

Enhanced Privacy-Preserving Decision Trees using Secure Multiparty Computation and Differential Privacy

Arisa Tajima¹ Wei Jiang² Virendra Marathe² Hamid Mozaffari²

¹University of Massachusetts Amherst, Amherst, MA, USA

`atajima@umass.edu`

²Oracle Labs – East, Burlington, MA, USA

`{wei.wj.jiang, virendra.marathe, hamid.mozaffari}@oracle.com`

Abstract

We address the problem of decision tree learning from data that may be distributed across multiple data owners while protecting the privacy of the training data and the model from each data owner. To protect the training data, Secure Multiparty Computation (MPC) provides a very promising privacy-preserving solution with none or minimum accuracy loss, but there is a potential privacy leakage on the training data during model prediction or inference. To this end, existing solutions use a hybrid framework that combines MPC and Differential Privacy (DP), where DP models are trained with MPC protocols. The existing approaches combine MPC and DP in a naive way and often lead to models with lower accuracy. In this work, to take the full advantage of MPC’s inherent security guarantee, we propose a novel way of utilizing both MPC and DP that can improve model accuracy while providing the same privacy guarantee. Our key design idea is to adopt MPC for building the entire model without leaking any intermediate results, and then carefully calibrated noise to enforce the DP guarantee (DP-noise) is only added at the leaf level to achieve the desired ϵ -DP. By doing so, less amount of noise is needed, and as a consequence, model accuracy can be significantly improved. We provide formal privacy proof of the proposed protocol and analyze the amount of required DP-noise. In addition, we implemented our protocol in a distributed environment, and our empirical results show that our approach can indeed improve model accuracy by up to 29% for the Adult dataset even with a small privacy budget of $\epsilon = 0.005$ comparing to the existing solution. However, our solution is computationally more expensive, and it trades off between accuracy and computation cost.

1 Introduction

Training machine learning (ML) models on distributed data while providing formal privacy guarantee regarding the training data brings a great value in various business settings. Without a privacy protection mechanism in place, data scientists may not be able to build any ML models. Suppose that data scientists want to build ML models on customer data from multiple organizations. Under the traditional and centralized solution, customer data need to be collected and managed by a central server. Data scientists interact with the server and perform the required computations to generate ML models. This centralized solution may not be always possible due to privacy risks and government regulations, such as the General Data Protection Regulation (GDPR) [30].

To address potential privacy risks, privacy-enhancing technologies have emerged as promising solutions that can help data scientists analyze customers’ data through the ML life-cycle without revealing private information in a distributed/federated environment. The common techniques include Federated Learning (FL) [21, 26, 43], Homomorphic Encryption (HE) [5, 14], Secure Multiparty Computation (MPC) [44] and Differential Privacy (DP) [9], where each of these techniques has pros and cons. In FL, while data never leave the data owners, the resulting model can be less accurate comparing to the centralized solution.

Additionally, intermediate model parameters may leak information about each party’s local training data. Recent works have studied secure training of decision trees under federated learning with HE and/or MPC. These techniques rely on cryptographic schemes and allow multiple parties to jointly train decision trees without revealing each party’s private data. Also, they can provide as good accuracy as the centralized solution. However, the output model may not necessarily be private and does not provide a formal DP guarantee. It is well known that models are vulnerable to membership inference attack [28, 37]. In other words, when models and/or prediction results are shared with data scientists or a third party, they can leak private information about individuals in the training data.

Independently, DP decision trees have extensively studied in the context of centralized learning [25] where decision trees are learned over training data by answering queries under DP. A resulting model that satisfies DP can be publicly released for inference without degrading privacy loss. Because noise is introduced during a model training to satisfy DP, there is a cost of accuracy loss. The accuracy is mainly affected by the total privacy budget, the number of queries required to build decision trees, and the sensitivity of the queries. Since the server, which produces the DP-trees, has access to the original data, the scheme does not work in a federated environment where the participating parties do not want to disclose their private training data.

In order to protect the training data during both training and inference, a hybrid framework was introduced in [41] that combines MPC and DP by training a centralized DP decision tree using MPC. However, the model accuracy provided by this approach is limited to what can be achieved by the centralized DP decision tree solutions [25, 41]. In addition, a synergy between MPC and DP is not well understood as the two techniques are often considered as orthogonal problems. In fact, we provide a novel way of combining MPC and DP that can provide a better accuracy than the existing approach under the same DP guarantee while offering data privacy and model confidentiality. Overall, the existing solutions do not satisfy one or more of the following properties:

- Data privacy: Each party’s private training data should not be disclosed during training and inference.
- Accuracy: The model should be as accurate as possible.
- Model confidentiality: This property is essential when the model has monetary value and disclosing the model may open the door for more effective inference attacks.
- Query privacy: Query content and prediction result should be hidden from the ML service providers since query content may contain sensitive information.

1.1 Problem Definition

We consider a cloud-assisted setting where data scientists build tree-based models using cloud computing nodes for clients who own private data. Prediction results are shared with data scientists. Our goal is to train a DP tree-based model as if the model were trained on a trusted server without participating parties sharing their private training data. The proposed protocol achieves all four properties defined previously. In order to maximize efficiency, we utilize a multi-server model where data owners/clients secretly share their training data with three independent computing servers. Model training is performed by the servers using secret shares, and the learned model is also secretly shared and stored at the servers. A user or data scientist secretly shares its query with the servers who perform the required computation to derive the prediction result in secret share format. The shares are sent to the user who then reconstructs the actual prediction result. To summarize, we consider the following three classes of entities:

- Clients C : Each client C_i holds private training data and secretly shares its data with the cloud servers.
- Servers S : The computing servers manage secretly shared data from the clients and collaboratively perform MPC-based training protocols. The servers also store the trained model and provide inference service to data scientists/users.
- Users U : A user issues a query to the servers and receives shares of the prediction result. A user may be affiliated with a client or is authorized to the inference service.

The proposed solution works for one or more clients. When there are multiple data owners, we assume that the data are either horizontally or vertically partitioned. The multi-server setting offers a high-level of scalability especially when a number of clients need to build an ML model on their aggregate data. Moreover, it works for either horizontal or vertical data partitioning scheme.

1.2 Threat Model

We consider the semi-honest model [15] where an adversary follows the protocol but may try to infer the other parties' private information based on its own input, output and the messages received during protocol execution. We also assume that the participating parties do not collude. On a high level, here we briefly emphasize how data privacy is achieved from the perspective of each entity group:

- **Clients:** Each client secretly shares its data with the servers and does not participate in any other computations related to model training and inference. Thus, it cannot learn any information about the other clients' training data.
- **Servers:** They only perform computations on secretly shared data, and the original data are never reconstructed at the servers. In addition, any intermediate results, accessible to the servers, are either secretly shared or randomized. As a consequence, the servers cannot learn any information about the training data, except for certain domain knowledge, e.g., the domain size of each attribute.
- **Users:** Because the prediction result is differentially private, the user cannot learn anything about the training data beyond what is allowed by the DP privacy budget.

1.3 Our Contribution

Adopting the multi-server setting, we design and implement MPC protocols for training tree-based ML models. Our protocols satisfy: data privacy, high model accuracy, and model confidentiality, and query privacy. Our solutions work for any number of clients and users, but it requires at least three independent computing servers to perform secret sharing based MPC computations. In summary, this work makes several significant contributions:

- We propose a novel way of combining MPC and DP for tree based ML models that requires less amount of noise and with provable DP guarantee.
- A formal privacy proof is provided to theoretically justify that our approach can achieve the same ϵ -DP by only adding noise at the leaf level.
- Inspired by the high level discussions given in [8], we develop concrete and fully secure ID3 tree generation and prediction protocols.
- Our protocols are implemented in a real distributed environment to generate extensive empirical data. The results are consistent with our design and theoretical analysis, e.g., up to 29% accuracy improvements on benchmarks like the Adult dataset.

The existing work most relevant to ours is discussed in Appendix A.

2 The Enhanced MPC+DP Framework

To completely hide all intermediate information about the data and keep the tree secret, every non-leaf node must have the same number of branches. This can be achieved by adding dummy values to each attribute so that every attribute $A_i \in \mathcal{A}$ has the same domain size of $w = \max_i |dom(A_i)|$. In other words, we will build a complete w -ary tree except for leaf nodes with $|dom(A_0)|$ labels where A_0 denotes the class attribute. Also, the stopping criterion only checks if the tree reaches the maximum depth.

Before executing the model generation protocol, each data owner needs to locally pre-process its data. We assume that data owners apply consistent pre-processing strategies, such as discretization, handling missing values, etc. In addition, data integration at the servers can be tricky. Under the vertical partition scheme, the data owners generate one-hot encodings for their attributes and then secretly share the encodings with the three computing servers. The servers need to store the shares in the same ordering which can be guaranteed easily even if they do not know the data schema and attribute domains. Detailed discussions on our proposed secure tree generation and prediction protocols are provided in Appendix B and Appendix C respectively.

3 Experimental Results

We empirically evaluate the model accuracy and efficiency of our proposed ODP+FMPC-ID3 protocol (Algorithm 1) comparing to the existing solution CDP+PMPC-ID3. We also implemented MPC-ID3 which offers identical accuracy as the non-secure ID3 and serves as the baseline for accuracy evaluation. All implementations were written in python, and we used the MPyC library [36] for basic MPC operations. In our experiments, we use the following three real-world datasets: Titanic [29], Heart [20] and Adult [3]. All datasets have binary class labels. For the utility evaluation, all experiments are performed over 5 different 80-20 train-test splits. We measure the accuracy and ROC AUC on the test set. To evaluate protocol efficiency, we run secure 3-party computations where all experiments are repeated at least 3 times on a 80-20 train-test split. We measure the runtime and the network traffic for training and inference. The network traffic refers to the total amount of messages each party received during the computation.

- **Accuracy Evaluation:** To evaluate accuracy, we mainly focus on three parameters: privacy loss ϵ , tree depth and DP mechanisms (Laplace vs. Exponential). The results are presented in Appendix D.
- **Efficiency Evaluation:** We show the runtime and the network traffic of training ID3 for our solution and the existing work under a fixed ϵ of 1.0. We varied the maximum tree depth from 2 to 5 for each dataset. The main results are given in Appendix E.

In summary, comparing to the existing solution, our proposed protocol provides between 22% to 29% accuracy improvement. However, our solution is computationally more expensive and incurs larger amount of network traffic.

4 Conclusion and Future Work

In this paper, we proposed a new approach for building a differentially private decision tree. Utilizing MPC to its max, the decision tree can be kept hidden even during model prediction. As a result, our solution achieves (1) data privacy, (2) high model accuracy, (3) model confidentiality and (4) query privacy. Furthermore, our solution can be naturally extended to random forest model.

The main drawback of our solution is its efficiency. As a future research direction, we will explore more efficient ways to construct a fully secure protocol for generating decision trees. One potential solution is to design customized MPC functionalities for our specific ML application. For example, since secure comparison has been used extensively in the sub-protocols, a more efficient secure comparison can greatly improve the overall protocol efficiency. In addition, the efficiency of the prediction protocol could be improved by directly using secret shares to represent attributes for the internal nodes instead of the vector indicator representation. However, the challenge is to develop an efficient and customized equality test or attribute selection sub-protocol. We will also investigate a tighter bound on DP guarantee of our random forest solution. By utilizing the additional randomness among the sub-samples, we may be able to derive a more precise DP bound for the random forest.

References

- [1] Mark Abspoel, Daniel Escudero, and Nikolaj Volgushev. Secure training of decision trees with continuous attributes. *Cryptology ePrint Archive*, 2020.
- [2] Borja Balle, Gilles Barthe, and Marco Gaboardi. Privacy amplification by subsampling: Tight analyses via couplings and divergences. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [3] Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: <https://doi.org/10.24432/C5XW20>.
- [4] Mariusz Bojarski, Anna Choromanska, Krzysztof Choromanski, and Yann LeCun. Differentially-and non-differentially-private random decision trees. *arXiv preprint arXiv:1410.6973*, 2014.
- [5] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. *ACM Trans. Comput. Theory*, 6(3), jul 2014.
- [6] Leo Breiman and JH Friedman. Random forests. *Classification and regression trees*, 1984.
- [7] Shorya Consul and SA William. Differentially private random forests for regression and classification. *Association for the Advancement of Artificial Intelligence*, 2021.
- [8] Sebastiaan De Hoogh, Berry Schoenmakers, Ping Chen, and Harm op den Akker. Practical secure decision tree learning in a tele-treatment application. In *Financial Cryptography and Data Security: 18th International Conference, FC 2014, Christ Church, Barbados, March 3-7, 2014, Revised Selected Papers 18*, pages 179–194. Springer, 2014.
- [9] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer, 2006.
- [10] Fabienne Eigner, Aniket Kate, Matteo Maffei, Francesca Pampaloni, and Ivan Pryvalov. Differentially private data aggregation with optimal utility. In *Proceedings of the 30th Annual Computer Security Applications Conference*, pages 316–325, 2014.
- [11] Sam Fletcher and Md Zahidul Islam. A differentially private random decision forest using reliable signal-to-noise ratios. In *AI 2015: Advances in Artificial Intelligence: 28th Australasian Joint Conference, Canberra, ACT, Australia, November 30–December 4, 2015, Proceedings 28*, pages 192–203. Springer, 2015.
- [12] Sam Fletcher and Md Zahidul Islam. Differentially private random decision forests using smooth sensitivity. *Expert systems with applications*, 78:16–31, 2017.
- [13] Arik Friedman and Assaf Schuster. Data mining with differential privacy. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 493–502, 2010.
- [14] Craig Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the Forty-First Annual ACM Symposium on Theory of Computing, STOC '09*, page 169–178, New York, NY, USA, 2009. Association for Computing Machinery.
- [15] Oded Goldreich. *The Foundations of Cryptography*, volume 2, chapter General Cryptographic Protocols. Cambridge University Press, 2004.
- [16] Koki Hamada, Dai Ikarashi, Ryo Kikuchi, and Koji Chida. Efficient decision tree training with new data structure for secure multi-party computation. *arXiv preprint arXiv:2112.12906*, 2021.

- [17] Naoise Holohan, Stefano Braghin, Pól Mac Aonghusa, and Killian Levacher. Diffprivlib: the IBM differential privacy library. *ArXiv e-prints*, 1907.02444 [cs.CR], July 2019.
- [18] Jun Hou, Qianmu Li, Shunmei Meng, Zhen Ni, Yini Chen, and Yaozong Liu. Dprf: a differential privacy protection random forest. *Ieee Access*, 7:130707–130720, 2019.
- [19] Geetha Jagannathan, Krishnan Pillaipakkamnatt, and Rebecca N Wright. A practical differentially private random decision tree classifier. In *2009 IEEE International Conference on Data Mining Workshops*, pages 114–121. IEEE, 2009.
- [20] Steinbrunn William Pfisterer Matthias Janosi, Andras and Robert Detrano. Heart Disease. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C52P4X>.
- [21] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- [22] Marcel Keller. Mp-spdz: A versatile framework for multi-party computation. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, CCS '20*, page 1575–1590, New York, NY, USA, 2020. Association for Computing Machinery.
- [23] Nhan Khanh Le, Yang Liu, Quang Minh Nguyen, Qingchen Liu, Fangzhou Liu, Quanwei Cai, and Sandra Hirche. Fedxgboost: Privacy-preserving xgboost for federated learning. *arXiv preprint arXiv:2106.10662*, 2021.
- [24] Qinbin Li, Zhaomin Wu, Zeyi Wen, and Bingsheng He. Privacy-preserving gradient boosting decision trees. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01):784–791, Apr. 2020.
- [25] Samuel Maddock, Graham Cormode, Tianhao Wang, Carsten Maple, and Somesh Jha. Federated boosted decision trees with differential privacy. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 2249–2263, 2022.
- [26] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In Aarti Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 20–22 Apr 2017.
- [27] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE, 2007.
- [28] Shagufta Mehnaz, Sayanton V. Dibbo, Ehsanul Kabir, Ninghui Li, and Elisa Bertino. Are your sensitive attributes private? novel model inversion attribute inference attacks on classification models. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 4579–4596, Boston, MA, August 2022. USENIX Association.
- [29] Vinicius Barbosa Paiva. The Complete Titanic Dataset. <https://www.kaggle.com/datasets/vinicius150987/titanic3>, 2020.
- [30] European Parliament and of the Council. Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (general data protection regulation), 2016.
- [31] Abhijit Patil and Sanjay Singh. Differential private random forest. In *2014 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, pages 2623–2630. IEEE, 2014.
- [32] Sikha Pentylala, Davis Railsback, Ricardo Maia, Rafael Dowsley, David Melanson, Anderson Nascimento, and Martine De Cock. Training differentially private models with secure multiparty computation. *arXiv preprint arXiv:2202.02625*, 2022.

- [33] Benny Pinkas, Thomas Schneider, Gil Segev, and Michael Zohner. Phasing: Private set intersection using permutation-based hashing. In *24th USENIX Security Symposium (USENIX Security 15)*, pages 515–530, Washington, D.C., August 2015. USENIX Association.
- [34] J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1:81–106, 1986.
- [35] J Ross Quinlan. Program for machine learning. *C4*. 5, 1993.
- [36] Berry Schoenmakers. MPyC: Multiparty Computation in Python. <https://www.win.tue.nl/~berry/mpyc/>, 2018.
- [37] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18, 2017.
- [38] K.K.A. Thissen. *Achieving Differential Privacy in Secure Multiparty Computation*. PhD thesis, Master’s Thesis, Technische Universiteit Eindhoven, Eindhoven, 2019.
- [39] Zhihua Tian, Rui Zhang, Xiaoyang Hou, Jian Liu, and Kui Ren. Federboost: Private federated learning for gbdt. *arXiv preprint arXiv:2011.02796*, 2020.
- [40] Rui Wang, Oğuzhan Ersoy, Hangyu Zhu, Yaochu Jin, and Kaitai Liang. Feverless: Fast and secure vertical federated learning based on xgboost for decentralized labels. *IEEE Transactions on Big Data*, 2022.
- [41] Yuncheng Wu, Shaofeng Cai, Xiaokui Xiao, Gang Chen, and Beng Chin Ooi. Privacy preserving vertical federated learning for tree-based models. *Proc. VLDB Endow.*, 13(12):2090–2103, jul 2020.
- [42] Bangzhou Xin, Wei Yang, Shaowei Wang, and Liusheng Huang. Differentially private greedy decision forest. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2672–2676. IEEE, 2019.
- [43] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Trans. Intell. Syst. Technol.*, 10(2), jan 2019.
- [44] Andrew C. Yao. How to generate and exchange secrets. In *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*, pages 162–167, 1986.
- [45] Samuel Yeom, Matt Fredrikson, and Somesh Jha. The unintended consequences of overfitting: Training data inference attacks. *CoRR*, abs/1709.01604, 2017.

A Related Work

Differentially private decision trees have been well studied in the centralized setting, including decision trees, random forests and XGBoost [4, 7, 11–13, 18, 24, 42]. However, in the centralized setting, a party who builds a model is assumed to have access to the entire training dataset. Orthogonally, MPC-based decision tree training has been studied where a model is learned on secret-shared data [1, 8, 16]. Hoogh et al. proposed MPC-based protocols for an ID3-based tree algorithm for cases of revealing and hiding the tree [8]. Recent works have focused on a secure CART-based tree algorithm while keeping the model secret [1, 16]. While the training data are protected, those protocols do not provide DP guarantees and thus are susceptible to membership inference attacks [37, 45].

The combination of MPC and DP for tree-based models is considered in several works, especially in the context of federated learning where DP is involved in various contexts. One common usage of DP is to replace cryptographic techniques for efficiency gain. This approach is computationally efficient, but it incurs a significant loss in utility due to the added noise to achieve local DP [23, 39].

Another usage is to train DP models in the federated setting to make them robust to membership inference attacks, which align with our motivation [25, 40, 41]. Maddock et al. provided a general framework of learning DP gradient boosted decision tree models in a federated setting [25]. Their framework uses secure aggregation-based MPC and DP. However, since the majority operations are performed locally at each party, DP techniques that require centralized evaluation, e.g., the Exponential mechanism, were not considered. This can be a limitation since many DP decision tree-based algorithms require the Exponential mechanism to provide a good accuracy. Our framework relies on a secret-sharing scheme and thus can completely simulate a centralized DP learning in a federated setting. Wu et al. briefly introduced how to train a centralized DP decision tree with an MPC framework in a federated setting [41]. Their approach of combining MPC and DP during tree generation process is very different from ours (adding noise at the leaf level and producing more accurate results), and it only works for vertically partitioned situation. Furthermore, they do not have an implementation and empirical evaluation of the hybrid framework of MPC and DP and thus efficiency of the framework was not clear.

Another essential and often overlooked aspect is that when analyzing DP guarantee, the existing work on DP-tree models neglects the fact that the tree structure itself can reveal information about the training data, and their DP analyses ignore the impact of the tree structure. Therefore, the existing solutions may not achieve the expected DP guarantee. As far as we know, there is no existing work that comprehensively analyzes the connection between MPC and DP for decision trees from the perspective of privacy, accuracy and efficiency. This is because MPC and DP are often considered orthogonal problems. However, we will show that the existing way of combining MPC and DP can result in sub-optimal performance. In addition, by taking advantage of the inherent privacy guarantee of an MPC protocol, we will theoretically demonstrate that adding noise to the leaf nodes (and less noise overall) is sufficient to achieve the same level of DP while simultaneously improving accuracy comparing to the existing approaches.

B The Tree Generation Protocol

We adopt the following notation conventions:

- $[x]$: secret shares of x , and we use bold lower case letter to represent a vector: $[\mathbf{x}_i]$ represents the shares of its i -th component of vector \mathbf{x} . We also use capital letter to represent a multi-dimensional data or a vector.
- D : the dataset represented as a set of one-hot encodings. D is three-dimensional: $D_{i,j,k}$ represents the k -th component of the j -th one-hot encoding of the i -th attribute. In other words, D_i is matrix and a collection of one-hot encodings for attribute A_i .
- T : a binary vector representing if a tuple is in a dataset.
- R : a set of attribute indices (excluding class attribute index 0) and \mathcal{R} a binary vector of size $|R|$ indicating the available attributes for the current recursive call.
- h : the tree depth is bounded by $1 \leq h \leq |\mathcal{A} - \{A_0\}|$.
- $[\mathbf{x}] \bullet [\mathbf{y}]$: represent a secure dot product between two secretly shared vectors.
- $[\mathbf{x}] \times [\mathbf{y}]$: represent secure component-wise multiplications between two secretly shared vectors.

The key steps of our tree generation protocol (Algorithm 1) include:

- **Compute the class partitions and frequencies** (Steps 5-8): derive a binary vector C_i for each class i . $C_{i,j} = 1$ implies the j -th tuple belongs to class i , and 0 otherwise. These vectors are used to derive class frequencies which are stored in τ . Each frequency is DP randomized.
- **Check the stopping condition and derive the majority class** (Steps 9-11): When h reaches 0, the protocol computes the majority class returned as the leaf label.

- **Compute the quality scores** (Steps 12-14): the C_i vectors derived from the previous steps are used by the MAX protocol (Algorithm 2). At step 14, the scores for the available attributes are kept unchanged, but the scores for the other attributes are set to 0 so that these attributes would not be chosen for the current node.
- **Select the best attribute and its encoding** (Steps 15-19): the arg max protocol returns the index of an attribute whose quality score is the maximum. At step 16, the index is converted into a vector representation using the Indicator protocol (see Algorithm 3 for detail). Then one hot encodings of A_k are retrieved and stored in X .
- **Partition the data and recursively build sub-trees** (Steps 20-22): At step 21, $[T] \times [X_i]$ produces a binary vector that represents the data partitioned based on $A_{k,i}$. $[\mathcal{R}] - [k]$ means that the chosen attributed is removed from further consideration. Then a sub-tree is built with each data partition with tree height reduced by one.

Note that the dummy values are represented as zero vectors. As a result, they have no effect on the MAX protocol, and the quality scores for all attributes are unaffected by the dummy values. Although the resulting tree model has branches generated from these dummy values, they do not interfere with or affect the result of model inference tasks according to our prediction protocol. Section C provides more details.

Quality metric. There are various quality metrics that determine the best split. While popular metrics are information gain or Gini index, some works explore the median or the max operator. In the context of DP, since the choice of quality metrics affects the model accuracy, the max operator, which has a lower sensitivity, has shown to have a better performance than the other metrics [13]. For this reason, we consider the max operator as the quality metric for choosing the best split attribute. The max operator corresponds to the misclassification rate by picking the class with the highest frequency, which has the sensitivity of 1:

$$q(D, A_i) = \sum_{j \in A_i} \max_{c \in A_0} (D_{A_i,j} \bullet D_{A_0,c}) \quad (1)$$

The MAX protocol presented in Algorithm 2 is our MPC implementation of the max operator according to Equation 1.

B.1 Complexity Analysis

For MPC protocols, the asymptotic complexity is often based on the number of secure multiplications (SM). Since secret sharing based secure addition does not require any communication, its complexity is negligible comparing to SM. First, we analyze the complexity for each sub-protocol, and let l represent the share size in bits:

- $\arg \max([x])$: this protocol can be implemented using $|[x]| - 1$ secure comparisons (SC), each of which has a complexity of $O(l)$. Thus, the complexity for $\arg \max$ is $O(l|[x]|)$.
- $\text{MAX}([D_i], [C], u, w) \rightarrow [\gamma_i]$: since C_i has a size of n , each secure dot product requires n SMs. There are u dot products. The max protocol (step 2 of Algorithm 2) can be implemented using $u - 1$ SCs. As a result, the total complexity for the MAX protocol is $O((n + l)wu)$.
- $\text{Indicator}(m, [k]) \rightarrow [k]$: the protocol requires m SCs, so its complexity is bounded by $O(lm)$.

Next we estimate the complexity of each major component:

- Steps 5-8: there are u component-wise secure multiplications with size n each, so that complexity is $O(nu)$.
- Steps 9-11: the size of τ is u ; thus, the complexity is bounded by $O(lu)$.

- Steps 12-13: since MAX is called m times, the complexity is $O(m(n+l)wu)$.
- Step 14: requires $O(m)$ SM operations.
- Step 15: the size of γ is m ; thus, the complexity is bounded by $O(lm)$.
- Step 16: the complexity is $O(lm)$.
- Steps 17-19: there are wn secure dot products of size m , so the complexity is bounded by $O(mnw)$.
- Steps 20-22: there are w recursive calls, and before each call, a secure component-wise vector multiplication is performed. As a result, the complexity is bounded by $O(nw)$ before performing the recursive calls.

It is clear that the complexities at steps 12-13 dominate the rest of the protocol before each recursive call. In addition, since the vector sizes do not change for all recursive calls, the total complexity is bounded by: $O(w^{h-1}m(n+l)wu)$, for $1 \leq h \leq m$.

B.2 Privacy/Security Analysis

As stated before, our solution achieves (1) data privacy, (2) high model accuracy, (3) model confidentiality and (4) query privacy. For high model accuracy, we will demonstrate it in two ways: theoretically prove less noise is needed to achieve the same level of DP, and empirically justify the accuracy improvement using three real world datasets. For the others, we need to prove the following:

- The training data are not disclosed when building the tree.
- The final tree model and user queries are never disclosed during prediction.
- The tree is ϵ -DP.

The first two are straightforward to prove since we adopt the standard MPC functionalities to implement the MPC-DP mechanisms and the MPC-decision tree generation. Intuitively speaking, nothing is disclosed during model generation and prediction, except that a user obtains the final prediction result. As long as the servers follow the protocol and the secret sharing scheme is secure, so does our protocol. On the other and, since our way of combining DP and MPC is new, we need to prove the resulting protocol satisfies ϵ -DP. The key idea is to show that when the tree structure is hidden, adding less amount of noise to the leaf nodes is sufficient to achieve the same level of DP comparing to the existing solution where DP noises are also added to the attribute selection process for the internal nodes. Let T_1 and T_2 are two decision trees with the following properties:

- T_1 : The tree structures are completely hidden including the leaf labels. In addition, before leaf labels are derived, noises are securely added to achieve ϵ -DP. During prediction, the users only learn the predicted result. The entire process can be achieved using MPC. (Note that T_1 may be more private than ϵ -DP since the leaf labels are hidden.)
- T_2 : The tree structures are known but noises are added during tree generation process to achieve ϵ -DP. More specifically, MPC techniques are adopted to compute the scores for each attribute. Using MPC, noises are securely added to the scores before selecting the best split or attribute. After the selection, the chosen attribute is disclosed. Then the process repeats to build the rest of the tree. For prediction purpose, the DP-tree can be shared with the users.

It is well-known that when an adversary has black-box access to a decision tree model, the adversary can reconstruct the tree with very high accuracy by using a number of queries and their prediction results. As a result, we may ask that even though T_1 is completely hidden from anyone, is it still possible for an adversary to reconstruct another tree \hat{T}_1 (with black-box access to T_1) that is less private than T_1 or T_2 ? The short

answer is no because when deriving T_1 , we followed the sequential composition theorem and all intermediate computations and results are implemented using MPC protocols. If \hat{T}_1 is less private than T_1 or T_2 , then either the composition theorem is incorrect or the MPC protocol leaks non-negligible information. A more formal analysis is given next.

B.2.1 T_1 is ϵ -Differentially Private

Here we prove that adding ϵ -DP noise at each leaf node independently guarantees T_1 is ϵ -DP. Let f be the mechanism that securely selects the best attribute for the internal nodes of T_1 :

- $f([\vec{D}], [\vec{\Lambda}]) \rightarrow [\vec{D}_l], [\vec{D}_r], [\vec{\Lambda}_l], [\vec{\Lambda}_r], [a], [t_a]$

Without loss of generality, we assume each attribute splits the data into two partitions: left and right denoted by \vec{D}_l and \vec{D}_r respectively. All \vec{D} , \vec{D}_l and \vec{D}_r are binary vectors with the same size bounded by the size of the dataset, and the $[\]$ notation indicates the vectors are secretly shared component-wise. $\vec{\Lambda}$, $\vec{\Lambda}_l$ and $\vec{\Lambda}_r$ are also binary vectors with the same size determined by the number of attributes. The chosen attribute is denoted by a with threshold t_a . All inputs and outputs of f are secretly shared, and the height of the tree is fixed to a public parameter h indicating the tree has $0, 1, \dots, h$ layers. \vec{D}_l and $\vec{\Lambda}_l$ (or \vec{D}_r and $\vec{\Lambda}_r$) are inputs to build the left (or right) branch of the sub-tree rooted at a . $[a]$ and $[t_a]$ are the actual output of each tree layer. Let g be a function for the leaf node:

- $g([\vec{D}]) \rightarrow [c]$

where \vec{D} represents a set of tuples at the current leaf node and c is the class label for this node. The class label is chosen using the Laplace mechanism with privacy budget ϵ and sensitivity 1. Both \vec{D} and c are secretly shared. Based on f and g , we define a layer function L as follows:

- $L^k = \langle f_1^k, \dots, f_{2^k}^k \rangle$, for $k \in \{0, h-1\}$.
- $L^h = \langle g_1, \dots, g_{2^h} \rangle$.

Each layer is a parallel composition of either the f or the g functions. Next we show f is 0-DP. Suppose \vec{D} and \vec{D}' represent two neighboring datasets with the same vector size. (Note that we can always make them the same size by adding a dummy entry to the smaller vector.) The f function produces the following outputs on the two neighboring datasets \vec{D} and \vec{D}' :

- $f([\vec{D}], [\vec{\Lambda}]) \rightarrow [\vec{D}_l], [\vec{D}_r], [\vec{\Lambda}_l], [\vec{\Lambda}_r], [a], [t_a]$
- $f([\vec{D}'], [\vec{\Lambda}']) \rightarrow [\vec{D}'_l], [\vec{D}'_r], [\vec{\Lambda}'_l], [\vec{\Lambda}'_r], [a'], [t_{a'}]$

The output of f are secret shares from the same domain \mathcal{F} . Since secret shares are uniformly random in \mathcal{F} , the adversary cannot tell the difference between the two outputs. That is:

- $\Pr(f([\vec{D}], [\vec{\Lambda}]) \in \mathcal{F}) = \Pr(f([\vec{D}'], [\vec{\Lambda}']) \in \mathcal{F})$

This implies that $\epsilon = 0$ and f is 0-DP. In addition, since the layer function L^k is a parallel composition of independent calls of f , the output distributions of L^k on two neighboring datasets are the same. Therefore, L^k is also 0-DP. The same analysis is applicable for g and L^h , and we conclude L^h is 0-DP. Let p be the prediction function:

- $p([T_1], [\tau]) \rightarrow c_\tau$

where $[T_1]$ indicates the secret shares of T_1 , a collection of outputs from each layer function L^k, L^h , and τ is a secretly shared user record. c_τ is the classification result only known to the user. Because the class labels are ϵ -DP, based on the robustness to post-processing property of DP, the prediction result is ϵ -DP, and so does T_1 .

Another important consequence of the above proof is that since noise is only added at the leaf level, the total amount of noise added to the tree in our approach is significantly less than the existing solution. The empirical evaluation is presented in Section 3.

Algorithm 1 ODP+FMPC-ID3($[D], [T], [\mathcal{R}], R, h, \epsilon \rightarrow \{[k], \{[tr_1], \dots, [tr_w]\}\}$)

```

1:  $u \leftarrow |[D_0]|$  ▷ the number of classes
2:  $m \leftarrow [R]$  ▷ the number of attributes
3:  $n \leftarrow |[D_{0,0}]|$  ▷ the number of tuples
4:  $w \leftarrow |[D_1]|$  ▷ the domain size of each attribute
5: for  $i \leftarrow 1$  to  $u$  do
6:    $[C_i] \leftarrow [D_{0,i}] \times [T]$  ▷ tuples belong to class  $i$ 
7:    $[\tau_i] \leftarrow \sum_{j=1}^n [C_{i,j}]$  ▷ the size of class  $i$ 
8:    $[\tau_i] \leftarrow \text{MPC-LM}([\tau_i], \epsilon, 1)$  ▷ DP class sizes
9: if  $h = 0$  then
10:   $[c_{i*}] \leftarrow \arg \max([\tau])$  ▷ the class with the largest size
11:  return  $\{[c_{i*}]\}$  ▷ the leaf label
12: for  $i \leftarrow 1$  to  $m$  do
13:   $[\gamma_i] \leftarrow \text{MAX}([D_i], [C], u, w)$  ▷ the  $i$ -th attribute's score
14:   $[\gamma] \leftarrow [\gamma] \times [\mathcal{R}]$  ▷ keep the scores of available attributes
15:   $[k] \leftarrow \arg \max([\gamma])$  ▷ the attribute index with the max score
16:   $[k] \leftarrow \text{Indicator}(m, [k])$  ▷ converting  $k$  into a vector format
17: for  $i \leftarrow 1$  to  $w$  do
18:   for  $j \leftarrow 1$  to  $n$  do
19:     $[X_{i,j}] \leftarrow [D_{1,i,j}, \dots, D_{m,i,j}] \bullet [k]$  ▷ encoding of  $A_k$ 
20: for  $i \leftarrow 1$  to  $w$  do
21:   $[T_i] \leftarrow [T] \times [X_i]$ 
22:   $[tr_i] \leftarrow \text{ODP+FMPC-ID3}([D], [T_i], [\mathcal{R}] - [k], R, h - 1, \epsilon)$ 
23: return  $\{[k], \{[tr_1], \dots, [tr_w]\}\}$ 

```

Algorithm 2 MAX($[D_i], [C], u, w \rightarrow [\gamma_i]$)

Require: $[D_i]$ is a set of w one-hot encodings of A_i . $[C]$ is a set of binary vectors, and $[C_j]$ indicates if a tuple belongs to class j . w defines the domain size of all attributes, and u indicates the number of classes.

```

1: for  $j \leftarrow 1$  to  $w$  do
2:   $[\tau_j] \leftarrow \max([D_{i,j}] \bullet [C_1], \dots, [D_{i,j}] \bullet [C_u])$ 
3:  $[\gamma_i] \leftarrow \sum_{j=1}^w [\tau_j]$ 
4: return  $[\gamma_i]$ 

```

Algorithm 3 Indicator($m, [k] \rightarrow [k]$)

Require: m specifies the number of attributes, and k is a positive integer in $\{1, \dots, m\}$.

```

1: for  $i \leftarrow 1$  to  $m$  do
2:   $[k_i] \leftarrow \text{Equal}(i, [k])$  ▷ check if  $k$  is equal to  $i$ 
3: return  $[k]$ 

```

C The Prediction Protocol

Let $\mathcal{A} = A_0, A_1, A_2, A_3, A_4, A_5$ where A_0 is the class attribute. Ignoring the class attribute, the rest can be represented as indicator vectors $\{10000, 01000, 00100, 00010, 00001\}$ from A_1 to A_5 respectively. Suppose that each attribute A_i has three values represented by $\{100, 010, 001\}$, corresponding to $\{a_{i1}, a_{i2}, a_{i3}\}$ respectively or the tree branches from left to right. Let $\Gamma = \{\Gamma_1, \Gamma_2\} = \{\Gamma_1, \{\Gamma_{21}, \Gamma_{22}, \Gamma_{23}\}\}$ represent a decision tree where Γ_1 is the root attribute and Γ_2 is a set of children nodes of Γ_1 . If Γ_1 is a leaf node, then it is set

to c , one of the class labels; otherwise, it is defined recursively. As an example, a partial tree is given in Figure 1. The complete tree is full 3-nary tree; that is, each internal node has three branches. Due to space limitation, we only present part of it. According to the figure, Γ_1 represents A_2 as the chosen root of Γ , then Γ_{21} , Γ_{22} and Γ_{23} correspond to the sub-trees rooted at A_1 , A_3 and A_5 respectively. Γ_1 is represented using A_2 's indicator vector 01000. Let a tuple $t = [2, 3, 3, 1, 1]$ represented as a 3-by-5 matrix:

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

t_i indicates the i -th row of the matrix, and $|t_i|$ gives the number of attributes. The number of rows in t is denoted by w specifying the domain size of each attribute. In this example, $|t_i| = 5$ and $w = 3$.

The prediction protocol is given by Protocol 4. It derives the class label of a secretly shared record t

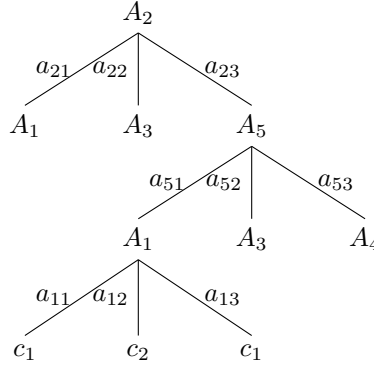


Figure 1: A partial decision tree

Algorithm 4 Predict($[t], [\Gamma]$) $\rightarrow [c]$

Require: $[t]$ secret shares of a tuple being classified, $[\Gamma]$ secret shares of the decision tree, and w denotes the number of attributes.

- 1: **if** $||[\Gamma]|| = 1$ **then**
 - 2: **return** $[\Gamma]$
 - 3: $[\Gamma_1], [\Gamma_2] \leftarrow [\Gamma]$
 - 4: **for** $i \leftarrow 1$ to w **do**
 - 5: $[f_i] \leftarrow [t_i] \bullet [\Gamma_1]$
 - 6: **for** $i \leftarrow 1$ to w **do**
 - 7: $[c_i] \leftarrow \text{Predict}([t], [\Gamma_{2i}])$
 - 8: $[c] \leftarrow [f] \bullet [c]$
 - 9: **return** $[c]$
-

recursively. The tree Γ is also secretly shared based on its smallest elements. For example, Γ_1 is secretly shared component-wise of A_2 's indicator vector at the top level of the tree. This recursively applies to the other levels. At the leaf level, the class label c is secretly shared. The tree structure, a complete w -nary tree, is preserved, so that we can easily check if Γ is a tree or a leaf node. At step 1 of the protocol. $||[\Gamma]|| = 1$ indicates that $[\Gamma]$ is a leaf node, and it is returned as the leaf label. Otherwise, $[\Gamma]$ is decomposed into two components $[\Gamma_1]$ and $[\Gamma_2]$. At steps 4-5, the protocol retrieves the correct attribute at the current root whose values serve as flags to indicate the right prediction path. At steps 6-7, the protocol makes w recursive calls on the sub-trees. At the end, the results from each sub-recursive calls are combine to produce the shares of

the class label. The key observation is that there is only one path in the tree whose associated f_i values are all one. Thus, the results from all other paths are 0, and the result from the correct path is preserved.

D Accuracy Evaluation

To evaluate accuracy, we mainly focus on three parameters: privacy loss ϵ , tree depth and DP mechanisms (Laplace vs. Exponential).

D.1 Varying Privacy Loss

We evaluate the accuracy and the ROC AUC of different variants of the ID3 algorithm. Figure 2 shows the relationship between model utility and the privacy loss ϵ with various datasets under the fixed tree depth. For Titanic, Heart, and Adult, the maximum depths are set to 4, 3, and 5, respectively. Regarding the privacy loss parameter, for Titanic and Heart, the values of ϵ are varied from 0.2 to 2.0. For Adult, a larger dataset, the values of ϵ are varied from 0.005 to 1.0. We did not use smaller ϵ values for Titanic and Heart because ϵ less than 0.2 would make the accuracy below 0.5 which causes the model useless.

The plots compare the accuracy of our proposed solution ODP-FMPC-ID3 against that of the existing solution CDP-PMPC-ID3 under the same ϵ -DP guarantee, benchmarking with MPC-ID3. Note that the accuracy of MPC-ID3 matches that of non-private centralized learning. As observed in Figure 2, the accuracy of our solution consistently performs better. For example, when $\epsilon = 0.2$, our solution produced an accuracy of 0.71 and ROC AUC of 0.71 for Heart. Whereas, under the same ϵ , the existing solution generated an accuracy of 0.575 and ROC AUC of 0.572. This amounts to 23.5% and 24.1% improvements respectively. For the Adult dataset, $\epsilon = 0.005$, our solution produced an accuracy of 0.755 and ROC AUC of 0.65. However, under the same ϵ , the existing solution only generated an accuracy of 0.585 and ROC AUC of 0.525. Thus, our solution amounts to 29.1% and 23.8% improvements. For Adult dataset, when ϵ is between 0.1 and 1.0, our approach even outperforms the non-private centralized learning. This would be the effect of the randomness introduced by DP, acting as regularization. Our solution performs better as ϵ gets smaller.

D.2 Varying Tree Depth

We also evaluate the model accuracy with different tree depths, which is shown in Figure 3. We measure ROC AUC under a fixed ϵ while varying the tree depths from 2 to 5. We use $\epsilon = 0.2, 0.2, 0.1$ for the three datasets Titanic, Heart, and Adult respectively. The performance trends are the same for other ϵ values and the accuracy metric.

In DP decision trees, the tree depth affects model accuracy. As a tree becomes deeper, training samples at each leaf node become very small and thus can hinder accuracy. Because DP leaf updates are based on the frequency of the class counts, smaller counts are susceptible to DP noise. As seen in Figure 3, since Titanic and Heart are smaller datasets, increasing the tree depth greatly lowers the model accuracy. However, for Adult dataset, since the data size is much bigger, under the chosen $\epsilon = 0.1$, increasing the tree depth from 2 to 5 does not degrade the accuracy. It simply requires a deeper tree to learn a model. On the other hand, this is not the case when using much smaller values of ϵ .

As also shown in Figure 3, our solution consistently provides a better accuracy for each of the maximum tree depths. In the existing solution, the privacy loss budget is split across different heights of tree. Thus, each node of the tree gets a budget of $\epsilon/(h + 1)$, where h is the maximum depth of the tree. Whereas, for our proposed approach, the full privacy budget ϵ is used for each leaf node. As a result, our solution is more robust to the maximum depth.

D.3 Comparison between Laplace and Exponential

In the previous results, we use the Laplace mechanism for DP leaf updates for both approaches. Here, we additionally show the model accuracy when using the Exponential mechanism. While Fletcher et al.

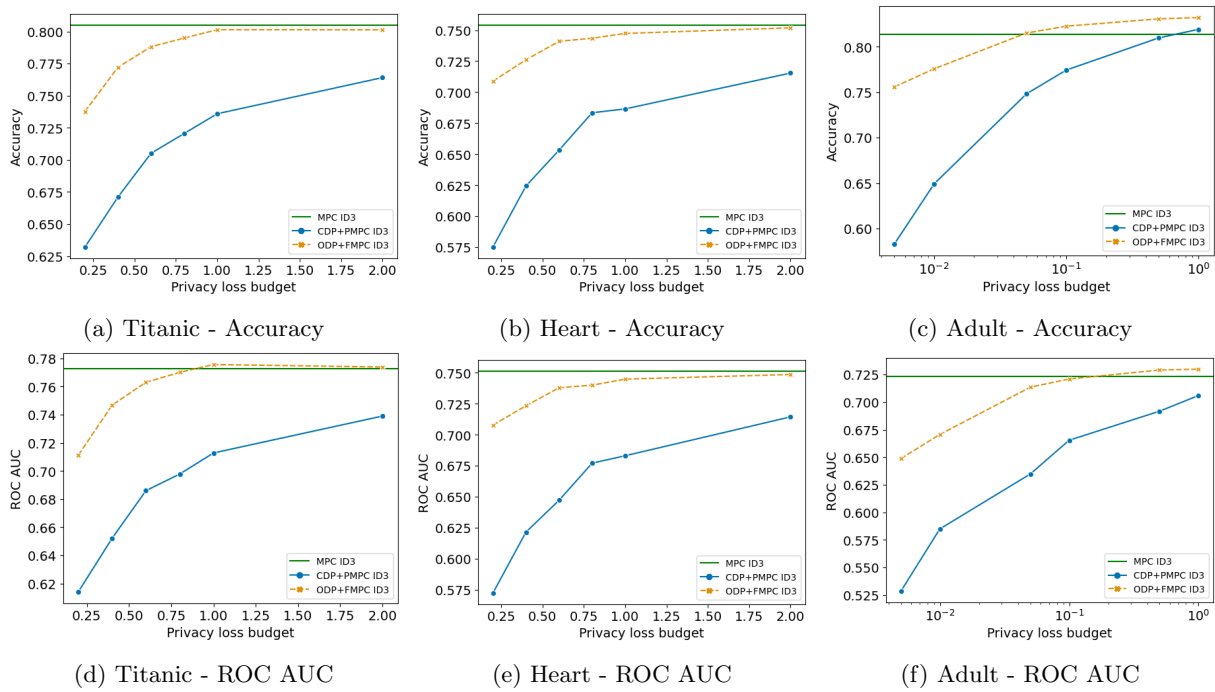


Figure 2: Model accuracy vs. privacy loss budget ϵ

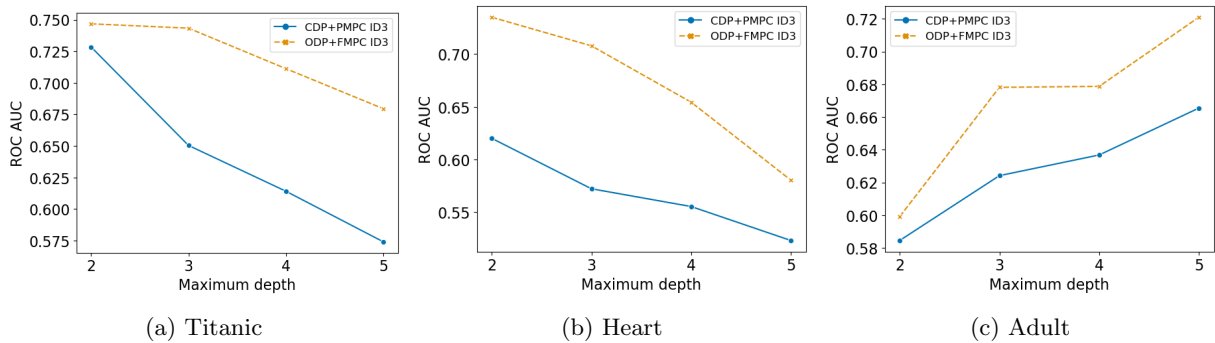


Figure 3: Model accuracy vs. tree depth

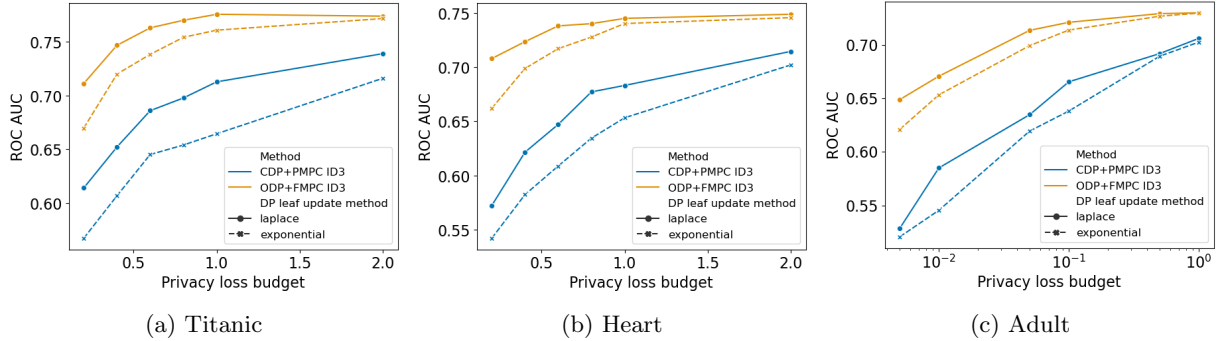


Figure 4: Accuracy comparison on different DP leaf update methods

suggested an advantage of the Exponential mechanism against the Laplace mechanism for multi-class random forest models [12], as far as we know, there is no consensus on which DP algorithm performs the best. In fact, recent works use the Laplace mechanism for leaf updates [7, 18].

Figure 4 compares the accuracy of the Laplace mechanism and that of the Exponential mechanism, varying privacy loss ϵ from 0.2 to 2.0 for Titanic and Heart and 0.005 to 1.0 for Adult. Similar to Figure 2, the maximum depth is fixed for each plot, where $h = 4, 3, 5$ for Titanic, Heart, and Adult respectively. For both approaches, the Laplace mechanism provides a better accuracy than the Exponential mechanism by up to 5%. While with large values of ϵ , there is no huge difference between the two DP mechanisms, with small values of ϵ , the Laplace mechanism consistently outperforms the Exponential mechanism.

E Efficiency Evaluation

We show the runtime and the network traffic of training ID3 for our solution and the existing work under a fixed ϵ of 1.0 in Figure 5. We varied the maximum tree depth from 2 to 5 for each dataset. For each approach, we have the two variants of the leaf update method, the Laplace mechanism-based update and the Exponential mechanism-based update.

Overall, our solution is computationally more expensive for both runtime and network traffic. It relies on the full MPC-based ID3 functionality, which is computationally more expensive than the partial MPC-based ID3 used in the existing work. Especially, our approach builds a completely oblivious tree where every non-leaf node has w branches, where w is the maximum domain size. Thus, the runtime and the network traffic greatly increase as the depth of the tree increases. This is the tradeoff between accuracy and efficiency, where a high accuracy observed in the enhanced framework was achieved by the computationally expensive full MPC protocol that can hide the model. Building efficient MPC-based decision trees are independent of our work, and it is still an active research area.

In addition, Figure 5 compares the efficiency of the MPC-based Laplace mechanism and the MPC-based Exponential mechanism. While there was no clear distinction for the existing approach, we observed that for our solution, the Laplace mechanism incurs less runtime and network traffic than the Exponential mechanism. The gap becomes very noticeable as we increase the maximum depth. The MPC-based Exponential mechanism requires a number of MPC operations, including a very costly secure exponentiation evaluation. The detailed cost analysis on MPC can be found in [10].

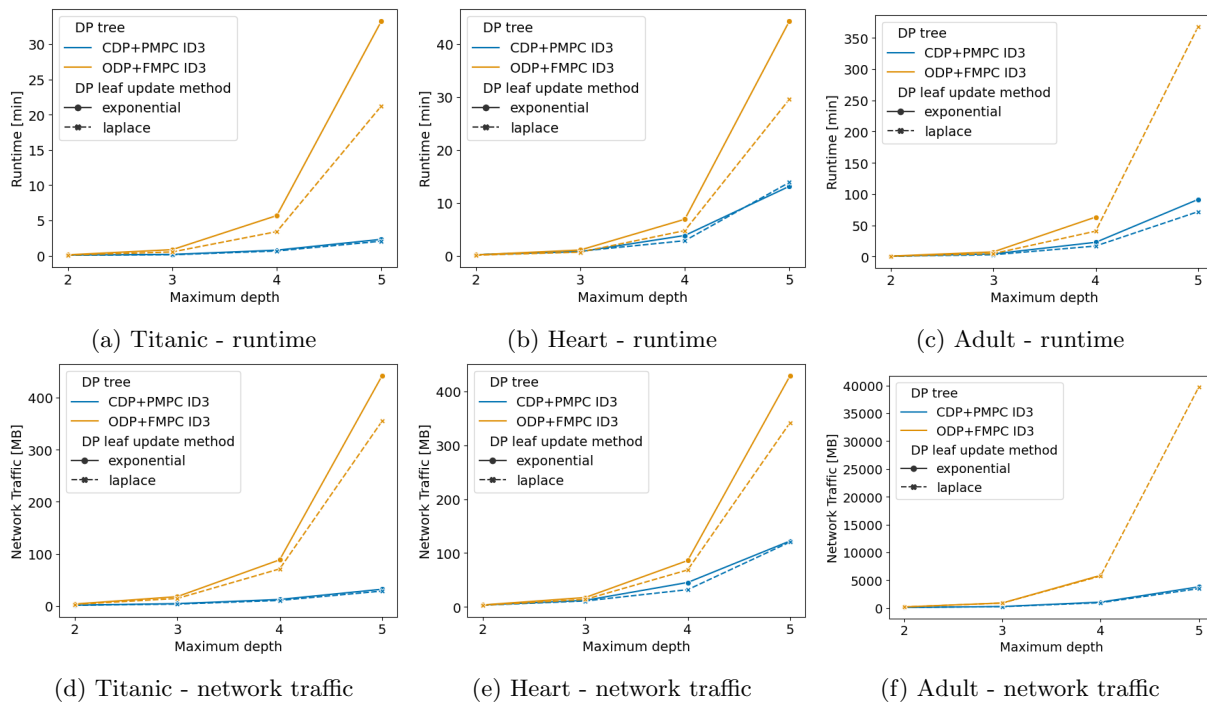


Figure 5: Efficiency vs. tree depth