



# 基于同态加密的医疗数据密文异常检测方法

李腾<sup>1</sup>, 方保坤<sup>2</sup>, 马卓<sup>1\*</sup>, 沈玉龙<sup>1</sup>, 马建峰<sup>1</sup>

1. 西安电子科技大学网络与信息安全学院, 西安 710071

2. 西安电子科技大学广州研究院, 广州 510555

\* 通信作者. E-mail: mazhuo@mail.xidian.edu.cn

收稿日期: 2022-05-29; 修回日期: 2022-07-22; 接受日期: 2022-11-07; 网络出版日期: 2023-07-06

国家自然科学基金 (批准号: 62272370, U21A20464)、青年人才托举工程 (批准号: 2022QNRC001)、111 计划 (批准号: B16037) 和中央高校基本科研业务费专项基金 (批准号: QTZX23071) 资助项目

**摘要** 为了避免电子健康数据 (electronic health records, EHRs) 在异常检测过程中泄漏患者信息和诊断结果等问题, 针对医院与患者敏感数据的隐私保护, 提出了一种基于 CKKS 全同态加密的 EHRs 异常检测隐私保护模型. 将医院以及患者的 EHRs 由 CKKS 算法实现浮点数同态加密, 设计一个基于密文比较算法的协议, 通过可信密钥服务器与第三方数据中心之间的通信建立密文态孤立森林模型. 并利用 CKKS 算法的 SIMD 技术, 实现密文数据在孤立森林模型上的异常检测, 最终返回密文结果. 理论分析和实验结果表明, 所提出的方案能够保证 EHRs 的隐私安全. 并且在不同的 EHRs 数据集上验证, 该模型优于传统的明文异常检测算法和同类型密文异常检测算法, 且在密文态上能够保持和明文孤立森林算法相近的检测效率, 有较好的异常检测效果.

**关键词** 同态加密, 孤立森林, 异常检测, 隐私保护, 密文比较

## 1 引言

随着新冠疫情的暴发和全球肆虐, 电子健康记录 (electronic health records, EHRs) 已成为大数据中不可忽略的一部分<sup>[1~4]</sup>. 目前 EHRs 的增长速度越来越快, 数量也越来越庞大. 在这些记录中, 往往包括患者的基本个人信息、病史数据、医疗检测信息、位置信息等. 此外, 随着人工智能技术的发展, 这些检测数据能够被应用于智慧医疗诊断, 例如利用异常检测模型<sup>[5~8]</sup>, 可以帮助患者从 EHRs 中尽早地检测出病理异常情况, 从而进行及时的治疗. 这类技术能够使医疗诊断更加迅速, 使患者足不出户就能够享受准确且个性化的医疗诊断, 对癌症的预防、新冠疫情的预测、心脑血管疾病的提前预防都有着重要的指导意义, 也是未来医疗领域的趋势<sup>[9]</sup>.

EHRs 中包含着许多患者的个人信息、诊断结果以及其他一些敏感数据, 通常存储在医疗机构的中央服务器或者可信数据中心中. 然而, 这种通用的架构并非安全, 这些数据一旦遭受外部攻击或内

**引用格式:** 李腾, 方保坤, 马卓, 等. 基于同态加密的医疗数据密文异常检测方法. 中国科学: 信息科学, 2023, 53: 1368–1391, doi: 10.1360/SSI-2022-0214  
Li T, Fang B K, Ma Z, et al. Homomorphic encryption-based ciphertext anomaly detection method for e-health records (in Chinese). Sci Sin Inform, 2023, 53: 1368–1391, doi: 10.1360/SSI-2022-0214

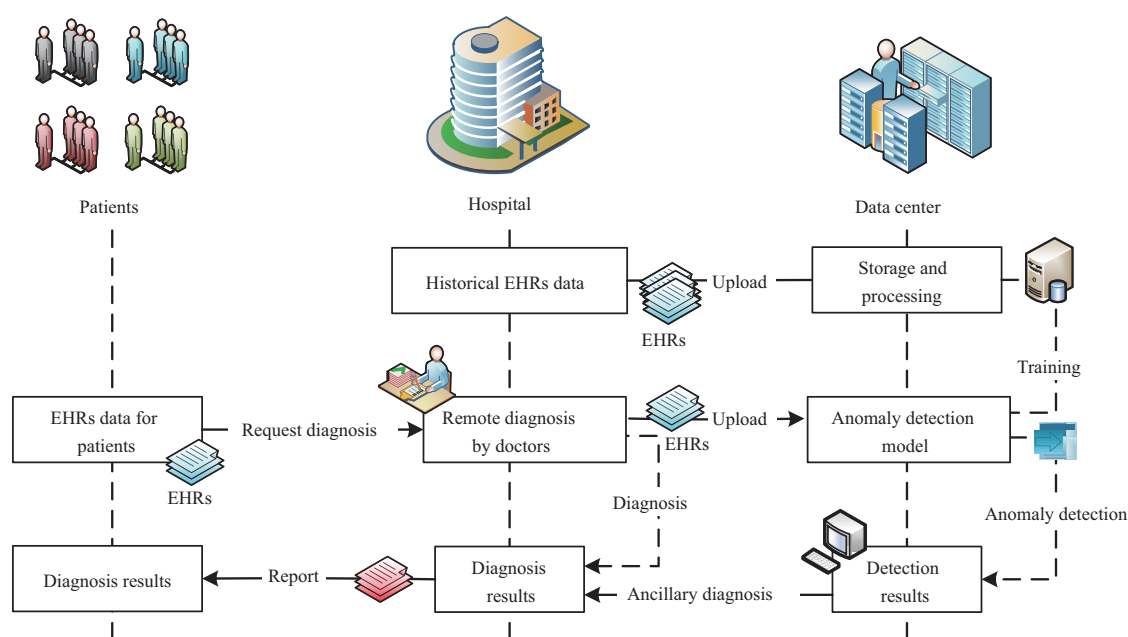


图 1 (网络版彩图) 传统 EHRs 应用场景架构

Figure 1 (Color online) Traditional EHRs application scenario

部人员误操作将会导致数据泄漏或丢失,造成不可挽回的损失<sup>[10~14]</sup>. 2021年10月,美国的公共卫生系统 Broward Health 遭受攻击,泄漏了超过 130 万人的患者信息,其中包括医疗数据、诊断数据、病情历史等. 因此如何在保护隐私的前提下,实现对 EHRs 的智能异常检测是一个大挑战,这也是未来智慧医疗的一大难题<sup>[15,16]</sup>.

面对这一挑战,目前的 EHRs 应用场景架构如图 1 所示,患者希望能够不泄漏自己的隐私敏感数据,获取自己的 EHRs 诊断结果,并且不希望泄漏自己的诊断结果. 同时,各医院拥有大量患者的 EHRs,其中包括医疗健康数据以及诊断结果,这些数据能够为异常检测模型提供丰富的训练集<sup>[6,17]</sup>,但医院又缺乏可以提供数据挖掘和数据分析技术的人员,因此需要第三方数据中心的协助. 虽然第三方机构能够帮助医院搭建异常检测模型、提供数据挖掘及分析服务,帮助患者实现在线诊断,但是其在数据的传输和存储的过程中,可能会由于安全隐患遭受攻击或者内部人员窃取导致敏感数据泄漏<sup>[18~20]</sup>. 针对上述场景的需求与挑战如表 1 所示,常见的隐私保护方法有: (1) K-匿名化技术<sup>[7,21,22]</sup>,通过删除或隐藏患者 EHRs 中的表示属性,使攻击者无法区分信息的归属者,但该方法却难以抵抗拥有背景知识的重识别攻击; (2) 引入噪声<sup>[22~25]</sup>,使用混淆和基于随机化的方法对 EHRs 引入噪声,使得攻击者无法获取准确的数据,但是这也会影响异常检测结果的精确性; (3) 密码学知识<sup>[26~30]</sup>,能够对患者的 EHRs 隐私数据进行加密保护,但是传统密码方法加密后的数据无法进行密文上的运算实现数据挖掘.

针对上述的场景与挑战,密码学家提出了全同态加密算法 (fully homomorphic encryption, FHE)<sup>[9,31~35]</sup>,该算法能够提供在密文环境下直接对加密的数据进行算术运算<sup>[36~38]</sup>,并且可以恢复出明文的运算结果,这能够使得 EHRs 在计算的过程中不需要对密文数据进行解密,避免了信息的泄漏和被窃取. 但是,现阶段的全同态加密技术限制性较多,其计算开销大,在一些问题上需要较长的运算时间<sup>[39~41]</sup>. 并且全同态加密只支持同态加法和同态乘法运算<sup>[42~44]</sup>. 此外,现阶段实现无监督机器学习

表 1 需求与挑战  
Table 1 Needs and challenges

Role	Needs	Challenge
Patients	Diagnostic results required	EHRs and diagnostic results are easy to leak
Hospital	Using EHRs to build anomaly detection model	Lack of data mining ability
Data center	Help hospitals to build anomaly detection models	Implement EHRs security and privacy computing

习的异常检测算法, 主要包括 (1) 基于距离的算法; (2) 基于密度的算法; (3) 基于聚类的算法. 但是这些方法通常都用于处理维度较小的小规模数据集, 因为这些算法涉及一些复杂函数, 导致计算量较高, 不适用于海量数据集. 且这些算法涉及的复杂函数与同态加密结合实现并不友好.

因此, 本文提出一种基于 CKKS 全同态加密的孤立森林 (isolation forest, IF) 异常检测模型<sup>[45, 46]</sup>, 对患者 EHRs 进行加密保护了数据隐私, 又能够对密文进行运算实现数据挖掘. 本文采用的孤立森林算法, 能够将训练和异常检测任务分布成多个孤立树 (isolation tree, iTree) 的子任务, 因此可以并行实现以优化同态加密计算时间. 该模型能够收集医院加密后的 EHRs, 并利用密文训练孤立森林异常检测模型, 此外还能够对患者加密的 EHRs 实现异常检测. 在此过程中确保了所有的数据皆以密文的形式保存和运算. 本文的主要贡献如下:

(1) 本文提出了基于 CKKS 全同态加密的异常检测隐私保护模型. 该模型针对 EHRs 提供了可靠的同态加密保护, 有效保护了数据的隐私安全. 首次将同态加密算法结合孤立森林算法实现异常检测, 能够对密文数据进行训练. 并且在数据中心的孤立森林异常检测算法能够同时进行多棵密文孤立树的构建, 利用数据中心的高算力多分布式的优势, 可以实现快速的异常检测.

(2) 为了实现密文数据在第三方数据中心的密文可计算性, 本文模型设计了 CKKS 中的 MaxIdx 算法, 该方案使数据中心与密钥服务器仅一次通信就能够实现对浮点数加密数据序列求最值序号的操作. 这个方案确保了孤立森林模型对密文数据的构建, 同时避免了数据泄漏导致隐私泄漏, 且具有较好的计算效率.

(3) 为了实现密文数据异常检测的高效性, 本文模型利用了 CKKS 算法的单指令多数据流技术 (single instruction multiple data, SIMD). 该技术可以对密文数据序列一次操作就进行多密文的同态计算操作. 在孤立森林的预测过程中, 利用 SIMD 技术可以无须进行通信与重加密, 密文数据在孤立树上的一次计算实现高度求解, 优化了检测效率.

(4) 本文设计了相关实验测试该模型的性能, 并通过安全性分析证明该模型能够实现隐私保护, 避免数据泄漏. 实验证实不同的 EHRs 数据集中, 实现密文异常检测, 最高能达到 95.6% 的异常检测准确率, 优于现阶段最好的密文异常检测算法 BGV-FCM, 且与明文孤立森林算法的异常检测能力相当. 在时间开销上, 一棵同态孤立树的构建仅为 312 s, 效率优于同类型密文异常检测算法.

## 2 相关工作

### 2.1 医疗数据隐私保护方法

现阶段差分隐私 (differential privacy, DP) 在保护医疗数据隐私方面有很强的理论保证, 其也被广泛应用于数据挖掘和机器学习当中. 2019 年 Du 等<sup>[7]</sup> 提出了基于差分隐私的鲁棒性异常检测, 将差分隐私应用于注入异常值的医疗数据集上训练的自动编码器网络, 但是需要对不同的数据集设置不同的  $\epsilon$  值, 以保证异常检测的精度. 2020 年 Fan 等<sup>[24]</sup> 针对云计算支持的数据中心, 提出了一种满足局部

差异隐私的分类算法. 通过加入拉普拉斯 (Laplace) 机制, 保证隐私保护的有效性和可靠性. 但是局部差分隐私在高维数据下还是难以保持高精度的异常检测. 2019 年 Hou 等<sup>[47]</sup> 将差分隐私应用于随机森林分类算法中, 提出了一种基于差分隐私的随机森林分类算法来保护数据分类过程中的隐私信息, 平衡了隐私预算以及分类精度. 但是总的来说, 差分隐私引入了噪声会导致数据挖掘性能下降. 2019 年 Alabdulatif 等<sup>[48]</sup> 提出了对医疗数据的分布式隐私保护检测, 但该方法依赖于高并发通信, 并且算法适用性弱. 此外, 2020 年 Li 等<sup>[12]</sup> 提出了一种基于安全多方计算 (secure multi-party computation, MPC) 与同态加密结合的方式进行自助医疗诊断方案. 2021 年 Reich 等<sup>[29]</sup> 提出了第一个基于安全多方计算的 EHRs 文本分类隐私保护方案, 从文本中提取特征, 使用逻辑回归和树集合进行后续的文本分类. 但总的来说安全多方计算严重依赖于通信, 且算法设计存在困难, 不能支持离线计算, 因此在 MPC 框架下建立模型在实际应用中通常并不容易.

## 2.2 基于同态加密的机器学习

现阶段基于同态加密的数据挖掘和机器学习技术研究逐渐兴起<sup>[49~52]</sup>. 在 2019 年 Alabdulatif 等<sup>[48, 53~55]</sup> 提出了基于 BGV 同态加密算法的异常检测, 该方案利用 FCM 的聚类方法实现异常检测, 但在该方案中, 使用 IEEE 754 标准实现整数表示浮点数, 增加了密文计算的复杂度, 效率较低<sup>[56]</sup>. 由于 BGV 算法只能加密整数型数据, 因此能够实现浮点数加密的 CKKS 算法, 开始得到广泛应用. 2020 年 Xu 等<sup>[57]</sup> 提出了基于 CKKS 同态加密算法的多分类的 Logistic 回归模型, 实现简单的二分类模型. 2021 年 Lv 等<sup>[58]</sup> 提出了基于 CKKS 同态加密的线性系统求解的方案, 利用 SIMD 技术实现了并行求解, 可以在机器学习中处理更复杂的矩阵运算, 2022 年 Lee 等<sup>[59]</sup> 在线性求解的基础上提出了基于 RNS-CKKS 算法的神经网络模型, 实现多分类模型. 但上述方案多用于有监督的机器学习算法, 这些算法仅涉及简单的密文加乘运算. 由于同态加密算法较难实现比较运算, 因此在树形结构的机器学习模型上的实现就较为困难<sup>[47, 60~62]</sup>. 2019 年 Aloufi 等<sup>[63]</sup> 提出了基于多密钥同态加密的随机森林评估实现多分类模型, 利用多密钥 BGV 算法 (multi-key BGV, MKBGV), 提出了多用户的随机森林评估协议, 结合安全多方的思想, 实现基于电路的比较运算的特点, 但是对于该方案多密钥协议如果参与方有一方脱机或离线就会导致评估难以进行. 2020 年 Huynh<sup>[64]</sup> 提出了 Cryptotree, 他利用 DNN 神经网络, 将随机森林模型转化为神经随机森林, 取代了密文比较的过程, 但是该方法增加了运算的时间复杂度.

## 2.3 密文数据比较方法

在一些机器学习算法中, 常会涉及一些比较算法, 常见的密文比较算法是基于安全多方的密文比较协议. 在 2019 年 Ichikawa 等<sup>[61]</sup> 提出了新的安全多方决策树分类协议, 通过隐藏输入向量和输出类, 实现不泄漏隐私就进行密文比较, 但是安全多方计算的方法也带来了指数级的通信复杂度. 而基于同态加密的密文比较算法研究中, 2017 年 Cheon 等<sup>[65]</sup> 提出 CKKS 同态加密算法, CKKS 支持针对实数或复数的浮点数加密, 适用于机器学习模型训练等不需要精确结果的场景. 但是 CKKS 加密方案并不支持浮点数据的比较. 因此 2019 年 Cheon 等<sup>[66]</sup> 提出了极小极大逼近的同态比较方案, 并证明了有理函数可以通过迭代算法进行计算. 由于这种迭代特性, 与以往的多项式逼近方法相比, 其计算复杂度达到对数级. 2020 年他们又提出了基于最优复杂性的同态比较方法, 该方法使基于 word-wise 加密的方案在密文比较运算上取得了接近 bit-wise 的效率<sup>[67]</sup>. 2021 年, Jia 等<sup>[68]</sup> 提出的基于 BGV 算法的 DBSCAN 聚类隐私保护方案中, 提出了三种整数编码的形式, 将浮点数据映射到整数型数据中, 并实现浮点数密文的比较, 但是该算法也需要频繁地进行服务器与用户端的通信, 会影响计算效率. 此外,

表 2 CKKS 方案符号含义表示

Table 2 List of CKKS notations

Symbol	Notation
$R$	Quotient ring
$q$	Modulus
$R_q$	The residue ring of $R$ modulo $q$ , $R_q = R/qR$
$L$	The initial noise budget
$l$	The noise budget
$p$	The noise of homomorphic multiplication
HWT( $h$ )	The distribution that chooses a polynomial uniformly from $R_{2L}$ , that it has exactly $h$ nonzero coefficients
$U(S)$	The uniform distribution of finite set $S$
DG( $\sigma^2$ )	The discrete Gaussian distribution of variance $\sigma^2$
pk	The public key
sk	The secret key
evk	The evaluation key
$m$	Plaintext
ct	Ciphertext

现阶段基于布尔运算的同态加密算法 TFHE 能够实现快速高效的 bootstrapping 操作<sup>[35]</sup>, 这些操作都是基于门电路运算, TFHE 能够快速执行单比特密文数据的比较<sup>[69, 70]</sup>. 但由于布尔电路的限制, TFHE 在整数和浮点数的比较上, 需要设计较复杂的布尔比较电路.

### 3 背景知识

#### 3.1 CKKS 同态加密算法

CKKS 同态加密提供了浮点数加密的近似算法<sup>[71]</sup>. 其明文空间和密文空间是基于整数多项式环结构构造的, 但样本数据往往是以向量形式出现的, 因此需要将其编码成多项式形式. 该算法在加密过程中应用一些“小”噪声来屏蔽消息. 噪声预算是在方案初始化时确定的, 对密文的计算会消耗这些预先分配的预算. 一旦噪声预算耗尽, 解密将返回错误的结果. 由于该方案具有处理实数的能力, 因此常应用于数据科学和机器学习. 以下是 CKKS 方案的一些细节:

CKKS 方案的噪声预算用参数  $L$  初始化. 密文的噪声预算表示为  $l$ . 对于每次执行同态乘法, 噪声预算被整数  $p$  减去, 即  $l = l - p$ . 当消息刚刚加密时,  $l = L$ . 当  $l < p$  时, 噪声预算被称为耗尽. 本节使用的符号含义表示如表 2 所示.

(1) KeyGen( $\sigma$ ). 设  $2^L$  是初始密文模, 设 HWT( $h$ ) 表示从  $R_{2L}$  中均匀选取一个多项式的分布, 它有  $h$  个非零系数. 从 HWT( $h$ ) 中抽取一个秘密值  $s$ , 从  $U(R_{2L})$  中抽取一个随机数  $a$ , 从 DG( $\sigma^2$ ) 中取一个误差值  $e$ . 则返回私钥为  $\text{sk}(\text{sk} \leftarrow (1, s))$ , 公钥为  $\text{pk}(\text{pk} \leftarrow (b, a) \subseteq R_{2L})$ ,  $b = -a \cdot s + e(\text{mod } L)$ . 然后取  $a' \leftarrow U(R_{2L})$ ,  $e' \leftarrow \text{DG}(\sigma^2)$  计算密钥为  $\text{evk}(\text{evk} \leftarrow (b', a'))$ ,  $b' = -a' + e' + L \cdot s^2(\text{mod } 2^L)$ .

(2) Encrypt(pk,  $m$ ). 对于明文  $m \subseteq R_{2L}$ , 取  $v \leftarrow U(R_{2L})$ ,  $e_0, e_1 \leftarrow \text{DG}(\sigma^2)$ , 则返回密文  $\text{ct} = -v \cdot \text{pk} + (m + e_0, e_1)(\text{mod } 2^L)$ .

(3) Decrypt(sk, ct). 对于密文  $\text{ct} = ((c_0, c_1), l) \subseteq R_{2^l}$ , 则返回明文  $c_0 + c_1 \cdot \text{sk}(\text{mod } 2^l)$ .

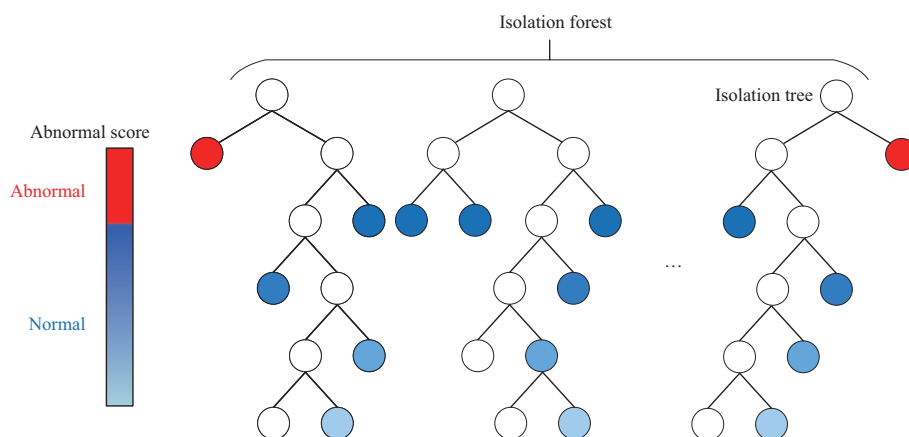


图 2 (网络版彩图) 孤立森林示意图

Figure 2 (Color online) Illustration of an isolation forest

(4)  $\text{Add}(ct_1, ct_2)$ . 对于密文  $ct_1 = ((c_0, c_1), l) \subseteq R_{2^l}$ ,  $ct_2 = ((c'_0, c'_1), l) \subseteq R_{2^l}$ , 则返回密文  $ct = (c_0, c_1) + (c'_0, c'_1) \pmod{2^l}$ .

(5)  $\text{Mult}(ct_1, ct_2)$ . 对于密文  $ct_1 = ((c_0, c_1), l) \subseteq R_{2^l}$ ,  $ct_2 = ((c'_0, c'_1), l) \subseteq R_{2^l}$ , 则  $(d_0, d_1, d_2) = (c_0 c'_0, c_1 c'_0 + c_0 c'_1, c_1 c'_1) \pmod{2^l}$ , 返回密文  $ct = (d_0, d_1) + [2^{-L} \cdot d_2 \cdot \text{evk} \pmod{2^l}]$ .

(6)  $\text{Rescale}(ct, p)$ . 对于密文  $ct = ((c_0, c_1), l) \subseteq R_{2^l}$ , 整数  $p \leq l$ , 密文重缩放后返回  $ct' = [2^{-p} \cdot (c_0, c_1)] \pmod{2^{l-p}}$ .

(7)  $\text{Rotate}(ct, r)$ . 对于密文向量  $ct$ , 返回  $ct$  向量向右移动  $r$  个位置的密文向量.

此外, CKKS 同态加密还支持 SIMD 操作, 也称为批量操作, 其将多组明文组合成向量, 加密成单一的密文向量. 然后这个密文向量就是一个可以批处理的密文向量. 利用 SIMD 操作, 就可以实现一次指令, 对密文向量中的所有密文进行相同的计算, 以实现同态加法和乘法运算的并行化. 并且结合上面的 Rotate 操作就可以实现密文向量中的旋转, 从而完成更复杂的向量运算, 因此 SIMD 技术可以降低计算的成本.

### 3.2 孤立森林

孤立森林 (IF) 算法是一种常用的无监督异常检测算法, 适用于高维数据集. 与其他异常检测算法不同, 孤立森林算法不再聚类正常样本点, 而是分离异常点. 在孤立森林中, 异常点被定义为“更容易被分离的点”, 可以理解为分布稀疏且远离稠密群体的点. 在特征空间中, 稀疏分布的区域表示在该区域发生事件的概率非常低, 因此在这些区域中的数据可以被视为异常. 孤立森林由多个孤立树 (iTree) 构成, 孤立树和二叉搜索树具有相同的结构. 在孤立树中, 数据集被递归地随机划分, 直到所有采样点被孤立或孤立树达到设置的高度. 在这种随机分割策略下, 异常点通常具有较短的路径, 如图 2 所示, 异常点被很快地分隔. 直观地说, 异常点在较低的叶子被分割, 在孤立树上的路径较短; 而正常点在较高的叶子被分割, 在孤立树上的路径较长.

孤立森林分为以下两个步骤.

(1) **生成孤立树**. 从训练组中取样, 建立孤立树, 再由孤立树组成孤立森林. 首先, 从训练数据中随机选取  $n$  个点作为子样本, 放入孤立树的根节点, 其次, 随机选择维度, 在当前节点数据范围内随机生成切割点  $p$ , 切割点在当前节点数据中指定维度的最大值和最小值之间生成; 然后, 利用  $p$  点将当前

节点数据空间划分为 2 个子空间: 将小于  $p$  的点放在当前节点的左分支上, 将大于或等于  $p$  的点放在节点的右分支上; 最后, 递归地在节点的左分支和右分支上执行两个步骤, 并继续构造新的叶节点, 直到叶节点上只有一个数据 (无法再分割), 或者树已生长到设置的高度。

(2) 计算异常得分. 在孤立树上记录测试集样本点的路径. 根据异常得分计算公式计算每个样本点的异常得分. 在获得孤立树后, 可以使用所有的孤立树计算异常分数. 对于每个样本, 需要综合计算每个树的结果, 异常得分由式 (1) 计算得出.

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \Rightarrow \begin{cases} s(x, n) \rightarrow 1, \text{ 异常,} \\ s(x, n) \rightarrow 0.5, \text{ 无明显异常,} \\ s(x, n) \rightarrow 0, \text{ 正常,} \end{cases} \quad (1)$$

$$c(n) = 2H(n-1) - \frac{2(n-1)}{n}, \quad (2)$$

$$H(i) = \ln(i) + \gamma (\gamma = 0.577215\dots), \quad (3)$$

其中  $h(x)$  表示样本  $x$  到达叶子节点的路径长度, 式 (2) 中  $c(n)$  表示一个包含  $n$  个样本的树的平均路径长度, 式 (3) 中  $H(i)$  表示调和级数.

## 4 系统模型

### 4.1 系统模型概述

本文提出基于 CKKS 全同态加密的孤立森林异常检测模型, 能够实现 EHRs 的隐私保护异常检测. 其主要分为两个场景: 第一个场景, 数据中心为医院服务, 是传统的两方模型. 由可信的密钥服务器, 向医院和第三方数据中心分发密钥, 医院将本地存储的 EHRs 进行同态加密后, 将密文发送给数据中心, 数据中心对加密后的密文数据进行运算实现孤立森林的异常检测模型训练并对密文进行异常检测. 最后数据中心可以将加密的孤立森林模型以及加密的异常检测结果返回给医院, 由医院解密出明文结果和模型. 第二个场景, 患者希望对自己 EHRs 进行隐私保护的异常检测. 由可信的密钥服务器, 向医院、患者和第三方数据中心分发密钥. 医院利用公钥对场景 1 中构建的模型加密后返回给数据中心, 患者将自己的 EHRs 加密发送给数据中心. 数据中心利用加密的模型对患者的加密数据进行异常检测. 并将结果以密文的形式返回给患者, 患者解密就可以诊断 EHRs 是否存在异常, 确保及时治疗.

针对该系统, 其两个应用场景主要细节如图 3 所示. 在场景 1 中, 首先, 密钥服务器与医院通信, 生成公钥  $pk$  与私钥  $sk$  并分发, 同时向第三方数据中心分发计算密钥  $evk$ . 其次, 医院将 EHRs 进行特征选取, 利用主成分分析算法 (principal component analysis, PCA) 简化数据维度, 降低数据的可解释性, 优化孤立森林算法的第一个随机过程, 将选取好特征的数据进行加密并发送给第三方数据中心. 在数据中心的, 需要执行密文比较运算, 通过计算可以得到若干较大值的序号, 并且与密钥服务器通信, 返回这些密文数据序号实现孤立树节点分离, 通过这种方式, 可以在数据中心建立密文态的孤立森林. 然后数据中心也可以对医院的 EHRs 进行异常检测. 利用 CKKS 的 SIMD 性质, 可以实现多密文同时比较, 计算出单个密文 EHRs 在各个孤立树中的高度, 并且数据中心无法获取这些加密的结果. 最后将每棵孤立树返回的加密结果平均计算, 就可以将加密结果返回给医院, 由医院将结果进行解密后, 进行异常得分的计算, 并与阈值进行比较, 就可以获得 EHRs 异常检测的结果. 并且数据中心

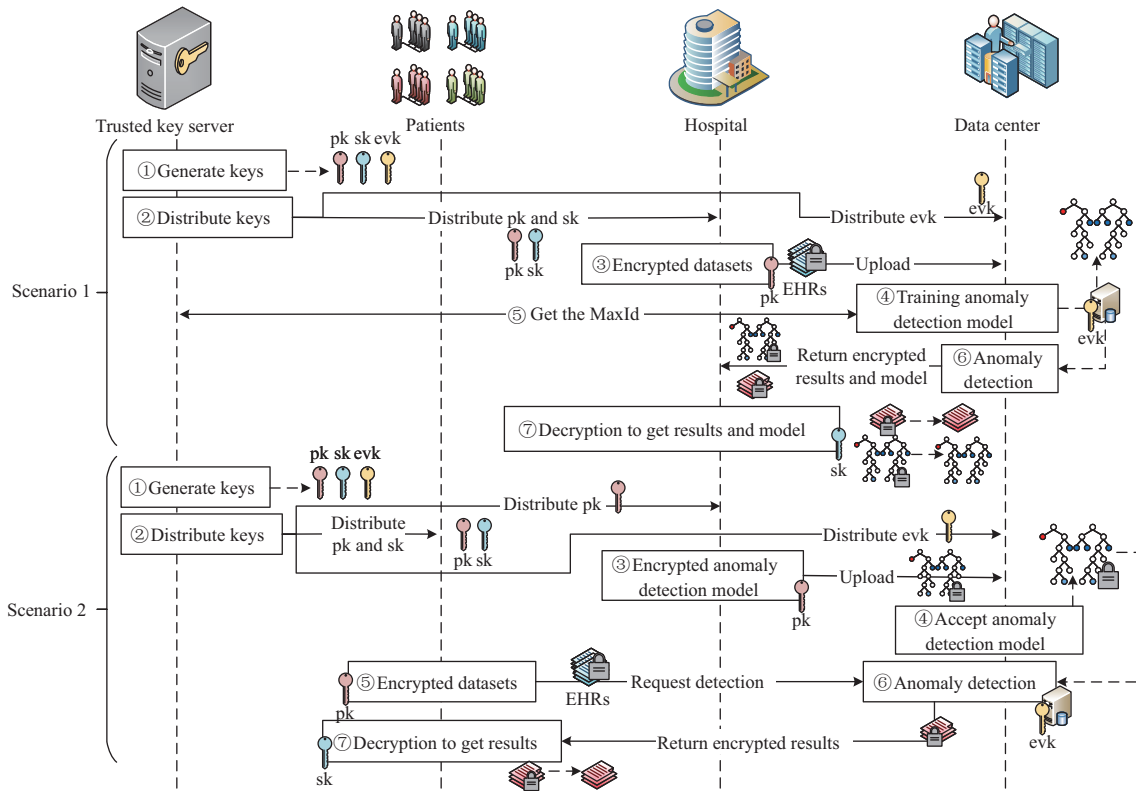


图 3 (网络版彩图) 系统框架图

Figure 3 (Color online) System framework of our proposed scheme

也可将加密的孤立森林模型发送回给医院, 以数据结构形式存储方便后续利用; 在场景 2 中, 首先密钥服务器与医院和患者进行通信, 生成公钥  $pk$  与私钥  $sk$  并分发, 同时向第三方数据中心分发计算密钥  $evk$ . 由医院方对场景 1 中生成的模型利用公钥  $pk$  进行加密后, 发送给数据中心. 同时患者也对自己的 EHRs 进行指定的 PCA 处理并加密以密文的形式发送给数据中心. 同理, 后者利用 CKKS 的 SIMD 性质, 可以实现多密文的同时比较, 计算出单个密文 EHRs 在各个孤立树中的高度, 每棵孤立树返回的加密结果平均计算后, 将加密结果返回患者, 由患者将加密结果进行解密后, 进行异常得分的计算, 并与阈值进行比较, 就可以获得 EHRs 异常检测的结果. 后续章节将对系统模型中的详细算法进行解释. 本节系统模型使用的符号含义表示如表 3 所示.

### 4.2 密文比较算法

在 CKKS 同态加密算法中, 并不具有密文比较的运算, 因此需要一些密文计算, 来近似地获取密文比较的结果 [72~74]. 在本文的密文孤立森林模型训练过程中, 需要将一组数据按大小随机分离, 这个问题可以转化为在一组数据中寻找若干个最大值的序号. 要在一组数据  $(a_1, a_2, \dots, a_n)$  中寻找最大值, 如式 (4) 所示:

$$\lim_{k \rightarrow \infty} \frac{a_i^k}{a_1^k + a_2^k + \dots + a_n^k} = \begin{cases} 1, & a_i \text{ 是最大值,} \\ 0, & a_i \text{ 非最大值.} \end{cases} \quad (4)$$

而在孤立森林的检测过程中, 也需要实现数据与孤立树的密文节点比较大小, 同理, 实现比较大小的



表 3 本文系统模型的符号含义表示  
Table 3 List of our proposed scheme notations

Symbol	Notation
$[\cdot]$	The ciphertext data
Inv	The algorithm to calculate the reciprocal of data $x$
Comp	The algorithm to compare the data $a$ and $b$
MaxIdx	The algorithm to calculate the sequence number of the max value in a data sequence
enc.iTree	Homomorphic encryption isolation tree
$d$	The number of loop iterations
inv	The reciprocal of data
MaxId	The sequence number of the max value in the data series
$\mathbf{x} = (x_1, x_2, \dots, x_n)$	Data sequence
$X_{n \times m}$	The matrix of $n \times m$
SplitAtt	The attribute for splitting each node in isolation tree
SplitValue	The Value corresponding to the attributes of each node in the isolation tree for segmentation
$N$	The node in isolation tree
$L, R$	The left child or right child node of isolation tree node $N$
SM	SM close to 1 means that the data finally falls in the left child of this node of the isolation tree, otherwise it falls in the right child
$L_m$	Each leaf node on the isolation tree
$l_L$	The height corresponding to each leaf node on the isolation tree

Comp 函数可以利用式 (5):

$$\text{Comp}(a, b) \approx \lim_{k \rightarrow \infty} \frac{a^k}{a^k + b^k} = \begin{cases} 1, & a > b, \\ \frac{1}{2}, & a = b, \\ 0, & a < b. \end{cases} \quad (5)$$

但是在同态运算中最基本的是加法和乘法, 难以实现上述除法运算, 因此, 文献 [65] 针对倒数 (Inv) 运算提出利用 Goldschmidt 算法. Goldschmidt 算法是由泰勒 (Taylor) 级数转换而来的, 其具体推导过程如式 (6) 所示:

$$\frac{1}{1+x} = 1 - x + x^2 - x^3 + \dots = (1-x)(1+x^2)(1+x^4)(1+x^8) \dots = \prod_{i=0}^{\infty} (1 + (-x)^{2^i}). \quad (6)$$

因此在  $0 < x < 2$  条件下,  $\frac{1}{x}$  的求解如式 (7) 所示:

$$\frac{1}{x} = \frac{1}{1+(x-1)} = \prod_{i=0}^{\infty} (1 + (1-x)^{2^i}) \approx \prod_{i=0}^d (1 + (1-x)^{2^i}), \quad 0 < x < 2. \quad (7)$$

该算法能够将倒数运算转化成有限次的乘法运算. 在 CKKS 密文中的同态 Inv 算法如算法 1 所示.

实现同态密文的 Inv 运算, 即利用其性质, 实现密文 Comp 运算, 密文的同态 Comp 算法如算法 2 所示, 可以求出两个密文数据的比较结果, 并且比较结果以密文的形式呈现. 同理, 利用 Inv 运算

**算法 1** 同态 Inv 算法**输入:** 密文  $\llbracket x \rrbracket$ ,  $x \in (0, 2)$ , 迭代次数  $d \in N$ ;**输出:** 密文  $\llbracket a_d \rrbracket$ ,  $a \approx \frac{1}{x}$ ;

```

1:  $a_0 \leftarrow 2 - x$ ;
2:  $b_0 \leftarrow 1 - x$ ;
3: for  $n = 0, 1, 2, \dots, d - 1$  do
4:    $b_{n+1} \leftarrow b_n^2$ ;
5:    $a_{n+1} \leftarrow a_n \cdot (1 + b_{n+1})$ ;
6: end for

```

**返回:**  $\llbracket a_d \rrbracket$ .**算法 2** 同态 Comp 算法**输入:** 密文数据  $\llbracket a \rrbracket$ ,  $\llbracket b \rrbracket$ ,  $a, b \in (\frac{1}{2}, \frac{3}{2})$ ,  $d, d', m, t \in N$ ;**输出:** 密文  $\llbracket a_t' \rrbracket$ ,  $a_t \approx 1 \Rightarrow a > b$ ,  $a_t \approx 0 \Rightarrow a < b$ ,  $a_t \approx \frac{1}{2} \Rightarrow a = b$ ;

```

1:  $\llbracket a_0 \rrbracket \leftarrow \frac{\llbracket a_0 \rrbracket}{2} \cdot \text{Inv}(\frac{\llbracket a \rrbracket + \llbracket b \rrbracket}{2}, d')$ ;
2:  $\llbracket b_0 \rrbracket \leftarrow 1 - \llbracket a_0 \rrbracket$ ;
3: for  $i = 0, 1, 2, \dots, t - 1$  do
4:    $\text{inv} \leftarrow \text{Inv}(\llbracket a_i^m \rrbracket + \llbracket b_i^m \rrbracket, d)$ ;
5:    $\llbracket a_{i+1} \rrbracket \leftarrow \llbracket a_i \rrbracket \cdot \text{inv}$ ;
6:    $\llbracket b_{i+1} \rrbracket \leftarrow 1 - \llbracket a_{i+1} \rrbracket \cdot \text{inv}$ ;
7: end for

```

**返回:**  $\llbracket a_t \rrbracket$ .

也能够实现多个密文数据组成的密文序列, 求序列中最大值的序号的 Max 算法如算法 3 所示. 但这两个算法的特点都是最后以密文的形式输出的. 对于 Comp 算法其输出是 0 或 1 的形式, 可以利用其布尔运算的性质进行后续运算. 但对于 Max 算法它只能返回最大值序号的密文, 在不解密的情况下, 难以利用, 因此针对这个问题, 本文利用 Max 算法设计与密钥服务器通信且返回最大值序号 MaxIdx 的算法 MaxIdx, 如算法 4 所示.

**算法 3** 同态 Max 算法**输入:** 密文序列  $(\llbracket x_1 \rrbracket, \llbracket x_2 \rrbracket, \dots, \llbracket x_n \rrbracket)$ ,  $x_i \in (\frac{1}{2}, \frac{3}{2})$ ,  $d, d', m, t \in N$ ;**输出:** 密文序列  $(\llbracket b_1 \rrbracket, \llbracket b_2 \rrbracket, \dots, \llbracket b_n \rrbracket)$ , 若  $a_i$  是最大值则  $b_i$  越接近 1;

```

1:  $\text{inv} \leftarrow \text{Inv}(\sum_{j=1}^n \llbracket x_j \rrbracket / n, d)$ ;
2: for  $j = 1, 2, \dots, n - 1$  do
3:    $\llbracket b_j \rrbracket \leftarrow \llbracket x_j \rrbracket / n \cdot \text{inv}$ ;
4: end for
5:  $b_n \leftarrow 1 - \sum_{k=1}^{n-1} \llbracket b_k \rrbracket$ ;
6: for  $i = 1, 2, \dots, t$  do
7:    $\text{inv} \leftarrow \text{Inv}(\sum_{j=1}^n \llbracket b_j \rrbracket^m, d)$ ;
8:   for  $j = 0, 1, 2, \dots, n - 1$  do
9:      $\llbracket b_j \rrbracket \leftarrow \llbracket b_j \rrbracket^m \cdot \text{inv}$ ;
10:  end for
11:    $\llbracket b_n \rrbracket \leftarrow 1 - \sum_{k=1}^{n-1} \llbracket b_k \rrbracket$ ;
12: end for

```

**返回:**  $(\llbracket b_1 \rrbracket, \llbracket b_2 \rrbracket, \dots, \llbracket b_n \rrbracket)$ .

综上, 本文列出了 MaxIdx 算法与其他的密文比较的对比, 如表 4 所示. 表中  $d$  表示同态 Inv 算

**算法 4** MaxIdx 算法**输入:** 密文序列 ( $\llbracket x_1 \rrbracket, \llbracket x_2 \rrbracket, \dots, \llbracket x_n \rrbracket$ );**输出:** 返回 MaxId 表示序列中最大值的序号;

- 1: 数据中心计算 ( $\llbracket b_1 \rrbracket, \llbracket b_2 \rrbracket, \dots, \llbracket b_n \rrbracket$ ) = Max( $\llbracket x_1 \rrbracket, \llbracket x_2 \rrbracket, \dots, \llbracket x_n \rrbracket$ );
- 2: 数据中心将 ( $\llbracket b_1 \rrbracket, \llbracket b_2 \rrbracket, \dots, \llbracket b_n \rrbracket$ ) 发送给密钥服务器;
- 3: 密钥服务器解密 ( $\llbracket b_1 \rrbracket, \llbracket b_2 \rrbracket, \dots, \llbracket b_n \rrbracket$ );
- 4: **if**  $b_i = 1$  **then**
- 5:     密钥服务器返回 MaxId =  $i$ ;
- 6: **end if**

**表 4** 不同密文比较算法对比**Table 4** Comparison of different ciphertext comparison algorithms

	Our MaxIdx scheme	BGV <sup>1</sup> scheme <sup>[19]</sup>	BGV <sup>2</sup> scheme <sup>[14]</sup>	TFHE scheme <sup>[30]</sup>	MPC <sup>[12]</sup>
Number of homomorphic operations	$2d + 4t + 2td + 1$	$2n$	$2^n$	$2^{\frac{n}{2}}$	–
Complexity	$O(N)$	$O(N)$	$O(N^2)$	$O(\log N)$	$O(N^2)$
Number of communications	1	$N$	$N$	1	$N$
Ciphertext comparison time (ms)	21.8	74.1	169.1	10.8	–

法的迭代次数,  $t$  表示同态 Comp 算法的迭代次数, 通常  $d < 5, t < 5$  可以保证 16 位精度;  $n$  表示布尔类型密文的位数;  $N$  表示密文序列的长度. 文献 [68] 提出了 BGV 编码 1, 该算法采用整数对编码, 用数对表示浮点数, 例如 1.5 用 (2, 3) 表示, 但无法表示高精度的浮点数, 并且求 MaxId 需要每次比较结束就与服务器通信. 文献 [63] 提出了 BGV 编码 2, 采用的是 IEEE 754 编码, 用 32 位编码表示浮点数, 但求 MaxId 同样需要  $N$  次交互, 并且需要迭代  $N^2$  次密文比较运算. 文献 [70] 改进了 TFHE 按位比较算法, 引入树结构实现并行的布尔比较, 能够实现快速的密文比较, 但 TFHE 算法依赖于布尔电路, 因此需要设置固定精度的加法器和比较器, 实现较为复杂. 安全多方比较同样基于位运算<sup>[61]</sup>, 但是其需要时刻保持通信, 性能急剧下降. 而 CKKS 同态加密能够对浮点数进行加密, 因此其在计算上会优于基于布尔电路的同态加密和多方安全比较. 并且本文的 MaxIdx 方案, 每次计算完最大值只需与服务器进行一次通信. 最后, 表 4 还列出了密文比较的时间, 其中本文的 MaxIdx 和 BGV 编码 2 比较算法, 都是基于 SIMD 技术将密文打包的, 其中每个密文向量中包含 5220 个密文数据, 因此单个密文比较时间是总时间/密文数量. 而 TFHE 无法利用 SIMD 技术, 因此每次比较运算只能针对单一密文数据.

**4.3 密文孤立森林模型**

本文对 CKKS 算法加密后的密文数据, 选用孤立森林算法实现异常检测. 孤立森林算法通过综合每一棵树上各个数据点的异常得分, 从而对点的异常判断做出综合评价, 这充分体现了集成学习的思想. 相较于 FCM, KNN 等基于距离的异常检测算法, 孤立森林具有更好的精度. 与诸如 XGboost 等集成学习方法相比, 孤立森林在精度相差不大的情况下具有更快的检测效果. 能够对大规模高维度的 EHRs 进行异常检测. 但是, 孤立森林在构建森林的过程中, 存在两个随机过程, 本文借助统计学习的方法消除这两个随机过程随机性带来的误差<sup>[75]</sup>.

在孤立森林的建立中, 第一次随机过程出现在数据特征的选取时, 然而并非正常样本和异常样本在每一个特征维度上都具有明显的区分性, 有可能在某些特征上正常样本与异常样本没有区分度. 因

此在未知数据标签的前提下挑选区分度明显的特征将有利于模型精度的提高. 第二个随机过程出现在特征维度上的随机值选取过程中, 这样将使得孤立树的建立过于离散, 如果在每一次选取过程中, 筛选出最可能异常的数, 也有利于精度的提高. 因此, 本文改进孤立森林算法, 简化两个随机过程, 实现高效的异常检测. 对第一个随机过程, 利用 PCA 算法, 将多维数据压缩处理, 降低数据维度的同时, 也将原始数据中重要的特征保留下来, 去除噪声和不重要的特征, 减少高纬度数据的随机性, 也避免数据被第三方数据中心或攻击者恢复出原始数据. 对第二个随机过程, 一个好的分割点应满足其左或右子节点有一个是数量较少的, 这样能使异常点在树上的路径更短, 提高异常点的异常得分, 因此在每一轮的随机分割过程中尽量使左右子树的数量相差较大. 结合密文比较算法, 可以在分割过程循环调用 MaxIdx, 直到满足左右节点数值差阈值. 密文孤立树的建立具体算法如算法 5 所示<sup>[76, 77]</sup>.

---

**算法 5** 同态 iTree 算法 enc.iTree
 

---

**输入:** 密文数据  $[[X_{n \times m}]] = ([[x_1]], [[x_2]], \dots, [[x_n]])$ , 当前树的高度  $e$ , 树的高度限制  $l$ ;

```

1: if  $e \geq l$  or  $|X| \leq 1$  then
2:   return enc_Leaf;
3: else
4:   令  $Q$  为密文数据  $X$  的一个特征列表;
5:   随机选择  $Q$  中的一个特征  $q$ ;
6:   随机选择一个分割比例  $p$ ,  $p$  在  $(0, 0.34)$  和  $(0.67, 1)$  之间;
7:   while  $\frac{|[[X]]|}{|[[X_r]]|} < p$  do
8:      $([[b_{1q}]], [[b_{2q}]], \dots, [[b_{nq}]]) = \text{MaxIdx}([[x_{1q}]], [[x_{2q}]], \dots, [[x_{nq}]])$ ;
9:     发送  $([[b_{1q}]], [[b_{2q}]], \dots, [[b_{nq}]])$  至密钥服务器解密;
10:    从密钥服务器接收  $(b_{1q}, b_{2q}, \dots, b_{nq})$ , 获取最大值的序号  $i$ ;
11:     $[[X]].\text{remove}([[x_i]])$ ;
12:     $[[X]]_r.\text{add}([[x_i]])$ ;
13:     $\text{iTree.SplitValue} \leftarrow [[x_{iq}]]$ ;
14:   end while
15: end if
16:  $[[X]]_l \leftarrow X$ ;
17:  $\text{iTree.SplitAtt} = q$ ;
18:  $\text{iTree.lchild} = \text{enc.iTree}([[X]]_l, e + 1, l)$ ;
19:  $\text{iTree.rchild} = \text{enc.iTree}([[X]]_r, e + 1, l)$ ;

```

**返回:** enc.iTree.

---

综上, 当加密 EHRs 传输到第三方数据中心后, 利用树的分散结构, 可以同时多次与密钥服务器通信获取序号, 这样就能快速同步地训练出密文态下的孤立森林模型. 整个密文孤立森林模型的流程图如图 4 所示. 此外针对训练出来的模型, 需要满足医院以及患者的异常检测要求, 利用 CKKS 同态加密算法的 SIMD 性质, 实现对加密数据的异常检测. 因此该算法结合密文比较算法, 针对每个数据点, 可以基于定义 1 实现其在 iTree 上的任一节点  $N$  的评估. SM 是节点通过比较运算获得的值, 它是或者 0 或者 1 的近似. 最后生成的向量如图 4 所示, 在哪个叶子上对应的叶子节点的向量值为 1.

**定义 1** SM 是孤立树节点与数据  $x$  通过 Comp 运算获得的值, 它总是或者 0 或者 1 的近似, 且以密文形式传递.

(1) 若节点  $L$  是节点  $N$  的左孩子, 则  $\text{SM}(L, x) = \text{SM}(N, x) \times \text{Comp}(N_{\text{SplitValue}}, x_{\text{SplitValue}})$ .

(2) 若节点  $R$  是节点  $N$  的右孩子, 则  $\text{SM}(R, x) = \text{SM}(N, x) \times \text{Comp}(N_{\text{SplitValue}}, x_{\text{SplitValue}})$ .

通过该定义, 就可以获得最终的每棵孤立树的叶子节点的 SM 值, 若该值接近于 1, 则表明该数据

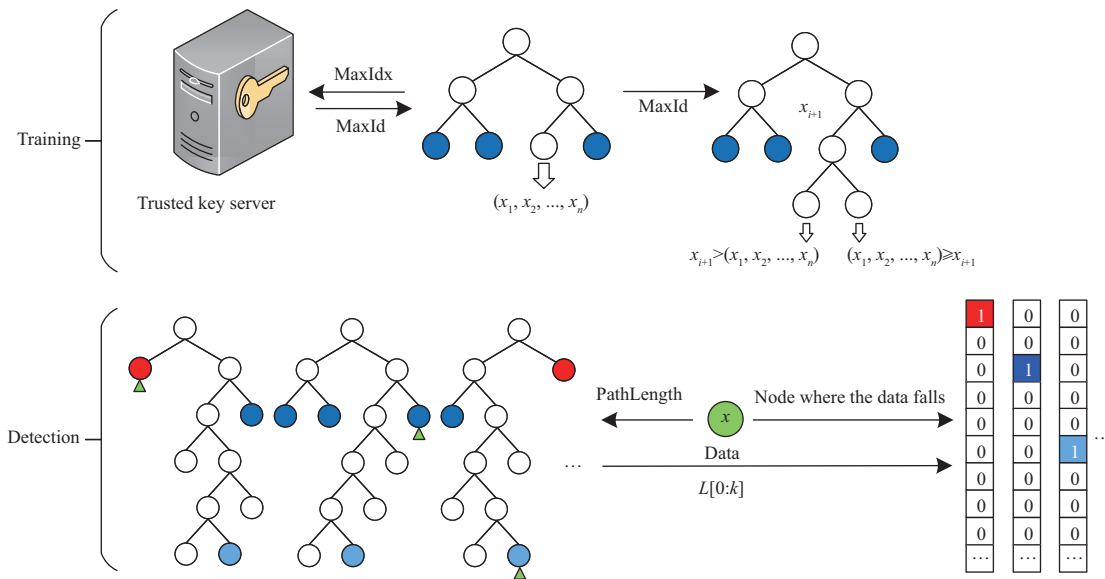


图 4 (网络版彩图) 密文孤立森林模型示意图

Figure 4 (Color online) Schematic diagram of the ciphertext isolation forest model

最后落在孤立树的这片叶子上, 否则为 0. 其中树上的每一片叶子节点记录了当前所处的高度值. 最后计算数据点在森林中每一棵树的高度的平均值就可以计算在孤立森林上的平均高度, 如式 (8) 所示,

$$h(x) = \frac{1}{n} \cdot \sum_n \left( \sum_{L_w} SM(L_m, x) \cdot l_{L_m} \right), \quad (8)$$

其中  $L_m$  表示孤立树上的每一个叶子节点,  $l_L$  表示叶子节点所对应的高度.

利用上述的定义与算法, 就可以实现数据在孤立树上的批量比较运算. 通过并行运算, 可以计算出数据点落在孤立树的具体叶子节点及具体高度, 具体算法如算法 6 所示. 在算法 6 中, 医院和患者利用 CKKS 的 SIMD 性质, 构造  $M^+$  和  $M^-$  两个模型参数加密向量, 每个向量都各自包含着密文在各个分割维度的数值, 以及孤立树上每个节点的 SplitValue 值, 各个数值在向量中的位置不同, 通过同态 Comp 算法和 Rotate 操作, 可以实现对向量中的每一个数值进行比较, 且这个过程不需要与密钥服务器进行通信. 利用 SIMD 可以一次性把所有数值逐个比较, 通过这些操作可以在最后得出密文数据在孤立树上每一个叶子节点的 SM 值, 且这些 SM 值中有且只有一个值为 1, 其表示数据最后落在这个叶子节点上, 其高度即为所求. 所以通过该算法, 可以在孤立森林中的每一棵孤立树定位数据的高度, 在求出平均高度后, 通过式 (9) 就可以将异常得分的计算结果以密文的形式返回给患者或者医院, 后者进行解密后, 对比异常得分的阈值就可以判断该 EHRs 是否存在异常.

$$[s([x], n)] = 2^{-\frac{E([h(x)])}{c(n)}}. \quad (9)$$

此外, 由于在孤立森林算法中每棵树都是随机选取样本进行孤立树建立的, 因此该模型可以并行运行以优化时间复杂度. 此外不同于常见的 K-Means, DBSCAN, FCM 等算法, 该模型不需要涉及有关距离和密度的指标, 因此可以大幅度地提升模型的计算速度, 减少系统的开销. 同时因为孤立森林的孤立树数量和采样数都是固定的, 因此孤立森林的时间复杂度是线性的, 且树的数量越多, 其算法越稳定. 其中关于本文方法的具体时间复杂度如表 5<sup>[68,78]</sup> 所示, 表中的  $N$  表示数据的多少,  $d$  表示密文比较的迭代次数,  $q$  表示密文的模数,  $t$  表示 K-Means 算法迭代的次数,  $k$  表示聚类簇数.

**算法 6** 同态 PathLength 算法

**输入:** 将 enc.iTree 填充为完全二叉树 (无子节点的节点将自身复制至子节点) 高度为  $l$ , 带有  $2^l - 1$  个节点和  $2^{l-1}$  个叶子节点, 密文数据  $[[x]]$ ;

**输出:** 密文  $L$  的前  $2^{l-1}$  个数据为每个叶子节点的  $SM(\cdot)$  值;

- 1: 假设  $N_1, N_2, \dots, N_{2^{l-1}}$  是 enc.iTree 中广度优先遍历访问的所有树节点,  $N_{1_{SV}}$  表示节点的 SplitValue 值,  $N_{1_{SA}}$  表示节点的 SplitAtt 值;
  - 2: 设  $\{[[x]], n\} = \underbrace{[[x]], [[x]], \dots, [[x]]}_{n \text{ 个}};$
  - 3:  $T^+ = \{N_{1_{SV}}, 2^{l-1}\} \{0, 2^{l-1}\} \{N_{2_{SV}}, 2^{l-2}\} \{0, 2^{l-2}\} \{N_{3_{SV}}, 2^{l-2}\} \{0, 2^{l-2}\} \dots \{N_{2^{l-1}_{SV}}, 1\} \{0, 1\};$
  - 4:  $T^- = \{0, 2^{l-1}\} \{N_{1_{SV}}, 2^{l-1}\} \{0, 2^{l-2}\} \{N_{2_{SV}}, 2^{l-2}\} \{0, 2^{l-2}\} \{N_{3_{SV}}, 2^{l-2}\} \dots \{0, 1\} \{N_{2^{l-1}_{SV}}, 1\};$
  - 5:  $V^+ = \{[[x]]_{N_{1_{SA}}, 2^{l-1}}\} \{0, 2^{l-1}\} \{[[x]]_{N_{2_{SA}}, 2^{l-2}}\} \{0, 2^{l-2}\} \{[[x]]_{N_{3_{SA}}, 2^{l-2}}\} \{0, 2^{l-2}\} \dots \{[[x]]_{N_{2^{l-1}_{SA}}, 1}\} \{0, 1\};$
  - 6:  $V^- = \{0, 2^{l-1}\} \{[[x]]_{N_{1_{SA}}, 2^{l-1}}\} \{0, 2^{l-2}\} \{[[x]]_{N_{2_{SA}}, 2^{l-2}}\} \{0, 2^{l-2}\} \{[[x]]_{N_{3_{SA}}, 2^{l-2}}\} \dots \{0, 1\} \{[[x]]_{N_{2^{l-1}_{SA}}, 1}\};$
  - 7:  $M^+ = T^+ + V^+;$
  - 8:  $M^- = T^- + V^-;$
  - 9:  $L = \text{Comp}(M^+, M^-);$
  - 10: **for**  $i = 0, 1, 2, \dots, \log_2 l - 1$  **do**
  - 11:      $L = L \times \text{Rotate}(L, -2^{l+i-1});$
  - 12: **end for**
- 返回:**  $L[0 : 2^{l-1}]$ .

**表 5** 时间复杂度对比  
**Table 5** Time complexity comparison

	Plaintext time complexity	Time complexity of ciphertext calculation (one time)	Total time complexity
Our scheme (training)	$O(1)$	$O(N \log N)$	$O(N^2 \log N)$
Our scheme (detection)	$O(N)$	$O(N)$	$O(N)$
DBSCAN <sup>[68]</sup>	$O(N^2)$	$O(dN \log q)$	$O(dN^2 \log q)$
K-Means <sup>[78]</sup>	$O(Nkt)$	$O(dN \log q)$	$O(N^3)$

**4.4 安全性分析**

本小节将对本文模型的安全性进行分析. 在本文的模型中, 将第三方数据中心设定为诚实且好奇的半诚实模型. 在选取的密码算法中, CKKS 同态加密已经被证明满足选择明文攻击下的不可区分性 (indistinguishability under chosen plaintext attack, IND-CPA)<sup>[79]</sup>, IND-CPA 等价于语义安全, 保证了 CKKS 算法更强的安全性. 此外, 针对 CKKS 算法的密文  $ct$  有

$$\begin{aligned}
 ct &= (c_0, c_1) \\
 &= (-v \cdot b + m + e_0, -v \cdot a + e_1) \bmod 2^L.
 \end{aligned}
 \tag{10}$$

对其进行解密算法后, 得到的明文消息  $\hat{m}$  是

$$\begin{aligned}
 \hat{m} &= (c_0 + c_1 \cdot sk) \\
 &= [-v \cdot (-a \cdot s) - v \cdot e + m + e_0] + (-v \cdot a \cdot s + e_1 \cdot s) \bmod 2^L \\
 &= (m - v \cdot e + e_0 + e_1 \cdot s) \bmod 2^L.
 \end{aligned}
 \tag{11}$$

由上式可知, 当误差  $e, e_0, e_1$  足够小时,  $\hat{m} \approx m$ , 因此即使 CKKS 解密后, 也无法恢复原始数据. 该模

型希望达到的安全目标是第三方数据中心和具有窃听能力的攻击者无法获取医院和患者的 EHRs 信息, 也无法通过同态孤立森林算法倒推出原始数据. 在以上设定的前提下, 下文将给出定理 1~3 并进行证明, 以说明本文方案的安全性.

**定理1** 第三方数据中心或具备窃听能力的敌手根据他们所获取的密文信息成功恢复用户原始数据的概率是可以忽略的.

**证明** 第三方数据中心或具备窃听能力的敌手在构建孤立森林模型时, 能够获得的数据包括密文  $X_{n \times m} = ([x_1], [x_2], \dots, [x_n])$ ,  $B = ([b_1], [b_2], \dots, [b_n])$  以及明文 MaxId; 从数据的组成来看, 除了 MaxId 之外, 其他数据都是通过 CKKS 同态加密所得到的密文, 第三方数据中心和敌手在不知道私钥的情况下成功获取医院的原始数据的概率等同于攻破 CKKS 算法. 由同态加密算法的语义安全可知, 攻破 CKKS 算法的概率是可以忽略的. 第三方数据中心和敌手能够获得的唯一明文数据是 MaxId, 而此明文信息只是表示最大值的序列号, 而数据在传送到第三方数据中心时, 已经经过 PCA 算法的处理, 因此该数据并不具有特殊的属性含义, 并不会泄漏原始数据; 而在异常检测的过程中, 能够获得的数据有密文  $X_{n \times m} = ([x_1], [x_2], \dots, [x_n])$  和  $[s]$ , 同理这些数据都是通过 CKKS 同态加密所得到的密文, 第三方数据中心和敌手在不知道私钥的情况下成功获取患者的原始数据和异常检测结果的概率等同于攻破 CKKS 算法.

**定理2** 数据中心或具备窃听能力的敌手获取同态孤立森林算法后, 能够成功获得用户原始数据的概率是可忽略的.

**证明** 数据中心或具有窃听能力的敌手可能通过窃听或攻击等方式获得在数据中心中执行的同态孤立森林算法的特定内容. 但由于孤立森林算法是一种无监督的机器学习算法, 因此不需要在数据中心中通过设置训练集来预先训练, 它可以直接对数据进行分割操作然后进行预测, 即可得到结果. 本文方案保留了孤立森林算法的这一特点, 直接作用于密文数据, 得到密文数据分割的结果, 并且不是固定的模型, 每个模型都是根据数据随机生成的, 每次生成的也都是不同的, 生成的每个节点孤立树都采用 CKKS 加密, 模型信息无法泄漏. 此外, 该算法基于密文数据大小的比较, 每一轮的采样和选择都是从多个维度中随机指定一个维度, 在当前数据中随机产生一个切割比例  $p$ , 随机循环多次 MaxIdx 操作,

$$([b_{1q}], [b_{2q}], \dots, [b_{nq}]) = \text{MaxIdx}([x_{1q}], [x_{2q}], \dots, [x_{nq}]). \quad (12)$$

同时, 由于 CKKS 算法同态运算的特性, 在每次乘法运算后, 都将执行重缩放操作, 如下所示:

$$ct' = [2^{-p} \cdot (c_0, c_1)] \pmod{2^{l-p}}. \quad (13)$$

这将使密文  $ct$  转化为一个新的密文  $ct'$ , 破解难度加大. 并且 CKKS 算法的同态运算具有近似结果的性质, 相当于在数据中加入噪声, 与真实数据存在偏差. 因此, 即使敌手获得算法过程, 也无法通过算法内容推断出原始数据信息.

**定理3** 假扮患者拥有私钥的敌手窃听获取医院传递的数据集后, 能够成功恢复医院原始数据的概率是可忽略的.

**证明** 敌手可以假扮成患者获取私钥, 然后窃听医院传递的数据进行窃听. 在场景 1 中, 密钥服务器仅向医院发送私钥, 在这个阶段敌手无法假扮患者获取医院的私钥. 因此敌手窃听医院原始数据的概率如定理 1 的证明所述, 相当于攻破 CKKS 算法的概率. 在场景 2 中敌手可以假扮患者获取密钥服务器分发的私钥. 但是在该场景中, 医院只通过公钥加密传递在场景 1 中所生成的加密模型, 并不传递原始数据. 由于孤立森林算法的特性, 孤立树上的每个节点的数据都是数据集中的随机数据的随

表 6 数据集属性  
Table 6 Dataset properties

Dataset	Number of instances	Number of dimensions	Abnormal proportion (%)	Number of dimensions after PCA
Vertebral <sup>[80]</sup>	256	6	12.5	5
Arrhythmia <sup>[45]</sup>	452	274	15	128
Breast Cancer <sup>[45]</sup>	683	9	35	6
Cardiotocography <sup>[81]</sup>	1831	21	9.6	15
Annthyroid <sup>[82]</sup>	7200	6	7.42	4
Mammography <sup>[83]</sup>	11183	6	2.32	3

机维度选取的, 并且节点不会记录随机数据的信息. 因此孤立树上的每一个并不会透露出是哪一个患者的哪一个 EHRs 数据. 由于 CKKS 满足 IND-CPA, 能够满足隐私数据的机密性, 同时 CKKS 算法的同态运算具有近似结果的性质, 相当于在数据中加入噪声点. 因此, 即使敌手利用私钥窃听获取了整个明文模型数据, 也无法通过模型内容推断出医院原始数据集的信息.

综上, 由于 EHRs 在传输到第三方数据中心前经过 PCA 算法处理, 数据难以解释原始特征, 此外孤立森林模型异常检测的过程中数据都以密文形式传递与运算, 且 CKKS 算法满足语义安全, 因此攻击者和第三方数据中心难以获取医院或患者的原始数据, 本文的模型经证明是安全的.

## 5 实验与分析

本文实验使用的配置为 Intel(R) Core(TM) i7-10700 CPU @ 2.90 GHz、16 GB 内存, 借助基于 python 开发的同态加密库 tenseal 实现对数据集的加密和同态比较运算. 在实验中选择异常检测中常用的医疗数据集, 其中包括脊椎 (Vertebral)、心律失常 (Arrhythmia)、乳腺癌 (Breast Cancer)、心电图 (Cardiotocography)、甲状腺疾病 (Annthyroid) 和乳房 x 光 (Mammography) 数据集<sup>[45, 80~83]</sup>. 表 6 给出了这些数据集对应的数量、维度以及异常比例. 本文将通过比较明文数据的异常检测效果和密文数据上的异常检测效果来验证该模型的准确性. 并且与现阶段表现最好的密文异常检测算法 BGV-FCM<sup>[54]</sup> 进行对比, 该异常检测算法采用 BGV 同态加密算法和模糊 C 均值聚类算法 (fuzzy C-means, FCM) 对密文数据进行异常检测. 实验中明文数据集和密文数据集在执行孤立森林算法时, 所用的参数  $t$  和  $\psi$  一致, 其中  $t = 100$  表示孤立森林中孤立树的数量,  $\psi = 256$  表示子采样的数量.

### 5.1 数据预处理

在本文提出的模型中, 首先需要由医院和患者在本地将数据进行 PCA 处理, 这样降低数据维度的同时, 也可以将原始数据中重要的特征保留下来, 去除掉噪声和不重要的特征, 并且避免数据特征泄漏. 因此, 在数据预处理中先对各个数据集进行 PCA 处理, 选取能够较好地表示原始数据集的多个维度, 这些数据集降低维度后的数量如表 6 所示. 再将处理后的数据进行 Min-Max 标准化处理, 使数据集中在能够执行密文比较运算的  $\frac{1}{2}$  与  $\frac{3}{2}$  之间.

然后, 利用同态加密库将 PCA 处理过的数据进行 CKKS 同态加密, 在本实验中选择的 CKKS 同态加密参数为 poly\_modulus\_degree (多项式模数) = 16384, coeff\_mod\_bit\_sizes (系数模数) = [60, 40, 40, 40, 40, 40, 40, 40, 60], global\_scale (缩放因子) =  $2^{40}$ , 经过测试在该参数下能够保证较小的比较运算误差以及更高的计算效率.



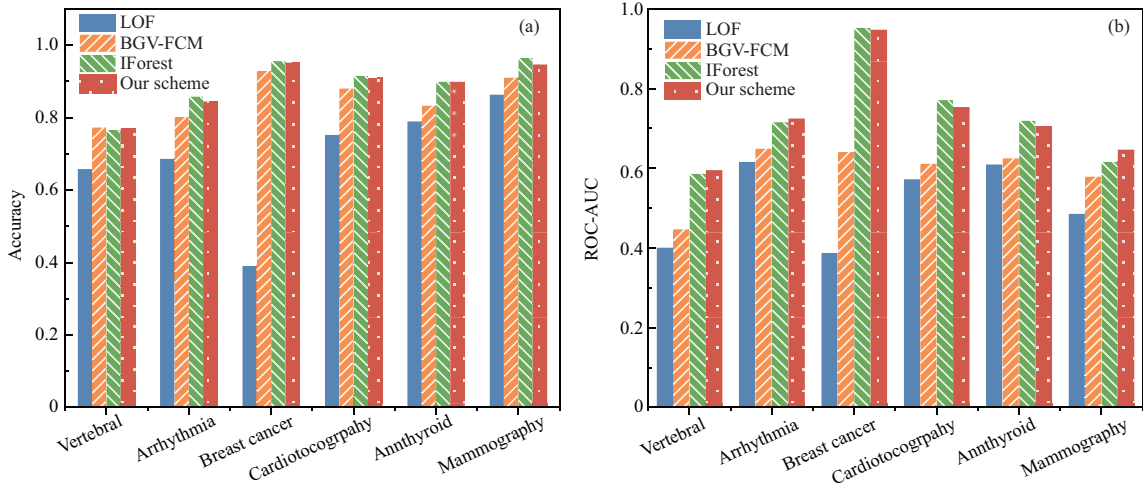


图 5 (网络版彩图) 4 种异常检测模型在不同数据集上的检测效果. (a) Accuracy; (b) ROC-AUC  
 Figure 5 (Color online) Detection effects of four anomaly detection models in different datasets. (a) Accuracy; (b) ROC-AUC

### 5.2 性能比较

实验阶段, 本文在同样的数据集中采取 LOF (local outlier factor)、孤立森林算法、密文 BGV-FCM 算法以及本文模型 4 种算法进行异常检测的对比实验. 其中明文的 LOF 算法是传统的基于密度的异常检测算法, 密文 BGV-FCM 采用的是基于聚类的 FCM 异常检测算法. 在本实验中, 主要验证孤立森林模型能够高效地处理 EHRs 的异常检测, 并且本文涉及的密文孤立森林模型能够在检测效果上与明文的孤立森林模型相当. 因此, 在对比实验中, 本小节将结合异常检测的准确度 (accuracy) 以及 ROC-AUC (area under the ROC curve) 值来进行比较. 其中 Accuracy 表示异常检测模型的准确性, ROC-AUC 表示 ROC 曲线下的面积, 该值经常用于二分类模型或异常检测模型, 由 TPR 和 FPR 两个参数共同决定了 ROC 曲线, 其中 TPR 和 FPR 如式 (14) 所示:

$$TPR = \frac{TP}{TP + FN}, \quad FPR = \frac{FP}{FP + TN}. \quad (14)$$

其中 TP 表示真阳性, FP 表示假阳性, TN 表示真阴性, FN 表示假阴性. ROC 曲线就是以 FPR 为横坐标, TPR 为纵坐标绘制的曲线. ROC-AUC 适用于每个类别之间的观察平衡, AUC 可以用于清楚地描述该异常检测模型的正确性, AUC 越接近 1 表明模型越准确, AUC 越接近 0.5 表明模型越倾向于随机分类. 由于孤立森林算法的随机采样性, 因此孤立森林算法和本文模型都采取多次训练和预测, 并求平均值.

实验首先验证了本文模型的正确性和有效性, 详细实验结果由图 5(a) 和 (b) 所示. 在 Breast Cancer, Cardiocography, Annthyroid 和 Mammography 数据集上, 本文模型分别可以达到 0.956, 0.916, 0.899 和 0.964 的 Accuracy, 在 ROC-AUC 指标上, 也可以达到 0.947, 0.754, 0.726 和 0.647. 甚至有些超过了明文的孤立森林算法效果. 此外, 在表现优异的数据集上还画出了 ROC 曲线图, 如图 6 所示, 相较于明文模型的 LOF 和密文模型 BGV-FCM, 本文模型都有较好的 ROC 曲线, 并且与孤立森林模型相近. 表明了密文状态下本文模型与明文下的孤立森林算法具有一定的可比性, 同时也证明了该模型的正确性, 由于本文模型能够在密文的条件下, 保证孤立森林算法的结构, 因此可以达到和明文相同的异常检测效果. 并且该模型和孤立森林算法在所有的数据集上都有比传统的 LOF 算

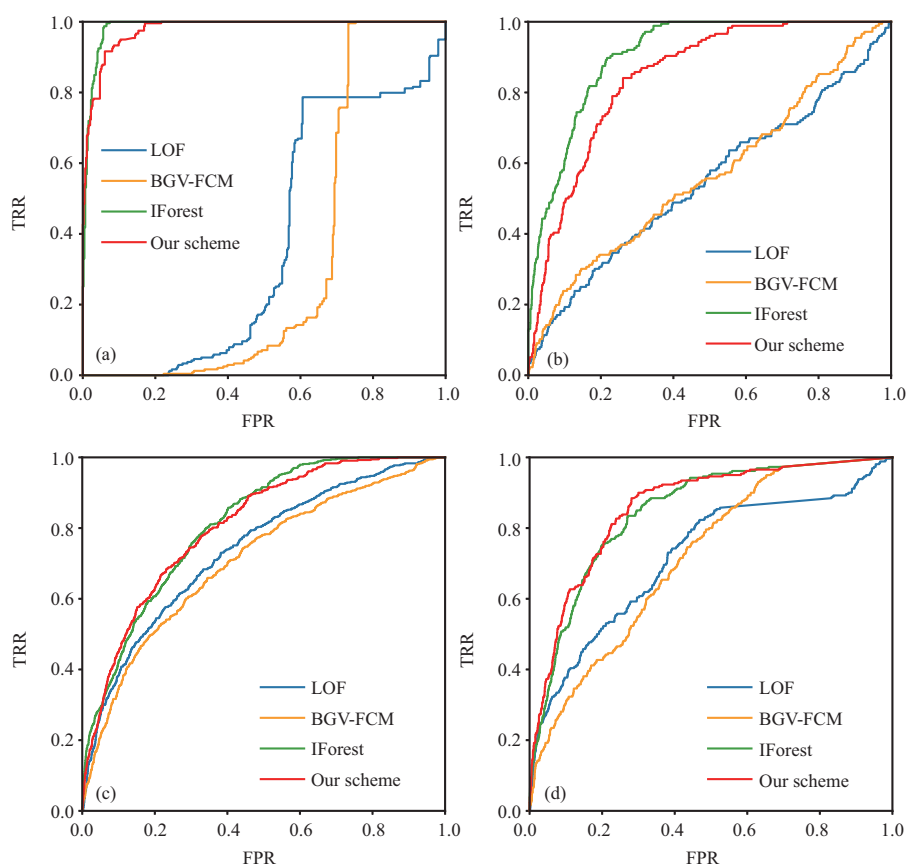


图 6 (网络版彩图) 4 种模型在不同数据集下的 ROC 曲线图, 本文模型与明文孤立森林算法有相近的检测性能, 且都优于明文的 LOF 和密文的 BGV-FCM 算法. (a) Breast Cancer; (b) Cardiotocography; (c) Anthyroid; (d) Mammography

Figure 6 (Color online) The ROC curves of the 4 models under different datasets, the model in this paper has similar detection performance to the plaintext isolation forest algorithm, and are better than the plaintext LOF and ciphertext BGV-FCM algorithm. (a) Breast Cancer; (b) Cardiotocography; (c) Anthyroid; (d) Mammography

法和 BGV-FCM 算法更好的异常检测效果, 这也验证了该模型的有效性并且超过了现有的密文异常检测方案. 可以看出该模型能够确保数据隐私安全的同时, 保证异常检测的有效进行.

此外, 针对在异常检测效果较好的 Breast Cancer 数据集, 本小节还针对该数据集实现明文下孤立森林算法和本文模型的异常检测可视化对比. 在数据预处理过程中, 将 Breast Cancer 数据集降维至 6 维, 其中数据的异常三维分布如图 7 所示. 针对这些数据进行异常检测后, 于维度 1 与维度 2、维度 3 与维度 4、维度 5 与维度 6 进行展示与效果分析, 绘制出异常数据和正常数据的二维分布的图像. 异常检测结果如图 8 和 9 所示. 通过对比可以看出, 两个模型方法实现的聚类效果大体一致, 都能够较好地将异常数据和正常数据进行聚类, 通过聚类效果的示意图可以看出, 两个模型只有个别数据的偏差, 因此可见明文孤立森林算法与本文模型在对数据的异常检测效果几乎相同, 并且从图的数据点分布可以看出, 该模型与明文孤立森林算法相同, 能够隔离出异于正常数据的异常点, 从而能够较好地划分出异常数据, 并且达到较好的效率, 但性能并没有下降.

最后, 在时间开销方面, 文献 [54] 中的方案利用 BGV 同态加密与 FCM 算法结合的方式对密文数据进行异常检测. 但在该方案中对数据进行逐比特加密, 导致同态计算时间较长, 对 304 条的密态

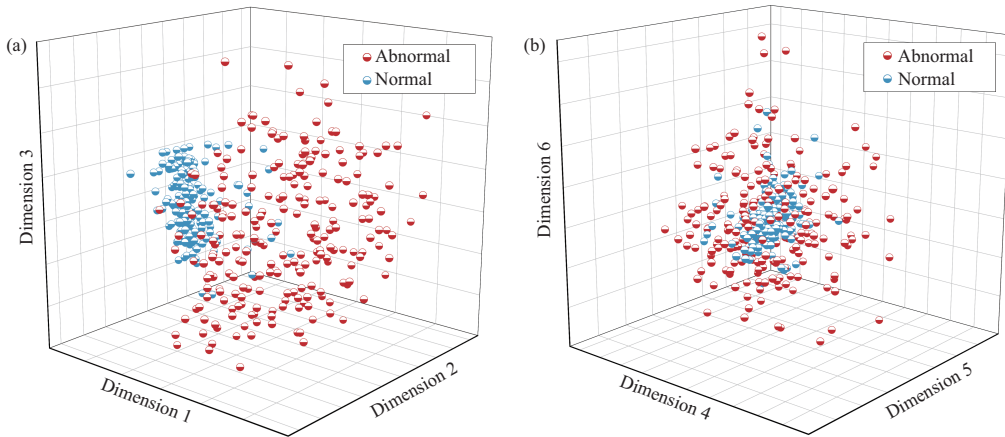


图 7 (网络版彩图) Breast Cancer 数据集经过 PCA 处理后不同维度下的正常和异常数据分布. (a) 异常分布 -1; (b) 异常分布 -2

Figure 7 (Color online) The normal and anomalous data distributions in different dimensions of the Breast Cancer dataset after PCA processing. (a) Anomaly distribution-1; (b) anomaly distribution-2

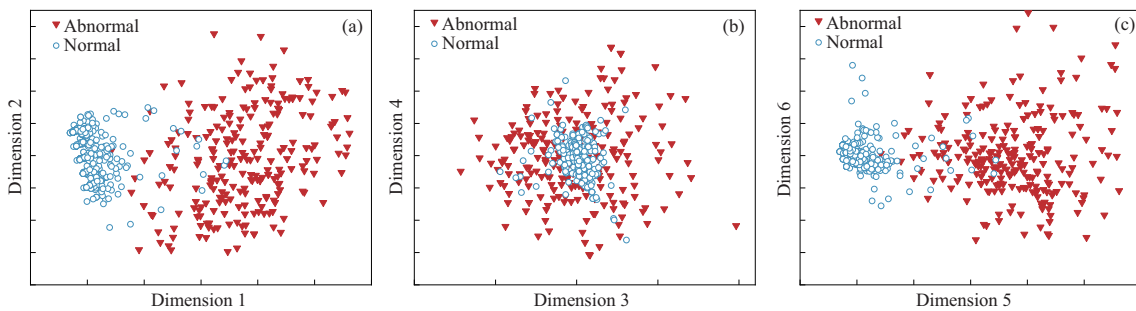


图 8 (网络版彩图) 明文孤立森林算法的异常检测效果. (a) 维度 1 和维度 2 下的异常检测效果; (b) 维度 3 和维度 4 下的异常检测效果; (c) 维度 5 和维度 6 下的异常检测效果

Figure 8 (Color online) The anomaly detection effect of plaintext isolation forest algorithm. (a) Anomaly detection effect under dimension 1 and dimension 2; (b) anomaly detection effect under dimension 3 and dimension 4; (c) anomaly detection effect under dimension 5 and dimension 6

75 维数据在单设备上的处理时间长达 952 s. 而本文的方案时间开销较低, 在处理数据集 Arrhythmia 时, 对 452 条的密态 128 维数据构造密文孤立森林模型, 在该算法中生成一棵采样 256 个数据的同态孤立树仅需 312 s. 在数据中心中利用并行化处理就可以在 312 s 内构建 100 棵同态孤立树. 并且 BGV-FCM 算法依托在云中的 Mapreduce 技术实现高性能并行计算, 而本模型只需依赖孤立森林的树独立性, 就可以在任意设备实现并行处理. 此外, 密文数据在孤立森林中进行检测时, 利用 CKKS 的 SIMD 性质也可以在 1 s 内完成在同态孤立树中的路径定位, 找到数据在孤立树中的高度.

综上, 经过实验与分析, 本文提出的模型在密文状态下与明文孤立森林算法具有几乎相同的异常检测效果, 在针对异常与正常数据分布不平衡的数据集上也能够有良好的表现. 并且该模型与孤立森林算法一样能够适用于分布式运行, 降低运行时间. 因此, 以上的实验均表明, 本文模型能够在密文下有效且准确地检测出异常值, 保障敏感数据的隐私安全.

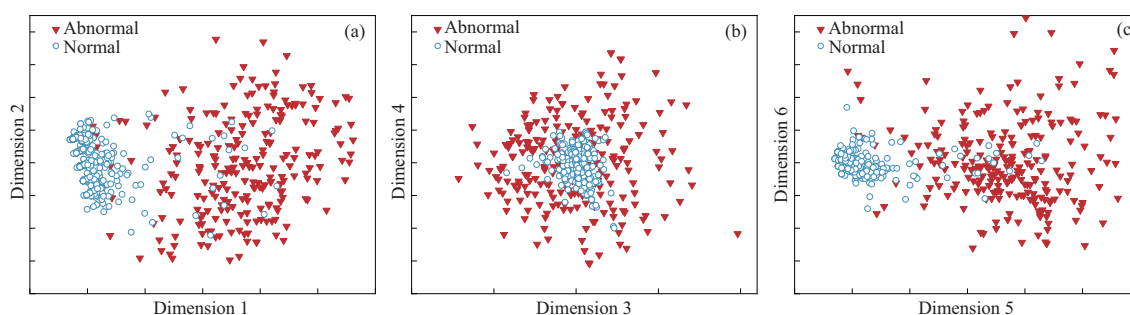


图 9 (网络版彩图) 本文模型的异常检测效果. (a) 维度 1 和维度 2 下的异常检测效果; (b) 维度 3 和维度 4 下的异常检测效果; (c) 维度 5 和维度 6 下的异常检测效果

**Figure 9** (Color online) The anomaly detection effect of the model in this paper. (a) Anomaly detection effect under dimension 1 and dimension 2; (b) anomaly detection effect under dimension 3 and dimension 4; (c) anomaly detection effect under dimension 5 and dimension 6

## 6 结束语

为了解决当前 EHRs 存在的数据泄漏和隐私保护问题, 本文提出了一种基于 CKKS 同态加密的 EHRs 异常检测的隐私保护模型. 该模型能够确保医院以及患者在进行 EHRs 数据的异常检测时, 避免数据泄漏. 针对 EHRs 的使用场景, 设计了由多方共同构建孤立森林的模型, 且该模型支持密文数据在孤立森林上进行节点比较, 并支持对密文数据进行异常检测. 通过实验证明本文的密文异常检测模型优于现阶段最好的密文检测方法, 且具有和明文孤立森林算法相同的异常检测效果, 在保证良好的异常检测效果的同时, 又能够提供可靠的数据安全性和隐私保护. 未来的研究将包括, 改进该模型, 减少通信次数, 以及提高同态计算的能力, 让其能够参与更广泛的数据挖掘与机器学习方法当中, 并提高检测效率.

## 参考文献

- 1 Su Y, Li Y, Zhang K, et al. A privacy-preserving public integrity check scheme for outsourced EHRs. *Inf Sci*, 2021, 542: 112–130
- 2 Zhang M W, Huang J J, Han L. Range-based multi-keyword searchable scheme with privacy protection in e-healthcare cloud systems. *J Software*, 2021, 32: 3266–3282 [张明武, 黄嘉骏, 韩亮. 医疗大数据隐私保护多关键词范围搜索方案. *软件学报*, 2021, 32: 3266–3282]
- 3 Carпов S, Nguyen T H, Sirdey R, et al. Practical privacy-preserving medical diagnosis using homomorphic encryption. In: *Proceedings of 2016 IEEE 9th International Conference on Cloud Computing (CLOUD)*, 2016. 593–599
- 4 Wang N W, Liu H Z, Xu C. Deep learning for the detection of COVID-19 using transfer learning and model integration. In: *Proceedings of 2020 IEEE 10th International Conference on Electronics Information and Emergency Communication (ICEIEC)*, 2020. 281–284
- 5 Pang G S, Shen C H, Cao L B, et al. Deep learning for anomaly detection. *ACM Comput Surv*, 2021, 54: 1–38
- 6 Nguyen T D, Marchal S, Miettinen M, et al. DioT: a federated self-learning anomaly detection system for IOT. In: *Proceedings of IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, 2019. 756–767
- 7 Du M, Jia R X, Song D. Robust anomaly detection and backdoor attack detection via differential privacy. 2019. ArXiv:191107116
- 8 Ni T G, Zhu J Q, Qu J, et al. Labeling privacy protection SVM using privileged information for COVID-19 diagnosis. *ACM Trans Internet Technol*, 2022, 22: 1–21
- 9 Gentry C. A fully homomorphic encryption scheme. Dissertation for Ph.D. Degree. Stanford: Stanford University, 2009

- 10 Pulido-Gaytan B, Tchernykh A, Cortés-Mendoza J M, et al. Privacy-preserving neural networks with Homomorphic encryption: challenges and opportunities. *Peer-to-Peer Netw Appl*, 2021, 14: 1666–1691
- 11 Lu G H, Duan C H, Zhou G H, et al. Privacy-preserving outlier detection with high efficiency over distributed datasets. In: *Proceedings of IEEE 40th International Conference on Computer Communications*, 2021. 1–10
- 12 Li D, Liao X F, Xiang T, et al. Privacy-preserving self-serviced medical diagnosis scheme based on secure multi-party computation. *Comput Security*, 2020, 90: 101701
- 13 Ren W, Tong X, Du J, et al. Privacy-preserving using homomorphic encryption in Mobile IoT systems. *Comput Commun*, 2021, 165: 105–111
- 14 Yonetani R, Boddeti V N, Kitani K M, et al. Privacy-preserving visual learning using doubly permuted homomorphic encryption. In: *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2017. 2059–2069
- 15 Giacomelli I, Jha S, Joye M, et al. Privacy-preserving ridge regression with only linearly-homomorphic encryption. In: *Proceedings of Applied Cryptography and Network Security (ACNS 2018)*, 2018. 243–261
- 16 Wang Y, He M X. CPDS: a cross-blockchain based privacy-preserving data sharing for electronic health records. In: *Proceedings of IEEE 6th International Conference on Cloud Computing and Big Data Analytics (ICCCBDA)*, 2021. 90–99
- 17 Ma Z R, Ma J F, Miao Y B, et al. Privacy-preserving and high-accurate outsourced disease predictor on random forest. *Inf Sci*, 2019, 496: 225–241
- 18 Liu G, Yan Z, Feng W, et al. SeDID: an SGX-enabled decentralized intrusion detection framework for network trust evaluation. *Inf Fusion*, 2021, 70: 100–114
- 19 Liu L, Chen R M, Liu X M, et al. Towards practical privacy-preserving decision tree training and evaluation in the cloud. *IEEE Trans Inform Forensic Secur*, 2020, 15: 2914–2929
- 20 Rahman M S, Khalil I, Atiquzzaman M, et al. Towards privacy preserving AI based composition framework in edge networks using fully homomorphic encryption. *Eng Appl Artif Intelligence*, 2020, 94: 103737
- 21 Kanwal T, Anjum A, Malik S U R, et al. A robust privacy preserving approach for electronic health records using multiple dataset with multiple sensitive attributes. *Comput Secur*, 2021, 105: 102224
- 22 Jiang B, Li J Q, Yue G H, et al. Differential privacy for industrial Internet of Things: opportunities, applications, and challenges. *IEEE Internet Things J*, 2021, 8: 10430–10451
- 23 Jia B, Zhang X S, Liu J W, et al. Blockchain-enabled federated learning data protection aggregation scheme with differential privacy and homomorphic encryption in IIoT. *IEEE Trans Ind Inf*, 2022, 18: 4049–4058
- 24 Fan W B, He J, Guo M J, et al. Privacy preserving classification on local differential privacy in data centers. *J Parallel Distrib Comput*, 2020, 135: 70–82
- 25 Ma J, Naas S-A, Sigg S, et al. Privacy-preserving federated learning based on multi-key homomorphic encryption. 2021. [ArXiv:210406824](https://arxiv.org/abs/210406824)
- 26 Ibrahim A, Mahmood B, Singhal M. A secure framework for sharing electronic health records over clouds. In: *Proceedings of 2016 IEEE International Conference on Serious Games and Applications for Health (SeGAH)*, 2016. 1–8
- 27 Lu Y, Zhu M H. Privacy preserving distributed optimization using homomorphic encryption. *Automatica*, 2018, 96: 314–325
- 28 Pang H P, Wang B C. Privacy-preserving association rule mining using homomorphic encryption in a multikey environment. *IEEE Syst J*, 2021, 15: 3131–3141
- 29 Reich D, Todoki A, Dowsley R, et al. Privacy-preserving classification of personal text messages with secure multi-party computation: an application to hate-speech detection. 2021. [ArXiv:190602325](https://arxiv.org/abs/190602325)
- 30 Zhang C K, Liu H D, Li Y. Time series discord discovery under multi-party privacy preserving. In: *Proceedings of IEEE Second International Conference on Data Science in Cyberspace (DSC)*, 2017. 467–474
- 31 Ducas L, Micciancio D. FHEW: bootstrapping homomorphic encryption in less than a second. In: *Proceedings of Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT 2015)*, 2015. 617–640
- 32 Brakerski Z, Gentry C, Vaikuntanathan V. (Leveled) fully homomorphic encryption without bootstrapping. *ACM Trans Comput Theor*, 2014, 6: 1–36
- 33 Phong L T, Aono Y, Hayashi T, et al. Privacy-preserving deep learning via additively homomorphic encryption.

- IEEE Trans Inform Forensic Secur, 2018, 13: 1333–1345
- 34 Fan J, Vercauteren F. Somewhat practical fully homomorphic encryption. Cryptology ePrint Archive, 2012. <https://eprint.iacr.org/2012/144>
- 35 Chillotti I, Gama N, Georgieva M, et al. TFHE: fast fully homomorphic encryption over the torus. J Cryptol, 2020, 33: 34–91
- 36 Lou Q, Feng B, Fox G C, et al. Glyph: fast and accurately training deep neural networks on encrypted data. 2020. ArXiv:191107101
- 37 Lou Q, Jiang L. HEMET: a homomorphic-encryption-friendly privacy-preserving mobile neural network architecture. In: Proceedings of International Conference on Machine Learning (PMLR), 2021. 7102–7110
- 38 Wood A, Najarian K, Kahrobaei D. Homomorphic encryption for machine learning in medicine and bioinformatics. ACM Comput Surv, 2021, 53: 1–35
- 39 Lu W J, Huang Z C, Hong C, et al. PEGASUS: bridging polynomial and non-polynomial evaluations in homomorphic encryption. In: Proceedings of IEEE Symposium on Security and Privacy (SP), 2021. 1057–1073
- 40 Iezzi M. Practical privacy-preserving data science with homomorphic encryption: an overview. In: Proceedings of 2020 IEEE International Conference on Big Data (Big Data), 2020. 3979–3988
- 41 Li J, Kuang X H, Lin S J, et al. Privacy preservation for machine learning training and classification based on homomorphic encryption schemes. Inf Sci, 2020, 526: 166–179
- 42 Fang H K, Qian Q. Privacy preserving machine learning with homomorphic encryption and federated learning. Future Internet, 2021, 13: 94
- 43 Takabi H, Hesamifard E, Ghasemi M. Privacy preserving multi-party machine learning with homomorphic encryption. In: Proceedings of the 29th Annual Conference on Neural Information Processing Systems (NIPS), 2016
- 44 Xiao X D, Wu T, Chen Y F, et al. Privacy-preserved approximate classification based on homomorphic encryption. Math Comput Appl, 2019, 24: 92
- 45 Liu F T, Ting K M, Zhou Z H. Isolation forest. In: Proceedings of the 8th IEEE International Conference on Data Mining, 2008. 413–422
- 46 Al-Emran M, Shaalan K, Hassanien A E, et al. Recent Advances in Intelligent Systems and Smart Applications. Berlin: Springer, 2021
- 47 Hou J, Li Q M, Meng S M, et al. DPRF: a differential privacy protection random forest. IEEE Access, 2019, 7: 130707–130720
- 48 Alabdulatif A, Khalil I, Yi X, et al. Secure edge of things for smart healthcare surveillance framework. IEEE Access, 2019, 7: 31010–31021
- 49 Cheon J H, Kim D, Kim Y, et al. Ensemble method for privacy-preserving logistic regression based on homomorphic encryption. IEEE Access, 2018, 6: 46938–46948
- 50 Meftah S, Tan B H M, Mun C F, et al. DOReN: toward efficient deep convolutional neural networks with fully homomorphic encryption. IEEE Trans Inform Forensic Secur, 2021, 16: 3740–3752
- 51 Kim S, Kim J, Koo D, et al. Efficient privacy-preserving matrix factorization via fully homomorphic encryption: extended abstract. In: Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security, 2016. 617–628
- 52 Hong S, Kim S, Choi J, et al. Efficient sorting of homomorphic encrypted data with  $k$ -way sorting network. IEEE Trans Inform Forensic Secur, 2021, 16: 4389–4404
- 53 Alabdulatif A, Kumarage H, Khalil I, et al. Privacy-preserving anomaly detection in cloud with lightweight homomorphic encryption. J Comput Syst Sci, 2017, 90: 28–45
- 54 Alabdulatif A, Khalil I, Kumarage H, et al. Privacy-preserving anomaly detection in the cloud for quality assured decision-making in smart cities. J Parallel Distrib Comput, 2019, 127: 209–223
- 55 Alabdulatif A, Khalil I, Zomaya A Y, et al. Fully homomorphic based privacy-preserving distributed expectation maximization on cloud. IEEE Trans Parallel Distrib Syst, 2020, 31: 2668–2681
- 56 Alabdulatif A, Khalil I, Yi X. Towards secure big data analytic for cloud-enabled applications with fully homomorphic encryption. J Parallel Distrib Comput, 2020, 137: 192–204
- 57 Xu X W, Cai B, Xiang H, et al. Research and implementation of secure multinomial classification logistic regression model based on homomorphic encryption. J Cryptol Res, 2019, 7: 179–186 [许心炜, 蔡斌, 向宏, 等. 基于同态加密

- 的多分类 Logistic 回归模型. 密码学报, 2019, 7: 179–186]
- 58 Lv Y, Wu W Y. Linear system solving scheme based on homo-morphic encryption. *Comput Sci*, 2022, 49: 338–345 [吕由, 吴文渊. 基于同态加密的线性系统求解方案. *计算机科学*, 2022, 49: 338–345]
- 59 Lee J W, Kang H, Lee Y, et al. Privacy-preserving machine learning with fully homomorphic encryption for deep neural network. *IEEE Access*, 2022, 10: 30039–30054
- 60 Froelicher D, Troncoso-Pastoriza J R, Sousa J S, et al. Drynx: decentralized, secure, verifiable system for statistical queries and machine learning on distributed datasets. *IEEE Trans Inform Forensic Secur*, 2020, 15: 3035–3050
- 61 Ichikawa A, Ogata W, Hamada K, et al. Efficient secure multi-party protocols for decision tree classification. In: *Proceedings of Australasian Conference on Information Security and Privacy*, 2019. 362–380
- 62 Itokazu K, Wang L H, Ozawa S. Outlier detection by privacy-preserving ensemble decision tree using homomorphic encryption. In: *Proceedings of International Joint Conference on Neural Networks (IJCNN)*, 2021. 1–7
- 63 Aloufi A, Hu P, Wong H W H, et al. Blindfolded evaluation of random forests with multi-key homomorphic encryption. *IEEE Trans Dependable Secure Comput*, 2021, 18: 1821–1835
- 64 Huynh D. Cryptotree: fast and accurate predictions on encrypted structured data. 2020. ArXiv:200608299
- 65 Cheon J H, Kim A, Kim M, et al. Homomorphic encryption for arithmetic of approximate numbers. In: *Proceedings of International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT 2017)*, 2017. 409–437
- 66 Cheon J H, Kim D, Kim D, et al. Numerical method for comparison on homomorphically encrypted numbers. In: *Proceedings of International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT 2019)*, 2019. 415–445
- 67 Cheon J H, Kim D, Kim D. Efficient homomorphic comparison methods with optimal complexity. In: *Proceedings of International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT 2020)*, 2020. 221–256
- 68 Jia C F, Li R Q, Wang Y F, et al. Privacy protection scheme of DBSCAN clustering based on homomorphic encryption. *J Commun*, 2021, 42: 1–11 [贾春福, 李瑞琪, 王雅飞. 基于同态加密的 DBSCAN 聚类隐私保护方案. *通信学报*, 2021, 42: 1–11]
- 69 Bourse F, Olivier S, Jacques T. Improved secure integer comparison via homomorphic encryption. In: *Proceedings of 2020 Cryptographers' Track at the RSA Conference*, 2020. 391–416
- 70 Chakraborty O, Zuber M. Efficient and accurate homomorphic comparisons. In: *Proceedings of the 10th Workshop on Encrypted Computing & Applied Homomorphic Cryptography*, 2022. 35–46
- 71 Sim J J, Chan F M, Chen S, et al. Achieving GWAS with homomorphic encryption. *BMC Med Genomics*, 2020, 13: 90
- 72 Lee E, Lee J W, No J S, et al. Minimax approximation of sign function by composite polynomial for homomorphic comparison. *IEEE Trans Dependable Secure Comput*, 2022, 19: 3711–3727
- 73 Marcano N J H, Moller M, Hansen S, et al. On fully homomorphic encryption for privacy-preserving deep learning. In: *Proceedings of 2019 IEEE Globecom Workshops (GC Wkshps)*, 2019. 1–6
- 74 Salem M, Taheri S, Yuan J S. Utilizing transfer learning and homomorphic encryption in a privacy preserving and secure biometric recognition system. *Computers*, 2018, 8: 3
- 75 Malik R, Singhal V, Gottfried B, et al. Vectorized secure evaluation of decision forests. In: *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, 2021. 1049–1063
- 76 Li J Y, Huang H. Faster secure data mining via distributed homomorphic encryption. In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020. 2706–2714
- 77 Sarpatwar K, Ratha N, Nandakumar K, et al. Privacy enhanced decision tree inference. In: *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020. 154–159
- 78 Wu W, Liu J, Wang H M, et al. Secure and efficient outsourced k-means clustering using fully homomorphic encryption with ciphertext packing technique. *IEEE Trans Knowl Data Eng*, 2020, 33: 3424–3437
- 79 Li B Y, Micciancio D. On the security of homomorphic encryption on approximate numbers. In: *Proceedings of Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT 2021)*, 2021. 648–677
- 80 Sathe S, Aggarwal C. LODS: local density meets spectral outlier detection. In: *Proceedings of the SIAM*

- International Conference on Data Mining, Society for Industrial and Applied Mathematics, 2016. 171–179
- 81 Heigl M, Anand K A, Urmann A, et al. On the improvement of the isolation forest algorithm for outlier detection with streaming data. *Electronics*, 2021, 10: 1534
- 82 Abe N, Zadrozny B, Langford J. Outlier detection by active learning. In: *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2006. 504–509
- 83 Aggarwal C C, Sathe S. Theoretical foundations and algorithms for outlier ensembles. *SIGKDD Explor Newsl*, 2015, 17: 24–47

## Homomorphic encryption-based ciphertext anomaly detection method for e-health records

Teng LI<sup>1</sup>, Baokun FANG<sup>2</sup>, Zhuo MA<sup>1\*</sup>, Yulong SHEN<sup>1</sup> & Jianfeng MA<sup>1</sup>

1. *School of Cyber Engineering, Xidian University, Xi'an 710071, China;*

2. *Guangzhou Institute of Technology, Xidian University, Guangzhou 510555, China*

\* Corresponding author. E-mail: mazhuo@mail.xidian.edu.cn

**Abstract** To avoid the leakage of patient information and diagnosis results in electronic health data (EHRs) in the process of anomaly detection, a privacy protection model for EHRs' anomaly detection based on CKKS fully homomorphic encryption of sensitive data of hospitals and patients is proposed. The EHRs of hospitals and patients are encrypted using the CKKS algorithm to achieve floating-point number homomorphic encryption. Then, a protocol based on the ciphertext comparison algorithm is designed to establish a ciphertext state isolation forest model through the communication between the trusted key server and the third-party data center. Using the SIMD technology of the CKKS algorithm, the anomaly detection of the ciphertext data on the isolation forest model is realized, and the ciphertext result is finally returned. The theoretical analysis and experimental results show that the proposed scheme can ensure EHR privacy and security. It is verified on the dataset that this model is superior to the traditional plaintext anomaly detection algorithm and the same type of ciphertext anomaly detection algorithm. The model can maintain detection efficiency similar to the plaintext isolation forest algorithm in the ciphertext state and has a good anomaly detection effect.

**Keywords** homomorphic encryption, isolation forest, anomaly detection, privacy protection, ciphertext comparison