

高效的决策树隐私分类服务协议

马立川^{1,2}, 彭佳怡^{1,2}, 裴庆祺^{1,2}, 朱浩瑾³

(1. 西安电子科技大学综合业务网理论及关键技术国家重点实验室, 陕西 西安 710071;

2. 陕西省区块链与安全计算重点实验室, 陕西 西安 710071; 3. 上海交通大学计算机学院, 上海 200240)

摘 要: 为了有效解决物联网大数据场景中的决策树隐私分类服务问题, 将决策树分类模型与安全多方计算技术相结合, 提出了一种高效的决策树隐私分类服务协议。该协议包括: 决策树分类模型混淆、基于布尔共享的隐私比较和基于不经意传输的隐私分类结果获取 3 个阶段。该协议能够同时保护服务提供商决策树分类模型参数及结构特征和用户需要进行分类的特征数据不被泄露。安全性分析表明, 所提决策树隐私分类服务协议能够抵抗“诚实好奇”的攻击者。将所提协议用于通过公开数据集得到的决策树分类模型, 以分类准确率和完成隐私分类服务的时间效率为指标与现有方法进行对比, 实验结果验证了所提出隐私分类服务协议的准确性和高效性。

关键词: 决策树; 隐私保护; 不经意传输; 安全多方计算

中图分类号: TN92

文献标识码: A

DOI: 10.11959/j.issn.1000-436x.2021149

Efficient privacy-preserving decision tree classification protocol

MA Lichuan^{1,2}, PENG Jiayi^{1,2}, PEI Qingqi^{1,2}, ZHU Haojin³

1. The State Key Laboratory of Integrated Services Networks, Xidian University, Xi'an 710071, China

2. Shaanxi Key Laboratory of Blockchain and Secure Computing, Xi'an 710071, China

3. The Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China

Abstract: To provide privacy-preserving decision tree classification services in the Internet of things (IoT) big data scenario, an efficient privacy-preserving decision tree classification protocol was proposed by adopting the secure multiparty computation framework into the classification model. The entire protocol consisted of three parts: the original decision tree model mixing, the Boolean share-based privacy-preserving comparing, and the 1-out-of- n oblivious transfer-based classification result obtaining. Via the proposed protocol, the service providers could protect the parameters of their decision tree models and the users were able to derive the classification result without exposing their privately hold data. Through a concrete security analysis, the proposed protocol was proved to be secure against semi-honest adversaries. By implementing the proposed protocol on various practical decision tree models from open datasets, the classification accuracy and the average time cost for completing one privacy-preserving classification service were evaluated. After compared with existing related works, the performance superiority of the proposed protocol is demonstrated.

Keywords: decision tree, privacy preserving, oblivious transfer, secure multiparty computation

收稿日期: 2021-04-01; 修回日期: 2021-06-15

基金项目: 国家重点研发计划基金资助项目 (No.2020YFB1807500); 国家自然科学基金资助项目 (No.61902292, No.61972453, No.62072355); 陕西省重点研发计划基金资助项目 (No.2021ZDLGY06-03, No.2019ZDLGY13-07, No.2019ZDLGY13-04); 中央高校基本科研业务费基金资助项目 (No.XJS201502)

Foundation Items: The National Key Research and Development Program of China (No.2020YFB1807500), The National Natural Science Foundation of China (No.61902292, No.61972453, No.62072355), The Key Research and Development Programs of Shaanxi (No.2021ZDLGY06-03, No.2019ZDLGY13-07, No.2019ZDLGY13-04), The Fundamental Research Funds for the Central Universities (No. XJS201502)

1 引言

随着信息化和网络化进程的加快以及嵌入式设备的普及，物联网（IoT, Internet of things）技术已经成为学术界和工业界的研究热点。作为联接网络空间和物理世界的“桥梁”，物联网已经在智能医疗、智慧城市、无人驾驶等与民生息息相关的领域扮演了越来越重要的角色^[1]。相关调研报告指出，2025年全球物联网设备的数量将会达到754.4亿^[2]。

数以亿计的物联网终端设备持续对其所处的环境状态进行捕捉，并源源不断地产生诸如日志、声音、视频等多样化的海量数据。然而，由于物联网设备是计算、通信、存储等资源受限的小型设备，其本身难以执行复杂的运算。因此，物联网终端产生的海量数据一般被上传到云计算中心，利用大数据分析技术对数据中蕴含的价值进行充分挖掘。在此背景下，便产生了“物联网大数据”的概念^[3]。

与此同时，能够从多样化数据中进行模式挖掘与特征提取的机器学习算法已经被成功地应用于语音视频分析、自然语言处理、趋势预测等领域，其已经构成了大数据分析技术的重要组成部分。其中，基于规则空间划分的决策树分类算法因其易于实现和高效性，已经成为机器学习中应用最广泛的分类算法之一^[4]。在物联网大数据中，往往采用“机器学习即服务”（MLaaS, machine learning as a service）的方式来对用户提供服务，即云数据中心将来自物联网终端设备的海量数据进行汇聚并进行训练得到最终的决策树分类模型，然后通过该模型对外提供分类服务^[5-6]。

然而，用户在以 MLaaS 的方式便利地使用决策树分类服务的同时，面临着严重的隐私泄露风险。一方面，服务提供商在提供决策树分类服务时，需要保护其所训练出来的分类模型不被泄露。另一方面，用户在请求分类服务时，一般需要向服务提供商提交其需要进行分类预测的数据，而这些数据往往包含用户的行为习惯、偏好、位置、收入等敏感信息，随着用户隐私保护的意识越来越强，在进行分类时需要兼顾用户的数据隐私。此外，各国颁布的隐私保护法规（如欧盟的《通用数据保护条例》，美国的《加利福尼亚消费者隐私法案》^[7]和我国的《中华人民共和国网络安全法》）严格要求服务提供商在提供服务时需要对用户的隐私信息进行保护^[7-9]。因此，在物联网大数据背景下，服务提

供商对用户提供服务时要求保护预测模型和用户的属性数据不被泄露，即服务提供商需要提供决策树隐私分类服务。

目前，为了实现决策树隐私分类服务，一般采用可搜索加密^[10]、同态加密^[11-12]以及安全多方计算^[13]等工具。然而，文献[10]所提基于可搜索加密的决策树隐私分类方法泄露了树形分类模型的整体结构。文献[11-12]基于同态加密所设计的方法虽然能够同时保护分类模型结构和用户数据不被泄露，但给服务提供商和用户带来了巨大的计算负担。文献[13]引入安全多方计算框架，将 Yao 混淆电路和不经意传输（OT, oblivious transfer）协议相结合，提升了决策树隐私分类的效率，但是其中涉及了多个按固定顺序依次计算的混淆电路，故在一定程度上限制了该方法的实际应用。

为了保证服务提供商能够提供决策树隐私分类服务，并克服现有工作的缺点，本文提出了一种面向物联网大数据的决策树隐私分类服务方法，进一步提升了决策树隐私服务分类方法的效率。本文具体的研究工作如下。

1) 提出了面向物联网大数据的决策树隐私分类服务系统模型，基于该模型，给出了威胁模型及安全定义。

2) 设计了一种高效的决策树隐私分类服务协议，其包括决策树分类模型混淆、基于布尔共享的隐私比较和基于不经意传输的隐私分类结果获取 3 个阶段。该协议能够保护服务提供商决策树分类模型参数及结构特征和用户需要进行分类的特征数据。

3) 通过安全性分析证明了所提决策树隐私分类服务能够抵抗“诚实好奇”的恶意攻击者。同时，将所提协议用于通过公开数据集得到的决策树分类模型，以分类准确率和完成隐私分类服务的效率为指标，与现有方法进行对比，验证了本文所提隐私分类服务协议的高效性。

2 相关研究工作

作为机器学习中的一种典型方法，决策树因其易于实现和分类性能高效被广泛应用于移动通信^[14]、智慧医疗诊断^[15]、网络安全防护^[16]等各个方面。决策树分类方法的工作原理如下。通过训练数据得到一种树形的分类模型，其包括内部节点和叶子节点。每个内部节点具有一个属性标签和阈值，叶子

节点则代表一个分类。在利用决策树模型进行分类时，需要从根节点开始，将对应节点属性标签的属性值与阈值进行比较，根据比较结果选择其相应的子节点，直至到达叶子节点，得到最终的分类结果。上述分类过程可以总结为：利用数据的属性值找到决策树中一条从根节点到叶子节点的路径，叶子节点所对应的分类为该条数据的最终分类结果。

随着用户隐私保护意识的增强以及世界各国法规对隐私信息保护的要求越来越严苛，在物联网大数据环境下使用机器学习模型提供分类服务时，需要同时保护分类模型及用户数据不被泄露^[17]。近年来，为了在兼顾隐私的同时实现决策树分类方法，一般引入可搜索加密、同态加密和安全多方计算等工具。其中，文献[10]将决策树中根据每一个内部节点所定义的阈值对决策树从根节点到叶子节点的路径进行编码，并将路径的编码与叶子节点所定义的类别建立映射，此时，可以将决策路径选取问题转化为以路径编码为关键词的搜索问题。然而，该方法泄露了决策树的整体结构，并且难以处理内部节点所定义的阈值为非整数的情况。文献[11]给出了包括决策树模型在内的多种隐私分类方法，其采用了全同态加密方法，给服务提供商和用户带来了巨大的计算负担。文献[12]对上述方法进行了改进，其方案仅需要利用加法同态加密即可。文献[11-12]的方法复杂度均取决于决策树内部节点的数量，当决策树规模变大时，便变得不实用。文献[13]引入

安全多方计算框架，将 Yao 混淆电路与不经意传输协议相结合，使决策树隐私分类服务的复杂度只与决策树的深度相关。虽然文献[13]所提方法在这些方法中性能最优，但是其每次迭代的需要引入多个混淆电路的计算，故其实用性仍受到限制。

因此，本文综合考虑了现有决策树隐私分类服务协议的优点，将决策树分类模型与安全多方计算框架相结合，提出了一种更加高效的决策树隐私分类服务协议。

3 系统模型

决策树隐私分类服务系统模型如图 1 所示。服务提供商在向用户提供决策树隐私分类服务时，考虑了云计算中典型的 Server-Client 模型。服务器(用 S 表示)位于云计算中心，其主要负责收集来自物联网设备的数据，并对数据进行标记。本文充分利用云数据中心的计算和存储能力，对所收集的数据进行训练，得到树形结构的决策树分类模型，并利用该模型为用户提供分类服务。用户(用 C 表示)可以向 S 提供用于分类的数据，经过计算后由 S 返回分类结果。

决策树的分类模型用 T 表示，其为树形结构，包括根节点、内部节点和叶子节点。使用集合 $V = \{v_1, \dots, v_m\}$ 表示 T 中除去叶子节点的所有节点集合， $|V| = m$ 。节点的标号从根节点开始，逐层从左往右依次标记， v_1 表示根节点。 T 中叶子节点的集合为 $Z = \{z_1, \dots, z_n\}$ ， $|Z| = n$ ，标记的顺序仍为从

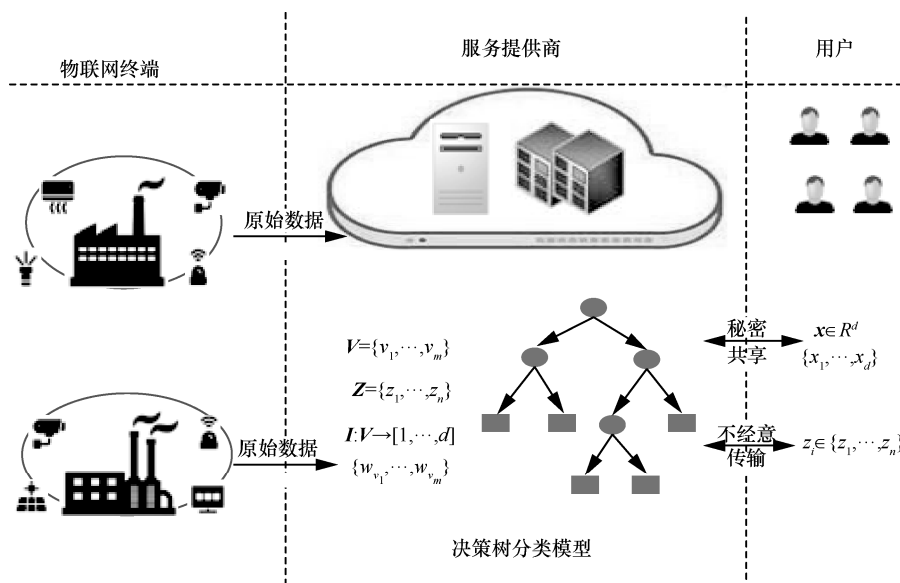


图 1 决策树隐私分类服务系统模型

左到右。此时，树 T 包含节点的数量为 $m+n$ 。分类模型 T 为 V 中的每个节点 v_k ($k=1, \dots, m$) 分配权重 w_k 、布尔函数 $f_k(x)=1$ ($x \leq w_k$)，以及标记函数 $I:V \rightarrow [1, \dots, d]$ ，此处， $I(v_k)$ 返回的是内部节点 v_k 所对应的属性序号，记为 i_k 。

假设用户的数据用 \mathbf{x} 表示，其具有 d 个属性，则 $\mathbf{x} \in R^d$ 。在不考虑隐私保护的前提下通过决策树分类模型 T 对 \mathbf{x} 进行分类时，首先从根节点 v_1 开始，计算 $f_1(x_{i_1})$ 并根据其取值确定接下来需要考虑的子节点，依次类推，最终到达叶子节点，该叶子节点所对应的类别即为 \mathbf{x} 的最终分类结果。此时，可以将 T 看作函数 $T:R^d \rightarrow Z(= \{z_1, \dots, z_n\})$ 。

与文献[12-13]中的假设条件不同，本文中考虑的决策树分类模型 T 不一定是一个深度 t 的完全二叉树。

考虑到隐私保护的要求，本文与大多数隐私保护相关工作的恶意模型假设相同，即服务提供商和用户均为“诚实好奇”的，其能够遵循协议的规定正确地完成任务，但试图从接收到的数据中推断另一方的原始输入。严格的“诚实好奇”模型下安全定义如下^[18]。

定义 1 令 π 表示一个协议。如果 π 能够在“诚实好奇”攻击者 A 存在的前提下安全计算指定函数 ε ，那么存在一个可信的模拟器 Sim ，其能够模拟协议 π 的运行过程， Sim 与 A 共同完成协议 π ，在此过程中， Sim 以产生随机比特串的方式模拟实际情况下另一方的输入，并且满足

$$\text{View}_{\pi, A, \varepsilon}^{\text{real}} \triangleq \text{View}_{\text{Sim}, \varepsilon}^{\text{ideal}} \quad (1)$$

其中， $\text{View}_{\pi, A, \varepsilon}^{\text{real}}$ 表示协议 π 在执行过程中攻击者 A 获得的全部数据， $\text{View}_{\text{Sim}, \varepsilon}^{\text{ideal}}$ 表示攻击者 A 在与 Sim 交互完成协议 π 所获得的全部数据，符号 \triangleq 表示计算不可区分。

定义 1 表明，如果协议 π 在“诚实好奇”攻击者存在的情况下是安全的，那么该攻击者在完成协议的过程中仅仅能获得输入和协议规定的输出，无法获取除此之外的任何信息。

在本文所提决策树隐私分类服务协议中，“诚实好奇”攻击者具有两层含义：1) 当拥有决策树分类模型的 S 为“诚实好奇”攻击者时，其基于完成隐私分类服务协议过程中所接收到的交互数据，试图推断用户进行分类的私密数据；2) 当 C 为攻击者时，则会基于交互数据去推断 S 所拥有

决策树分类模型的相关信息。本文第 5 节将综合考虑上述 2 种情况，给出所提出隐私分类协议的安全性证明。

4 决策树隐私分类协议设计

本节首先给出了所提决策树隐私分类协议的概述和基本工作原理；其次，分别从决策树模型变换、分类路径确定及分类结果获取 3 个方面详细地介绍了该协议每个步骤的实现细节。

4.1 协议概述

本文所提协议聚焦在决策树分类环节，服务提供商拥有决策树分类模型 T ，用户拥有属性数据 \mathbf{x} 。在不考虑隐私保护的情况下，服务提供商将 \mathbf{x} 输入模型 T ，得到一个分类结果 z_x 并将其返回给用户。然而，在隐私保护的前提下，服务提供商不能将 T 泄露给用户，而用户则不希望将 \mathbf{x} 和 z_x 泄露给服务提供商。为此，在决策树隐私分类协议设计时，需要同时保护服务提供商的分类模型 T 以及用户的数据 \mathbf{x} 及分类结果 z_x 。

在使用决策树分类模型 T 对 \mathbf{x} 进行分类时，需要从 T 的根节点 v_1 开始，得到根节点所对应的属性序号 i_1 ，然后 w_1 与 x_{i_1} 进行比较，根据比较结果选择 v_1 的左子节点或右子节点继续同样的步骤，直至到达叶子节点，叶子节点所对应的分类即为 \mathbf{x} 的分类结果。为了在上述过程中保护分类模型 T 、用户数据 \mathbf{x} 和分类结果 z_x 不被泄露，需要满足以下条件。

1) 对 T 所定义的内部节点阈值比较顺序地进行混淆，使用户无法通过进行比较操作的属性值顺序推断出 T 除叶子节点以外的树形结构。

2) 对需要进行比较操作的阈值及用户数据的属性值进行保护，使用户无法推断出内部节点所对应的阈值且服务提供商无法推断出用户数据。

3) 对分类结果 z_x 进行保护，使服务提供商无法获取用户数据 \mathbf{x} 所对应的分类结果。

为了满足上述 3 个针对决策树分类的隐私保护条件，本文所提树隐私分类协议由决策树分类路径混淆、基于布尔共享的隐私比较和基于不经意传输的分类结果获取三部分构成。本文用到的数学符号以含义如表 1 所示。

4.2 决策树分类模型混淆

通过决策树模型进行分类时，可以从根节点到叶子节点的一条路径称为分类路径，每一条分类

路径对应一个唯一的分类结果。文献[12-13]将决策树扩展为一个深度 t 的完全二叉树，每一条分类路径包含 $d-1$ 个根节点和一个叶子节点。如果在每个非叶子节点，用 0 表示属性值小于阈值，1 表示其他情况。那么，可以将 T 看作函数 $T: \{0,1\}^{t-1} \rightarrow \{z_1, \dots, z_n\}$ 。由于在 T 中包含了 m 个内部节点，文献[12]将 $\{0,1\}^{t-1}$ 扩展到 $\{0,1\}^m$ ，其中每个比特对应一个内部节点的比较结果，此时分类模型变为 $T: \{0,1\}^m \rightarrow \{z_1, \dots, z_n\}$ 。

表 1 数学符号及含义

参数	含义
S	拥有决策树分类模型的云服务器
C	请求决策树隐私分类的用户
T	决策树分类模型
\mathbf{x}	用户数据
z_x	\mathbf{x} 所对应的分类结果
$I: V \rightarrow [1, \dots, d]$	标记函数
\mathbf{IV}_0	原始决策树分类模型内部节点标号序列
\mathbf{IX}_0	与 \mathbf{IV}_0 对应的属性标号序列
\mathbf{W}_0	与 \mathbf{IV}_0 对应的内部节点阈值序列
δ_r	定义在 \mathbf{IV}_0 上的随机置换
\mathbf{IV}	\mathbf{IV}_0 经过函数 δ_r 作用后的随机置换
\mathbf{LX}	与 \mathbf{IV}_r 对应的属性标号序列
\mathbf{W}_r	与 \mathbf{IV}_r 对应的内部节点阈值序列
$[x]$	x 的二进制表示
$[x]_s, [x]_c$	$[x]$ 的布尔共享，且 $[x] = [x]_s \oplus [x]_c$
\oplus	“异或”运算
\otimes	“与”运算
$(\cdot)_{OT}$	1-out-of- n 不经意传输过程

然而，在获取长度为 m 的比特串作为输入时，首先，记决策树分类模型内部节点的标号序列为 $\mathbf{IV}_0 = \{1, \dots, m\}$ ，其中将根节点的序号标为 1；然后，按照广度优先搜索的原则逐层从左到右依次对内部节点进行标号。由标号序列 \mathbf{IV}_0 所确定的内部节点序列记为 \mathbf{V}_0 ，那么将 \mathbf{V}_0 中每个内部节点所对应的属性标号序列和阈值序列分别记为 \mathbf{IX}_0 和 \mathbf{W}_0 ，其中 $\mathbf{IX}_0 = \{I(v_{0,k}): k=1, \dots, m\}$ ， $\mathbf{W}_0 = \{w(v_{0,k}): k=1, \dots, m\}$ 。此时，决策树分类模型 T 由 \mathbf{IV}_0 、 \mathbf{IX}_0 和 \mathbf{W}_0 唯一确定，即可以看作函数 $T[\mathbf{IV}_0, \mathbf{IX}_0, \mathbf{W}_0]: \mathbf{x} \in R^d \rightarrow \{z_1, \dots, z_n\}$ 。

服务提供商在提供决策树分类服务时，如果直

接将 \mathbf{IX}_0 发送给用户，则会暴露分类模型 T 的树形结构。为了保护 \mathbf{IX}_0 ，文献[12]采用了树变换的方法，然而此方法仍会暴露根节点所对应的数据属性标号。本文直接采用随机置换的方法，通过 \mathbf{IV}_0 的随机置换对 \mathbf{IX}_0 进行混淆，从而保护树形结构信息不被泄露。

定义函数 δ_r 为 \mathbf{IV}_0 的随机置换，即 $\delta_r: \mathbf{IV}_0 \rightarrow \mathbf{IV}_r$ ，由 \mathbf{IV}_r 所确定的内部节点序列表示为 \mathbf{V}_r ，那么由 \mathbf{V}_r 中内部节点所确定的属性标号序列 $\mathbf{IX}_r = \{I(v_{r,k}): k=1, \dots, m\}$ 。此时，通过作用在 \mathbf{IV}_0 上的随机置换函数 δ_r 将 $T[\mathbf{IV}_0, \mathbf{IX}_0, \mathbf{W}_0]$ 进行混淆得到新的决策树分类模型 $T[\mathbf{IV}_r, \mathbf{IX}_r, \mathbf{W}_r]$ 。该过程如算法 1 所示。

算法 1 决策树分类模型混淆算法

输入 $\mathbf{IV}_r, \mathbf{IX}_r, \mathbf{W}_r$ 为空集

输出 经过随机置换函数 δ_r 作用后的决策树分类模型

1) 根据训练数据得到原始的决策树分类模型

$$T[\mathbf{IV}_0, \mathbf{IX}_0, \mathbf{W}_0]: \mathbf{x} \in R^d \rightarrow \{z_1, \dots, z_n\}$$

2) for $j=1:m$

3) 从 \mathbf{IV}_0 中随机选取一个元素 $iv_{0,j}$

$$4) \mathbf{IV}_r = \mathbf{IV}_r \cup \{iv_{0,j}\}$$

$$5) \mathbf{IX}_r = \mathbf{IX}_r \cup \{I(v_{iv_{0,j}})\}$$

$$6) \mathbf{W}_r = \mathbf{W}_r \cup \{w_{iv_{0,j}}\}$$

7) end for

8) 输出经过随机置换函数 δ_r 作用后的决策树分类模型 $T[\mathbf{IV}_r, \mathbf{IX}_r, \mathbf{W}_r]$

对于任意的用户数据 $\mathbf{x} \in R^d$ ，利用原始分类模型 $T[\mathbf{IV}_0, \mathbf{IX}_0, \mathbf{W}_0]$ 进行分类时，可以将 \mathbf{x} 映射为 $\sigma_x \in \{0,1\}^m$ ，此时，定义函数 $\phi_0: \sigma \in \{0,1\}^m \rightarrow \{1, \dots, n\}$ 表示决策路径 σ 与分类标号之间的映射。而利用经过混淆后的决策树分类模型 $T[\mathbf{IV}_r, \mathbf{IX}_r, \mathbf{W}_r]$ 进行分类时， \mathbf{x} 被 ϕ_r 映射为 $\sigma_x \in \{0,1\}^m$ ，其可以看作 σ_x 在函数 δ_r 作用下的一个置换。用户在请求分类服务后， ϕ_r 与 \mathbf{IX}_r 可以由服务提供商发送给请求用户。

值得注意的是，通过该方法得到的 \mathbf{IX}_r ，攻击者能够猜对的概率为 $1/(m!)$ ；而文献[12]中作者通过决策树变换方法得到 \mathbf{IV}_r ，攻击者能够猜对的概率为 $1/2^m$ ，并且攻击者总是能够获取根节点所对应的数据属性标号。

4.3 基于布尔共享的隐私比较

用户 C 提交隐私分类服务请求后，服务提供商 S 将 ϕ_r 与 \mathbf{IX}_r 发送给用户 C。接下来，用户 C 将根据 \mathbf{IX}_r 所确定的属性标号，选择对应的属性值与服务提供商拥有的阈值序列 \mathbf{W}_r 中对应的阈值进行比较，进而确定最终的决策路径 $\sigma_{rx} \in \{0,1\}^m$ ，随后可以通过公开的函数 ϕ_r 得到数据 \mathbf{x} 所对应的类别标号。

在上述过程中，对于 \mathbf{IX}_r 中的任意属性标号 $\tau_j (j=1, \dots, m)$ ，需要将 x_{τ_j} 与对应的 w_j 进行比较，如果 $x_{\tau_j} < w_j$ ， $\sigma_{rx,j} = 1$ ；否则， $\sigma_{rx,j} = 0$ 。此时，用户 C 拥有 x_{τ_j} ，服务提供商 S 拥有 w_j 。为了满足隐私保护要求，在进行比较的时候，不能向对方泄露 x_{τ_j} 和 w_j 的具体数值，并且只有用户 C 能够获取

最终的比较结果。为了实现隐私比较，文献[11]应用了全同态加密^[19]的方法得到最终比较结果；文献[12]使用基于加法同态加密的比较方法^[20-21]，减少了使用全同态加密时服务提供商和用户的计算负担。然而，无论是基于全同态加密还是加法同态加密，其给服务提供商和用户带来的计算负担仍然很大。为进一步提升隐私比较的效率，本文采用了基于布尔共享的隐私比较方法，其基本思路将需要比较的属性值和阈值转化为定长（长度设为 l ）的比特串并产生对应的布尔共享，按照经典的 GMW（Goldreich-Micali-Wigderson）协议^[22]确定完成比较操作的布尔电路，确定每个参与方（S 或 C）所要进行的运算。在整个过程中，没有涉及复杂的密文运算，故可以提升隐私比较的效率。

在实现基于布尔共享的隐私比较时，对于任意 $j(=1, \dots, m)$ ，用户 C 将 x_{τ_j} 转化为长度为 l 的二进制表示 $[x_{\tau_j}]_c$ ，然后随机产生长度为 l 的比特串 $[x_{\tau_j}]_c$ ，并令

$$[x_{\tau_j}]_s = [x_{\tau_j}]_c \oplus [x_{\tau_j}] \quad (2)$$

此时， $[x_{\tau_j}]_c$ 和 $[x_{\tau_j}]_s$ 便构成了 $[x_{\tau_j}]$ 的布尔共享，其中，用户 C 将 $[x_{\tau_j}]_c$ 保留，并将 $[x_{\tau_j}]_s$ 发送给服务提供商 S。按照同样的方式，S 得到 $[w_j]_s$ 、 $[w_j]_c$ 和 $[w_j]_c$ ，并将 $[w_j]_c$ 发送给用户 C。

接下来，需要确定实现比较操作的布尔电路，如图 2 所示，本文所提方法采用了 CMP 布尔电路^[23]。可以看到，在实现比较操作时，涉及了 2 种运算：“异或”运算 \oplus 和“与”运算 \otimes 。

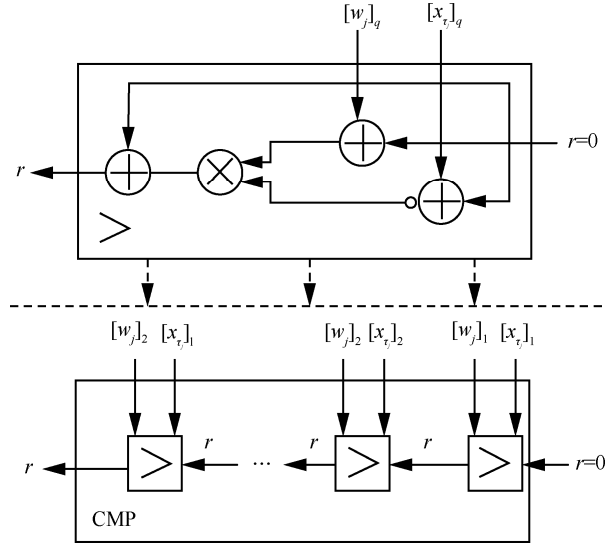


图 2 CMP 布尔电路

对于任意的 $\alpha, \beta \in \{0,1\}$ ，假设 $\alpha = \alpha_s \oplus \alpha_c$ ， $\beta = \beta_s \oplus \beta_c$ ，则有

$$[\alpha \otimes \beta]_{S(\text{or } C)} = \alpha_{S(\text{or } C)} \otimes \beta_{S(\text{or } C)} \quad (3)$$

要求 $\alpha \otimes \beta$ 的布尔共享，则需要借助能够预先计算得到的三元组 $\langle a_c, b_c, g_c \rangle$ 和 $\langle a_s, b_s, g_s \rangle$ 使^[24]

$$g_c \oplus g_s = (a_c \oplus a_s) \otimes (b_c \oplus b_s) \quad (4)$$

此时，按照文献[24]给出的步骤得到 $[\alpha \otimes \beta]_s$ 和 $[\alpha \otimes \beta]_c$ 。

利用上述的布尔电路 CMP 以及布尔共享的前提下进行异或运算和与运算，算法 2 给出了 S 和 C 基于其所拥有的布尔共享实现隐私比较的过程。

算法 2 基于布尔共享的隐私比较

输入 服务提供商 S 输入 $[x_{\tau_j}]_s$ 、 $[w_j]_s$ ，辅助比特 r_s ；用户 C 输入 $[x_{\tau_j}]_c$ 、 $[w_j]_c$ 及 r_c ，其中， $r_s \oplus r_c = 0$ 。

1) for $q = 1:l$

2) S 计算

$$r_s = (r_s \oplus [x_{\tau_j}]_{s,q}) \otimes \overline{(r_s \oplus [w_j]_{s,q})} \oplus r_s$$

3) C 计算

$$r_c = (r_c \oplus [x_{\tau_j}]_{c,q}) \otimes (r_c \oplus [w_j]_{c,q}) \oplus r_c$$

4) end for

5) S 将 r_s 发送给 C

6) C 计算 $\sigma_{rx,j} = r_s \oplus r_c$

将算法 2 执行 m 次即可得到 $\sigma_{rx} \in \{0,1\}^m$ 。值得注意的是,虽然基于布尔共享的隐私比较需要进行多次,但是不同的隐私比较操作相互独立,故可以用并行的方式完成该 m 对数的比较,进一步提升实现隐私比较的效率。此时,用户 C 便可以根据公开的 ϕ_r 和 σ_{rx} 得到 \mathbf{x} 所对应的叶子节点标号。

4.4 基于不经意传输的隐私分类结果获取

经过基于布尔共享的隐私比较之后,用户 C 获得了数据 \mathbf{x} 所对应的叶子节点标号,记为 γ 。而对于服务提供商 S 而言,叶子节点集合 $\mathbf{Z} = \{z_1, \dots, z_n\}$ 中的每个叶子节点对应一个类别,假设 $z_j (j=1, \dots, n)$ 为一个长度为 λ 的比特串,即 $z_j \in \{0,1\}^\lambda$ 。在获取最终的分类结果时,C 希望在 S 无法知晓 γ 的前提下获取 z_γ , S 则希望 C 只能得到 z_γ 而无法获取其余叶子节点所对应的类别信息。

此时, C 从 S 处获取隐私分类结果的过程为典型的 1-out-of- n 不经意传输过程,记为 $(\binom{n}{1})\text{OT}_\lambda$, 其中 S 的输入为 $\{z_1, \dots, z_n\}$, C 的输入为 γ 。该过程与文献[12]中用户获取最终分类结果的方法保持一致。

目前,实现 $(\binom{n}{1})\text{OT}_\lambda$ 的最优方案为 Pinkas 等^[25]提出的方法,其基本原理是通过引入能够快速实现的伪随机函数将 $(\binom{n}{1})\text{OT}_\lambda$ 转化为 $\lceil \ln n \rceil \times (\binom{n}{1})\text{OT}_k$, 其中 $\lceil \cdot \rceil$ 表示向上取整操作。在后面的仿真实验中,本文也采用该方法对上述过程进行实现。具体的实现算法如算法 3 所示。

算法 3 基于不经意传输的隐私分类结果获取

输入 服务提供商 S 输入叶子节点集合 $\mathbf{Z} = \{z_1, \dots, z_n\}$; 用户 C 输入对应于数据 \mathbf{x} 的叶子节点标号 γ

- 1) S 确定选取伪随机函数 F , 并准备 $L (= \lceil \ln n \rceil)$ 对密钥 $(K_1^0, K_1^1), \dots, (K_L^0, K_L^1)$
- 2) for $I = 1 : n$
- 3) S 将 I 表示成二进制形式 (i_1, \dots, i_L)
- 4) S 计算 $Y_I = z_I \oplus (\bigoplus_{j=1}^L F_{K_j^{i_j}}(I))$
- 5) end for
- 6) C 将 γ 转化为其二进制形式 $(\gamma_1, \dots, \gamma_L)$
- 7) S 与 C 执行 L 次 $(\binom{1}{1})\text{OT}_k$, 使 C 获取 $(K_1^{\gamma_1}, \dots, K_1^{\gamma_L})$
- 8) S 将 (Y_1, \dots, Y_n) 发送至 C
- 9) C 计算得到最终的分类结果 z_γ 为

$$z_\gamma = Y_\gamma \oplus (\bigoplus_{j=1}^L F_{K_j^{\gamma_j}}(I))$$

至此,决策树分类模型混淆、基于布尔共享的隐私比较和基于不经意传输的隐私分类结果获取便构成了本文所提的面向物联网大数据的决策树隐私分类协议。

5 性能分析

本节首先通过严格的安全性分析证明了所提决策树隐私分类服务协议在抵抗“诚实好奇”攻击者的安全性。其次,通过设置实验,分别对所提协议的分类准确率和实现效率进行验证。此处,在验证分类准确率时,将本文协议的分类准确率与明文下的分类准确率进行对比,对比结果表明,两者的分类准确率保持一致,该结果进一步验证了所提分类协议的正确性。在估计实现效率时,以完成一次决策树分类服务所需的时间为指标,将本文所提方法与文献[12-13]中的方法进行比较,实验结果表明,无论是对小型决策树分类模型还是内部节点和深度比较大的决策树分类模型,本文所提方法均优于其余 2 种方法。

5.1 安全性分析

根据定义 1 所给出的安全定义,如果本文所提方案对于“诚实好奇”攻击者而言是安全的,那么在整个隐私分类服务实现的过程中,服务提供商或者用户仅能获取其协议所规定的的数据,无法获取任何与其原始输入相关的任意信息。在进行安全性分析时,本文考虑了 2 种情况,即服务提供商和用户分别变为“诚实好奇”攻击者的情形。

当服务提供商变为攻击者时, $\text{View}_{\pi,S}^{\text{real}}$ 表示服务提供商在执行隐私分类协议时所能获得的信息。根据整个隐私分类服务协议的工作流程, $\text{View}_{\pi,S,1}^{\text{real}}$ 、 $\text{View}_{\pi,S,2}^{\text{real}}$ 和 $\text{View}_{\pi,S,3}^{\text{real}}$ 分别表示决策树分类模型混淆、基于布尔共享的隐私比较和基于不经意传输协议的隐私分类结果获取 3 个阶段服务提供商所能获取的数据。那么有

$$\text{View}_{\pi,S}^{\text{real}} = \{\text{View}_{\pi,S,1}^{\text{real}}, \text{View}_{\pi,S,2}^{\text{real}}, \text{View}_{\pi,S,3}^{\text{real}}\} \quad (5)$$

服务提供商通过训练数据得到原始的决策树分类模型 $T_0 = T[\mathbf{IV}_0, \mathbf{IX}_0, \mathbf{W}_0]$ 后,进入决策树分类模型混淆阶段,得到 $T_r = [\mathbf{IV}_r, \mathbf{IX}_r, \mathbf{W}_r]$, 此时不存在与用户的交互,故可以得到

$$\text{View}_{\pi,S,1}^{\text{real}} \triangleq \{T_0, T_r\} \quad (6)$$

假设存在一个可信的模拟器 Sim_S ，其能够以随机产生的方式模拟用户的输入，并与服务提供商进行交互完成隐私分类服务协议。此时，用 $\text{View}_{\pi,S}^{\text{sim}}$ 表示在该过程中服务提供商所能够获取的数据，与协议真实的实现过程类似，引入 $\text{View}_{\pi,S,1}^{\text{sim}}$ 、 $\text{View}_{\pi,S,2}^{\text{sim}}$ 和 $\text{View}_{\pi,S,3}^{\text{sim}}$ 使

$$\text{View}_{\pi,S}^{\text{sim}} = \{\text{View}_{\pi,S,1}^{\text{sim}}, \text{View}_{\pi,S,2}^{\text{sim}}, \text{View}_{\pi,S,3}^{\text{sim}}\} \quad (7)$$

由于在决策树分类模型混淆阶段不存在与用户的交互，故此时不存在与 Sim_S 的交互，因此有

$$\text{View}_{\pi,S,1}^{\text{sim}} \triangleq \text{View}_{\pi,S,1}^{\text{real}} = \{T_0, T_r\} \quad (8)$$

同时，根据基于布尔共享的隐私比较以及 1-out-of- n 不经意传输协议的安全性，可以得到 $\text{View}_{\pi,S,2}^{\text{sim}} \triangleq \text{View}_{\pi,S,2}^{\text{real}}$ 和 $\text{View}_{\pi,S,3}^{\text{sim}} \triangleq \text{View}_{\pi,S,3}^{\text{real}}$ ，因此

$$\text{View}_{\pi,S}^{\text{sim}} \triangleq \text{View}_{\pi,S}^{\text{real}} \quad (9)$$

基于定义 1 给出的安全性定义，本文所提的协议能够很好地抵制服务提供商变为“诚实好奇”恶意攻击者的情形。类似地，可以证明在用户变为“诚实好奇”恶意攻击者，本文所提隐私分类服务协议仍然是安全的。综上所述，本文所提协议能够很好地抵制“诚实好奇”的恶意攻击者。

5.2 分类准确率验证及时间效率对比

安全性分析证明，本文所提决策树隐私分类服务协议能够很好地抵抗“诚实好奇”模型下的恶意攻击者。本节通过实验验证本文所提隐私分类协议的分类准确率和实现效率。

本节实验通过 C++ 实现，代码运行于装有 Ubuntu 18.04 的虚拟机上，该虚拟机的内存和硬盘容量分别为 16 GB 和 50 GB，处理器个数为 6。在决策树分类模型混淆算法的实现过程中，由于混淆的本质是对原始的内部节点序列进行随机置换，因此实验中利用标准的 AES-128 对称加密算法来保证混淆过程的安全性。在实现基于布尔共享的隐私比较时，其安全性主要基于随机产生长度为 l 的比特串以实现数据以及决策树内部节点阈值的布尔共享，为此，具体的实现过程中调用了开源的 Openssl 库来产生满足条件的随机数。在实现基于不经意传输的隐私分类结果获取时，采用了基于文献[25]所设计的 $\text{OT}_{\text{Extension}}$ 开源框架^[26]来实现高效的 (ϵ) OT_2 协议，该框架中所涉及的大数运算、对称密码以哈希运算则基于 Openssl 库和 GMP 库来实

现，以确保其安全性。此外，所有的数据均用长度为 64 的比特串表示，前 48 bit 表示整数位，后 16 bit 表示小数位。

为了得到真实的决策树分类模型，本文采用了与文献[12-13]相同的数据集，包括来自加州大学欧文分校机器学习标准测试数据集 UCI 的 Nursery、Breast-cancer、Housing、Credit-screening 和 Spambase，以及来自 PhysioBank 的 ECG 数据集。利用 Python 编程调用 sklearn 库中的 DictVectorizer 模块对数据集进行特征提取，并使用 tree 模块得到对应于不同数据集的决策树分类模型，对协议的分类正确性和时间效率进行验证。

在验证所提模型的分类准确率时，首先针对明文状态下训练得到的决策树分类模型，对测试数据进行分类，得到其分类准确率。然后，基于本文所设计的隐私分类协议和已经得到的决策树分类模型，使用相同的测试数据集，得到利用本文协议的数据分类准确率。分类准确率结果是对不同数据集分别进行 50 次实验得到分类准确率的平均值，如表 2 所示。通过表 2 的结果可以看出，通过本文协议得到的分类准确率与明文状态下直接使用分类模型得到的分类准确率保持一致。这主要是由于本文协议包括决策树分类模型、基于布尔共享的隐私比较和基于不经意传输的隐私分类结果获取 3 个阶段，每个阶段的处理均不影响决策树分类模型与分类结果之间的一一对应关系。因此，本文所提决策树隐私分类服务协议的分类准确率得以验证。

表 2 分类准确率对比

数据集	分类准确率	
	明文状态	本文协议
Breast-cancer	0.894	0.894
Nursery	0.936	0.936
ECG	0.886	0.886
Credit-screening	0.877	0.877
Housing	1.000	1.000
Spambase	0.990	0.990

下面对本文协议的实现效率进行评估。本文以完成隐私分类服务的时间效率为指标，将所提协议分别与文献[12-13]方法进行对比，验证本文协议的高效性。值得注意的是，通过数据集 Housing 和 Spambase 训练得到的决策树分类模型中所包含的内部节点数

量及其深度远大于其余 4 个数据集，即完成一次隐私分类服务所需的时间要远大于其余分类模型。

本文协议与文献[12]方法进行对比得到的结果如图 3 所示。由于通过 Breast-cancer、Nursery、ECG 和 Credit-screening 训练得到的决策树分类模型规模比较小，故 2 种方法均能很快地完成一次隐私分类服务，本文协议的时间效率优于文献[12]方法。而对于数据集 Housing 和 Spambase 而言，经过训练得到的决策树分类模型规模比较大，无论是本文协议还是文献[12]方法，完成一次隐私分类服务所需的时间均大量增加，但是本文协议仍能将其运行时间控制在 0.5 s 左右，而文献[12]方法在决策树规模变大时完成隐私分类服务所需的时间急剧增加，这主要是因为在进行隐私保护下的比较运算时，文献[12]引入了基于同态加密的比较方案，同态加密在处理算术运算时效率较高，而进行比较运算时则需要进行复杂的密文运算，使该运算的时间开销较大。当决策树分类模型规模比较大时，其存在大量的内部节点，使完成隐私分类服务所需的比较运算数量较大，故文献[12]方法在分类模型规模变大时所需的时间大量增加。而本文协议在进行比较运算时使用布尔共享的思路，比较运算效率较高，因此，即使在决策树分类模型规模变大时，本文协议仍能高效地完成隐私分类服务。

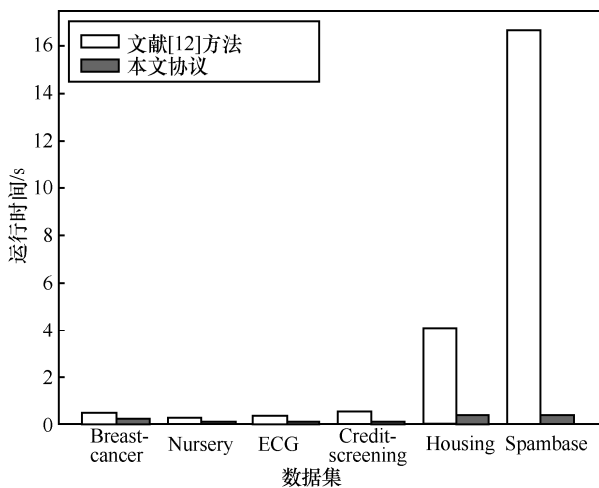


图 3 本文协议与文献[12]方法的时间效率对比

图 4 给出了本文协议与文献[13]方法在完成一次隐私分类服务时的时间效率对比。实验结果表明，无论是通过 Breast-cancer、Nursery、ECG 和 Credit-screening 训练得到规模较小的决策树分类模型，还是通过 Housing 和 Spambase 得到规模较大的分类模型，2 种方法均能够高效地完成隐私分类

服务。虽然文献[13]方法时间复杂度只与决策树分类模型的深度有关，但是其每次迭代引入了大量的混淆电路生成与解析，而本文协议只涉及相对简单的比较运算和最后的 1-out-of-n OT 协议，故本文协议的时间效率根据所选的数据集实现隐私分类服务时全面优于文献[13]方法。

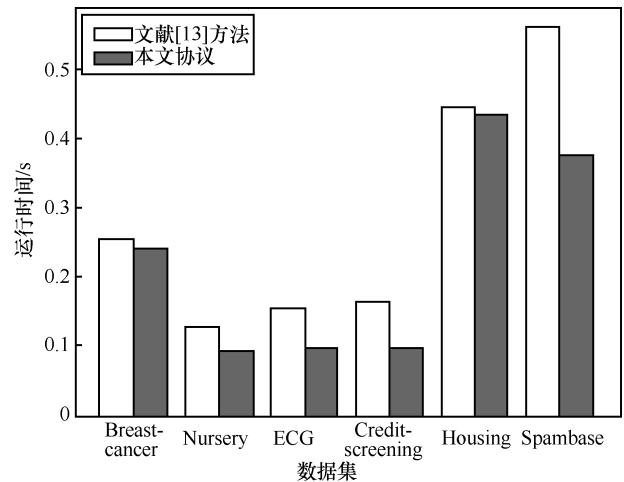


图 4 本文协议与文献[13]方法的时间效率对比

综上所述，本文协议不仅能够很好地抵制“诚实好奇”模型下的恶意攻击者，还能准确并高效地提供决策树隐私分类服务。

6 结束语

本文面向物联网大数据场景，兼顾越来越迫切的隐私保护需求，将决策树分类模型与安全多方计算框架相结合，设计了一种高效的决策树隐私分类服务协议。该协议包括：决策树分类模型混淆、基于布尔共享的隐私比较和基于不经意传输的隐私分类结果获取 3 个阶段。通过该协议，能够同时保护服务提供商决策树分类模型参数及结构特征和用户需要进行分类的特征数据。安全性分析表明，本文所提决策树隐私分类服务协议能够抵抗“诚实好奇”的恶意攻击者。同时，将本文协议应用于通过 Breast-cancer、Nursery、ECG、Credit-screening、Housing 及 Spambase 等 6 个公开数据集得到的决策树分类模型，以分类准确率和完成单次隐私分类服务的平均时间为指标，验证了该决策树隐私分类服务协议的准确性和高效性。本文的研究能够在不同的决策树分类场景中，兼顾隐私保护需求的前提下，进一步提高物联网大数据场景中的决策树隐私分类服务效率。

参考文献：

- [1] NESHENKO N, BOU-HARB E, CRICHIGNO J, et al. Demystifying IoT security: an exhaustive survey on IoT vulnerabilities and a first empirical look on Internet-scale IoT exploitations[J]. IEEE Communications Surveys & Tutorials, 2019, 21(3): 2702-2733.
- [2] CVITIĆ I, PERAKOVIĆ D, PERIŠA M, et al. Novel approach for detection of IoT generated DDoS traffic[J]. Wireless Networks, 2021, 27(3): 1573-1586.
- [3] MOHAMMADI M, AL-FUQAHA A, SOROUR S, et al. Deep learning for IoT big data and streaming analytics: a survey[J]. IEEE Communications Surveys & Tutorials, 2018, 20(4): 2923-2960.
- [4] 陈立南, 刘阳, 马严, 等. 基于统计的高效决策树分组分类算法[J]. 通信学报, 2014, 35(Z1): 58-64.
CHEN L N, LIU Y, MA Y, et al. Efficient-cutting packet classification algorithm based on the statistical decision tree[J]. Journal on Communications, 2014, 35(Z1): 58-64.
- [5] LI T, LI X, ZHONG X Y, et al. Communication-efficient outsourced privacy-preserving classification service using trusted processor[J]. Information Sciences, 2019, 505: 473-486.
- [6] PINHEIRO A J, DE M BEZERRA J, BURGARDT C A P, et al. Identifying IoT devices and events based on packet length from encrypted traffic[J]. Computer Communications, 2019, 144: 8-17.
- [7] STALLINGS W. Handling of personal information and deidentified, aggregated, and pseudonymized information under the California consumer privacy act[J]. IEEE Security & Privacy, 2020, 18(1): 61-64.
- [8] SHASTRI S, WASSERMAN M, CHIDAMBARAM V. GDPR anti-patterns[J]. Communications of the ACM, 2021, 64(2): 59-65.
- [9] QI A M, SHAO G S, ZHENG W T. Assessing China's cybersecurity law[J]. Computer Law & Security Review, 2018, 34(6): 1342-1354.
- [10] LIANG J W, QIN Z, XIAO S, et al. Efficient and secure decision tree classification for cloud-assisted online diagnosis services[J]. IEEE Transactions on Dependable and Secure Computing, 2021, 18(4): 1632-1644.
- [11] BOST R, POPA R A, TU S, et al. Machine learning classification over encrypted data[C]//Proceedings 2015 Network and Distributed System Security Symposium. VA: Internet Society, 2015: 4324-4325.
- [12] WU D J, FENG T, NAEHRIG M, et al. Privately evaluating decision trees and random forests[J]. Proceedings on Privacy Enhancing Technologies, 2016, 2016(4): 335-355.
- [13] TUENO A, KERSCHBAUM F, KATZENBEISSER S. Private evaluation of decision trees using sublinear cost[J]. Proceedings on Privacy Enhancing Technologies, 2019, 2019(1): 266-286.
- [14] 刘琳岚, 高声荣, 舒坚. 基于随机森林的链路质量预测[J]. 通信学报, 2019, 40(4): 202-211.
LIU L L, GAO S R, SHU J. Link quality prediction based on random forest[J]. Journal on Communications, 2019, 40(4): 202-211.
- [15] CHEN X, ZHU C C, YIN J. Ensemble of decision tree reveals potential miRNA-disease associations[J]. PLoS Computational Biology, 2019, 15(7): 1-24.
- [16] 李佳, 云晓春, 李书豪, 等. 基于混合结构深度神经网络的 HTTP 恶意流量检测方法[J]. 通信学报, 2019, 40(1): 24-33.
LI J, YUN X C, LI S H, et al. HTTP malicious traffic detection method based on hybrid structure deep neural network[J]. Journal on Communications, 2019, 40(1): 24-33.
- [17] LI F H, LI H, NIU B, et al. Privacy computing: concept, computing framework, and future development trends[J]. Engineering, 2019, 5(6): 1179-1192.
- [18] RAN C. Security and composition of multiparty cryptographic protocols[J]. Journal of Cryptology, 2000, 13(1): 143-202.
- [19] BRAKERSKI Z, GENTRY C, VAIKUNTANATHAN V. Fully homomorphic encryption without bootstrapping[J]. ACM Transactions on Computation Theory, 2014, 6(3): 1-36.
- [20] DAMGARD I, GEISLERR M, KROIGAARD M. Efficient and secure comparison for on-line auctions[C]//2007 Australasian conference on information security and privacy. Berlin: Springer, 2007: 416-430.
- [21] ERKIN Z, FRANZ M, GUAJARDO J, et al. Privacy-preserving face recognition[C]//International Symposium on Privacy Enhancing Technologies Symposium. Berlin: Springer, 2009: 235-253.
- [22] GOLDREICH O, MICALI S, WIGDERSON A. How to play any mental game, or a completeness theorem for protocols with honest majority[C]//Proceedings of the 19th Symposium on Theory of Computing. New York: ACM Press, 2019: 307-328.
- [23] KOLESNIKOV V, SADEGHI A R, SCHNEIDER T. Improved garbled circuit building blocks and applications to auctions and computing minima[C]//International Conference on Cryptology and Network Security. Berlin: Springer, 2009: 1-20.
- [24] BEAVER D. Efficient multiparty protocols using circuit randomization[C]//Annual International Cryptology Conference. Berlin: Springer, 1991: 420-432.
- [25] NAOR M, PINKAS B. Computationally secure oblivious transfer[J]. Journal of Cryptology, 2005, 18(1): 1-35.
- [26] ASHAROV G, LINDELL Y, SCHNEIDER T, et al. More efficient oblivious transfer and extensions for faster secure computation[C]//2013 ACM SIGSAC conference on Computer & communications security. New York: ACM Press, 2013: 535-548.

[作者简介]



马立川 (1988-), 男, 山东潍坊人, 博士, 西安电子科技大学讲师, 主要研究方向为信任管理机制、隐私保护、边缘计算安全等。



彭佳怡 (1997-), 女, 陕西西安人, 西安电子科技大学硕士生, 主要研究方向为隐私保护、安全多方计算等。

裴庆祺 (1975-), 男, 广西玉林人, 博士, 西安电子科技大学教授, 主要研究方向为认知网络、物联网与边缘计算安全、无线网络物理层安全、区块链技术、分布式协同攻防技术等。

朱浩瑾 (1980-), 男, 湖北武穴人, 博士, 上海交通大学教授, 主要研究方向为车联网安全、移动网络安全与隐私保护等。