varnothing(N)^{-1} \bmod N \right\rceil.$$

Paillier 满足加法同态和数乘同态. 具体来说, Paillier 的加法同态可表示为  $Dec_{sk}(c_1 \times c_2) = m_1 + m_2$ ; Paillier 的数乘同态可表示为  $Dec_{sk}(c_1^k) = k \times m_1$ .

## 4 问题定义

本文提出面向征信数据安全共享的 SVM 训练机制. 该机制利用多方共享在联盟链上的加密征信数据, 协同多个数据提供方共同完成对于密态征信数据的计算, 以满足某种征信活动, 比如训练基于 SVM 的信用评价模型. 本节围绕系统模型、系统面临的潜在威胁以及安全目标三方面, 进行了问题定义.

### 4.1 系统模型

围绕面向征信数据安全共享的 SVM 训练机制设计系统. 如图 1 所示, 系统由三类实体组成: 征信机构、联盟链平台、数据请求方. 这三个实体以联盟链为核心实现数据的安全计算.

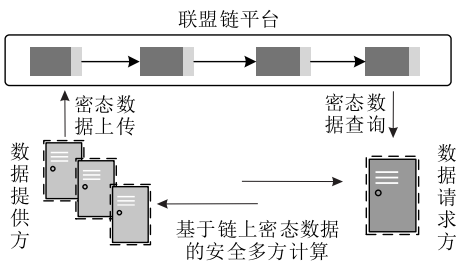


图 1 系统实体关系图

(1) 征信机构(数据提供方). 征信数据由征信机构提供, 征信机构主要包括政府征信机构和企业征信机构. 在系统中, 征信机构将以区块链节点的身份加入到联盟链中.

(2) 联盟链平台. 相比于公有链和私有链, 联盟链更适合征信场景. 数据提供方的征信数据在经过预处理后, 将被以交易的形式存储在共享账本中.

(3) 数据请求方(模型训练方). 在该机制中, 数据请求方同时是模型训练方. 数据请求方从联盟链

上获取到密态数据后, 为了满足征信请求, 需要联合数据提供方共同完成对链上密态数据的计算. 数据请求方可以被任意一个征信机构运行的区块链节点所替代.

### 4.2 潜在威胁

系统的主要威胁来自三个方面:

(1) 网络窃听者. 攻击者通过网络窃听或者数据拦截等方式非法获取到传输过程中的数据, 来达到某种目的.

(2) 半诚实参与者. 是指本文系统中的征信机构, 在整个模型的训练过程中, 数据提供方和模型训练方均可能是半诚实参与者. 这些参与者想要得到最终结果, 所以会按照约定正确执行每一步. 同时基于某种目的, 又想窃取到其他参与者的数据, 通过保留计算过程中的中间数据, 对其他参与者的数据进行推断并窥探.

(3) 恶意攻击者. 这一类攻击者可以通过任意方法恶意对方案中的数据进行非法盗取. 这种攻击行为不在本文的考虑范围之内.

由于参与共享的数据提供方以及模型训练方在加入联盟链之前需要进行认证, 因而本系统更加关注数据共享过程中联盟链成员对系统的威胁, 所以本系统重点防范的是第一类和第二类威胁, 暂不考虑第三种威胁.

### 4.3 安全目标

本系统有三个安全目标: 一是在系统各实体通信期间, 网络窃听者无法从通信流量中窃取到通信双方或多方的征信数据隐私; 二是保证任何征信机构都无法从密态数据在计算期间产生的中间结果和最后结果中推测出原始的征信数据; 三是各机构共享的征信数据在存储期间不被恶意窃取及篡改.

为了实现第一个安全目标, 在传输过程中, 经过加密之后的征信数据能有效防止网络窃听者的攻击. 本系统的数据加密将由征信机构在数据共享之前进行.

为了实现第二个目标, 系统在数据共享时, 使用基于安全多方计算的支持向量机算法来保证隐私.

为了实现第三个安全目标, 令共享的征信数据能可靠存储, 系统借助区块链平台以去中心化的存储方式, 保证加密数据在机构间公开同时又不可篡改. 由于存储的数据都经过加密, 因此数据隐私性也得到保证.

综上所述, 文中介绍的系统能够有效保护各个征信机构的数据隐私安全.

## 5 信用评价模型训练系统

为了验证本文提出机制的可用性以及性能, 本节建立了基于该机制的信用评价模型训练系统. 在该机制的安全多方计算部分实现了安全的 SVM 方法, 用于训练信用评价模型. 本节从系统工作流程、功能模块构成、各个模块具体介绍、关键算法四个方面详细展示了信用评价模型训练系统.

### 5.1 系统工作流程

如图 2 所示, 系统运行的流程如下:

(1) 系统基于联盟链开发, 每一个征信机构都需要以区块链节点的形式加入联盟链.

(2) 各个征信机构共享的数据应具备统一的数据格式, 因此征信数据在写入区块链共享账本之前需要进行预处理, 修改格式不规范的数据, 丢弃内容不完整的数据, 对格式标准的数据进行加密.

(3) 征信机构通过各自运行的联盟链节点调用联盟链上智能合约, 将预处理后的征信数据以交易的形式记录到区块链中.

(4) 记录在公开的共享账本上的征信数据可以被联盟链内节点查看. 通过同步包含交易信息的区块到节点本地, 任意机构都可以获取其他机构共享的密态征信数据. 除此之外, 各机构也可以调用智能合约查询链上的数据.

(5) 模型训练方利用共享来的密态征信数据, 使用基于安全多方计算的 SVM 方法, 与数据提供方协同训练模型.

(6) 最终训练完成, 模型训练方得到模型.

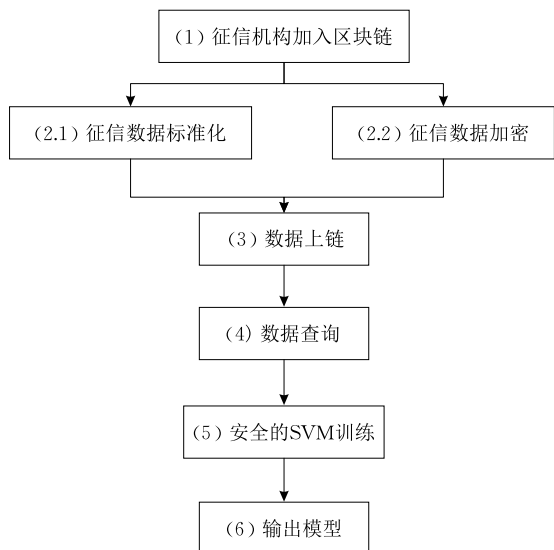


图 2 系统运行流程图

### 5.2 功能模块介绍

#### 5.2.1 功能模块组成

如图 3 所示, 系统主要由三个功能模块组成: 征信数据预处理模块、安全多方计算模块、联盟链模块. 其中, 数据预处理模块运行在链下, 数据格式规范化和加密工作在征信机构本地进行. 考虑到系统的运行效率, 需要计算资源和频繁交互的安全多方计算模块同样运行在链下. 联盟链模块运行在链上, 链上运行的智能合约实现上传和查询密态征信数据的功能. 通过调用链上的智能合约, 能够实现链上链下模块的交互.

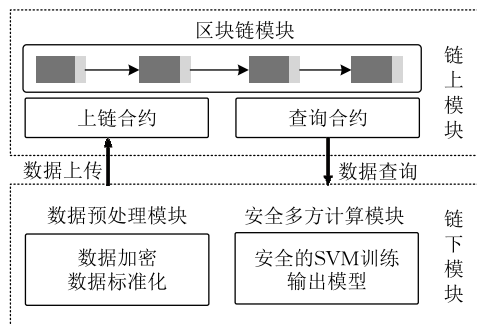


图 3 系统功能模块图

#### 5.2.2 数据预处理模块

数据预处理模块需要具备两个功能: 统一征信数据格式、数据加密.

多个征信机构之间共享的征信数据需要定义统一的数据格式, 不同机构持有的征信数据格式存在差异, 除了需要统一数据维度的含义、数量、次序之外, 在存储上也需要使用统一的类型、长度, 消除格式上的差异将方便信用数据的计算. 另一方面, 综合各方机构的征信数据内容, 可以得出一个全面的数据格式. 对于维度缺失、缺少标签等存在问题的征信数据将直接丢弃, 不参与数据共享.

为了保证在共享的过程中, 原始征信数据不被其他机构直接复制, 不被网络窃听器窃取有效信息, 征信将在本地进行加密. 本文的方案中选择 Paillier 作为加密征信数据的算法.

#### 5.2.3 征信数据计算模块

该模块的工作主要由信用评价模型训练方负责. 该模块解决的问题是如何协调征信各方对加密征信数据开展计算, 实现安全的 SVM 方法, 训练出信用评价模型. 在安全的 SVM 算法中, 模型的训练涉及到数据的加解密计算以及同态计算, 因此 SVM 的训练在链下完成, 智能合约不参与复杂的模型计算过程. 同时, 由于安全的 SVM 算法基于安全多方

计算,模型在训练过程中需要数据提供方和模型训练方的共同参与.在模型训练结束之后,模型训练方利用数据提供方的密态数据训练得到高质量的模型,同时对于数据提供方而言,数据在共享的过程中隐私安全得到充分的保障.该部分使用到的关于数据计算的算法在 5.3 节有详细的介绍.

#### 5.2.4 联盟链模块

征信机构加密后的征信数据上传到区块链上,数据请求方从区块链上获取加密的征信数据.

如图 4 所示,区块链共享账本由一个个按照时间顺序连接在一起的区块组成,每一个区块中包含若干个交易,征信数据按照固定格式存储在交易中.数据的上链和查询由分别由上链合约 *uploadData* 和查询合约 *requireData* 完成.两个合约编写完成之后,部署在区块链上,供征信机构调用.

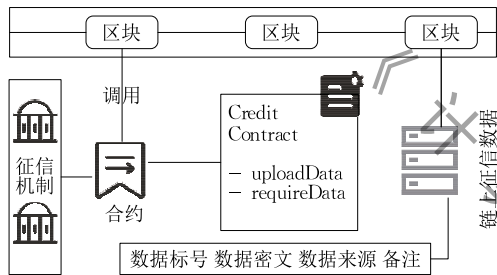


图 4 联盟链模块运行示意图

#### 过程 1. 上链合约.

输入:密态数据 *encData*; 备注 *dataRemarks*

1.  $data.ownerAddress \leftarrow$  从 *msg.sender* 获得合约调用者地址,即数据来源
2.  $data.encData \leftarrow$  从输入中获得数据密文
3.  $data.serialNum \leftarrow$  从合约中自动生成数据标号
4.  $data.remarks \leftarrow$  从输入中获得数据备注
5. 触发事件 *dataUploaded*

输出:链上征信数据 *data*

图 4 中的链上征信数据包括四个字段:数据标号、数据密文、数据来源、备注.每一条数据在区块链中都有唯一的标识,数据密文、备注来自于征信机构调用合约时的输入,数据来源是征信机构所运行节点调用合约时使用的地址,该地址与征信机构一一对应.在上链合约中,建立结构体存储链上征信数据.如过程 1 所示,上链合约的主要工作是对结构体的四个成员变量进行赋值,并设置事件将数据标号返回给上传数据的征信机构.

### 5.3 关键算法

#### 5.3.1 SVM<sup>[38]</sup>中的安全操作

SVM 的模型是一个划分超平面  $y = \mathbf{w}^T x + b$ ,  $(x_i, y_i) \in D$ , 有:  $\mathbf{w}^T x_i + b \geq 1, y_i = +1, \mathbf{w}^T x_i + b \leq -1,$

$y_i = -1$ . SVM 的基本型是:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (1)$$

$$\text{s. t. } y_i (\mathbf{w}^T x_i + b) \geq 1, i = 1, 2, \dots, m$$

系统采用的计算方案使用线性核函数,选择随机梯度下降作为求解式(1)的优化算法,基于随机梯度下降的 SVM 优化算法简单且高效. SVM 利用梯度下降求最优解时的目标函数如式(2)所示.令  $\beta = (\mathbf{w}, b)$ ,  $\nabla_{t+1} = \lambda \mathbf{w}_t - c \sum_{i=1}^m \max\{0, 1 - (\mathbf{w}_t x_i + b_t)\}$ . 其中  $\lambda$  表示学习率,由算法执行人员设定.  $c$  是误分类的惩罚项,通常取值为  $\frac{1}{m}$ .

$$\beta^{t+1} = \beta^t - \lambda \nabla_{t+1} \quad (2)$$

式(2)包括一些基本操作,如多项式乘法和正整数比较.为了安全地训练 SVM 模型,在这一节中,我们介绍如何利用 Paillier 的加法和数乘同态设计在 SVM 训练算法中的安全操作.

#### (1) 安全多项式乘法

利用 Paillier 的同态性质,我们可以轻松地得到安全加法和安全减法.基于 Paillier 的安全加法表示为:  $[m_1 + m_2] = [m_1] \times [m_2] \pmod{N^2}$ . 基于 Paillier 的安全减法表示为:  $[m_1 - m_2] = [m_1] \times [(m_2)^{-1}] \pmod{N^2}$ . 利用安全加法和安全减法我们可以很自然地得到安全多项式乘法:  $[am_1 + bm_2] = [m_1]^a \times [m_2]^b \pmod{N^2}$ .

通过 Paillier 构造的安全加法、安全减法和安全多项式乘法的安全性依赖于 Paillier 的计算不可区分.

#### (2) 安全比较

在我们方案中的安全比较是一个密文  $[m]$  与常数 1 的比较(求解  $\max\{0, 1 - (\mathbf{w}_t x_i + b_t)\}$ ). 对于参与协议的模型训练方 C 和数据提供方  $P_i$ , 输入所隐含的信息,任何一方都不能得到任何其它信息.我们的安全比较协议如算法 1 所示.

#### 算法 1. 安全比较算法.

$P_i$  输入:  $(SK_{P_i}, PK_{P_i})$

C 输入:  $[a]_c$

1. C 随机选择正整数  $r_1, r_2$  和  $r_3$ , 其中  $|r_3 - r_2| < r_1$
2. C 通过安全多项式乘法计算  $[ar_1 + r_2]$  和  $[r_1 + r_3]$  后发送给  $P_i$
3.  $P_i$  解密并比较  $(ar_1 + r_2)$  和  $(r_1 + r_3)$  的大小,得到比较结果
4. 当且仅当  $(ar_1 + r_2) > (r_1 + r_3), a > 1$ ; 否则  $a \leq 1$
5. RETURN ( $a < 1$ ) 给 C

**正确性分析:**  $|r_3 - r_2| < r_1 \leftrightarrow \frac{|r_3 - r_2|}{r_1} < 1$ .

$(ar_1 + r_2) = (r_1 + r_3) \leftrightarrow (a-1) = \frac{r_3 - r_2}{r_1}$ . 因为  $a$  是整数, 如果  $(ar_1 + r_2) > (r_1 + r_3)$ , 我们可以得到  $(a-1) \geq 1 \rightarrow a \geq 1$ , 否则  $a < 1$ .

**安全性分析:** 当面临半诚实敌手时, 我们遵循两个常用的定义: 安全的两方计算<sup>[39]</sup> 和模块化顺序组合<sup>[40]</sup>, 满足安全两方计算的协议可视为安全. 首先证明安全比较算法的安全性.

在安全比较算法中存在数据提供方  $P_i$  和模型训练方  $C$  两个实体, 它们共同计算函数  $F: F([a]_{P_i}, 1, PK_{P_i}, SK_{P_i}) = (\phi, (a < 1))$ . 实体  $C$  能够掌握的信息包括经过  $P_i$  加密的密态数据  $[a]$ 、 $P_i$  的公钥和自己生成的三个随机数. 由于  $C$  不掌握  $P_i$  的私钥, 除了上述信息之外, 不再得知其他的信息. 实体  $P_i$  掌握的信息包括  $(ar_1 + r_2)$ 、 $(r_1 + r_3)$ 、自己的公钥和私钥.  $P_i$  的模拟器运行生成理想状态下能够掌握的信息. 最终, 能够发现  $P_i$  实际掌握的信息和模拟器生成的理想的信息在分布上相同, 因此在统计上不可区分, 满足安全两方计算的统计安全性要求.

### 5.3.2 安全的 SVM 训练算法

我们的安全目标是在任何情况下都必须保证每个数据提供者提供的训练数据集的安全. 假设有  $n$  个数据提供方  $P_i$ , 和一个模型训练方  $C$ . 算法 2 详细介绍了我们的安全 SVM 训练算法, 其中  $\{XY_i^{Enc}, X_i^{Enc}, Y_i^{Enc}\} \in D_i^{Enc}$ .

#### 算法 2. 安全 SVM 训练算法.

每个  $P_i$  输入:  $D_i^{Enc}$ , 公钥  $PK_i$

$C$  输出: 支持向量机模型

1.  $C$  初始化模型参数  $\omega$
2. WHILE 未达预设精度 DO
3. FOR  $i=1$  TO  $n$  DO
4.  $P_i$  发送  $[D_i]$  给  $C$
5.  $C$  从所有  $[D_i]$  中随机选择  $m$  个  $x_j$  对梯度进行更新
6. FOR  $j=1$  TO  $m$  DO
7.  $C$  通过安全多项式乘法计算  $[y_j(\omega x_j - b)]$
8.  $C$  通过安全比较算法比较  $[y_j(\omega x_j - b)]$  和 1 的大小
9.  $C$  通过安全多项式乘法更新梯度  $[\nabla_{i+1}]$
10.  $C$  随机选择随机数  $r$ , 并通过安全多项式乘法计算  $[\nabla_{i+1} + r]$
11.  $P_i$  解密后返回  $(\nabla_{i+1} + r)$  给  $C$
12.  $C$  减去随机数得到  $\nabla_{i+1}$ , 更新模型参数  $\omega$
13. RETURN 支持向量机模型参数  $\omega$  给  $C$

**安全性分析:** 安全多项式乘法和安全比较是以

模块化的方式设计的, 因此执行算法 2 调用相应的模块即可. 所使用的模块算法是安全的, 算法 2 的安全性由模块化组合定理实现<sup>[34]</sup>.

算法 2 中参与的实体包括多个数据提供者和一个模型训练者. 每个数据提供者同模型训练者执行相同的算法操作, 因此只需证明一个数据提供者和模型训练者之间的算法操作满足安全性要求即可. 该算法计算的函数  $F$  为  $F(D_{P_i}, PK_i) = (\phi, (\omega, b))$ .

数据提供方能够掌握的信息包括明文数据集、每一轮带扰动的参数、自己的公私钥. 由于模型的扰动值随机生成, 因此模型的训练方无法提取出原始参数信息.

模型训练方掌握的信息包括密态数据、密态的计算中间值、安全比较结果、明文模型参数, 以及数据提供方的公钥. 由于模型训练方不知道数据提供方的私钥, 因而无法使用除暴力破解之外的任何方式, 从安全两方比较结果以及密态计算中间值中推测出任何关于明文数据集的信息.

## 6 实验验证及对比分析

区块链、加密算法以及安全多方计算的结合使用充分保证了信用评价系统的安全性. 使用 Solidity 编写智能合约, 实现了数据上链、查询等功能, 并将合约部署在联盟链 FISCO BCOS 平台上. 每一个征信机构对应一个联盟链节点, 节点可通过控制台调用智能合约.

在此基础上, 本节针对系统的数据计算模块开展了实验, 分析安全的 SVM 方法(算法 2)是否可行, 从而判断在隐私安全得到保障下信用评价系统的可用性和性能. 我们使用 Java 分别实现了常规的非隐私保护的 SVM 方法以及算法 2. 验证实验分成三部分: 第一部分将使用非隐私保护的 SVM 方法在未加密的训练集上训练信用模型; 第二部分, 训练集经加密, 通过本文系统使用的安全多方计算下安全的 SVM 方法训练信用评价模型. 实验设置了多个评价指标, 用于对比两部分实验的结果从而分析安全的 SVM 方法的性能. 第三部分在第二部分实验的基础之上, 针对数据提供方(安全多方计算的参与方)对系统扩展性的影响, 进行实验验证.

为了证明本文提出的安全的 SVM 方法在训练时间上的优势, 设置实验与文献[27]提出的基于部分同态加密的安全 PEGASOS 方案进行对比分析. 本文的方案同该方案均对基于梯度下降的 SVM 方



法进行了隐私保护,具备可比性。

## 6.1 实验准备

### 6.1.1 实验环境

实验代码在一台 PC 上运行,在该 PC 上模拟参与安全多方计算的征信机构,包括数据提供方和模型训练方. PC 的配置为:4-core Intel i7 (i7-3770 64 bit) processor at 3.40GHz and 8GB RAM.

### 6.1.2 训练集

实验主要使用三个不同的源自真实情况的数据集,前两个是征信领域常见的数据集,分别为 Australian Credit Approval 和 German Credit Data,第三个数据集用于验证当数据集规模增大时,模型训练时长的变化,因此选用样本数量较大的数据集 Electrical Grid Stability Simulated Data Data Set (EGSSD Data Set). 三个数据集均公开在 UCI 机器学习数据集资源网站上. 每一个数据集将按比例分成两个部分,一部分作为训练集训练模型,另一部分作为测试集测试模型的性能. 表 1 是训练集的基本情况.

表 1 训练集基本信息

训练集名称	案例数	属性数
Australian Credit Approval	690	14
German Credit Data	999	24
EGSSD Data Set	10 000	14

### 6.1.3 加密设置

浮点数处理:Paillier 的操作位于整数空间,但模型的实际训练过程中有大量的浮点数运算. 因此,为了处理加密的实值数据,在训练之前,数据集的实值通过格式转换成整数的表示形式. 根据国际标准 IEEE 754,任意一个二进制浮点数  $D$  可以表示为  $D = (-1)^s \times M \times 2^E$  (1),  $s$  表示符号位,  $M$  表示有效数字,  $E$  表示指数位. 另一方面,将明文映射到密文上进行计算需保证计算结果仍处于该系统的明文空间内,所以我们需要对加密系统的密钥长度进行设计,避免溢出的可能. 经过对可能的中间结果进行分析,我们将 Paillier 的密钥  $N$  设置为 2048 bits.

### 6.1.4 训练参数设置

训练模型时所使用的参数包括:训练的最大迭代次数为 1500 次,学习率为 0.05,阈值设置为 0.00001,每轮迭代中从训练数据集里选取的样例数为 15.

## 6.2 实验结果分析

### 6.2.1 评价指标

对于实验结果,我们将从时间成本、准确率两个

指标对实验结果进行性能分析:时间成本指的是模型训练方训练模型所花费的时间. 准确率表示训练出的信用评价模型在测试集上正确分类的结果所占所有结果的比例.

### 6.2.2 时间成本

多次实验后,统计得到非隐私保护的 SVM 方法和安全的 SVM 方法在两个训练集下训练模型的耗时. 在非隐私保护的 SVM 方法下,使用 Australia credit approval 训练集训练模型的时间均在 200 ms 左右,使用 German credit data 训练集训练模型的时间也在 200 ms 上下浮动. 结果表明非隐私保护的 SVM 方法能够在短时间内训练出模型. 在安全的 SVM 方法下,由于训练模型过程中涉及到加密数据(大整数)的运算,相比于非隐私保护 SVM 方法,模型训练时间大幅增加. 使用 Australia credit approval 训练集训练模型的时间均在 1 500 000 ms 左右,使用 German credit data 训练集训练模型的时间在 2 000 000 ms 上下浮动.

由于使用安全的 SVM 方法涉及到多个参与方共同计算,为了验证参与方数量对模型训练耗时的影响,因此另外开展实验测试在参与方数量分别为 1、2、3、4、5 时,在相同训练集上训练模型所用的时间成本,实验结果如图 5 所示. 在安全的 SVM 方法下,倘若不同数量的参与方最终组合得到的数据集相同,参与方数量的增加不会给计算耗时带来显著的改变.

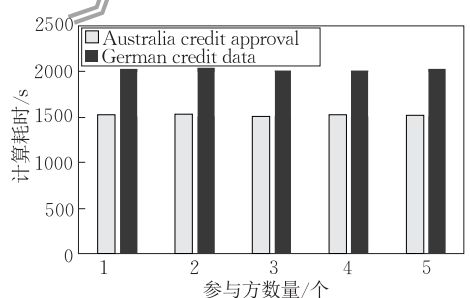


图 5 时间成本统计图

对于安全的 SVM 算法而言,模型的训练时间主要取决于模型训练的迭代次数. 因此,当数据集规模增加时,如果保证模型分类准确率不明显下降的同时,模型训练的迭代次数在小幅度增加或者不变,那么模型的训练时间将不会有较大幅度的变化. 为了测试数据集规模对安全的 SVM 方法训练时间的影响,我们分别从数据集 EGSSD 挑选 1000, 2000, 3000, 4000 条样本作为实验数据集进行实验. 实验结果如表 2 所示.

表 2 准确率统计表

数据集规模	训练时间/ms	准确率/%
1000	619372	80.1
2000	573736	79.5
3000	603712	78.8
4000	622092	78.6

随着训练集规模的增加,在 2000 次迭代的训练条件下,所获取的模型准确率不存在较大范围的波动.这得益于基于随机梯度下降方法的 SVM 面对大规模数据集时所具备的优势.

### 6.2.3 准确率

实验将数据集中大部分的数据用于训练模型,剩下的小部分作为测试集测试模型的准确率.

为了对比常规条件下的 SVM 方法和安全的 SVM 方法的分类准确率,在两个数据集上分别进行了 30 组实验,统计两种方法下模型分类的平均准确率,此时统计得到的平均准确率已趋向稳定.实验结果统计如表 3 所示,相比于常规条件下基于明文数据集的 SVM 方法,使用安全的 SVM 方法训练出的评价模型基本不会造成准确率损失.

表 3 准确率统计表

训练集名称	安全的 SVM/%	SVM/%
Australia credit approval	80.3	80.7
German credit data	75.6	76.2

同样的,由于使用安全的 SVM 方法涉及到多个参与方共同计算,为了验证参与方数量对模型准确率的影响,因此测试在参与方的数量分别为 1、2、3、4、5 时训练模型的准确率,在安全的 SVM 方法下,倘若不同数量的参与方最终组合得到的数据集相同,参与方数量的增加不会给准确率带来显著的改变.实验结果如图 6 所示.

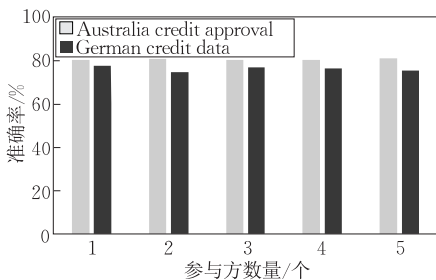


图 6 准确率统计图

### 6.3 同类方法对比分析

以上实验表明本文提出的方案能够在可接受时间内训练出准确率无损失的 SVM 模型,同时当计算参与方的数量增加时,方案具备扩展性.为了进一步分析方案在计算效率上的优势,我们同文献[27]

中安全的基于 PEGASOS 的 SVM 方案 (SPSVM) 进行对比.在引入隐私保护方案后,SVM 训练时间主要由许多耗时计算操作和参与方之间的交互产生.因此,通过同类的耗时计算操作数量和交互次数的统计,可以从理论分析上对比得出两种方案在训练时间上的优劣.除此之外,从实际的训练过程中,也对两种方法的耗时进行了实验统计.

首先,对两个方案的同类耗时计算操作的数量进行统计.在本文提出的安全的 SVM 方案中,训练过程中的耗时计算操作主要包含两种类型:安全多项式乘法和安全两方比较.在文献[27]提出的 SPSVM 中,耗时计算操作主要包括密态向量的欧拉距离计算、密态向量的内积计算,密态数据的比较以及安全的位移动计算等.为了方便统计对比两种方法所需的耗时操作,我们对两种方法涉及到的各种计算操作进一步划分.由于两个方案均基于部分同态加密,因此两个方案中涉及的计算操作均由下面四种基本操作组成,包括:加密、解密、同态数乘、同态加法.

统计结果如表 4 所示,由于 SPSVM 方法在一轮迭代中仅挑选一条数据,因此本文的安全 SVM 方法在一轮迭代中,设置挑选的样本数量为 1.其中  $d$  代表数据集的维度, $l$  代表 SPSVM 在训练过程中需要完成比较的整数其二进制格式的位数.

表 4 耗时计算对比表

	安全的 SVM	SPSVM <sup>[27]</sup>
同态数乘	$2d$	$3m+14l-7$
同态加法	$3d$	$7(m+6l)-25$
加密	2	$3(m+7l-4)$
解密	$d$	$2(m+7l-4)$

通过对比发现,在一轮迭代过程中,本文提出的安全的 SVM 方案中四种基本耗时操作数量以及参与方之间的交互次数远小于文献[27]中的 SPSVM 方案.在 SPSVM 方案的运行过程中,存在大量两个密态数据相乘的操作,而部分同态加密方案无法直接满足此类的同态乘法,因此训练方需要认证方的协助通过同态加法和同态数乘构造同态乘法,此过程导致了训练过程中大量的耗时操作和频繁的交互.反观本文提出的方案,我们构造的基于部分同态加密方案的安全多项式乘法和安全比较算法能够满足训练过程中所有的操作需求,无需通过复杂的计算实现同态乘法,因而在耗时操作的数量上和交互次数上保持很低.耗时操作数量越少代表模型训练时间越短,由此分析结果能够表明在训练时间开销方面,我们提出的方案具备明显的优势.

本文实验验证了模型训练的时间开销,并对核矩阵的计算和训练过程中的耗时分开统计.一方面,由于本文提出的安全的 SVM 方法采用线性核函数,无需计算核矩阵,而 SPSVM 方法采用高斯核函数,需要提前计算好核矩阵,因此我们按照文献[27]中的描述,将矩阵预先计算好,仅统计模型训练时间.此时,实验数据集采用完整的 Australian Credit Approval 数据集.在迭代次数设置为 100 次的条件下,SPSVM 方法的训练时间与本文提出的安全的 SVM 训练时间对比结果如表 5 所示.

表 5 耗时计算对比表

安全的 SVM/ms	SPSVM <sup>[27]</sup> /ms
45 863	1 054 738

实验结果显示,在假设核矩阵计算好的条件下,SPSVM 方法模型的训练时间显著高于我们提出的安全的 SVM 方法,本文方法的训练耗时不到 SPSVM 方法耗时的 5%.

另一方面,为了进一步证明本文提出的方法相较于 SPSVM 具备优势,开展实验统计 SPSVM 计算核矩阵的耗时.由于核矩阵的计算过程中涉及的欧式距离计算属于耗时操作,我们仅在小规模数据集进行实验.其中数据集中的样本取自 Australian Credit Approval,实验样本数量分别为 10, 20, ..., 90, 100. 实验结果如图 7 所示.

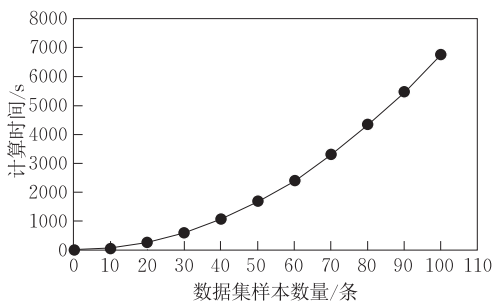


图 7 SPSVM 方法核矩阵计算耗时

由于矩阵中每一个元素的计算都涉及两个密态向量欧式距离的计算,因此计算耗时随着数据集规模的增长迅速增长.显然,在实际的应用场景中,此方法的可用性低.

除此之外,SPSVM 在训练过程中需引入第三方保证算法的顺利执行,因而存在潜在的安全性问题.相比之下,我们的方案完全去中心化,运行过程中只存在数据提供方和模型训练方两种角色,因此能够避免引入第三方带来的问题,更适用现实应用场景.

## 7 结 论

区块链技术的产生给安全的数据共享提供了许多新的解决思路,安全多方计算能够保证在各参与方不泄露数据隐私的条件下共同完成对数据的计算.基于两种技术的优势特征,本文提出了一种面向征信数据安全共享的 SVM 训练机制,满足征信机构之间的数据共享对安全性的要求.并在数据隐私不泄露的前提下,面向信用评价、金融分析、数据建模等征信应用提供了安全的数据共享和计算方案,一定程度上能够推进征信体系的完善.本文设计并实现了基于安全多方计算的安全 SVM 方法,在不引入可信第三方的条件下,模型训练者协同数据提供方训练信用评价模型,并保证训练过程中的数据隐私安全.文中从威胁模型、安全目标、算法流程等角度对机制的安全性进行了全面的分析.同时对系统的架构设计和功能模块进行了详细的阐述.最后展开实验验证信用评价模型训练方案的可用性,以此验证系统的可用性和性能.实验结果表明在安全性得到充分保证的情况下,密态数据依旧可用,同时模型的准确率没有受到影响.该机制的提出和实现为征信机构间数据的安全共享提供具备实用价值的解决思路.在未来工作中,一方面,我们将继续对方案的计算开销和通信开销进行优化,进一步降低模型训练所需的时间.另一方面,基于该方案的机器学习方法将从 SVM 扩展到其他机器学习算法,满足不同行业在数据共享时对隐私保护的需求.

致 谢 感谢向本文提出宝贵建议的审稿专家!

## 参 考 文 献

- [1] Li Jun-Li. A Study on the Construction of Individual Credit Investigation System and Application in China. Beijing: China Social Sciences Publishing House, 2010 (in Chinese)  
(李俊丽. 中国个人征信体系的构建与应用研究. 北京: 中国社会科学出版社, 2010)
- [2] Zhang Ying. Research on the Construction of China's Credit Information System. Beijing: Capital University of Economics and Business, 2005 (in Chinese)  
(张颖. 构建中国征信体系研究. 北京: 首都经济贸易大学, 2005)
- [3] Shi J, Zhang S, Qiu L. Credit scoring by feature-weighted support vector machines. Journal of Zhejiang University Science C: Computer & Electronics, 2013, 14(3): 197-204

- [4] Bellotti T, Crook J. Support vector machines for credit scoring and discovery of significant features. *Expert Systems with Applications*, 2009, 36(2): 3302-3308
- [5] Harris T. Credit scoring using the clustered support vector machine. *Expert Systems with Applications*, 2015, 42(2): 741-750
- [6] Maldonado S, Pérez J, Bravo C. Cost-based feature selection for support vector machines: An application in credit scoring. *European Journal of Operational Research*, 2017, 261(2): 656-665
- [7] Wang Y, Wang S, Lai K K. A new fuzzy support vector machine to evaluate credit risk. *IEEE Transactions on Fuzzy Systems*, 2005, 13(6): 820-831
- [8] Zhou L, Lai K K, Yu L. Least squares support vector machines ensemble models for credit scoring. *Expert Systems with Applications*, 2010, 37(1): 127-133
- [9] Xiao Q, Xu K, Wang D, et al. TCP performance over mobile networks in high-speed mobility scenarios//Proceedings of the International Conference on Network Protocols. Raleigh, USA, 2014: 281-286
- [10] Li L, Xu K, Wang D, et al. A measurement study on TCP behaviors in HSPA+ networks on high-speed rails//Proceedings of the IEEE Conference on Computer Communications. Hong Kong, China, 2015: 2731-2739
- [11] Yang Dong, Xu Xin-Yu. Thoughts on the development of Baihang credit. *China Information Security*, 2018, (2): 59-62(in Chinese)  
(杨东, 徐信予. 百行征信发展的几点思考. *中国信息安全*, 2018, (2): 59-62)
- [12] Yang Dong, Xu Xin-Yu. Two changes required by the Baihang credit. *Finance Economy*, 2018, (6): 22-23(in Chinese)  
(杨东, 徐信予. 百行征信需要的两个转变. *金融经济(市场版)*, 2018, (6): 22-23)
- [13] Shen M, Tang X, Zhu L, et al. Privacy-preserving support vector machine training over blockchain-based encrypted IoT data in smart cities. *IEEE Internet of Things Journal*, 2019, 6(5): 7702-7712
- [14] Shen M, Ma B, Zhu L, et al. Secure phrase search for intelligent processing of encrypted data in cloud-based IoT. *IEEE Internet of Things Journal*, 2018, 6(2): 1998-2008
- [15] Andrychowicz M, Dziembowski S, Malinowski D, et al. Secure Multiparty Computations on Bitcoin. New York; ACM, 2016: 59, 76-84
- [16] Xue Teng-Fei, Fu Qun-Chao, Wang Cong, et al. A medical data sharing model via blockchain. *Acta Automatica Sinica*, 2017, 43(9): 1555-1562(in Chinese)  
(薛腾飞, 傅群超, 王枞等. 基于区块链的医疗数据共享模型研究. *自动化学报*, 2017, 43(9): 1555-1562)
- [17] Fan K, Wang S, Ren Y, et al. MedBlock: Efficient and secure medical data sharing via blockchain. *Journal of Medical Systems*, 2018, 42(8): 1-11
- [18] Xia Q, Sifah E B, Asamoah K O, et al. MeDShare: Trustless medical data sharing among cloud service providers via blockchain. *IEEE Access*, 2017, 5(2017): 14757-14767
- [19] Shen M, Deng Y, Zhu L, et al. Privacy-preserving image retrieval for medical iot systems: A blockchain-based approach. *IEEE Network*, 2019, 33(5): 27-33
- [20] Aitzhan N Z, Svetinovic D. Security and privacy in decentralized energy trading through multi-signatures, blockchain and anonymous messaging streams. *IEEE Transactions on Dependable and Secure Computing*, 2018, 15(5): 840-852
- [21] Wu Zhen-Quan, Liang Yu-Hui, Kang Jia-Wen, et al. Secure data storage and sharing system based on consortium blockchain in smart grid. *Journal of Computer Applications*, 2017, 37(10): 2742-2747(in Chinese)  
(吴振铨, 梁宇辉, 康嘉文等. 基于联盟区块链的智能电网数据安全存储与共享系统. *计算机应用*, 2017, 37(10): 2742-2747)
- [22] Shi Ming-Sheng. Analysis of the application of blockchain technology in credit information industry. *Credit Reference*, 2018, (1): 20-24(in Chinese)  
(时明生. 区块链技术在征信业的应用探析. *征信*, 2018, (1): 20-24)
- [23] Ta Lin, Li Meng-Gang. An analysis of the prospects for application of blockchain technology in Internet financial credit. *Journal of Northeastern University (Social Science)*, 2018, 20(5): 466-474(in Chinese)  
(塔琳, 李孟刚. 区块链在互联网金融征信领域的应用前景探析. *东北大学学报(社会科学版)*, 2018, 20(5): 466-474)
- [24] Dong Xiang-Qian, Guo Bing, Shen Yan, et al. An efficient and secure decentralizing data sharing model. *Chinese Journal of Computers*, 2018, 41(5): 1021-1036(in Chinese)  
(董祥千, 郭兵, 沈艳等. 一种高效安全的去中心化数据共享模型. *计算机学报*, 2018, 41(5): 1021-1036)
- [25] Yue X, Wang H, Jin D, et al. Healthcare data gateways: Found healthcare intelligence on blockchain with novel privacy risk control. *Journal of Medical Systems*, 2016, 40(10): 1-8
- [26] Abadi M, et al. Deep learning with differential privacy//Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. Vienna, Austria, 2016: 308-318
- [27] Gonzalez-Serrano F-J, Navia-Vazquez A, Amor-Martin A. Training support vector machines with privacy-protected data. *The Journal of the Pattern Recognition Society*, 2017, 72(2017): 93-107
- [28] De Cock M, Dowsley R, Horst C, et al. Efficient and private scoring of decision trees, support vector machines and logistic regression models based on pre-computation. *IEEE Transactions on Dependable and Secure Computing*, 2017, 16(2): 217-230
- [29] Bost R, Popa R A, Tu S, et al. Machine learning classification over encrypted data//Proceedings of the Network and Distributed System Security Symposium. San Diego, USA, 2015: 4325
- [30] Li H, Liu J, Xu K, et al. Understanding video propagation in online social networks//Proceedings of the International Workshop on Quality of Service. Coimbra, Portugal, 2012: 1-9
- [31] Saito T. Bitcoin: A search-theoretic approach. *International Journal of Innovation in the Digital Economy*, 2015, 6(2): 52-71
- [32] Gao Feng, Mao Hong-Liang, Wu Zhen, et al. Lightweight transaction tracing technology for Bitcoin. *Chinese Journal of Computers*, 2018, 41(5): 989-1004(in Chinese)  
(高峰, 毛洪亮, 吴震等. 轻量级比特币交易溯源机制. *计算机学报*, 2018, 41(5): 989-1004)

- [33] Zhu Lie-Huang, Gao Feng, Shen Meng, et al. Survey on privacy preserving technique for blockchain technology. *Journal of Computer Research and Development*, 2017, 54(10): 2170-2186(in Chinese)  
(祝烈煌, 高峰, 沈蒙等. 区块链隐私保护研究综述. *计算机研究与发展*, 2017, 54(10): 2170-2186)
- [34] Huang Jun-Fei, Liu Jie. Survey on blockchain research. *Journal of Beijing University of Posts and Telecommunications*, 2018, 41(2): 1-8(in Chinese)  
(黄俊飞, 刘杰. 区块链技术研究综述. *北京邮电大学学报*, 2018, 41(2): 1-8)
- [35] Yuan Yong, Wang Fei-Yue. Blockchain: The state of the art and future trends. *Acta Automatica Sinica*, 2016, 42(4): 481-494(in Chinese)  
(袁勇, 王飞跃. 区块链技术发展现状与展望. *自动化学报*, 2016, 42(4): 481-494)
- [36] Katz J, Lindell Y. *Introduction to Modern Cryptography*. 2nd Edition. Boca Roca: Chapman and Hall/CRC, 2014
- [37] Paillier P. Public-key cryptosystems based on composite degree residuosity classes//*Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques*. Berlin, German, 1999: 223-238
- [38] Cortes C, Vapnik V. Support-vector networks. *Machine Learning*, 1995, 20(3): 273-297
- [39] Goldreich O. *Foundations of Cryptography: Volume 2, Basic Applications*. Cambridge, United Kingdom: Cambridge University Press, 2009
- [40] Canetti R. Security and composition of multiparty cryptographic protocols. *Journal of Cryptology*, 2000, 13(1): 143-202



**SHEN Meng**, Ph. D., associate professor. His research interests include privacy-preserving algorithms in cloud computing, blockchain technology and blockchain application.

**ZHANG Jie**, M. S. candidate. His research interests include blockchain application and data privacy-preserving.

**ZHU Lie-Huang**, Ph. D., professor. His research interests

include cryptography, network and information security.

**XU Ke**, Ph. D., professor. His research interests include network security and trustiness, blockchain applications.

**ZHANG Kai-Xiang**, B. S. His research interest is blockchain technology.

**LI Hui-Zhong**, M. S. His research interest is blockchain technology.

**TANG Xiang-Yun**, Ph. D. candidate. Her research interest is applied cryptography.

## Background

The problem that this manuscript focuses on is how to provide multiple data providers to assist the model trainer in training the SVM model under the condition that the data privacy is not leaked. Since credit data has the characteristics of high privacy, easy replication, and high value, the secure sharing of credit data has always been a difficult problem. There is still no ideal data sharing scheme to solve the problems. Blockchain technology has been used in data sharing schemes in medical and other industries in recent years. In most researches, data is stored in encrypted or plaintext on the blockchain. The result of the data sharing is that the data requester obtains the original plaintext data, and then the data requester use the plain data to train a SVM model. Obviously, such data sharing scheme cannot meet the requirement of security in the data sharing process in the credit data field as well.

Many current researches involve privacy protection in the process of machine learning, and the methods currently used are mainly around algorithms such as homomorphic encryption and differential privacy protection. However, on one hand, these solutions cannot achieve a balance between security and availability, and on the other hand, it is necessary to introduce a trusted third party to ensure the successful completion of the training process.

The solution proposed in this article strikes a balance

between security and usability, while avoiding the introduction of trusted third parties. Based on the scheme, a secure credit evaluation model training system for credit investigation scenarios is designed. Compared with other schemes, the secure SVM method based on secure multi-party computing does not need to introduce a third party, and in the collaborative training process, there is no need for frequent communication between data providers and data trainers. This paper conducts experiments on Australian Credit Approval and German Credit Data respectively. Experiments show that the security and usability of the mechanism can be guaranteed. In an acceptable time, the mechanism can train a model whose accuracy can be the same to the model under unprotected conditions.

This work is supported by the National Key Research and Development Program of China(No. 2018YFB0803405), the National Natural Science Foundation of China (Nos. 61902039, 61872041, 61932016), the Natural Science Foundation of Beijing, China (No. 4192050), and the CCF-Tencent Open Fund WeBank Special Funding. Prior to this study, our research team had achieved some results in the privacy protection during the machine learning process and the multi-source medical image retrieval scheme based on blockchain. We also research on the application of blockchain to cross-domain authentication and incentives scheme during the data sharing process.