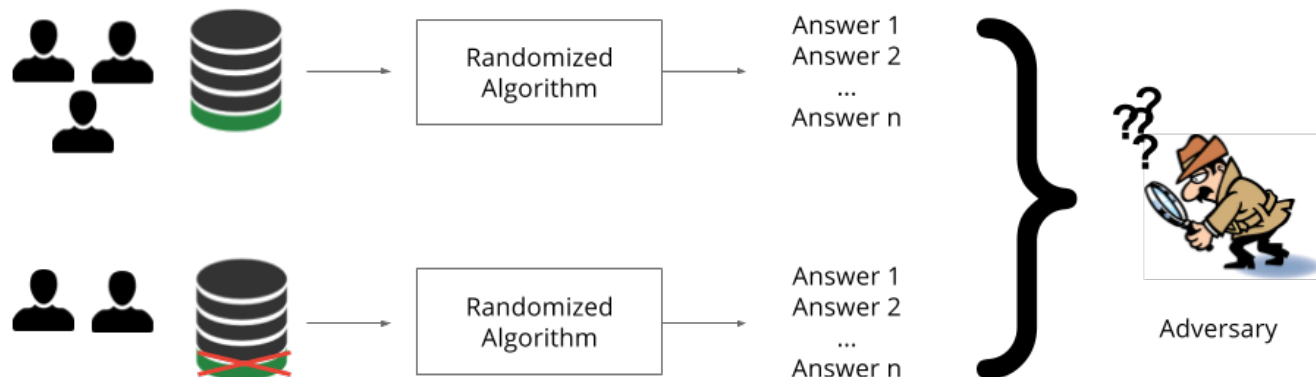


Privacy-Preserving Gradient Boosting Decision Trees

Kobe University Engineering ES5
M2 201T264T
Fuki Yamamoto

Background

- ✓ Recently, privacy concerns are growing
 - (e.g., inference of original data from statistics, leakage of sensitive info etc.)
 - Related legislation (e.g. General Data Protection Regulation (GDPR))
 - Balancing privacy protection and data utilization is required
 - Privacy-Preserving Data Mining
- ✓ Differential privacy (DP) provides privacy-preserving statistics
 - (Other privacy-preserving techniques protect only the computational process)



Background

- ✓ Previous DP algorithms in GBDT suffer from serious accuracy loss (GBDT is a ML model used in various domain)
- ✓ In this study, the accuracy loss is improved by focusing on **bounding of sensitivity** and **privacy budget allocation**

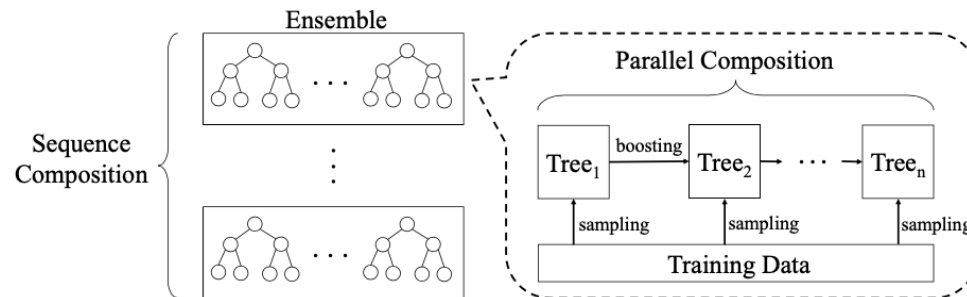
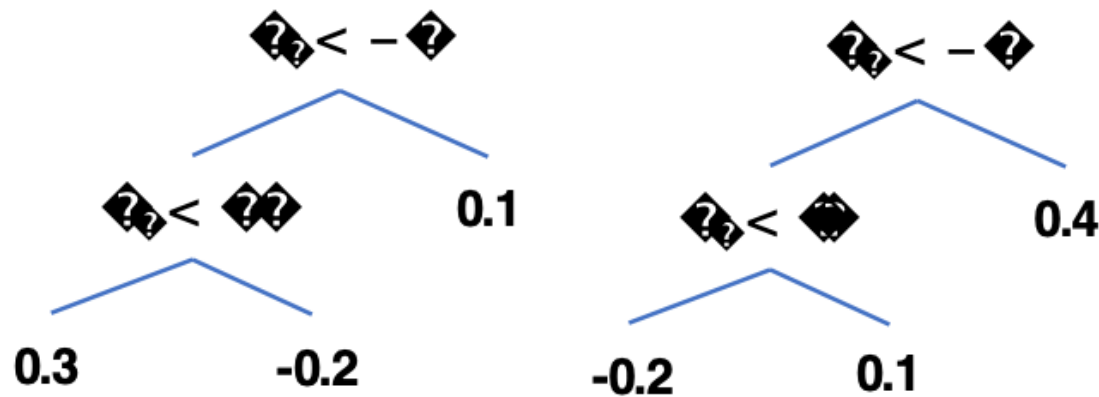


Figure 1: The two-level boosting design of DPBoost

Gradient Boosting Decision Trees (GBDT)

- ✓ Ensemble model with decision tree as the base model
- ✓ Training by gradient boosting, which improves the cost function by sequentially building the base model
- ✓ High efficiency and predictive performance, useful in various fields



Red: $0.3 + 0.4 = 0.7 > 0 \Rightarrow \text{Yes}$

Purple: $0.1 - 0.2 = -0.1 < 0 \Rightarrow \text{No}$

GBDT: Training

Cost function (for t -th tree)

$$\mathcal{L}^{(t)} = \sum_{i=1}^n \left(g_i f_t(x_i) + \frac{1}{2} f_t^2(x_i) \right) + \Omega(f_t)$$

Leaf Weight

$$V(I) = - \frac{\sum_{i \in I} g_i}{|I| + \lambda}$$

Splitting

$$G(I_L, I_R) = \frac{\left(\sum_{i \in I_L} g_i \right)^2}{|I_L| + \lambda} + \frac{\left(\sum_{i \in I_R} g_i \right)^2}{|I_R| + \lambda}$$

I_L, I_R : instance sets of
left and right nodes

$$I: I_L \cup I_R$$

$$g_i: \partial_{\hat{y}} l(y_i, \hat{y})$$

l : loss function (square error)

\hat{y} : prediction
by previous trees

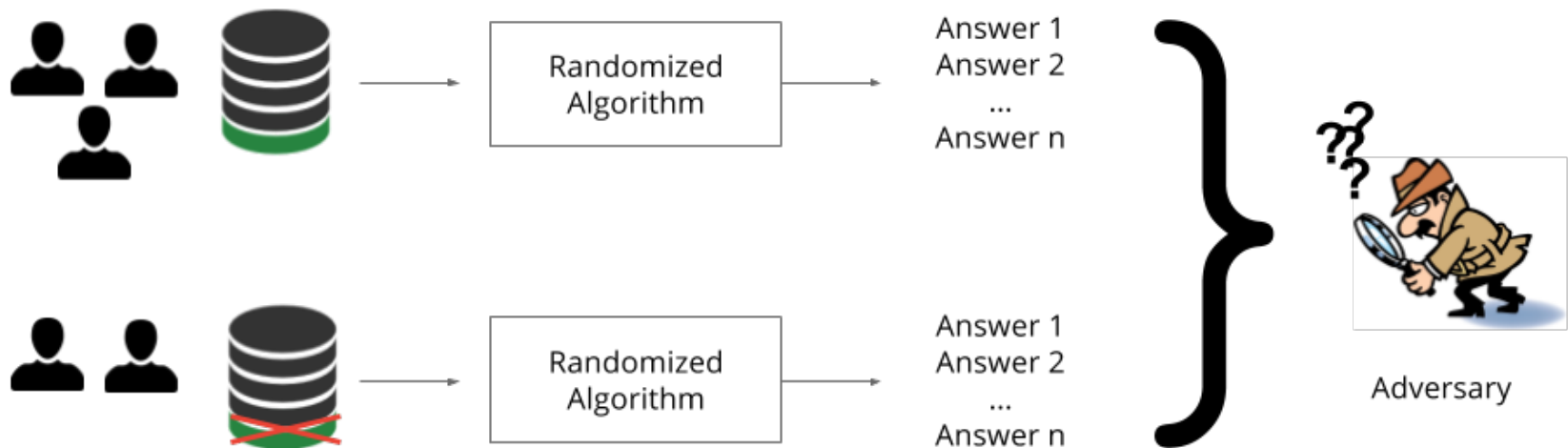
$f_t(\cdot)$: output of t -th tree

$\Omega(\cdot)$: regularization term

λ : regularization parameter

Differential Privacy (DP)

- ✓ Powerful and quantitative **privacy-preserving technology**
- ✓ An attacker with arbitrary background knowledge
(i.e., information about all the data contained in the dataset) is assumed
- ✓ Random noise is added to the statistics to obscure individual data



DP: Formula

- ✓ Adding random noise (F) makes it difficult to discriminate between neighboring datasets ($D_1, D_2 \in D$)
- ✓ The smaller ϵ is, the stronger the privacy is protected
- ✓ Effective also against an attacker with information on D

ϵ -Differential Privacy For $D_1, D_2 \in D$,

$$\Pr[F(D_1) \in O] \leq e^\epsilon \cdot \Pr[F(D_2) \in O]$$

F : randomized function, ϵ : privacy budget,

O : any output of f , D : entire dataset,

D_1, D_2 : datasets that differ in a single record ($D_1, D_2 \in D$)

DP: Sensitivity

- ✓ Dataset dependency of the output defined for each f
- ✓ The lower sensitivity, the smaller the required noise for DP

$$\Delta f = \max_{D_1, D_2 \in D} \left\| f(D_1) - f(D_2) \right\|_1$$

Δf : sensitivity for f , f : original function, D : entire dataset,
 D_1, D_2 : datasets that differ in a single record ($D_1, D_2 \in D$)

DP: Mechanisms

Laplace Mechanism: (for a query of **continuous** values)

Add a random number generated from the Laplace distribution

$$F(D) = f(D) + \text{Lap}(0, \frac{\Delta f}{\epsilon})$$

Exponential Mechanism: (for a query of **discrete** values)

Convert to a probabilistic algorithm based on a gain function

$$F(D, u) = \text{choose } r \in R \text{ with probability } \propto \exp(\frac{\epsilon u(D, r)}{2\Delta u})$$

F : mechanisms, f : original function, ϵ : privacy budget,

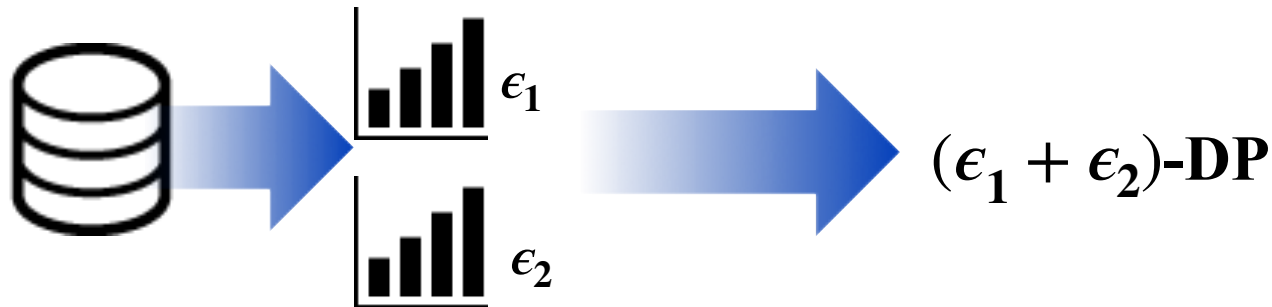
Lap : Laplace distribution, u : gain function $(D \times R) \rightarrow \mathbb{R}$

DP: Composition

Sequential Composition

(A series of functions $f = \{f_1, \dots, f_m\}$ is executed on the **same dataset**)

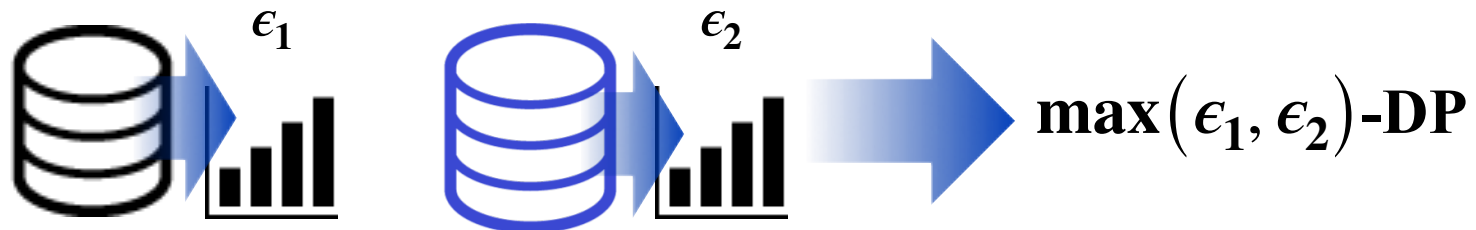
If f_i provides ϵ_i -DP, then f provides $\sum_{i=1}^m \epsilon_i$ -DP



Parallel Composition

(A series of functions $f = \{f_1, \dots, f_m\}$ is performed on **disjoint datasets**)

If f_i provides ϵ_i -DP, then f provides $\max(\epsilon_1, \dots, \epsilon_m)$ -DP



Related Works

Liu et al.'s method

Simply apply sequential composition

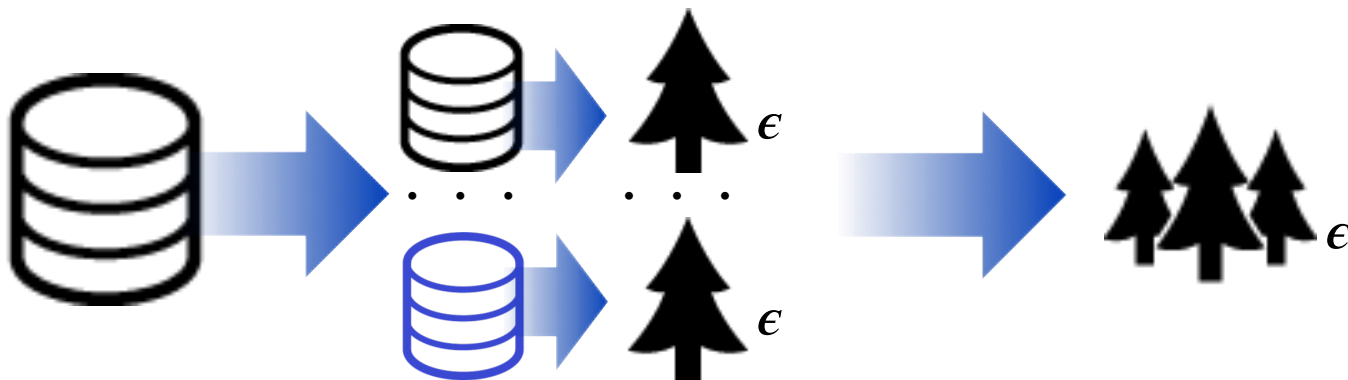
→ Scale of noise increases in proportion to the number of trees



Zhao et al's method

Simply apply parallel composition

→ Data samples for each tree decrease in proportion to the number of trees



Abstract

Title: Privacy-Preserving Gradient Boosting Decision Trees
proposed **DPBoost**, a DP framework that improves the accuracy loss in GBDT

Contributions:

✓ Sensitivity Bounds

- More tightly and accurately bounded sensitivity
- Gradient-based Data Filtering (GDF)
- Geometric Leaf Clipping (GLC)

✓ Privacy Budget Allocations

- Ensemble of Ensembles (EoE)

Tighter Sensitivity Bounds

Conventionally,

sensitivity was bounded by estimating the output range of the function

$$\Delta G \leq 2 \max \left| \frac{\left(\sum_{i \in I_L} g_i \right)^2}{|I_L| + \lambda} + \frac{\left(\sum_{i \in I_R} g_i \right)^2}{|I_R| + \lambda} \right|, \Delta V \leq 2 \max \left| \frac{\sum_{i \in I} g_i}{|I| + \lambda} \right|$$

→ **Proportional** to the number of data samples

DPBoost

bound sensitivity more accurately based on definition

$$\Delta f = \max_{D_1, D_2 \in D} \left\| f(D_1) - f(D_2) \right\|_1$$

$$\Rightarrow \Delta G \leq \frac{3\lambda + 2}{(\lambda + 1)(\lambda + 2)} g^{*2}, \Delta V \leq \frac{g^*}{1 + \lambda} \quad (g^* = \max_{i \in D} |g_i|)$$

→ **Not proportional** to the number of data samples

Gradient-based Data Filtering (GDF)

- ✓ Sensitivity can be bounded more tightly by filtering out data with large g_i
→ Filter out data with gradients larger than g_l^*

$$g_l^* = \max_{y_p \in [-1, 1]} \left\| \frac{\partial l(y_p, y)}{\partial y} \Big|_{y=0} \right\| \implies \Delta G \leq \frac{3\lambda + 2}{(\lambda + 1)(\lambda + 2)} g_l^{*2}, \quad \Delta V \leq \frac{g_l^*}{1 + \lambda}$$

- ✓ g_l^* depends only on the loss function l , and $g_l^* = 1$ when $l = \text{MSE}$
- ✓ Table 3 shows that GDF filters only a small number of data

Table 3: The number of training instances with/without gradient-based data filtering

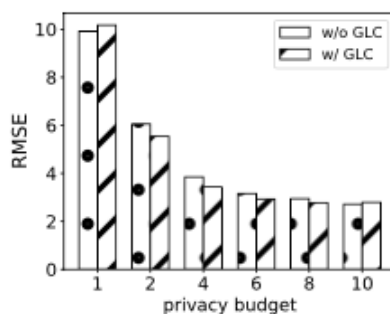
	w/ GDF	w/o GDF	filtered ratio
abalone	3340	3292	1.44%
YearPredictionMSD	370902	340989	8.06%
sklearn_reg	800000	799999	0

Geometric Leaf Clipping (GLC)

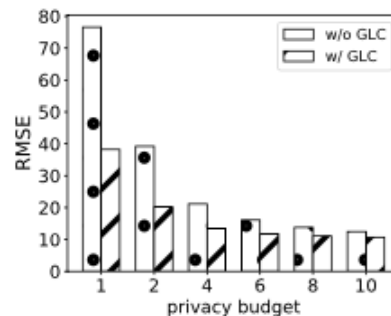
- ✓ Sensitivity can be bounded more tightly by clipping leaf weight $V(I)$
 - However, simple clipping causes a severe accuracy loss
 - Clipping according to gradient decay in GBDT learning
 - However, deriving the exact decay pattern is difficult
- ✓ When only one data corresponds to each leaf,

$$|V_t| \leq g_l^* (1 - \eta)^{t-1} \implies \Delta V \leq \min\left(\frac{g_l^*}{1 + \lambda}, 2g_l^* (1 - \eta)^{t-1}\right)$$

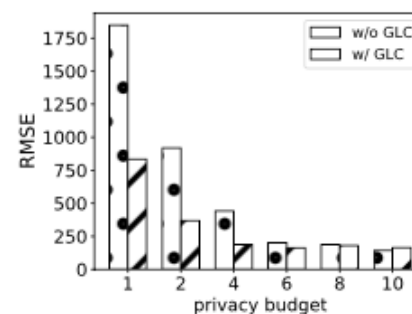
- ✓ The following figure shows that GLC leads to improved accuracy



(a) abalone



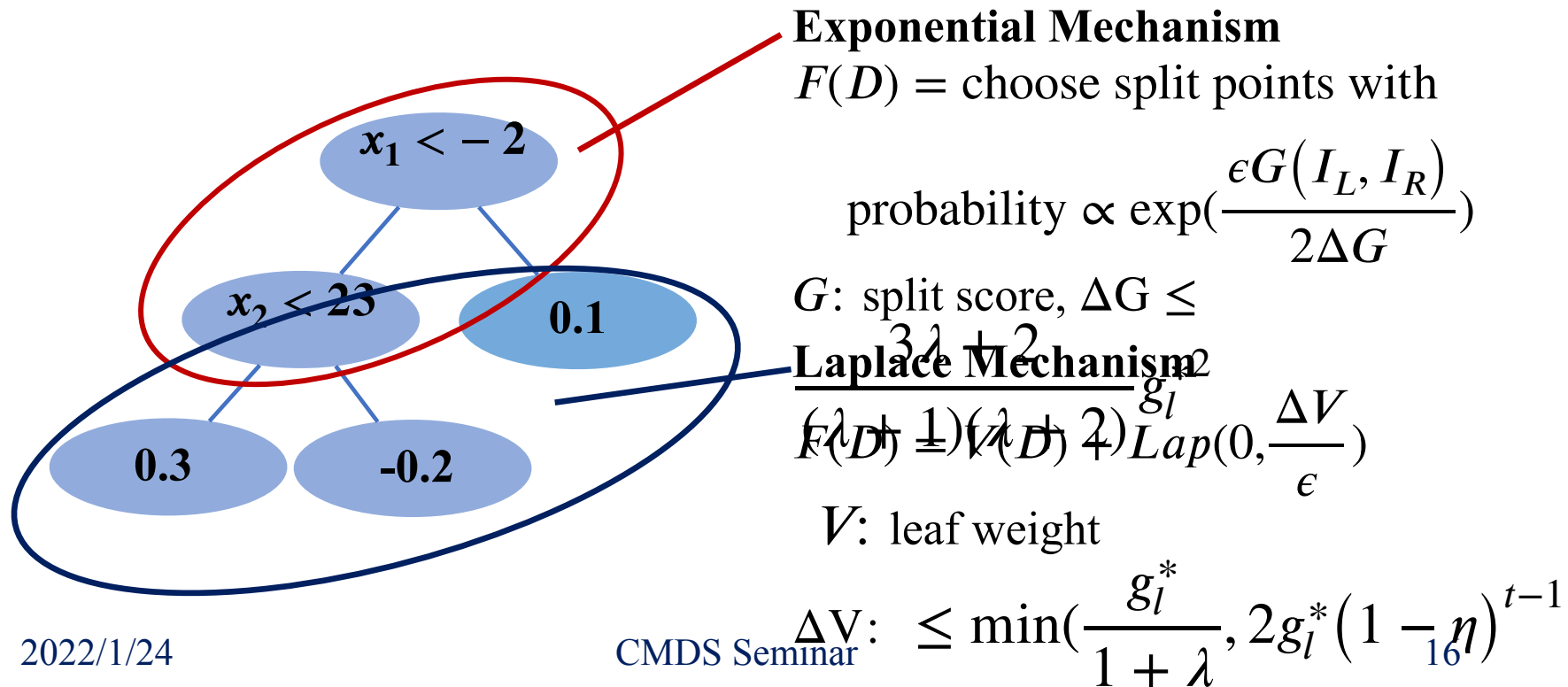
(b) YearPrediction



(c) synthetic_reg

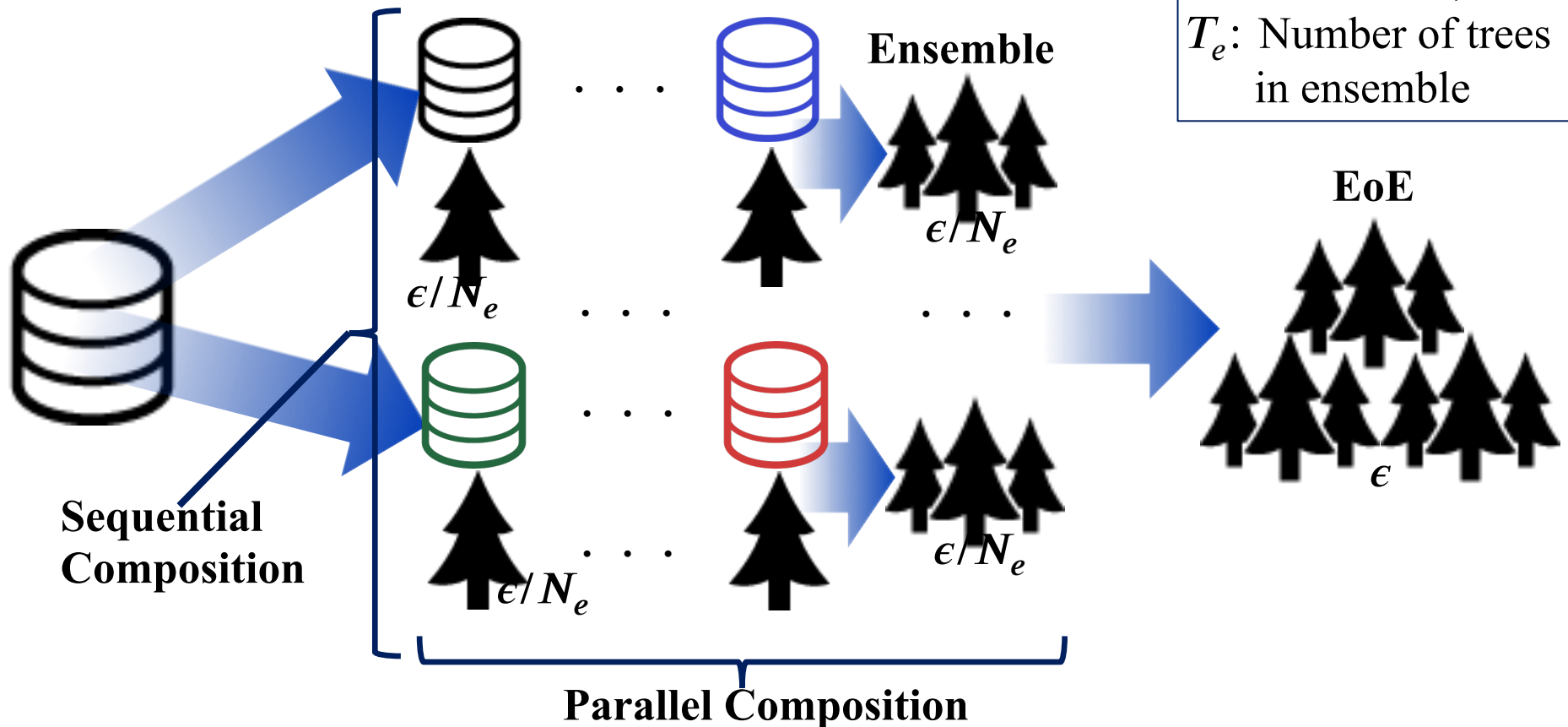
Privacy Budget Allocations for A Single Tree

- ✓ Introducing methods similar to those used in previous studies
- ✓ Applying exponential mechanism for splitting,
Laplace mechanism for leaf weight
- ✓ Allocating $\epsilon/2$ to splitting and $\epsilon/2$ to leaf weights



Privacy Budget Allocations Across Trees

- ✓ Applying parallel composition prevents the reduction of noise
- ✓ Ensemble of ensembles (EoE) reduces the accuracy loss caused by parallel composition



Privacy Budget Allocations Across Trees

- ✓ Applying parallel composition prevents the reduction of noise
- ✓ Ensemble of ensembles (EoE) reduces the accuracy loss caused by parallel composition

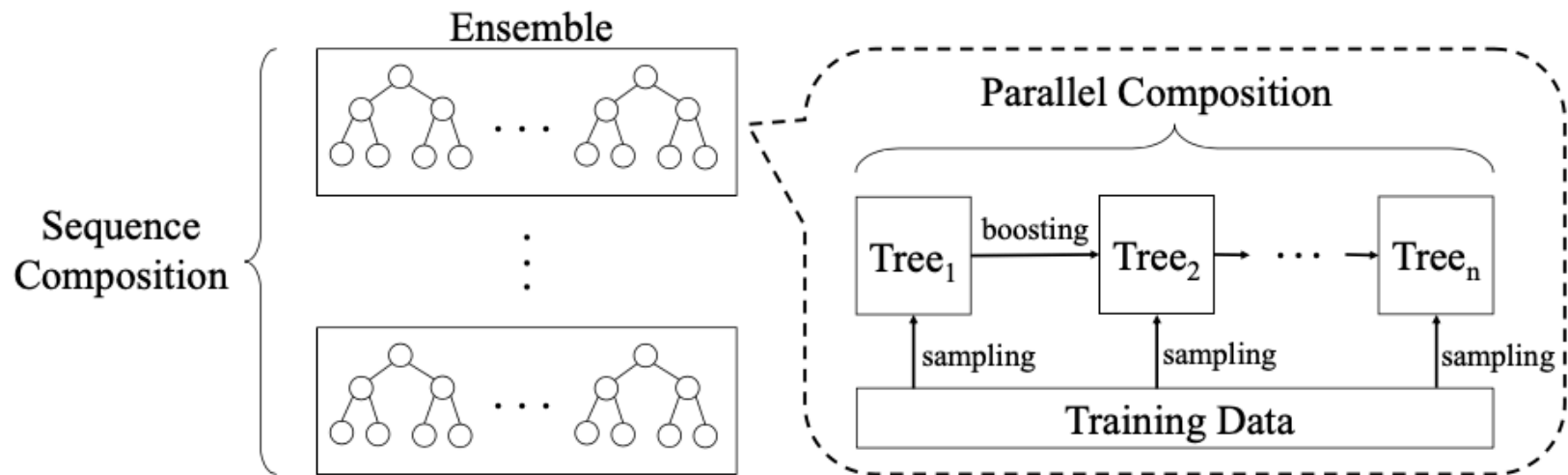


Figure 1: The two-level boosting design of DPBoost

Experimental Setups

- ✓ Use following public datasets
- ✓ Compare DPBoost with the following three methods:
 - **PARA** (Zhao et al.'s approach, 2018)
 - **SEQ** (Liu et al.'s approach, 2018)
 - **NP** (Train GBDTs without privacy concerns)
- ✓ Tighter sensitivity (proposed in this study) is applied to all methods

Table 1: Datasets used in the experiments.

datasets	#data	#features	task
adult	32,561	123	classification
real-sim	72,309	20,958	
covtype	581,012	54	
susy	5,000,000	18	
cod-rna	59,535	8	
webdata	49,749	300	
synthetic_cls	1,000,000	400	
abalone	4,177	8	regression
YearPredictionMSD	463,715	90	
synthetic_reg	1,000,000	400	

Experiments: Accuracy for different ϵ

- ✓ Set N_e to 1 in DPBoost and T to 50 for all approaches
→ Validation of **effects of GDF, GLC**
- ✓ DPBoost always outperformed the previous studies in small ϵ

N_e : Number of ensembles,
 T : Total number of trees
of trees

N_e : number of ensembles, T : total number of trees

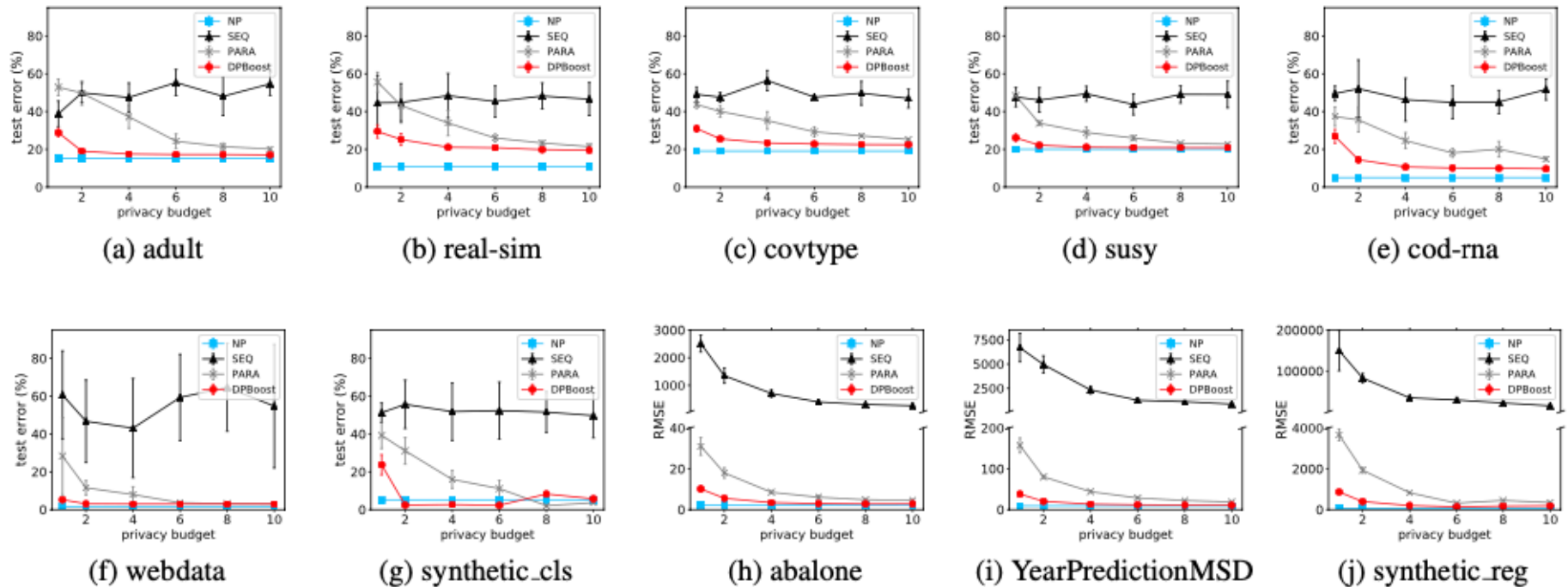


Figure 2: Comparison of the test errors/RMSE given different total privacy budgets. The number of trees is set to 50.

Experiments: Accuracy for different iterations

- ✓ Set N_e to 20 in DPBoost and T to 20 ~ 1000, $\epsilon = 100$ for all approaches
→ Validation of **effects of EoE**
- ✓ Unlike PARA and SEQ, DPBoost's accuracy was improved by increasing T

N_e : Number of ensembles,
 T : Total Number of trees

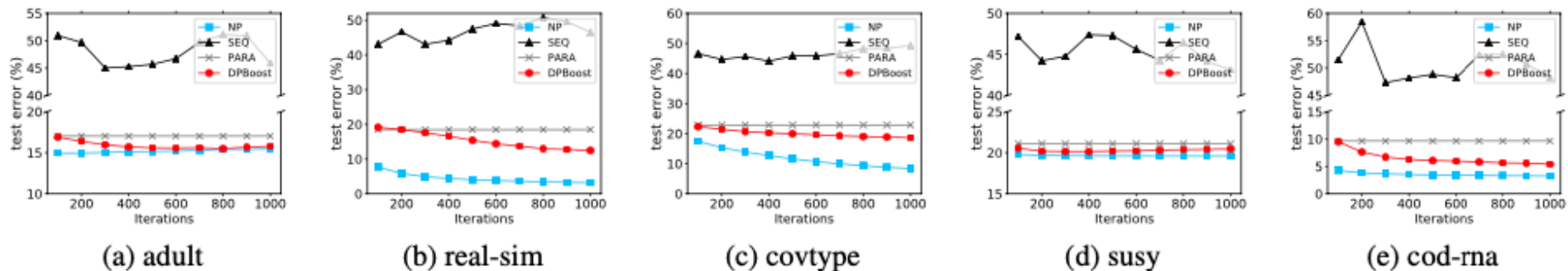


Figure 3: Comparison of test error convergence. The number of trees is set to 1000.

Experiments: Training Time

- ✓ Exponential mechanism (split point computation) is costly
- ✓ Time overhead increases in proportion to the size of the dataset

Table 2: Training time per tree (second) of DPBoost and NP.

datasets	DPBoost	NP
adult	0.019	0.007
real-sim	2.97	0.82
covtype	0.085	0.044
SUSY	0.38	0.32
cod-rna	0.016	0.009
webdata	0.032	0.013
synthetic_cls	1.00	0.36
abalone	2.95	2.85
YearPrediction	0.38	0.12
synthetic_reg	0.96	0.36

Conclusions

DPBoost:

- ✓ bounded sensitivity more strictly (GDF, GLC)
- ✓ allocated privacy budget more appropriately (EoE)
- ✓ achieved higher accuracy than previous studies

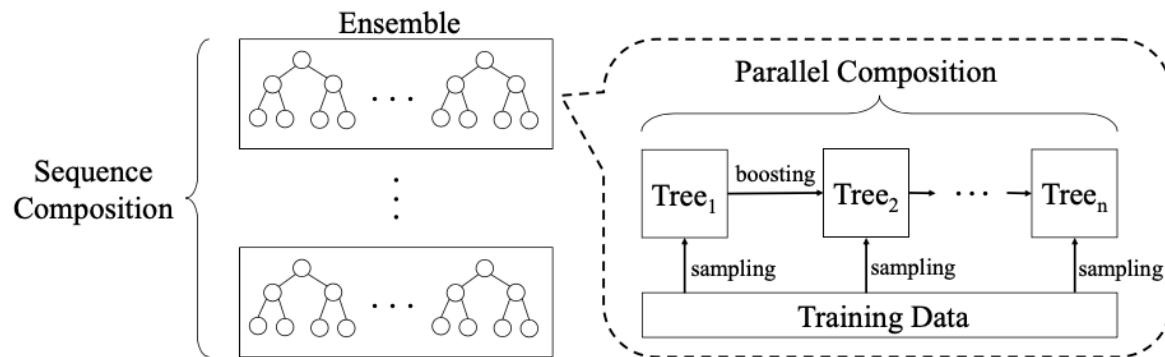


Figure 1: The two-level boosting design of DPBoost

Privacy Budget Allocations for A Single Tree

Algorithm 1: TrainSingleTree: Train a differentially private decision tree

Input: I : training data, $Depth_{max}$: maximum depth

Input: ε_t : privacy budget

```
1  $\varepsilon_{leaf} \leftarrow \frac{\varepsilon_t}{2}$  // privacy budget for leaf nodes
2  $\varepsilon_{nleaf} \leftarrow \frac{\varepsilon_t}{2Depth_{max}}$  // privacy budget for internal nodes
3 Perform gradient-based data filtering on dataset  $I$ .
4 for  $depth = 1$  to  $Depth_{max}$  do
5   for each node in current depth do
6     for each split value  $i$  do
7       Compute gain  $G_i$  according to Equation (3).
8        $P_i \leftarrow \exp(\frac{\varepsilon_{nleaf} G_i}{2\Delta G})$ 
9       /* Apply exponential mechanism */
10      Choose a value  $s$  with probability  $(P_s / \sum_i P_i)$ .
11      Split current node by feature value  $s$ .
12 for each leaf node  $i$  do
13   Compute leaf value  $V_i$  according to Equation (4).
14   Perform geometric leaf clipping on  $V_i$ .
15   /* Apply Laplace mechanism */
16    $V_i \leftarrow V_i + Lap(0, \Delta V / \varepsilon_{nleaf})$ 
```

Output: A ε_t -differentially private decision tree

Privacy Budget Allocations Across Trees

Algorithm 2: Train differentially private GBDTs

Input: \mathcal{D} : Dataset, $Depth_{max}$: maximum depth

Input: ε : privacy budget, λ : regularization parameter

Input: T : total number of trees, l : loss function

Input: T_e : number of trees in an ensemble

```
1  $N_e \leftarrow \lceil T/T_e \rceil$  // the number of ensembles
2  $\varepsilon_e \leftarrow \varepsilon/N_e$  // privacy budget for each tree
3 for  $t = 1$  to  $T$  do
4   Update gradients of all training instances on loss  $l$ .
5    $t_e \leftarrow t \bmod T_e$ 
6   if  $t_e == 1$  then
7      $I \leftarrow \mathcal{D}$  // initialize the dataset for an ensemble
8   Randomly pick  $(\frac{|\mathcal{D}|\eta(1-\eta)^{t_e}}{1-(1-\eta)^{T_e}})$  instances from  $I$  to
     constitute the subset  $I_t$ .
9    $I \leftarrow I - I_t$ 
10  TrainSingleTree (dataset =  $I_t$ ,
11    maximum depth =  $Depth_{max}$ ,
12    privacy budget =  $\varepsilon_e$ ,
13     $\Delta G = \frac{3\lambda+2}{(\lambda+1)(\lambda+2)} g_l^{*2}$ ,
14     $\Delta V = \min(\frac{g_l^*}{1+\lambda}, 2g_l^*(1-\eta)^{t-1})$ ).
```

Output: ε -differentially private GBDTs

Applications of DP and GBDT
