

PriPL-Tree: Accurate Range Query for Arbitrary Distribution under Local Differential Privacy

Leixia Wang
Renmin University of China
leixiawang@ruc.edu.cn

Haibo Hu
Hong Kong Polytechnic University
haibo.hu@polyu.edu.hk

Qingqing Ye
Hong Kong Polytechnic University
qqing.ye@polyu.edu.hk

Xiaofeng Meng*
Renmin University of China
xfmeng@ruc.edu.cn

ABSTRACT

Answering range queries in the context of Local Differential Privacy (LDP) is a widely studied problem in Online Analytical Processing (OLAP). Existing LDP solutions for range queries all assume a uniform data distribution within each domain partition, which may not align with real-world scenarios where data is skewed, resulting in inaccurate estimates. To address this problem, we introduce PriPL-Tree, a novel data structure that combines hierarchical tree structures with piecewise linear (PL) functions to answer range queries for arbitrary distribution. PriPL-Tree precisely models the underlying data distribution with a few line segments, leading to more accurate results for range queries. Furthermore, we extend it to multi-dimensional cases with novel data-aware adaptive grids. These grids leverage the insights from marginal distributions obtained through PriPL-Trees to partition the grids adaptively, adapting the density of underlying distributions. Our extensive experiments on both real and synthetic datasets demonstrate the effectiveness and superiority of PriPL-Tree over state-of-the-art solutions in answering range queries across arbitrary data distributions.

PVLDB Reference Format:

Leixia Wang, Qingqing Ye, Haibo Hu, and Xiaofeng Meng. PriPL-Tree: Accurate Range Query for Arbitrary Distribution under Local Differential Privacy. PVLDB, 14(1): XXX-XXX, 2020.
doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/LeixiaWang/PriPLT>.

1 INTRODUCTION

With an increasing amount of personal information being collected by third-party entities (a.k.a., data collectors), the protection of individual privacy is garnering much more attention. Local Differential Privacy (LDP) has emerged as a rigorous privacy-preserving standard widely employed in academia and industries [6, 11, 31]. Under LDP, users only need to submit perturbed values, preserving

the privacy of their raw data. The data collector collects these noisy values and invests effort in estimating various statistics to support data analysis tasks.

Range queries, as a prevalent query type, have been extensively studied in LDP, where the data collector aims to estimate the frequency of specific ranges within a domain. To support these queries, existing solutions generally construct hierarchical trees [4, 7, 33, 35] or grids [37, 40] over the whole domain, and then estimate frequencies of the partitioned subdomains (i.e., nodes in trees or cells in grids). To answer a range query, the frequencies of those nodes or cells covered by the given range will be summed up. When some subdomains are just partially covered, the data within them are assumed to be uniformly distributed, so that the corresponding frequency can be estimated based on the overlapped query proportion. However, the data we indexed typically exhibits various data distributions rather than uniform in reality. It is inevitable to introduce non-uniform errors by existing methods, leading to inaccurate responses.

Fig. 1 shows an example of this non-uniform estimation error. Given a distribution depicted as the black curve, Fig. 1(a) partitions the domain in a coarse-grained manner, resulting in a large non-uniform error. Fig. 1(b) partitions the domain in a fine-grained manner to reduce the non-uniform error, but incurs significant aggregated LDP noise error due to a large number of bins.

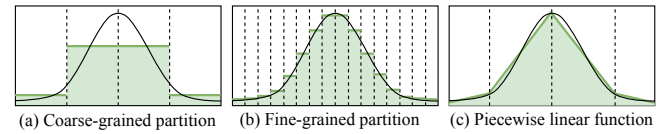


Figure 1: An Illustration on Non-uniform Errors

To tackle this challenge, we propose an innovative solution that employs a **piecewise linear (PL) function** to model the underlying data distribution instead of relying on a uniform assumption. By partitioning the data domain into several intervals and approximating the data distribution within each interval with a line segment, even complex data distributions can be well-approximated with a few parameters [30]. As shown in Fig. 1(c), the PL function accurately approximates the data distribution with a few segments (represented as frequency-slope pairs), which alleviates both non-uniform error and LDP noise error significantly.

Building upon the PL function, we introduce the **Private Piecewise Linear Tree (PriPL-Tree)**. In this tree, each leaf node corresponds to a segment and stores the corresponding frequency and

*Corresponding author: Xiaofeng Meng.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX

slope, while each non-leaf node represents the interval formed by combining the intervals of child nodes and stores associated frequencies. Compared to traditional hierarchical trees, the PriPL-Tree offers several significant advantages: (1) within a node with the same interval, it provides a more accurate fit to the underlying data distribution than the uniform assumption. (2) A few segments are sufficient to model the distribution, resulting in fewer leaf nodes and a lower tree height. This reduction in height is crucial in LDP as it facilitates the allocation of users or privacy budget among fewer tree layers (a necessary step to meet the privacy requirements of LDP), thereby mitigating noise errors in frequency estimation for each node. (3) The number of parameters in the tree depends only on the shape of the data distribution, not the domain size, enabling it to adapt to large domains with a more concise structure and a more accurate result.

However, constructing the PriPL-Tree under LDP is non-trivial because the data remains invisible to the data collector. To address this, we employ a three-phase approach. Firstly, we allocate a portion of users to estimate the histogram of the data, providing a rough view of the underlying data distribution. Using this noisy histogram, we fit the PL function and construct the PriPL-Tree with an adaptive node reduction strategy to optimize its errors for answering range queries. Next, we allocate the remaining users to estimate the frequencies of each node in the PriPL-Tree and perform a post-processing step to maintain the non-negativity and consistency of the frequencies of nodes.

Besides handling 1-D range queries, we have further extended PriPL-Tree for multi-dimensional scenarios by incorporating 2-D adaptive grids. These adaptive grids are also data-aware, having dense or sparse cell partitions depending on the density of the data distribution, and can be constructed based on insights from the marginal distributions obtained through PriPL-Trees. By leveraging both the marginal distributions from 1-D PriPL-Trees and the 2-D grids, we can answer λ -D range queries ($\lambda > 1$) using a maximum entropy approach [33, 37, 40]. To summarize, our contributions are:

- **Innovative PriPL-Tree:** We design PriPL-Tree, a novel data structure that models the underlying data distribution using a piecewise linear (PL) function instead of relying on a uniform data assumption in LDP. In this way, PriPL-Tree can answer range queries for arbitrary data distribution accurately.
- **Adaptive data-aware Grids:** By leveraging marginal distributions revealed by the PriPL-Trees, we design adaptive grids tailored to the density of underlying data distributions, which serves as the building block for answering multi-dimensional range queries effectively.
- **Extensive Experimental Evaluation:** We conduct comprehensive experimental studies on both real and synthetic datasets, which confirms the effectiveness and superiority of our proposed methods. Compared to existing work, our method achieves one order of magnitude on accuracy improvement.

In the remainder of this paper, we introduce LDP and analyze existing range query methods in Section 2. Our primary method, PriPL-Tree, is proposed in Section 3, extended with adaptive grids for multi-dimensional cases in Section 4. We evaluate them in Section 5. Finally, we review related works in Section 6 and conclude our paper in Section 7.

2 PRELIMINARIES

2.1 Local Differential Privacy (LDP)

In the context of data collection, LDP provides a mechanism \mathcal{R} that enables users to perturb their data v before sharing it with an untrusted data collector. By ensuring that the resulting perturbed data $\mathcal{R}(v)$ satisfies ϵ -LDP, the data collector cannot distinguish a value v from any other possible value v' with high confidence, thus safeguarding the privacy. A higher level of privacy is achieved when a smaller value of ϵ is employed.

Definition 2.1 (ϵ -Local Differential Privacy (ϵ -LDP) [9]). A perturbation mechanism \mathcal{R} satisfies ϵ -LDP ($\epsilon > 0$) iff for any pair of input data $v, v' \in D$ and any output z of \mathcal{R} , we have

$$\Pr[\mathcal{R}(v) = z] \leq e^\epsilon \Pr[\mathcal{R}(v') = z].$$

We introduce two state-of-the-art LDP mechanisms for fundamental frequency and numerical distribution estimation, respectively, both ensuring ϵ -LDP.

Optimal Unary Encoding Mechanism (OUE) [34] is the state-of-the-art frequency estimation mechanism, involving three critical steps — encoding, perturbation, and aggregation. In the encoding step, each user u_i ($i \leq N$) encodes his value $v_i \in D$ into a one-hot bit vector $\mathbf{B} = \{0, 1\}^{|D|}$, where the i -th position is set to 1, and others are set to 0. In the perturbation step, each user perturbs each bit in \mathbf{B} separately. The original bit “1” is retained with a probability of $p = 1/2$, while the bit “0” is flipped to “1” with a probability of $q = 1/(e^\epsilon + 1)$. All N users report their perturbed values to the data collector, who then counts the occurrences of each v as n'_v and calibrates it to an unbiased estimation of the frequency of v , i.e., $\hat{f}_v = (n'_v - Nq)/N(p - q)$. This mechanism achieves an optimized estimation variance of $\text{Var}(\hat{f}_v) \approx 4e^\epsilon / (N \cdot (e^\epsilon - 1)^2)$, denoted as σ^2 .

Square Wave Mechanism (SW) [23] is for numerical distribution estimation, involving perturbation and aggregation steps. In the perturbation step, each user perturbs his value $v_i \in D$ into v'_i within a domain with size $d + 2b$. Specifically, with a larger probability $p = e^\epsilon / ((2b+1)e^\epsilon + d - 1)$, he perturbs v_i to a value v'_i within $|v - v'_i| < b$; with a smaller probability $q = 1 / ((2b+1)e^\epsilon + d - 1)$, he perturbs it to other values. Maximizing the mutual information of input and output distributions yields the optimal $b = \left\lfloor \frac{\epsilon e^\epsilon - e^\epsilon + 1}{2e^\epsilon(e^\epsilon - 1 - \epsilon)} d \right\rfloor$. Upon receiving the perturbed data, the data collector estimates the distribution using the expectation maximization (EM) algorithm or the EM algorithm with the smoothing step (EMS). SW with EM can capture spiky distributions effectively, while SW with EMS provides a more accurate estimation by smoothing the LDP noise. We denote these two results as $\hat{\mathbf{F}}^{\text{EM}}$ and $\hat{\mathbf{F}}^{\text{EMS}}$, respectively.

2.2 Problem Definition

Consider N users and each user u_i ($i \leq N$) owns a private record containing m private values on attributes (A_1, A_2, \dots, A_m) . We use the notation $[x]$ to represent the set $\{1, 2, \dots, x\}$. Each attribute A_j ($j \in [m]$) has a public domain D_j , and each user u_i 's record across m attributes is denoted by $\mathbf{v}_i = (v_i^1, v_i^2, \dots, v_i^m)$, where $v_i^j \in D_j$ represents the attribute A_j 's value of user i . For convenience, we assume $D_j = [0, d_i]$ for continuous data and $D_j = [d_j]$ for discrete data. And we abuse v_i to denote the user u_i 's value of the attribute in 1-dimensional (a.k.a., 1-D) data.

Table 1: Notations

Symbols	Description
N	The total number of users
A_j	The j -th attribute
D_j, d_j	The attribute A_j 's domain D_j of size d_j
m	The number of private attributes in data
λ	The number of attributes involved in a range query
ϵ	Privacy budget
n_k	The node n_k in the PriPL-Tree
f_k, β_k	The frequency and the slope in n_k
I_k, s_{k-1}, s_k	The interval I_k of node n_k including $ I_k $ bucketized values between two breakpoints s_{k-1} and s_k .
α	User allocation ratio in phase 1 for PriPL-Tree
σ^2	The variance of OUE with N users and a privacy budget of ϵ .

The λ -dimensional (a.k.a., λ -D) range query is performed on a set of private attributes $\Phi \subseteq \{A_j | j \leq m\}$, where $\lambda = |\Phi| \leq m$. Let $[l_j, r_j]$ denote the specified range for the attribute $A_j \in \Phi$. The λ -D range query returns the frequency of records where all queried attribute values v_i^j ($j \in \Phi$) are within these specified ranges. Formally,

$$Q(\cap_{A_j \in \Phi} [l_j, r_j]) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\cap_{A_j \in \Phi} \{l_j \leq v_i^j \leq r_j\}},$$

where $\mathbb{1}$ is an indicator function that outputs 1 if the predicate is true and 0 otherwise.

Our goal is to let the untrusted data collector answer the range query $Q(\cap_{A_j \in \Phi} [l_j, r_j])$ while ensuring individual privacy under ϵ -LDP. Extensive research has been conducted on this problem in the context of LDP, as we reviewed below. The notations are summarized in Table 1.

2.3 Existing Methods

The hierarchical tree (HT) is the primary data structure for 1-D range queries in LDP. It hierarchically decomposes the entire domain into disjoint sub-domains (a.k.a., intervals), constructing a B -ary tree. Each node in the tree represents an interval and stores an estimated interval frequency. Non-leaf nodes aggregate the sum of frequencies from their B child nodes. As such, range queries can be answered by summing up the frequencies of a few nodes in the tree with disjoint intervals rather than summing up all individual bins' frequencies within the range as in a histogram, reducing the accumulated noise error. For example, considering a domain $D = [0, 16]$, we can either uniformly partition the domain into 16 bins to compute the histogram or construct a complete binary tree on the histogram. In both cases, the data within each bin is assumed to be uniformly distributed. When responding to a range query $Q([0, 5])$ from the tree, only the estimated frequencies of nodes with intervals $[0, 4]$ and $[4, 5]$ need to be added instead of the frequencies of 5 bins (i.e., $[0, 1]$, $[1, 2]$, $[2, 3]$, $[3, 4]$ and $[4, 5]$) in the histogram. To build this tree under LDP, the privacy budget or users are allocated among layers to estimate nodes' frequencies. To further reduce the error of range queries, a lot of optimization methods have been developed, including Haar transforming of data [4], optimizing branch number for the tree [4, 35], customizing branch number for nodes [33] and merging nodes with low frequencies [7].

Beyond 1-D, the HT can be extended to multi-dimensional cases [7, 35]. However, finely partitioning users or the privacy budget among layers and dimension combinations would increase the noise

error. To overcome the curse of dimensionality, grid-based methods are typically employed in λ -D ($\lambda > 1$) range queries. Given m private attributes, they estimate the frequencies of m 1-D grids for individual attributes and $\binom{m}{2}$ 2-D grids for attribute pairs. Users or privacy budgets are allocated among these grids, where the data in each cell is still assumed to be uniformly distributed. Based on these, a λ -D range query can be estimated from these 1-D and 2-D grids through the maximum entropy [44] or weighted updating [40] algorithms. To reduce the estimation error of range queries, Yang et al. [40] optimized the granularity for both 1-D and 2-D grids, and Wang et al. [37] further employed prefix-sum (PS) cubes.

2.4 Observations and Challenges

Drawing from current research on range queries in LDP, we summarize two key observations that guide our approach and identify a significant challenge. First, we outline the observations:

(1) **Tree vs. Grid:** Tree-based methods allocate users (or privacy budget) across multiple layers and dimensions, whereas grid-based methods allocate only among dimensions. Considering the significant noise from a few users or a small privacy budget, tree-based methods are preferable for 1-D range queries, while grid-based methods are more suitable for λ -D ($\lambda > 1$) range queries [33].

(2) **User Allocation vs. Privacy Budget Allocation:** To achieve ϵ -LDP, allocation of users or privacy budget is necessary among the layers of trees and the grids. Generally, user allocation is preferable in the LDP setting as it introduces less estimation error (i.e., the estimation variance) than privacy budget allocation [4, 7, 33, 35, 40].

Then, we present the significant challenge: **the unrealistic uniform assumptions**. Under LDP guarantee, as the data distribution is invisible to the data collector, existing works uniformly decompose the domain and/or assume a uniform data distribution in each decomposed sub-domain [4, 7, 33, 35, 37, 40]. However, real-world applications often involve data following various distributions (e.g., Gaussian or Zipf) rather than being uniformly distributed. This assumption inevitably leads to non-uniform errors or suboptimal estimation results.

In this work, we address uniform assumptions on the domain and data and correspondingly provide enhanced estimation accuracy of range queries in LDP. In what follows, we first present a solution for 1-D range queries in Section 3 and then extend it to a multi-dimensional setting in Section 4.

3 PRIPL-TREE: PRIVATE PIECEWISE LINEAR TREE

In this section, we present PriPL-Tree, a data structure combining piecewise linear functions and hierarchical trees to address the uniform assumption and accurately respond to 1-D range queries.

3.1 Design Rationale

The piecewise linear (PL) function is capable of approximating the underlying data distribution with only a few parameters, enabling us to not rely on uniform distribution assumptions. For instance, given a Gaussian distribution shown in Fig. 2(a), we can approximate it using 4 segments with 8 parameters. In this case,

the entire domain is divided into 4 intervals, and data in each interval $I_k = [s_{k-1}, s_k]$ ($1 \leq k \leq 4$) is fitted with a linear function defined by two parameters, i.e., β_k (slope) and f_k (sum of frequencies of all points in the interval). The linear expression is given by $y = \beta_k x + b_k$, where $b_k = f_k/|I_k| - \beta_k(|I_k| + 2s_{k-1} - 1)/2$ and $|I_k|$ represents the interval size.

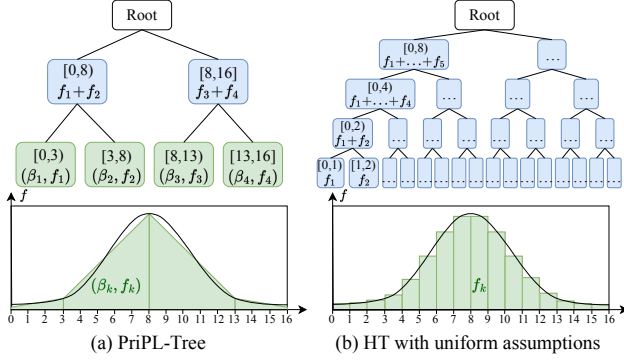


Figure 2: An Example of PriPL-Tree and HT

To facilitate range query processing, we integrate the PL function with a hierarchical tree structure, proposing the **Private Piecewise Linear Tree (PriPL-Tree)**. Each leaf node represents a segment (corresponding to an interval) of the PL function and stores its slope β_k and frequency f_k . Each non-leaf node represents an interval and only stores the interval frequency. As conventional hierarchical trees, the parent node stores the sum of its child nodes' frequencies. We count the layers of the tree starting from 0 at the top.

We provide an example in Fig. 2 to illustrate the PriPL-Tree, comparing it to the conventional hierarchical tree (HT) with uniform assumptions. Obviously, within the same interval, the PL function can provide a more accurate approximation of the underlying distribution than assuming uniform data distribution. As such, the PriPL-Tree captures the underlying distribution using significantly fewer leaf nodes (4 in PriPL-Tree vs. 16 in HT) and correspondingly less layers in the tree (2 in PriPL-Tree vs. 4 in HT). In the context of LDP, a smaller layer number implies that each layer in the tree can be allocated with more users, resulting in less estimation error due to the law of large numbers. Furthermore, the construction of the PriPL-Tree depends solely on the distribution of the underlying data, as opposed to HT which relies on the domain size. In HT, modeling data with a large domain size requires a higher tree or a coarser granularity for leaf nodes, increasing noise errors or non-uniform errors. The PriPL-Tree is well-suited for large domain-sized scenarios while reducing estimation errors for range queries.

However, constructing an effective PriPL-Tree in LDP settings is challenging due to the invisible data distribution. To address this, we first employ some users to collect a noisy histogram using LDP mechanisms, gaining insight into the underlying data distribution. Next, we fit PL functions based on this noisy histogram and use the remaining users to construct the entire tree. By considering frequency constraints, we can further optimize the frequencies and slopes in the tree, thereby supporting accurate range queries across arbitrary distributions. In the following text, we outline a three-phase workflow based on this idea and detail the methods for each phase in separate sections.

3.2 Workflow of PriPL-Tree

The workflow of the PriPL-Tree involves three phases: Private PL Fitting, PriPL-Tree Construction, and PriPL-Tree Refinement. To illustrate this workflow, we provide an example in Fig. 3.

Phase 1: Private Piecewise Linear (PL) Fitting. To gain a fundamental insight into the data distribution, we employ a proportion α of users to execute SW protocols with privacy budget ϵ , collecting a noisy histogram \hat{F}^H . For brevity, we use $[d]$ to denote the bucketized domain of the histogram. Then, we fit the PL function over it, as presented in Section 3.3. During PL fitting, we address two key issues: interval partitioning and segment fitting. Interval partitioning involves determining the number of intervals K and deciding $K+1$ appropriate breakpoints $\{s_0, s_1, \dots, s_K\}$. The derived interval is denoted as $I_k = [s_{k-1}, s_k]$ for $1 \leq k < K$ and $I_K = [s_{K-1}, s_K]$ for the last interval. To mapping to the histogram, we can also mark $[s_{k-1}, s_k]$ as $[s_{k-1}, s_k - 1]$. Segment fitting focuses on fitting the parameters of the PL function, obtaining the slope β_k ($1 \leq k \leq K$) of the line segment for each interval. Here, we do not concern the fitted intercept parameter of segments but record the slope parameter β_k and the estimated interval frequency, i.e., the sum of frequencies of values in each segment, denoted as $\hat{f}_k = \sum_{v \in S_k} \hat{f}_v^H$ ($1 \leq k \leq K$). These two parameters can be further optimized by utilizing the collected interval frequencies in subsequent phases.

Phase 2: PriPL-Tree Construction. Based on the PL function, the PriPL-Tree structure will be dynamically constructed in this phase. Each leaf node corresponds to a fitted segment in sequence, for example, n_1 to interval I_1 and n_2 to interval I_2 , as shown in Fig.3 (a) and (b). Each non-leaf node represents an interval encompassing its children and has a non-uniform branch number (i.e., fan-out). This flexible structure is designed to minimize average error in responding to range queries.

Based on this structure, we allocate the remaining $N \times (1 - \alpha)$ users to nodes and estimate their frequencies. Because the intervals of nodes along each path from the root to the leaves are overlapping, each user will be randomly allocated to one node per path. As a result, the total number of users along each path is $N \times (1 - \alpha)$. For each individual user, they are assigned to multiple nodes with non-intersecting intervals, and the union of these intervals covers the entire domain. After being informed of these intervals, users can encode their values into one-hot vectors to employ the OUE mechanism with a privacy budget of ϵ for frequency estimation. For example, if a user's value v is covered by nodes $\{n_3, n_7\}$ and he receives the intervals of nodes (n_6, n_3, n_4, n_5) , he can encode his value as $(0, 1, 0, 0)$ and apply OUE. By aggregating all users' perturbed values for corresponding nodes using OUE, we derive each node's frequency \hat{f}_k , forming a preliminary PriPL-Tree, as depicted in Fig.3 (b). Each leaf node in this tree has two frequencies: \hat{f}_k , estimated during private PL fitting in phase 1, and \tilde{f}_k , estimated by OUE in this phase. The detailed process of tree construction and user allocation is described in Section 3.4.

Phase 3: PriPL-Tree Refinement. In the current PriPL-Tree, there are several frequency inconsistencies: (1) the estimated frequency of values or intervals may be beyond the actual range of $[0, 1]$, (2) the frequency of a parent node may differ from the frequency sum of its child nodes, (3) two different frequencies occur

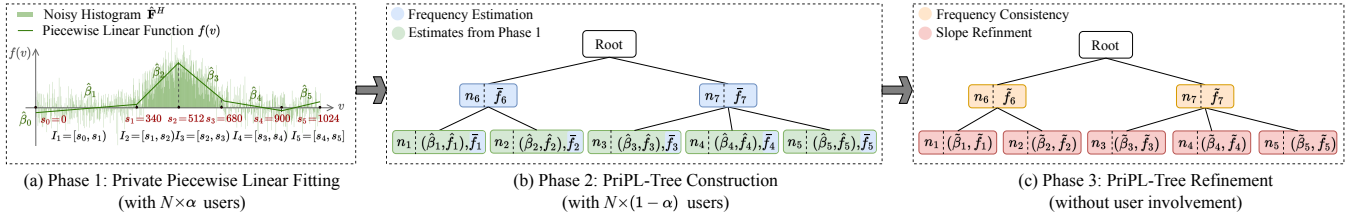


Figure 3: Workflow of PriPL-Tree

at leaf nodes. To address these issues, we propose a post-processing method in Section 3.5, yielding an optimized PriPL tree as shown in Fig.3 (c). It has consistent frequencies \tilde{f} across all nodes and optimized slopes $\tilde{\beta}_k$ for leaf nodes. By now, a well-estimated PriPL-Tree is ready to respond to range queries.

Response to 1-D Range Query $Q([l, r])$. Given a 1-D range query $Q([l, r])$, the response is obtained by summing the frequencies \tilde{f}_k of nodes n_k that are fully within the range $[l, r]$ but whose parents are not, as well as frequencies from parts of leaf nodes that overlap but are not completely within $[l, r]$. For example, when querying $Q([200, 1024])$ in Fig. 3 (c), we aggregate the frequencies \tilde{f}_2 of node n_2 , \tilde{f}_7 of node n_7 , and the frequency of the sub-range $[200, 340]$ (i.e., $[200, 339]$) within node n_1 . All these nodes and their corresponding frequencies can be derived by traversing the PriPL-Tree from top to bottom. Let $[l_{\text{sub}}, r_{\text{sub}}]$ denote the intersecting range of $[l, r]$ with the interval I_k of a leaf node n_k ; its frequency $Q([l_{\text{sub}}, r_{\text{sub}}])$ can be computed using Eq. (1). The detailed computation process is shown in Appendix A.1 in [32].

$$Q([l_{\text{sub}}, r_{\text{sub}}]) = (r_{\text{sub}} - l_{\text{sub}} + 1) \cdot \left(\tilde{\beta}_k \left(\frac{l_{\text{sub}} + r_{\text{sub}} + 1 - |I_k|}{2} - s_{k-1} \right) + \frac{\tilde{f}_k}{|I_k|} \right) \quad (1)$$

3.3 Private PL Fitting

As a fundamental data model, the PL function has been extensively studied in stream compression [2, 10, 19, 24, 39] and learned index [12, 13, 21, 22] applications. In these contexts, the original data distribution is available, and the PL model can be learned using heuristic algorithms with a specified restriction on the model's prediction error. However, our task poses a key challenge as we aim to use a noisy histogram to fit an unknown distribution while achieving an optimal error. To address this challenge, we carefully design the following segment fitting and interval partitioning steps to learn a PL model on the data distribution. The pseudocode is provided in Alg. 1.

3.3.1 Segment Fitting. Let's start with a simple case with K partitioned intervals, and we aim to optimize the PL function by minimizing the squared error between the fitted and the noisy values. To alleviate the impact of LDP noise, we assume the PL function is continuous. This allows us to model the entire noisy histogram as a whole, leveraging all the information in the histograms to fit each line segment, rather than utilizing only a subset of data located in individual intervals, where these data may be covered by LDP noise. Moreover, this assumption is practical even for a non-continuous distribution, where the data around the breakpoints can be approximated by a continuous function with a sharp line connecting two breakpoints. Note that this approximation would produce only a minor error, especially for bucketized values in histograms.

We model the continuous PL function in Eq.(2) and provide a graph with $K=5$ segments in Fig.3 (a). Each s_k ($0 < k < K$) represents the breakpoint between two segments, while s_0 and s_K mark the domain's endpoints. Let β_0 denote the intercept of the first line segment. For each line segment in the interval I_k , its slope is denoted by β_k , and its linear expression appears in the k -th row of Eq.(2).

$$f(v) = \begin{cases} \beta_0 + \beta_1(v - s_0), & s_0 \leq v < s_1 \\ \beta_0 + \beta_1(s_1 - s_0) + \beta_2(v - s_1), & s_1 \leq v < s_2 \\ \dots & \dots \\ \beta_0 + \sum_{k=1}^{K-1} \beta_k(s_k - s_{k-1}) + \beta_K(v - s_{K-1}), & s_{K-1} \leq v \leq s_K \end{cases} \quad (2)$$

Given the noisy histogram with frequency \hat{f}_v^H for $v \in [d]$, our objective is to minimize the squared error between the PL function $f(v)$ and the observed frequencies, i.e., $\min_{v \in [d]} (f(v) - \hat{f}_v^H)^2$.

For ease of optimization, we express the problem using matrices. Let $\mathbb{1}$ represent an indicator function, which is 1 when its predicate is met and 0 otherwise. We denote the frequencies in the noisy histogram by the vector $\hat{\mathbf{F}}^H = [\hat{f}_1^H \ \hat{f}_2^H \ \dots \ \hat{f}_d^H]^T$, and the parameters of the PL function by $\mathbf{B} = [\beta_0 \ \beta_1 \ \dots \ \beta_K]^T$. We define two matrices, $\mathbf{X}_{d \times (K+1)}$ and $\mathbf{A}_{d \times (K+1)}$, as described in Eq.(3) and Eq.(4) respectively. In both matrices, the v -th ($v \in [d]$) row corresponds to the expression of $f(v)$ by Eq.(2). For $k > 0$, the k -th column in \mathbf{X} refers to the term $(v - s_{k-1})$, and the k -th column in \mathbf{A} refers to the term $(s_k - s_{k-1})$.

$$\mathbf{X} = \begin{bmatrix} 1 & (1 - s_0)\mathbb{1}_{s_0 \leq 1 < s_1} & \dots & (1 - s_{K-1})\mathbb{1}_{s_{K-1} \leq 1 \leq s_K} \\ 1 & (2 - s_0)\mathbb{1}_{s_0 \leq 2 < s_1} & \dots & (2 - s_{K-1})\mathbb{1}_{s_{K-1} \leq 2 \leq s_K} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & (d - s_0)\mathbb{1}_{s_0 \leq d < s_1} & \dots & (d - s_{K-1})\mathbb{1}_{s_{K-1} \leq d \leq s_K} \end{bmatrix} \quad (3)$$

$$\mathbf{A} = \begin{bmatrix} 0 & (s_1 - s_0)\mathbb{1}_{1 > s_1} & \dots & (s_{K-1} - s_{K-2})\mathbb{1}_{1 > s_{K-1}} & 0 \\ 0 & (s_1 - s_0)\mathbb{1}_{2 > s_1} & \dots & (s_{K-1} - s_{K-2})\mathbb{1}_{2 > s_{K-1}} & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & (s_1 - s_0)\mathbb{1}_{d > s_1} & \dots & (s_{K-1} - s_{K-2})\mathbb{1}_{d > s_{K-1}} & 0 \end{bmatrix} \quad (4)$$

By calculating $(\mathbf{X} + \mathbf{A}) \cdot \mathbf{B}$, we derive the PL fitted frequencies for all values in $[d]$, i.e., $[f(1) \ f(2) \ \dots \ f(d)]^T$. Consequently, we reformulate the optimization problem using matrices as follows:

$$\min((\mathbf{X} + \mathbf{A}) \cdot \mathbf{B} - \hat{\mathbf{F}}^H)^T \cdot ((\mathbf{X} + \mathbf{A}) \cdot \mathbf{B} - \hat{\mathbf{F}}^H). \quad (5)$$

Let $L(\mathbf{B}) = ((\mathbf{X} + \mathbf{A}) \cdot \mathbf{B} - \hat{\mathbf{F}}^H)^T \cdot ((\mathbf{X} + \mathbf{A}) \cdot \mathbf{B} - \hat{\mathbf{F}}^H)$ denote the loss function. By setting $\partial L(\mathbf{B}) / \partial \mathbf{B} = 0$, we derive the closed-form solution for $\hat{\mathbf{B}}$, where \hat{B}_0 represents the estimated intercept $\hat{\beta}_0$ and \hat{B}_k ($k > 0$) represents the estimated slope $\hat{\beta}_k$ of the k -th segment.

$$\hat{\mathbf{B}} = ((\mathbf{X} + \mathbf{A})^T (\mathbf{X} + \mathbf{A}))^{-1} (\mathbf{X} + \mathbf{A})^T \hat{\mathbf{F}}^H. \quad (6)$$

Algorithm 1: Private PL Fitting

Input: Noisy histograms $\hat{\mathbf{F}}^{\text{EM}}$ and $\hat{\mathbf{F}}^{\text{EMS}}$.
Output: A PL function with adaptive K segments.

- 1 Set segment number $K=1$ and breakpoints $\mathbf{S}=(0, d-1)$;
- 2 Initialize search space $\Theta = \{s|s \in [d-1]\}$;
- 3 **for** $i \in \{1, 2\}$ **do**
- 4 **if** $i == 1$ **then** $\hat{\mathbf{F}}^{\text{H}} \leftarrow \hat{\mathbf{F}}^{\text{EM}}$ **else** $\hat{\mathbf{F}}^{\text{H}} \leftarrow \hat{\mathbf{F}}^{\text{EMS}}$;
- 5 **while** RSS not converges or $K \leq (K_{\max} \times i/2)$ **do**
- 6 $K \leftarrow K + 1$;
- 7 **foreach** Breakpoint candidate s in Θ **do**
- 8 Initialize $\mathbf{X}_{d \times (K+1)}$ and $\mathbf{A}_{d \times (K+1)}$ with \mathbf{S} and s ;
- 9 Calculate FL parameters $\hat{\mathbf{B}}$ with Eq.(6);
- 10 Calculate $\text{RSS}_k = \sum_{v \in S_k} (\hat{f}_v^{\text{H}} - f(v))^2$ for each interval S_k , record total $\text{RSS} = \sum_{1 \leq k \leq K} \text{RSS}_k$;
- 11 Append s^* to \mathbf{S} , where $s^* \in \Theta$ produces the minimal RSS;
- 12 Set $\Theta = \{s|s \in I_{k^*}\}$, where I_{k^*} has the maximum RSS_{k^*} and its frequency $\hat{f}_k = \sum_{v \in I_k} \hat{f}_v^{\text{EMS}} > \sigma\sqrt{(1-\alpha)}$;
- 13 **return** A PL function with breakpoints \mathbf{S} , slopes $\hat{\mathbf{B}}$ and frequencies $\hat{\mathbf{F}}$;

3.3.2 Interval Partitioning. Building upon segment fitting, we propose a greedy method to search the breakpoints one by one, achieving approximately optimal interval partitions. As described in Alg. 1, during each searching (i.e., each iteration in lines 5~11), we traverse all candidate breakpoint s in the search space Θ , fit segments based on it and existing breakpoints \mathbf{S} , and find the best breakpoint s^* that leads to the PL fitted segments with the minimal residual sum of squares (RSS). The search space initially contains all possible candidates in the whole domain. In subsequent iterations, we select values in the interval I_{k^*} as the new search space. This interval I_{k^*} has the maximum RSS for its fitted line segment, indicating that it requires further splitting to fit the underlying distribution more accurately, and its frequency \hat{f}_k should be greater than $\sigma\sqrt{(1-\alpha)}$, ensuring that will not be overwhelmed by the OUE noise in node frequency estimation in the next phase. This iteration continues until the maximum number of segments is reached (we set $K_{\max} = 32$ in experiments) or RSS converges, i.e., the ratio of total RSS between two consecutive iterations approaches 1, signifying no further gain in increasing segments to fit the underlying distribution.

Additionally, we propose two strategies to improve interval partitioning for effectiveness and efficiency. These are briefly described below, with detailed explanations provided in Appendix A.

(1) *Twice Partitioning Strategy for Effectiveness:* To effectively fit both jagged and smooth distributions, we sequentially perform interval partitioning on two distributions, $\hat{\mathbf{F}}^{\text{EM}}$ and $\hat{\mathbf{F}}^{\text{EMS}}$, as outlined in line 3 of Alg. 1. These distributions, derived from SW using EM and EMS, respectively, represent an initially calibrated (typically jagged) distribution and a further smoothed one with reduced noise error [23]. After this process, we can derive the necessary partitions to depict both types of distributions. Notably, these two estimations from SW require only one perturbation for each user, thus not increasing the privacy budget.

(2) *Search Acceleration Strategy for Efficiency:* To accelerate the search process, we propose a multi-granular search strategy instead of traversing all possible breakpoints during each search in current line 6 in Alg. 1. Given the granularity factor ϕ , which limits the

Algorithm 2: PriPL-Tree Construction

Input: K segments and $N(1-\alpha)$ users
Output: PriPL-Tree \mathcal{T}

/ Tree Structure Construction */*

- 1 Construct a basic balanced binary tree \mathcal{T} ;
- 2 **for** node n_k in postorder traversal **do**
- 3 Compute Err for \mathcal{T} with n_k and Err' w/o n_k using Eq.(7);
- 4 **if** $\text{Err}' < \text{Err}$ **then** Remove n_k ;
- 5 */* User Allocation */*
- 6 Assign all available users $U_0 = \{u_1, u_2, \dots, u_{N(1-\alpha)}\}$ to the root;
- 7 **for** node n_k in level order traversal **do**
- 8 **if** n_k is root **then** Set $U'_k = \phi$;
- 9 **else**
- 10 Compute h_k , the height of the subtree rooted at n_k ;
- 11 Set U'_k as $\lceil \frac{|U_k|}{h_k} \rceil$ randomly sampled users from U_k ;
- 12 Allocate U'_k to n_k for frequency estimation;
- 13 **for** node $n_c \in \text{children}(n_k)$ **do** Assign $U_c = U_k - U'_k$ to n_c ;
- 14 */* Node Frequency Estimation */*
- 15 **for** each user $u_i \in U_0$ **do**
- 16 Assign intervals of nodes $\mathcal{N}_i = \{n_k | u_i \in U'_k\}$ to u_i ;
- 17 Collect OUE perturbed vectors with size $|\mathcal{N}_i|$;
- 18 Estimate frequency \hat{f}_k for each node n_k based on OUE;
- 19 **return** PriPL-Tree \mathcal{T} ;

maximum number of candidate breakpoints during each search, we initially explore the space Θ using a step of $\lceil |\Theta|/\phi \rceil$ to identify an optimal breakpoint s^* . Subsequently, we narrow the search space to $[s^* - \lceil |\Theta|/\phi \rceil, s^* + \lceil |\Theta|/\phi \rceil]$ and search it with a finer step of $\lceil |\Theta|/\phi^2 \rceil$. We repeat this process until the step size reduces to 1. A relatively smaller ϕ ($\phi < d$) accelerates the search while increasing the probability of encountering local optima. Since breakpoints inherently represent different local optima partitioning the domain, ϕ primarily influences the order of breakpoint discovery rather than the final interval partitions. For clarity, this strategy is not included in Alg. 1, but it can replace line 3 in the algorithm.

3.4 PriPL-Tree Construction

Based on the derived partitioned intervals, we construct the PriPL-Tree, focusing on user allocation and tree structure construction and providing the pseudocode in Alg. 2.

3.4.1 User Allocation. Given a potentially unbalanced PriPL-Tree \mathcal{T} , we explore user allocation strategies. Nodes along each path from the root to the leaves have overlapping intervals and related frequencies, prompting us to allocate users to nodes along paths [33] rather than by layers [4, 7, 35]. The allocation process is detailed in lines 7~16 of Alg. 2. Initially, all unallocated users are assigned to the root (line 5). Because the root has a constant frequency of 1 and requires no estimation, it is allocated no users (line 7) and merely passes the user set to its children. For each non-root node n_k that we traversed in level order (line 6), it has inherited the unallocated user set U_k from its parent (line 12), uniformly samples $1/h_k$ of these users for itself, marked as U'_k (lines 10~11), and passes the remaining users $U_k - U'_k$ to its children (line 12). Here, h_k represents the height of the subtree rooted at n_k , ensuring uniform allocation

along the longest path. During this process, all child nodes of n_k will receive the same to-be-allocated user group $U_k - U'_k$ since they do not overlap in their intervals.

An example of user allocation is shown in Fig. 4, where colored rectangles above each node represent the users randomly assigned. Along each path, such as “ $n_0 - n_7 - n_8 - n_3$ ” in Fig. 4 (a), the total number of users is N' . For each user, like the one represented by the yellow rectangle, he will participate in frequency estimations for multiple nodes, such as $\{n_6, n_3, n_4, n_5\}$. The intervals of these nodes do not intersect and collectively cover the entire domain.

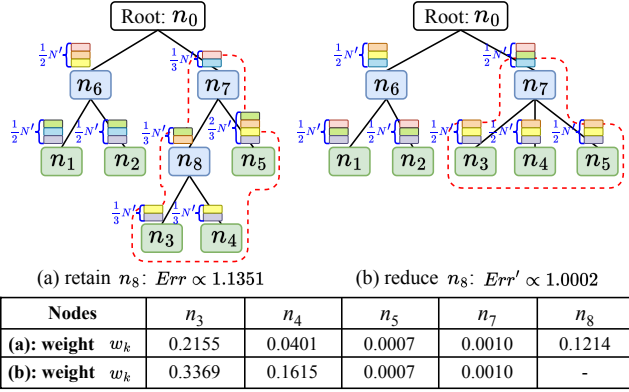


Figure 4: An Example of Tree Construction ($N' = N(1 - \alpha)$)

3.4.2 Tree Structure Construction. Initially, we can construct a basic balanced binary tree, as illustrated in Fig. 4 (a). For optimization, we perform *adaptive node reduction*, where the reduction of a node refers to removing it from the tree and linking its child nodes to its parent. As an example, we reduce the node n_8 in Fig. 4 (a), producing the tree in Fig. 4 (b). We examine all non-leaf nodes through postorder traversal, adaptively determining whether it should be reduced to minimize the average error in response to range queries.

The average error for all possible queries can be evaluated by accumulating the noise error from nodes within query ranges and the PL fitting error from intersecting leaf nodes. Since the PL fitting error primarily depends on the results from phase 1, we focus on the noise error of nodes, as formalized in the left part of Eq.(7). Here, \mathcal{N} denotes all nodes in the PriPL-tree, w_k denotes the probability of node n_k being involved in queries, and α_k denotes the ratio of allocated users for node n_k . As the variance of OUE employed in node frequency estimation is inversely proportional to its allocated user ratio, we simplify Err to the right part in Eq.(7). To compute weight w_k , we assume all queries arrive with equal probabilities. According to literature [26], w_k can be calculated by Eq. (8), where $[l_k, r_k]$ ($[l_p, r_p]$) denotes the interval of node n_k (its parent n_p), and the total number of possible queries is $(d + 1)d/2$.

$$Err = \sum_{n_k \in \mathcal{N}} w_k \cdot Var(\tilde{f}_k) \propto \sum_{n_k \in \mathcal{N}} w_k / \alpha_k \quad (7)$$

$$w_k = (l_k \cdot (d - r_k + 1) - l_p \cdot (d - r_p + 1)) / ((d + 1)d/2) \quad (8)$$

We provide an example in Fig. 4 to calculate errors and decide whether to reduce n_8 . Considering that the existence of a non-leaf node only influences user allocation and weight calculation for its ancestor and descendant nodes, as framed by a red dash line in Fig. 4, we compare the accumulated error from these nodes rather than the entire tree. Finally, we reduce n_8 for a smaller error.

3.5 PriPL-Tree Refinement

To address the frequency inconsistencies outlined in Section 3.2, we perform frequency and slope refinements as follows.

3.5.1 Frequency Refinement. As we know, several methods can address these inconsistencies individually, such as Norm-Sub for issue (1) [36] and the constrained inference method for issue (2) [14, 26]. However, applying one method may lead to the emergence of another inconsistency. AHEAD [7] employs these two techniques iteratively to solve issues (1) and (2), but this is inefficient and often inaccurate. To simultaneously address the three inconsistency issues, we devise an optimized constrained inference method comprising two steps: *weighted averaging* and *frequency consistency*, each requiring only a single tree traversal.

In the *weighted averaging* step, we traverse nodes from leaves to root, minimizing frequency variance for each node and addressing inconsistency issue (3). For a leaf node n_k with two frequency estimations, \hat{f}_k (from node frequency estimation) and \tilde{f}_k (from the noisy histogram estimation), we update its frequency to $\hat{f}_k = \theta \hat{f}_k + (1 - \theta) \tilde{f}_k$, where $\theta = Var(\tilde{f}_k) / (Var(\tilde{f}_k) + Var(\hat{f}_k))$, achieving the minimal variance $Var(\hat{f}_k) = Var(\tilde{f}_k) Var(\hat{f}_k) / (Var(\tilde{f}_k) + Var(\hat{f}_k))$. For non-leaf node n_k , we can similarly update its frequency to \hat{f}_k to minimize variance, using its frequency \hat{f}_k and the sum of child nodes' frequencies $\sum_{n_c \in child(n_k)} \hat{f}_c$.

In the *frequency consistency* step, we update frequencies from the root to leaves to address inconsistency issues (1) and (2). The root's frequency is fixed at 1, i.e., $\hat{f}_{root} = 1$. Given a parent node n_p with optimized frequency $\hat{f}_p \geq 0$ and its child nodes' to-be-optimized frequencies $\{\hat{f}_c | n_c \in child(p)\}$, we define the optimization problem in Eq.(9). Let D_+ be the set of child nodes with positive updated frequencies, and D_0 be those with zero updated frequencies. According to KKT conditions, the optimal frequency is $\hat{f}_c = \hat{f}_p + (\hat{f}_p - \sum_{n_j \in D_+} \hat{f}_j) / |D_+|$ if $n_c \in D_+$, and $\hat{f}_c = 0$ if $n_c \in D_0$.

$$\begin{aligned} & \min \sum_{n_c \in child(n_p)} (\hat{f}_c - \tilde{f}_c)^2 \\ \text{s.t. } & \sum_{n_c \in child(n_p)} \hat{f}_c = \hat{f}_p \text{ and } \forall n_c \in child(n_p), \hat{f}_c \geq 0 \end{aligned} \quad (9)$$

3.5.2 Slope Refinement. To meet frequency constraints within nodes, i.e., resolving inconsistency issue (1), we refine slopes based on optimal node frequencies. For node n_k with frequency \hat{f}_k , we must guarantee non-negativity at both endpoints of the interval I_k . Denoting the endpoints of I_k as l_k and r_k , we require $f(l_k) = \hat{f}_k / |I_k| - \hat{\beta}_k (|I_k| - 1) / 2 \geq 0$ and $f(r_k) = \hat{f}_k / |I_k| + \hat{\beta}_k (|I_k| - 1) / 2 \geq 0$. This establishes an effective range $[-C_k, C_k]$ for its slope, where $C_k = 2\hat{f}_k / |I_k| (|I_k| - 1)$. We then update $\hat{\beta}_k$ using Eq.(10), minimizing the error $(\hat{\beta}_k - \hat{\beta}_k)^2$ and ensuring all fitted frequencies in n_k are non-negative.

$$\hat{\beta}_k = \begin{cases} -C_k, & \hat{\beta}_k < -C_k \\ \hat{\beta}_k, & -C_k \leq \hat{\beta}_k \leq C_k \\ C_k, & \hat{\beta}_k > C_k \end{cases} \quad (10)$$

3.6 Privacy and Performance Analysis

In this subsection, we analyze the privacy guarantee, estimation error, and space and time complexity of PriPL-Tree for range queries.

3.6.1 Privacy Analysis. During the PriPL-Tree estimation, we collect and estimate frequencies twice from users via SW and OUE, each employing non-overlapping subsets of users and utilizing the full privacy budget ϵ . The PL-fitting, PriPL-Tree construction (excluding node estimation), and refinement are post-processing steps over these collected data. As such, our PriPL-Tree satisfies ϵ -LDP.

3.6.2 Error Analysis. When answering range queries from PriPL-Tree, errors arise from two sources: *noise and sampling error* and *PL estimation error*.

Noise and sampling error arise from the frequency estimation via LDP mechanisms using a subset of users. In the private PL fitting phase (phase 1), the estimated frequency \hat{f}_k of leaf node n_k , derived via the SW mechanism [23], exhibits bias depending on the data distribution and has a square error of $O(\frac{|I_k|}{N\alpha\epsilon^2})$ [8]. In the PriPL-Tree construction phase (phase 2), the node frequency \tilde{f}_k , estimated via the OUE mechanism [34], is unbiased and has variance $\text{Var}(\tilde{f}_k) = \frac{4e^\epsilon}{\alpha_k N \cdot (e^\epsilon - 1)^2} = O(\frac{1}{\alpha_k N \epsilon^2})$, where α_k represents the proportion of users allocated to node n_k . During the PriPL-Tree refinement (phase 3), we aggregate all these frequency estimates to minimize the variance of each node's frequency, leading to error bounds presented in Theorem 3.1. For the precise square error, which is necessary for subsequent multi-dimensional grids consistency refinement (i.e., Section 4.3), we introduce a numerical method. We treat the updated frequency after refinement as a weighted average of frequency estimates of all nodes from the first two phases. By accounting for specific weights, we can derive accurate errors. Due to space constraints, we elaborate this numerical method and the proof of Theorem 3.1 in Appendix C in [32].

THEOREM 3.1. *Given a PriPL-Tree with at most K segments (corresponding to K leaf nodes), the error variance of frequencies after weight averaging in refinement (phase 3) is $O\left(\frac{K \cdot \log K}{(1-\alpha) \cdot (K+1) \cdot N \cdot \epsilon^2}\right)$ for non-leaf nodes and $O\left(\frac{\log K}{(1-\alpha) \cdot N \cdot \epsilon^2}\right)$ for leaf nodes. After frequency consistency, these variance is capped at $O\left(\frac{K \cdot \log K}{(1-\alpha) \cdot N \cdot \epsilon^2}\right)$.*

PL estimation error arises when estimating the frequency of the sub-range $[l_{\text{sub}}, r_{\text{sub}}]$ within a leaf node n_k under a linear assumption. For each PL estimated frequency $f(v)$ of value v in this subrange, its square error can be approximated by $\mathbb{E}((f(v) - f_v)^2) = \mathbb{E}((f(v) - \hat{f}_v^H) + (\hat{f}_v^H - f_v))^2 \leq 2(\mathbb{E}(f(v) - \hat{f}_v^H)^2 + \mathbb{E}(\hat{f}_v^H - f_v)^2)$. The first term represents the square error of the PL function fitting the noisy histogram, while the second term represents the noise error of the noisy histogram. The magnitude of this error depends on the noisy histogram's distortion degree, the segment number, and the actual data distribution. It tends to be small for smooth distributions and large for jagged distributions. Ultimately, the total error of $Q([l_{\text{sub}}, r_{\text{sub}}])$ accumulates the error of values in it, leading to a result proportional to the range size square $O((r_{\text{sub}} - l_{\text{sub}} + 1)^2)$.

3.6.3 Space and Time Complexity Analysis. Assuming PriPL-Tree's maximum segment number, i.e., the number of leaf nodes, is K , the space complexity includes the size of PriPL-Tree, $O(K)$, and the size of parameter matrices \mathbf{X} and \mathbf{A} for private PL fitting, $O(K \cdot d)$, which is $O(K \cdot d)$ in total. The time complexity of PriPL-Tree involves two parts — the construction time and the query time. Overall, the construction time complexity mainly arises from private user data

aggregation, frequency estimation, and private PL fitting during phases 1 and 2, totaling $O(N \cdot K + d \cdot T + d \cdot \log d \cdot K^3)$. Here, T denotes the number of iterations for EM and EMS in SW mechanism during distribution estimation in Phase 1. The query time complexity, proportional to the tree height, is $O(\log_2 K)$. Due to space limitations, we provide a detailed analysis of the time complexity for each phase in Appendix E.1 of [32].

4 EXTENSION WITH ADAPTIVE GRIDS FOR MULTI-DIMENSIONAL QUERIES

Building upon the insights gained from existing works, as summarized in Section 2.4, we employ a combination of 1-D PriPL-Trees and 2-D grids to handle multi-dimensional range queries. In contrast to the uniformly constructed 2-D grids proposed by HDG [40], we introduce data-aware adaptive grids. These grids leverage the accurate 1-D marginal distributions depicted by PriPL-Trees, dynamically partitioning the entire domain into dense or sparse cells to adapt data distributions' density. As a result, they offer a more precise representation of 2-D data distribution and respond more accurately to range queries.

In this section, we first outline the workflow for multi-dimensional range queries and then present the core methods of *adaptive grid construction* and *consistency refinement*. Due to space limitations, we provide estimation error and runtime complexity in Appendix D and Appendix E.2 in our technical report [32], respectively.

4.1 Workflow of Multi-dimensional Cases

Taking inspiration from the idea of leveraging low-dimensional estimation for high-dimensional scenarios, we present the extension of PriPL-Tree for multi-dimensional range queries as follows, with a workflow diagram provided in Appendix B of [32].

Step 1: Estimating m 1-D Histograms using PriPL-Tree. Initially, we allocate half of the users to estimate 1-D marginal distributions. For each attribute A_i , we employ $N/2m$ users to estimate the PriPL-Tree and generate histograms \tilde{f}_i using the PL function within the PriPL-Tree. These histograms enable us to depict each attribute's underlying data distribution effectively.

Step 2: Estimating $\binom{m}{2}$ 2-D Adaptive Grids. For each attribute pair $\langle A_i, A_j \rangle$ ($i, j \in [m]$ and $i \neq j$), we construct a 2-D grid, as presented in Section 4.2. During each grid's construction, we partition the domain of each dimension into non-uniform intervals based on the marginal distribution, thereby forming 2-D grids that are denser in high-frequency regions and sparser in low-frequency regions, as depicted in Fig. 5. After construction, we assign $N/2\binom{m}{2}$ users to estimate the frequencies of cells in each grid using the OUE mechanism with a privacy budget of ϵ .

Step 3: Refining Consistency Between Grids and PriPL-Trees. Due to the noise introduced by the LDP mechanism, frequency inconsistencies for one attribute may arise among grids and PriPL-Trees, and some frequencies could be negative. To address these issues, we propose a post-processing method, detailed in Section 4.3, optimizing the frequencies of grids and adjusting both the frequencies and slopes of nodes in PriPL-Trees.

Step 4: Answering Range Queries. Leveraging these data structures, for 1-D queries, we can directly utilize PriPL-Trees to respond. For a 2-D query involving attributes $\langle A_i, A_j \rangle$, we answer it

by a response matrix with $d_i \times d_j$ values. This matrix represents the 2-D data distribution and is estimated based on the 1-D histograms and 2-D adaptive grids by utilizing the maximum entropy algorithm or weighted update as described in [33, 40]. The critical difference between us and [33, 40] is that they enforce the frequency sum of a sub-region equal to the frequency of its corresponding 1-D cell, while we enforce each frequency of the marginal distribution of the matrix to be equal to the frequency of the 1-D histogram. This fine-grained consistency between the response matrix and the histogram fully exploits the slope information in PriPL-Trees, yielding a more accurate distribution. Further, for λ -D ($\lambda > 2$) queries, we estimate with a response matrix with 2^λ values based on the associated 2-D queries as in [33, 40].

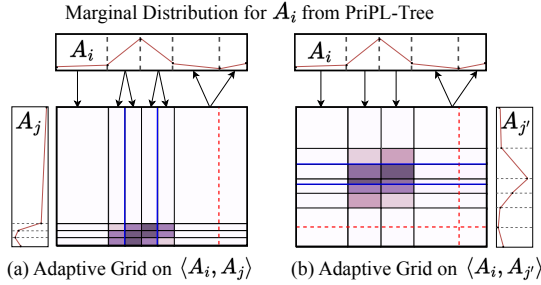


Figure 5: Examples of Adaptive Grids (Black solid lines represent partitions inherited from PriPL-Trees; blue solid lines indicate newly added partitions; red dashed lines indicate deleted partitions.)

4.2 Adaptive 2-D Grid Partitioning

The 2-D grid depicts the underlying data distribution by assuming uniform frequency distribution within each cell. To enhance its data depiction capability, we partition the grid according to the data density. In data-dense areas, densely partitioned cells provide a more accurate distribution representation, reducing reliance on uniform assumptions. Conversely, in data-sparse areas, low cell frequencies may be overwhelmed by OUE noise, diminishing their effectiveness. Therefore, in such areas, sparsely partitioned cells covering larger areas with relatively higher frequencies are preferable. Building on this concept, we introduce adaptive partitioning for 2-D grids, as illustrated in Fig. 5. For the attribute pair $\langle A_i, A_j \rangle$, we initially partition its domain according to the leaf node partitions in the PriPL-Trees, creating an initial grid G with $g_i \times g_j$ cells. We then dynamically adjust the partition lines, adding lines in high-frequency regions and removing them in low-frequency regions along each dimension, as indicated by the blue and red lines in Fig. 5. The goal of these adaptive adjustments is to minimize the squared error Err_G for range queries within this grid, while ensuring each cell’s frequency exceeds the standard deviation of the OUE noise, i.e., $\sqrt{2\sigma^2 \cdot \binom{m}{2}}$.

We introduce the squared error Err_G as below and detail the algorithm for adaptive partitioning in Appendix B.2 [32] due to space limitations. For a cell c in grid G , when fully covered by a query rectangle, it incurs a squared noise and sampling error calculated as $2\sigma^2 \binom{m}{2}$. If the cell intersects with a query rectangle, it incurs an estimation error proportional to the intersecting area o_c and intersecting frequency f_c . When adjusting the partition line on attribute

A_i , we can approximate this frequency by its marginal frequency $\tilde{f}_{\pi_i(c)}$ on A_i , where $\pi_i(\cdot)$ denotes the projection of the cell index on A_i . The squared estimation error is thus $\eta \cdot o_c \cdot \tilde{f}_{\pi_i(c)}^2$, where η is a constant and typically set at 0.03 for accurate estimations across various datasets according to our experiments. For a range query Q selecting a portion r of the area of a grid with $g_i \times g_j$ cells, the total squared error combines noise and sampling errors from $rg_i g_j$ cells and estimation errors proportional to r times the marginal frequency squares, totaling $2rg_i g_j \sigma^2 \binom{m}{2} + r\eta \sum_{c \in \pi_j(G)} o_c^2 \tilde{f}_c^2$.

4.3 Consistency Refinement

For partitions of A_j between the PriPL-Tree and the grid for $\langle A_i, A_j \rangle$, we observe several one-to-many relationships, as shown by arrows in Fig. 5. We treat each one-to-many relationship as a tree and can apply our optimized constrained inference method, detailed in Section 3.5, to ensure their frequency consistency and non-negativity. In a global view, each attributes A_i is linked to $m - 1$ grids, and each grid on $\langle A_i, A_j \rangle$ associates with two attributes. Applying the method above straightforwardly to update A_i and each related grid sequentially is challenging in maintaining global consistency. For instance, resolving consistency between A_i and $\langle A_i, A_j' \rangle$ might reintroduce inconsistencies between A_i and $\langle A_i, A_j \rangle$ that were previously resolved. Therefore, we first update all 1-D attributes, i.e., the leaf node frequencies in PriPL-Trees, by applying improved constrained inference sequentially across their $m - 1$ corresponding grids. Leveraging these updated 1-D frequencies, we then update the grids using the frequency consistency operation (i.e., the second step in the improved constrained inference). When multiple marginal cells in grids correspond to a single leaf node, we directly use the frequency of this leaf node to update the relevant cells. Conversely, if a marginal cell in grids corresponds to multiple leaf nodes, we update the cells’ frequencies using the sum of frequencies from these leaf nodes. For a grid on $\langle A_i, A_j \rangle$, to prevent reintroducing inconsistencies with A_i after aligning with A_j , we alternately update it with A_i and A_j until convergence.

Beyond ensuring consistency, by utilizing the updated frequencies of each attribute, we can further refine the PriPL-Tree to enhance accuracy for 1-D range queries. The slopes of leaf nodes can be updated based on these frequencies using the method described in Section 3.5, and non-leaf nodes can be updated by aggregating the frequencies of their children.

5 EVALUATION

In this section, we evaluate the performance of PriPL-Tree and its extension for both 1-D and multi-D range queries.

5.1 Experimental Setting

Competitors. We compare our methods with state-of-the-art techniques for range queries in LDP, including DHT [4], AHEAD [7], PrivNUD [33] for 1-D cases, and HDG [40], AHEAD, PrivNUD, PRISM [37] for multi-D cases. The baselines of hierarchical tree HH [4] and HIO [35] are beyond our scope since they have been demonstrated as inferior to DHT and AHEAD in 1-D cases [4, 7] and HDG in multi-D cases [40]. For a fair comparison, we implement these methods using their provided codes and parameters specified in their original papers.

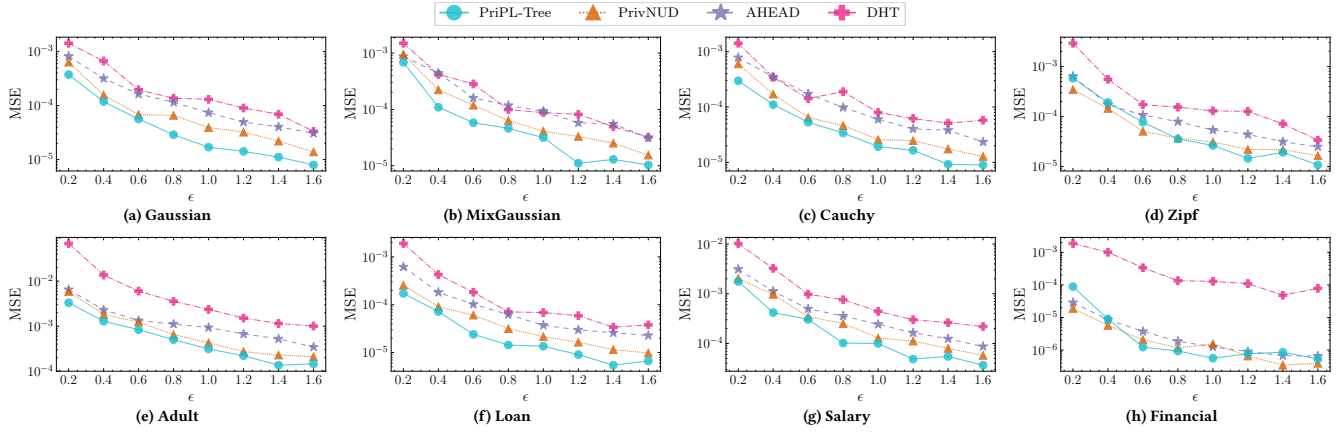


Figure 6: Evaluation for 1-D Range Queries with Varying Privacy Budget ϵ

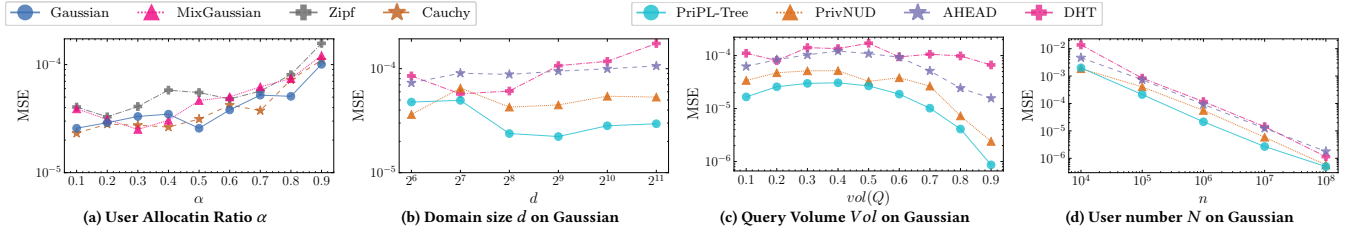


Figure 7: Evaluation for 1-D Range Queries with Varying Parameters

Datasets. We use four synthetic (Gaussian, MixGaussian, Cauchy, Zipf) and four real-world datasets (Adult [1], Loan [16], Salary [17], and Financial [18]), each with 5 dimensions. On each dimension, the synthetic datasets Gaussian, Cauchy, and Zipf sample data from $Gaussian(0, 1)$, $Cauchy(0, 1)$, and $Zipf(1.1)$ distributions, respectively. The MixGaussian dataset follows a mixture distribution of two Gaussian distributions, $Gaussian(0, 0.5)$ and $Gaussian(3, 0.8)$. For applying LDP mechanisms for frequency estimation and aligning with the existing works [4, 7, 33, 37, 40], all datasets are bucketized into domain [1024] for 1-D evaluations and [256] for multi-D evaluations, except for some attributes in the real-world dataset with discrete values and original domain sizes less than our specified ones. We provide statistics of these datasets in Tab. 2 where the mean and variance are for the default attribute in 1-D scenarios. A more detailed description is provided in Appendix F of [32].

Metrics. We employ the mean square error (MSE) [7, 33] to quantify the deviation between the estimated (\tilde{f}_Q) and actual (f_Q) answers to range queries Q , denoted as $MSE(Q) = \sum_{Q \in \mathcal{Q}} (f_Q - \tilde{f}_Q)^2 / |Q|$. During each evaluation, we test 1,000 randomly generated range queries with a specified query volume and report the final MSE by an average of 20 repeats of the experiment.

Default Settings. By default, we use a user allocation ratio $\alpha = 0.2$ and an acceleration granularity factor $\phi = 127$ in the PriPL-Tree workflow. For experiments, we set the privacy budget $\epsilon = 0.8$ and the query volume $vol(Q) = 0.5$, where query volume represents the ratio of the query range size to the domain size on each attribute. All experiments were conducted using Python 3.11 on a Linux server equipped with an Intel® Xeon® Gold 5218 CPU at 2.3GHz and 96GB of memory.

Table 2: Summary of Datasets

Dataset	U.# ^a	L.# ^b	Mean	Var.	Dataset	U.#	L.#	Mean	Var.
Gaussian	10 ⁶	0	155.0	775.34	Adult	32,561	4	30.36	336.88
MixGaussian	10 ⁶	0	135.12	2516.62	Loan	148,045	3	48.67	1620.19
Cauchy	10 ⁶	5	159.37	609.50	Salary	2,013,799	2	33.59	515.88
Zipf	10 ⁶	5	24.40	2517.47	Financial	6,362,620	5	0.23	2.65

^a U.#: The number of users (i.e., samples).

^b L.#: The number of leptokurtic attributes with a kurtosis exceeding 3.

5.2 1-D Experimental Results

In this subsection, we evaluate the performance of PriPL-Tree and its competitors (DHT, AHEAD, and PrivNUD) for 1-D range queries and analyze the impact of data and query parameters on them.

Overall Performance. We evaluate the performance of PriPL-Tree against three competitors across varying privacy budgets on synthetic and real-world datasets in Fig. 6. Our PriPL-Tree consistently outperforms its competitors on most continuous distributions, such as Gaussian, MixGaussian, Cauchy, and those found in the Adult, Loan, and Salary datasets. It significantly reduces MSEs by about 12.1% to 66.6%, averaging a 37.4% reduction across different privacy settings. In extremely leptokurtic distributions, like those in Zipf and Financial datasets, PriPL-Tree’s performance matches that of the leading competitor, PrivNUD. In these distributions, only a few values have significant frequencies, and our PriPL-Tree almost degenerates into an optimized hierarchical tree, similar to PrivNUD. It efficiently segregates high-frequency buckets into individual leaf nodes and merges those with negligible frequencies into a single node, with both frequency and slope nearing zero.

Impact of User Allocation Ratio α . In Fig. 7 (a), we assess the impact of the user allocation ratio α for phase 1 of the PriPL-Tree across four synthetic datasets. MSE remains stable for $\alpha \leq 0.5$

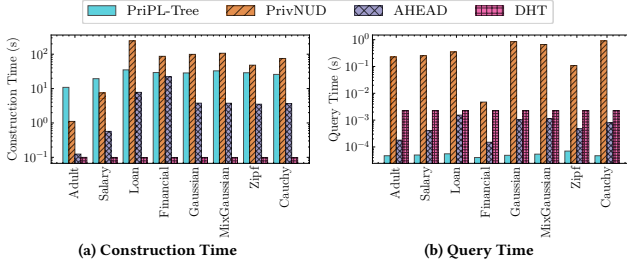


Figure 8: Runtime Evaluation under Different Datasets

and slightly increases for $\alpha > 0.5$, suggesting that fewer users are adequate for accurate PL fitting of underlying distributions. Consequently, we empirically set $\alpha = 0.2$ in our experiments.

Impact of Domain Size d . In Fig. 7 (b), we explore the impact of domain size d on the 1-D Gaussian dataset, demonstrating PriPL-Tree’s superiority, particularly in large domains. Notably, PriPL-Tree performs less effectively in very small domains, where coarse bucketizing diminishes the histogram’s ability to represent underlying distributions accurately, resulting in suboptimal PL functions and inferior outcomes.

Impact of Query Volume $vol(Q)$. In Fig. 7 (c), we assess the impact of query volume on a Gaussian dataset. PriPL-Tree consistently records the lowest MSE. All methods display an MSE that increases initially and then decreases, peaking around $vol(Q) = 0.5$. As detailed in Section 3.6.2, the error for range queries correlates with the query range size and the frequency of intersections between query ranges and leaf nodes. Below $vol(Q) = 0.5$, MSE increases primarily due to the expanding query range size. Above $vol(Q) = 0.5$, MSE decreases as intersections occur more frequently at the domain’s margins, where frequencies are lower and nearing 0, resulting in fewer PL fitting errors.

Impact of User Number N . In Fig. 7 (d), we examine the impact of user numbers on MSEs using a Gaussian dataset, noting that MSEs decrease as user numbers increase, which aligns with the law of large numbers. However, the advantage of PriPL-Tree diminishes with very small (e.g., 10^4) or very large (e.g., 10^8) user numbers. Insufficient users introduce excessive noise in LDP estimation, compromising PL parameter accuracy. On the other hand, a larger user pool mitigates LDP noise, enabling even simple hierarchical trees to provide accurate estimations.

Runtime Comparison: In Fig. 8, we explore the runtime of our method compared to competitors across various datasets. Given that our segment number K is typically small, resulting in fewer nodes and a lower tree height, our construction time remains competitive, generally under half a minute, even with additional PL fitting. On average, our construction time of 27.2s is shorter than the average for our three competitors, which is 29.8s. Additionally, our average query time is significantly lower, at around 50 μ s, while our competitors’ times remain in the millisecond range.

5.3 Multi-D Experimental Results

For multi-dimensional scenarios, we evaluate the performance of our method, PriPL-Tree with adaptive grids, against four competitors (HDG, AHEAD, PrivNUD, and PRISM). Typically, the standard

experimental setup is evaluating 2-D queries on 5-D datasets, as all high-dimensional query results are derived from these 2-D queries.

Overall Performance. We evaluate the performance of our PriPL-Tree method against competitors under various privacy budgets on 5-D synthetic and real-world datasets, as detailed in Fig. 9. Utilizing adaptive grids, the PriPL-Tree method achieves the lowest MSEs in most cases. Specifically, it achieves an average of 47.9% lower MSE than state-of-the-art solutions on real-world datasets and 23.7% lower MSE on synthetic datasets. The extent of improvement varies with the characteristics of different data distributions. Particularly in the Gaussian, MixGaussian, Loan, and Salary datasets, where most attributes are not leptokurtic, our PriPL-Tree method reduces the MSE by 10.6% to 81.9%, averaging a reduction of 56.7%. Conversely, for the other datasets, including Cauchy, Zipf, Adult, and Financial, which feature predominantly leptokurtic distributions where only a few values have significant frequencies, our PriPL-Tree almost degenerates to an optimized hierarchical tree. Consequently, the improvement becomes modest in these datasets, averaging a 14.9% reduction in MSE.

Impact of Data Dimension m . In Fig. 10 (a), we evaluate the performance of our PriPL-Tree method and competitors across varying data dimensions on the Gaussian dataset with a covariance of 0.6. PriPL-Tree consistently shows the lowest MSEs across different dimensions. As expected, all MSEs increase with the dimension m as users are distributed among more parts for estimation. Notably, two data points for AHEAD at $m \in \{25, 30\}$ are missing in the figure because their experiments exceed our server’s 96GB memory capacity, according to their open-source code. Despite these omissions, the available data points are sufficient to demonstrate that AHEAD’s performance is inferior to ours.

Impact of Query Dimension λ . In Fig. 10 (b), we assess the performance of our PriPL-Tree and other methods with varying query dimensions on a Gaussian dataset with $Cov = 0.6$. PriPL-Tree consistently records the lowest MSE, particularly noticeable in 1-D queries. During these experiments, we construct data structures across all five dimensions, with each 1-D range query selecting a dimension at random. This setup requires dividing users among $5 + \binom{5}{2} = 15$ parts, leading to fewer users per dimension and generally poorer 1-D estimations for competitors like HDG, AHEAD, PrivNUD, and PRISM. In contrast, our robust PriPL-Tree, enhanced by the consistency refinement operation in phase 3 for multi-dimensional range queries, effectively improves accuracy. This phase updates frequencies and slopes in PriPL-Trees using all related 2-D adaptive grids.

Impact of Attribute Correlation. In Figs. 10 (c) and (d), we examine the impact of attribute correlation on range queries using Gaussian and MixGaussian datasets. We use covariance to represent attribute correlation. The results reveal that PriPL-Tree consistently achieves the lowest MSEs, especially in datasets with high attribute correlations. This highlights our method’s superior capability to capture underlying distributions with adaptive 2-D grids, unlike competitors that rely on uniform grids.

6 RELATED WORK

In this section, we review related works in both central differential privacy (DP) and LDP scenarios.

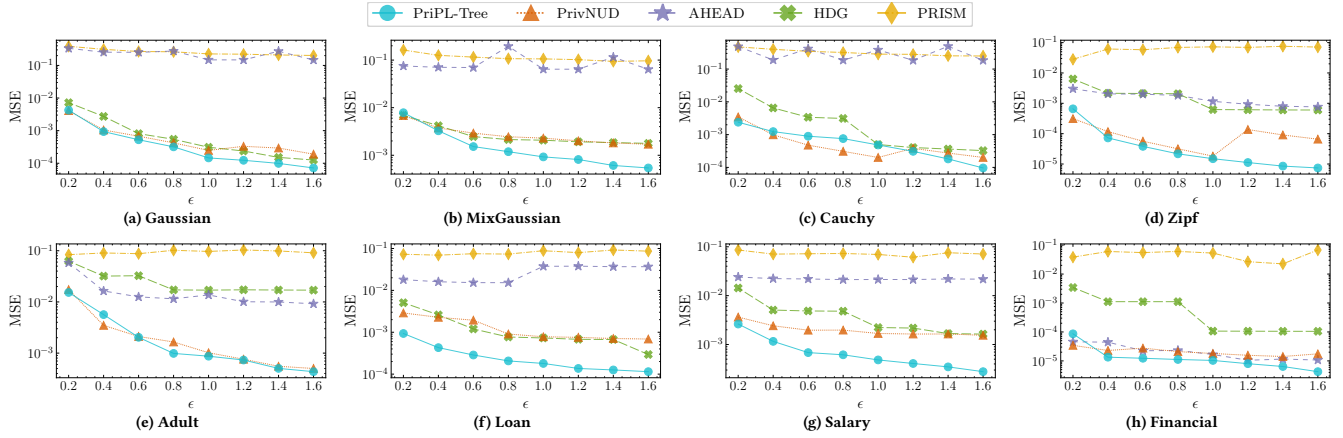


Figure 9: Evaluation for 2-D Range Queries on 5-D Datasets with Varying Privacy Budget ϵ

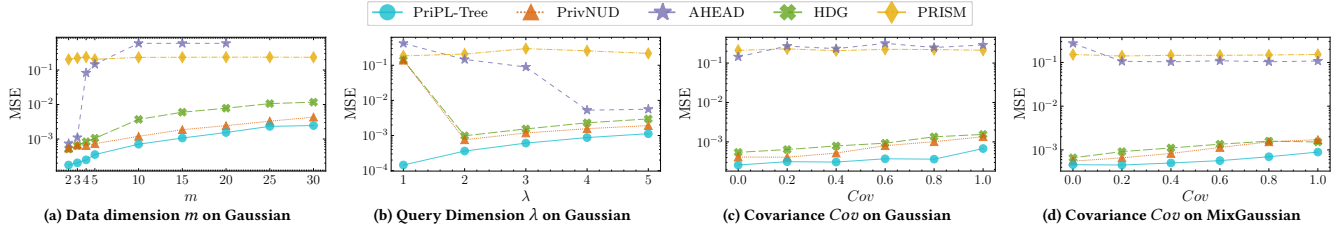


Figure 10: Evaluation for Multi-D Range Queries

Range Query under DP: In DP scenarios, the far-reaching hierarchical tree, and the constrained inference method, were first proposed by Hay et al. [14] and later analyzed and optimized by Qardaji et al. [27]. Various optimizations have since been proposed to mitigate noise errors on trees: Xiao et al. [38] enhanced trees using Haar wavelet transforms; Cormode et al. [5] proposed a geometric privacy budget allocation method; Li et al. [20] optimized the non-uniform domain partitioning and privacy budget allocation based on data distribution and query workloads; Zhang et al. [43] proposed PrivTree (i.e., a Quad-tree) with optimized node decomposition; Huang et al. [15] employed a balanced box-decomposition tree (BBD-tree) for counting arbitrarily shaped geometric ranges. Beyond hierarchical trees, Qardaji et al. [27] presented an adaptive-grid method with optimized grid granularity and declared that grid structures are more suitable for high-dimensional range queries compared to trees. Recently, Zeighami et al. [42] introduced a model-driven approach, where a model learns noisy answers from multiple 2-D range count queries and predicts results without relying on complex indexes.

Range Query under LDP: In this context, we summarize existing methods according to tree-based and grid-based methods. In tree-based methods, HH [4] and HIO [35] proposed the basic hierarchical tree in LDP almost simultaneously and optimize the tree’s fan-out (i.e., branching factor). As improvements, AHEAD [7] merged intervals with low frequencies; PrivNUD [33] customized the fan-out for each node; DHT [4] optimized this tree via Haar wavelet transformation as in [38]. In grid-based methods, HDG [40] proposed the state-of-the-art hybrid dimensional grids, and PRISM [37] replaced simple grids with prefix-sum (PS) cubes. Additionally,

there are other research topics covering range questions. McKenna et al. [25] proposed a general matrix mechanism for linear queries in LDP, which can also be applied to answer range queries. And Ye et al. [41] explored PrivKVM* for range-based estimation in key-value datasets.

Marginal Release under LDP: As a relevant problem to this work, we also review marginal release methods in LDP. Cormode et al. [3] introduced a Fourier transform-based method for private marginal release. Ren et al. [29] explored multi-dimensional joint distribution estimation using the Expectation-Maximization (EM) algorithm and Lasso regression. A more advanced method is CALM, proposed by Zhang et al. [44], which extends the idea of PriView [28] from central DP to LDP and reconstructs high-dimensional marginals using low-dimensional estimations. This idea has been widely adopted in multi-dimensional range queries.

7 CONCLUSION

In this paper, we propose the PriPL-Tree to accurately answer range queries on arbitrary data distributions. The key idea is to approximate the underlying distribution using piecewise linear functions, which alleviates both non-uniform error and LDP noise error. We further extend this with adaptive grids to handle multi-dimensional cases, where the grids dynamically adjust to the data density, thus more accurately modeling the 2-D distribution and improving accuracy for multi-dimensional range queries. Extensive experiments on both real and synthetic datasets demonstrate the effectiveness and superiority of PriPL-Tree over state-of-the-art solutions. For future work, we will explore automatic and data-aware machine learning models to further enhance estimation in LDP scenarios.

REFERENCES

- [1] Barry Becker and Ronny Kohavi. 1996. Adult. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5XW20>.
- [2] Chiranjeev Buragohain, Nisheeth Shrivastava, and Subhash Suri. 2006. Space efficient streaming algorithms for the maximum error histogram. In *2007 IEEE 23rd International Conference on Data Engineering*. IEEE, 1026–1035.
- [3] Graham Cormode, Tejas Kulkarni, and Divesh Srivastava. 2018. Marginal release under local differential privacy. In *Proceedings of the 2018 International Conference on Management of Data*. 131–146.
- [4] Graham Cormode, Tejas Kulkarni, and Divesh Srivastava. 2019. Answering Range Queries Under Local Differential Privacy. *Proc. VLDB Endow.* 12, 10 (2019), 1126–1138. <https://doi.org/10.14778/3339490.3339496>
- [5] Graham Cormode, Cecilia M. Procopiuc, Divesh Srivastava, Entong Shen, and Ting Yu. 2012. Differentially Private Spatial Decompositions. In *IEEE 28th International Conference on Data Engineering (ICDE 2012)*, Washington, DC, USA (Arlington, Virginia), 1–5 April, 2012, Anastasios Kementsietsidis and Marcos Antonio Vaz Salles (Eds.). IEEE Computer Society, 20–31. <https://doi.org/10.1109/ICDE.2012.16>
- [6] Bolin Ding, Janardhan Kulkarni, and Sergey Yekhanin. 2017. Collecting telemetry data privately. *Advances in Neural Information Processing Systems* 30 (2017).
- [7] Linkang Du, Zhikun Zhang, Shaojie Bai, Changchang Liu, Shouling Ji, Peng Cheng, and Jiming Chen. 2021. AHEAD: Adaptive Hierarchical Decomposition for Range Query under Local Differential Privacy. In *CCS '21: 2021 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, Republic of Korea, November 15 - 19, 2021*, Yongdae Kim, Jong Kim, Giovanni Vigna, and Elaine Shi (Eds.). ACM, 1266–1288. <https://doi.org/10.1145/3460120.3485668>
- [8] Jiawei Duan, Qingqing Ye, and Haibo Hu. 2022. Utility analysis and enhancement of LDP mechanisms in high-dimensional space. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 407–419.
- [9] John C Duchi, Michael I Jordan, and Martin J Wainwright. 2018. Minimax optimal procedures for locally private estimation. *J. Amer. Statist. Assoc.* 113, 521 (2018), 182–201.
- [10] Hazem Elmeleegy, Ahmed K Elmagarmid, Emmanuel Cecchet, Walid G Aref, and Willy Zwaenepoel. 2009. Online piece-wise linear approximation of numerical streams with precision guarantees. *Proceedings of the VLDB Endowment* 2, 1 (2009), 145–156.
- [11] Ulfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. 2014. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security*. 1054–1067.
- [12] Paolo Ferragina and Giorgio Vinciguerra. 2020. The PGM-index: a fully-dynamic compressed learned index with provable worst-case bounds. *Proceedings of the VLDB Endowment* 13, 8 (2020), 1162–1175.
- [13] Alex Galakatos, Michael Markovitch, Carsten Binnig, Rodrigo Fonseca, and Tim Kraska. 2019. Fiting-tree: A data-aware index structure. In *Proceedings of the 2019 International Conference on Management of Data*. 1189–1206.
- [14] Michael Hay, Vibhor Rastogi, Gerome Miklau, and Dan Suciu. 2010. Boosting the accuracy of differentially private histograms through consistency. *Proceedings of the VLDB Endowment* 3, 1-2 (2010), 1021–1032.
- [15] Ziyue Huang and Ke Yi. 2021. Approximate Range Counting Under Differential Privacy. In *37th International Symposium on Computational Geometry (SoCG 2021)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik.
- [16] Kaggle. 2007. All Lending Club loan data. <https://www.kaggle.com/datasets/wordsforthewise/lending-club>.
- [17] Kaggle. 2017. SF Salaries. <https://www.kaggle.com/datasets/kaggle/sf-salaries>.
- [18] Kaggle. 2017. Synthetic Financial Datasets For Fraud Detection. <https://www.kaggle.com/datasets/elaali/paysim1>.
- [19] Eamonn Keogh, Selina Chu, David Hart, and Michael Pazzani. 2001. An online algorithm for segmenting time series. In *Proceedings 2001 IEEE international conference on data mining*. IEEE, 289–296.
- [20] Chao Li, Michael Hay, Gerome Miklau, and Yue Wang. 2014. A Data- and Workload-Aware Query Answering Algorithm for Range Queries Under Differential Privacy. *Proc. VLDB Endow.* 7, 5 (2014), 341–352. <https://doi.org/10.14778/2732269.2732271>
- [21] Pengfei Li, Yu Hua, Jingnan Jia, and Pengfei Zuo. 2021. FINEdex: a fine-grained learned index scheme for scalable and concurrent memory systems. *Proceedings of the VLDB Endowment* 15, 2 (2021), 321–334.
- [22] Pengfei Li, Hua Lu, Qian Zheng, Long Yang, and Gang Pan. 2020. LISA: A learned index structure for spatial data. In *Proceedings of the 2020 ACM SIGMOD international conference on management of data*. 2119–2133.
- [23] Zitao Li, Tianhao Wang, Milan Lopuhaa-Zwakenberg, Ninghui Li, and Boris Skoric. 2020. Estimating numerical distributions under local differential privacy. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 621–635.
- [24] Xiaoyan Liu, Zhenjiang Lin, and Huaqing Wang. 2008. Novel online methods for time series segmentation. *IEEE Transactions on Knowledge and Data Engineering* 20, 12 (2008), 1616–1626.
- [25] Ryan McKenna, Raj Kumar Maity, Arya Mazumdar, and Gerome Miklau. 2020. A workload-adaptive mechanism for linear queries under local differential privacy. *Proc. VLDB Endow.* 13, 11 (2020), 1905–1918. <http://www.vldb.org/pvldb/vol13/p1905-mckenna.pdf>
- [26] Wahbeh Qardaji, Weining Yang, and Ninghui Li. 2013. Understanding hierarchical methods for differentially private histograms. *Proceedings of the VLDB Endowment* 6, 14 (2013), 1954–1965.
- [27] Wahbeh H. Qardaji, Weining Yang, and Ninghui Li. 2013. Differentially private grids for geospatial data. In *29th IEEE International Conference on Data Engineering, ICDE 2013, Brisbane, Australia, April 8-12, 2013*, Christian S. Jensen, Christopher M. Jermaine, and Xiaofang Zhou (Eds.). IEEE Computer Society, 757–768. <https://doi.org/10.1109/ICDE.2013.6544872>
- [28] Wahbeh H. Qardaji, Weining Yang, and Ninghui Li. 2014. PriView: practical differentially private release of marginal contingency tables. In *International Conference on Management of Data, SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014*, Curtis E. Dyreson, Feifei Li, and M. Tamer Özsu (Eds.). ACM, 1435–1446. <https://doi.org/10.1145/2588555.2588575>
- [29] Xuebin Ren, Chia-Mu Yu, Weiren Yu, Shusen Yang, Xinyu Yang, Julie A. McCann, and Philip S. Yu. 2018. LoPub: High-Dimensional Crowdsourced Data Publication With Local Differential Privacy. *IEEE Trans. Inf. Forensics Secur.* 13, 9 (2018), 2151–2166. <https://doi.org/10.1109/TIFS.2018.2812146>
- [30] Cosma Shalizi. 2015. Lecture 1: Optimal Prediction (with Refreshers). Website. <https://www.stat.cmu.edu/~cshalizi/mreg/15/lectures/01/lecture-01.pdf>.
- [31] Apple Differential Privacy Team. 2017. Learning with Privacy at Scale. (2017). <https://machinelearning.apple.com/research/learning-with-privacy-at-scale>
- [32] Leixia Wang, Qingqing Ye, Haibo Hu, and Xiaofeng Meng. 2024. PriPL-Tree: Accurate Range Query for Arbitrary Distribution under Local Differential Privacy. (2024). <https://github.com/LeixiaWang/PriPLT>
- [33] Ning Wang, Yaohua Wang, Zhigang Wang, Jie Nie, Zhiqiang Wei, Peng Tang, Yu Gu, and Ge Yu. 2023. PrivNUD: Effective Range Query Processing under Local Differential Privacy. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 2660–2672.
- [34] Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. 2017. Locally differentially private protocols for frequency estimation. In *26th USENIX Security Symposium (USENIX Security 17)*. 729–745.
- [35] Tianhao Wang, Bolin Ding, Jingren Zhou, Cheng Hong, Zhicong Huang, Ninghui Li, and Somesh Jha. 2019. Answering Multi-Dimensional Analytical Queries under Local Differential Privacy. In *Proceedings of the 2019 International Conference on Management of Data, SIGMOD Conference 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019*, Peter A. Boncz, Stefan Manegold, Anastasia Ailamaki, Amol Deshpande, and Tim Kraska (Eds.). ACM, 159–176. <https://doi.org/10.1145/3299869.3319891>
- [36] Tianhao Wang, Milan Lopuhaa-Zwakenberg, Zitao Li, Boris Skoric, and Ninghui Li. 2020. Locally Differentially Private Frequency Estimation with Consistency. In *NDSS'20: Proceedings of the NDSS Symposium*.
- [37] Yufei Wang and Xiang Cheng. 2022. PRISM: Prefix-Sum based Range Queries Processing Method under Local Differential Privacy. In *38th IEEE International Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9-12, 2022*. IEEE, 433–445. <https://doi.org/10.1109/ICDE53745.2022.00037>
- [38] Xiaokui Xiao, Guozhang Wang, and Johannes Gehrke. 2010. Differential privacy via wavelet transforms. *IEEE Transactions on knowledge and data engineering* 23, 8 (2010), 1200–1214.
- [39] Qing Xie, Chaoyi Pang, Xiaofang Zhou, Xiangliang Zhang, and Ke Deng. 2014. Maximum error-bounded piecewise linear representation for online stream approximation. *The VLDB journal* 23 (2014), 915–937.
- [40] Jianyu Yang, Tianhao Wang, Ninghui Li, Xiang Cheng, and Sen Su. 2020. Answering Multi-Dimensional Range Queries under Local Differential Privacy. *Proc. VLDB Endow.* 14, 3 (2020), 378–390. <https://doi.org/10.5555/3430915.3442436>
- [41] Qingqing Ye, Haibo Hu, Xiaofeng Meng, Huadi Zheng, Kai Huang, Chengfang Fang, and Jie Shi. 2023. PrivKVM*: Revisiting Key-Value Statistics Estimation With Local Differential Privacy. *IEEE Trans. Dependable Secur. Comput.* 20, 1 (2023), 17–35. <https://doi.org/10.1109/TDSC.2021.3107512>
- [42] Sepanta Zeighami, Ritesh Ahuja, Gabriel Ghinita, and Cyrus Shahabi. 2022. A Neural Database for Differentially Private Spatial Range Queries. *Proc. VLDB Endow.* 15, 5 (2022), 1066–1078. <https://www.vldb.org/pvldb/vol15/p1066-zeighami.pdf>
- [43] Jun Zhang, Xiaokui Xiao, and Xing Xie. 2016. PrivTree: A Differentially Private Algorithm for Hierarchical Decompositions. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, Fatma Özcan, Georgia Koutrika, and Sam Madden (Eds.). ACM, 155–170. <https://doi.org/10.1145/2882903.2882928>
- [44] Zhikun Zhang, Tianhao Wang, Ninghui Li, Shibo He, and Jiming Chen. 2018. CALM: Consistent Adaptive Local Marginal for Marginal Release under Local Differential Privacy. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018, Toronto, ON, Canada, October 15-19, 2018*, David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang (Eds.). ACM, 212–229. <https://doi.org/10.1145/3243734.3243742>

A TECHNICAL DETAILS FOR 1-D PRIPL-TREE

In this section, we detail the computation of Eq.(1) for responding to 1-D range queries and provide examples to illustrate interval partitioning and its two included strategies.

A.1 Response to 1-D Range Query

In this subsection, we focus on computing $Q([l_{\text{sub}}, r_{\text{sub}}])$ within node n_k . Given the interval $I_k = [s_{k-1}, s_k]$, slope $\tilde{\beta}_k$, and node frequency \tilde{f}_k of node n_k , we first calculate the corresponding linear function $f = \tilde{\beta}_k \cdot v + b_k$, where b_k is the intercept. This function is represented by the green line in Fig. 11. Because $\tilde{f}_k = \sum_{v \in I_k} (\tilde{\beta}_k \cdot v + b_k)$, we deduce that $b_k = \frac{\tilde{f}_k}{|I_k|} - \frac{\tilde{\beta}_k \cdot (2s_{k-1} + |I_k| - 1)}{2}$. Next, we compute $Q([l_{\text{sub}}, r_{\text{sub}}])$, which represents the frequency of the red-hatched region in Fig. 11:

$$\begin{aligned} Q([l_{\text{sub}}, r_{\text{sub}}]) &= \sum_{v \in [l_{\text{sub}}, r_{\text{sub}}]} (\tilde{\beta}_k \cdot v + b_k) \\ &= \tilde{\beta}_k \cdot \frac{(l_{\text{sub}} + r_{\text{sub}}) \cdot (r_{\text{sub}} - l_{\text{sub}} + 1)}{2} + b_k \cdot (r_{\text{sub}} - l_{\text{sub}} + 1) \\ &= (r_{\text{sub}} - l_{\text{sub}} + 1) \cdot \left(\tilde{\beta}_k \left(\frac{l_{\text{sub}} + r_{\text{sub}} + 1 - |I_k|}{2} - s_{k-1} \right) + \frac{\tilde{f}_k}{|I_k|} \right). \end{aligned}$$

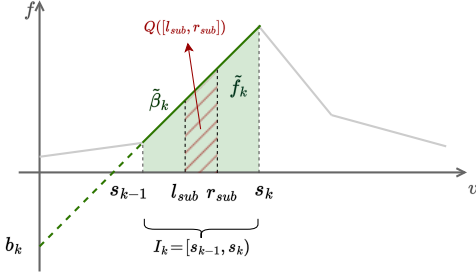


Figure 11: An Example of $Q([l_{\text{sub}}, r_{\text{sub}}])$ within node n_k .

A.2 Interval Partitioning

In this subsection, we illustrate the basic workflow of interval partitioning in Fig 12 (a). During interval partitioning, we have presented two strategies — twice partitioning and search acceleration — to enhance effectiveness and efficiency, respectively, in Section 3.6.2. Here, we provide additional examples of these two strategies for clarity and demonstrate their benefits.

A.2.1 Twice Partitioning Strategy. To illustrate the twice partitioning strategy, we present its results on the Cauchy and Salary datasets in Fig. 13. The blue solid lines, representing \hat{F}^{EM} derived by EM in SW, are jagged and noisy, while the red solid lines, representing \hat{F}^{EMS} derived by EMS in SW, are smoother but tend to flatten peak features. Initial partitions based on \hat{F}^{EM} are marked by dashed blue lines and help identify critical peaks indicated by yellow stars. Subsequent partitions on \hat{F}^{EMS} are shown with red dashed lines, identifying critical turning points marked by black

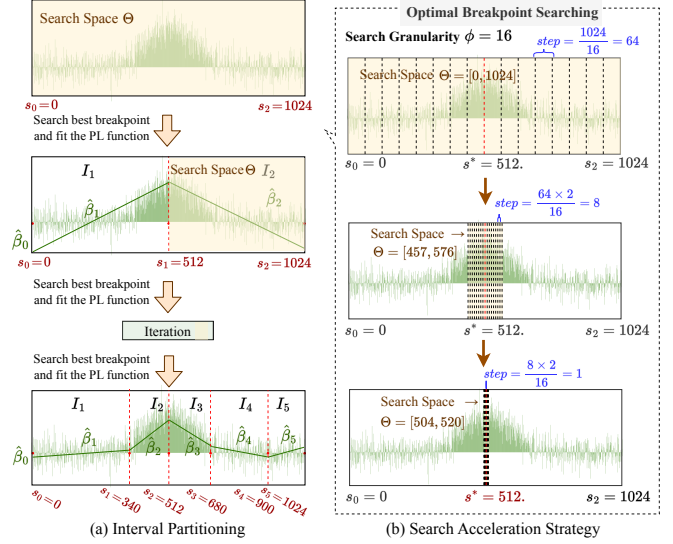


Figure 12: An Example of Interval Partitioning.

stars. In combination, this strategy facilitates the discovery of the most crucial partitions for precise PL fitting.

It is important to note that while \hat{F}^{EMS} does not directly provide accurate slopes at this phase, they can be refined during the PriPL-Tree refinement phase, deriving piecewise linear lines denoted by the black solid lines in Fig. 13. Additionally, this strategy does not significantly increase the time complexity of interval partitioning. Partitioning over the jagged \hat{F}^{EM} typically converges quickly, and the partitioning over the smoother \hat{F}^{EMS} , building on initial partitioning results, also converge swiftly. The total number of partitions remains within the maximum segment limit K_{max} .

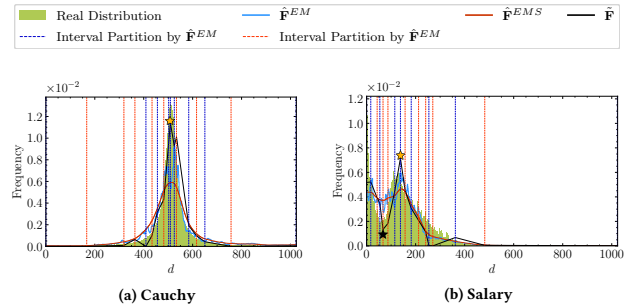


Figure 13: An Example of Twice Partitioning with $\epsilon = 0.2$

A.2.2 Search Acceleration Strategy. In Fig.12 (b), we illustrate the search acceleration strategy, which speeds up the initial breakpoint search during interval partitioning. Using a granularity factor of $\phi = 16$, we reduce the search step from 64 to 1, limiting the breakpoint candidates to no more than 16 per search. Black dashed lines mark the candidates, with the red dashed line indicating the optimal breakpoint. This method, compared to the basic one traversing all domain values, reduces the times of PL fitting for breakpoints from $O(d)$ to $O(\phi \cdot \log_{\phi} d)$. Although this acceleration might lead

to a local rather than global optimal breakpoint, it has minimal impact on final results since each breakpoint in a piecewise linear function naturally represents a local optimum to partition the domain. Moreover, because PL fitting occurs on a noisy histogram, a global optimal breakpoint might not accurately represent the actual distribution. Using a multi-granular search strategy that increases randomness can help mitigate this issue.

To validate the search acceleration strategy, we evaluated both the accelerated PriPL-Tree construction time and the MSE of queries on corresponding trees across four synthetic and four real-world datasets in Fig. 14. A granularity factor ϕ of 1024, meaning all values in the domain are traversed, represents no acceleration. The green bar indicates the accelerated private PL fitting time, which increases with ϕ . The red line shows the MSE, fluctuating irregularly with increases in ϕ , yet the variation between the highest and lowest MSEs remained within a two-fold difference. Observationally, a larger ϕ in the range of [128, 256] often provides better utility while effectively speeding up the process.

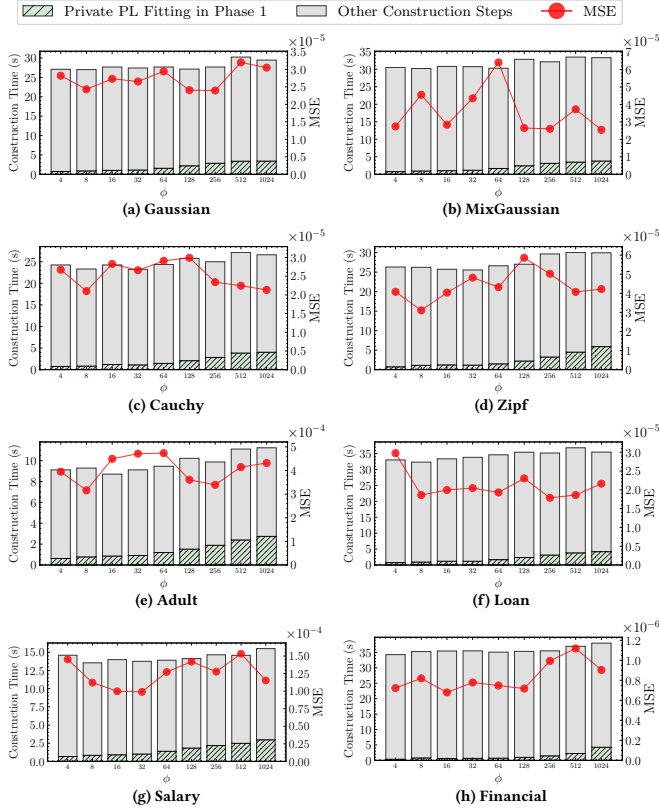


Figure 14: Construction Time and MSE for PriPL-Tree with Search Acceleration Strategy across Different Granularity Factors ϕ

B TECHNICAL DETAILS FOR MULTI-DIMENSIONAL RANGE QUERIES

In this section, we illustrate the entire workflow and detail the algorithm for adaptive 2-D grid partitioning.

Algorithm 3: Adaptive 2-D Grid Partitioning

Input: Attribute pair $\langle A_i, A_j \rangle$ and their marginal histograms
Output: Grid G with adaptive partitions

- 1 Initialize a Grid G with size $g_i \times g_j$;
- 2 Set $S = \{I_1^i, I_2^i, \dots, I_{g_i}^i\} \cup \{I_1^j, I_2^j, \dots, I_{g_j}^j\}$ and $M = \phi$;
- 3 **while** $S \neq \phi$ **do**
- 4 Select marginal cell c_I ($I \in S$) with maximal frequency \tilde{f}_I ;
- 5 Consider splitting c_I into c_{I_1} and c_{I_2} with $\tilde{f}_{c_{I_1}} \approx \tilde{f}_{c_{I_2}}$;
- 6 **if** $\tilde{f}_I > \bar{\sigma}$ and $Err_G(S) > Err_G(S - \{I\} \cup \{I_1, I_2\})$ **then**
- 7 Split c_I and update $S \leftarrow S - \{I\} \cup \{I_1, I_2\}$;
- 8 **else**
- 9 Delete I from S and add it to M ;
- 10 **while** $M \neq \phi$ **do**
- 11 Select marginal cell c_I ($I \in M$) with minimal frequency;
- 12 **if** the index I' of its neighbor exists in M **then**
- 13 Consider merging c_I and $c_{I'}$ into c_{I_3} ;
- 14 **if** $\tilde{f}_{I_3} \leq \bar{\sigma}$ or $Err_G(M) > Err_G(M - \{I, I'\} \cup \{I_3\})$ **then**
- 15 Merging c_I and $c_{I'}$, update $M \leftarrow M - \{I, I'\} \cup \{I_3\}$;
- 16 Delete I' from M ;
- 17 Delete I from M ;

B.1 The Workflow

We provide an example of the workflow in Fig. 15 to illustrate this process further.

B.2 The Algorithm For Adaptive Grid Partitioning

Building on the concept of adaptive partitioning described in Section 4.2, we detail the algorithm in Alg. 3 and provide an example in step 2 of Fig. 5. For the attribute pair $\langle A_i, A_j \rangle$, we initially partition its domain based on the leaf node partitions in the PriPL-Trees, creating a grid G with $g_i \times g_j$ cells (as shown in line 1). We then dynamically adjust partition lines along each dimension for adaptive partitioning, where adding a line splits a marginal cell (as shown in line 5) and removing a line merges marginal cells (as shown in line 13). In the algorithm, let I index marginal cells on both A_i and A_j ; we consistently select the most significant cells for splitting (high-frequency cells) or merging (low-frequency cells). This adjustment of marginal cells adheres to two principles: ensuring each cell's frequency exceeds the standard deviation of the OUE noise $\bar{\sigma} = \sqrt{2\sigma^2 \cdot \binom{m}{2}}$, and minimizing the squared error Err_G for range queries, as dictated by the splitting and merging conditions in lines 6 and 14.

C ERROR ANALYSIS ON 1-D PRIPL-TREE

In this section, we prove the asymptotic bound for noise and sampling error in PriPL-Tree and present a numerical method to accurately compute this error.

C.1 Proof and Analysis of Theorem 3.1

For convenience, we restate the theorem as follows.

THEOREM 3.1. *Given a PriPL-Tree with at most K segments (corresponding to K leaf nodes), the error variance of frequencies after*

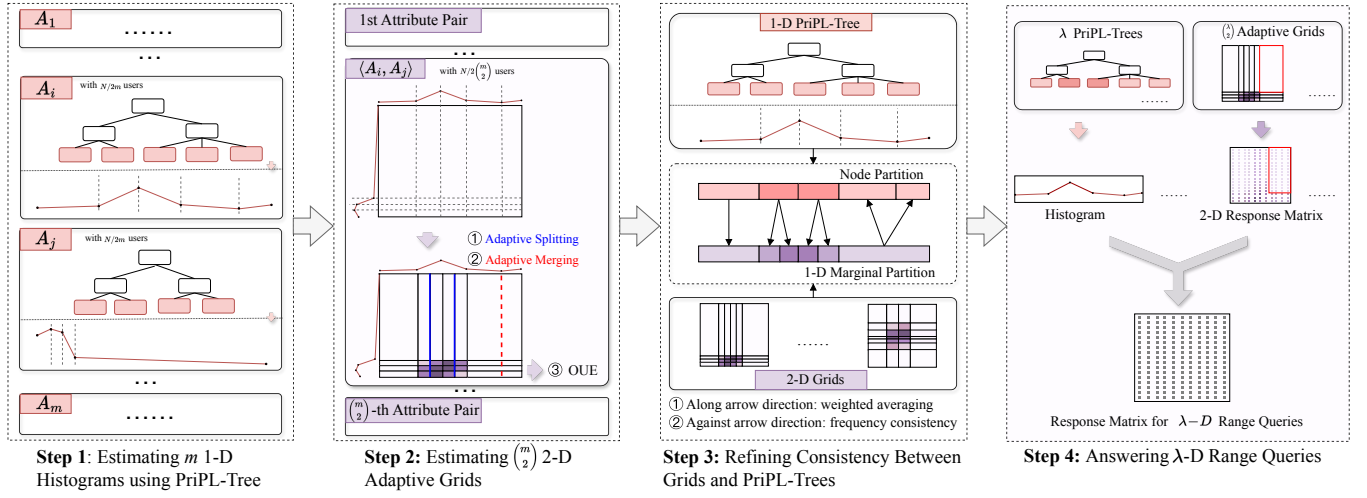


Figure 15: Workflow of Multi-dimensional Range Queries

weight averaging in refinement (phase 3) is $O\left(\frac{K \cdot \log K}{(1-\alpha) \cdot (K+1) \cdot N \cdot \epsilon^2}\right)$ for non-leaf nodes and $O\left(\frac{\log K}{(1-\alpha) \cdot N \cdot \epsilon^2}\right)$ for leaf nodes. After frequency consistency, these variance is capped at $O\left(\frac{K \log K}{(1-\alpha) N \epsilon^2}\right)$.

C.1.1 Analysis. It is important to note that after the frequency consistency step, some nodes may experience an increase in variance from $O\left(\frac{\log K}{(1-\alpha) N \epsilon^2}\right)$ to $O\left(\frac{K \log K}{(1-\alpha) N \epsilon^2}\right)$. However, most nodes exhibit reduced variances, which can be verified by computing the precise variances numerically. This phenomenon is common to any tree structure with non-uniform branchings, such as PrivNUD [33] and AHEAD [7], which exhibit heteroscedastic frequencies across nodes. In particular, in cases where the tree features uniform branching at each layer and all leaf nodes are at the same depth, each layer's nodes will be allocated an equal number of users, resulting in uniform variance among their estimated frequencies. According to the Gauss-Markov theorem [14], the variance post-frequency consistency in such scenarios can be estimated as $O\left(\frac{\log K}{(1-\alpha) N \epsilon^2}\right)$.

C.1.2 Proof. We prove the Theorem 3.1 as follows.

PROOF. For convenience, we assume the PriPL-Tree has a maximum height of $h_{\max} \leq \log_2 K$. Each node utilizes α_k users, with $\alpha_k \geq \frac{1-\alpha}{h_{\max}} \geq \frac{1-\alpha}{\log_2 K}$, and each non-leaf node has b_k branches with $1 \leq b_k \leq K$.

First, we prove the error bound after the weight averaging step.

For leaf nodes n_k , there are two frequencies: \hat{f}_k derived from the SW mechanism using αN users and \tilde{f}_k from the OUE mechanism using α_k users. By weighted averaging, the variance of the updated frequency \hat{f}_k is given by

$$\begin{aligned} \text{Var}(\hat{f}_k) &= \frac{\text{Var}(\hat{f}_k) \cdot \text{Var}(\tilde{f}_k)}{\text{Var}(\hat{f}_k) + \text{Var}(\tilde{f}_k)} \\ &= \text{Var}(\tilde{f}_k) - \frac{\text{Var}^2(\tilde{f}_k)}{\text{Var}(\hat{f}_k) + \text{Var}(\tilde{f}_k)} \\ &\leq \text{Var}(\tilde{f}_k) \end{aligned}$$

$$\begin{aligned} &= \frac{4\epsilon^\epsilon}{N \cdot \alpha_k \cdot (\epsilon^\epsilon - 1)^2} \\ &= O\left(\frac{1}{N \alpha_k \epsilon^2}\right) \\ &= O\left(\frac{\log K}{(1-\alpha) N \epsilon^2}\right). \end{aligned}$$

For the non-leaf node n_p with b_p child nodes, its frequency is updated based on its own and its children's frequencies $\{\hat{f}_c | n_c \in \text{child}(n_p)\}$. Before computing the updated variance, we present two preliminary conclusions:

(1) For a child node n_c , its updated frequency variance $\text{Var}(\hat{f}_c)$ is less than or equal to the variance of the parent node's original frequency, i.e., $\text{Var}(\tilde{f}_p)$. Assuming N' users are allocated to the subtree rooted at n_p and h_p denotes the subtree's maximum height. Since node n_p estimates its frequency \tilde{f}_p with at most N'/h_p users and its children are allocated at least N'/h_p users, the child node's frequency variance $\text{Var}(\hat{f}_c)$ does not exceed that of n_p , i.e., $\text{Var}(\hat{f}_c) \leq \text{Var}(\tilde{f}_p)$. Due to weighted averaging, the updated variance $\text{Var}(\hat{f}_c) \leq \text{Var}(\tilde{f}_c)$, implying $\text{Var}(\hat{f}_c) \leq \text{Var}(\tilde{f}_p)$.

(2) Estimates for different child nodes $\{\hat{f}_c | n_c \in \text{child}(n_p)\}$ of node n_p are independent, as each aggregates frequencies from its respective subtree, with no overlap among these subtrees.

As such, the variance of n_p for non-leaf nodes can be deduced as follows:

$$\begin{aligned} \text{Var}(\hat{f}_p) &= \frac{\text{Var}(\tilde{f}_p) \cdot \text{Var}\left(\sum_{n_c \in \text{child}(n_p)} \hat{f}_c\right)}{\text{Var}(\tilde{f}_p) + \text{Var}\left(\sum_{n_c \in \text{child}(n_p)} \hat{f}_c\right)} \\ &= \frac{\text{Var}(\tilde{f}_p) \cdot \left(\sum_{n_c \in \text{child}(n_p)} \text{Var}(\hat{f}_c)\right)}{\text{Var}(\tilde{f}_p) + \sum_{n_c \in \text{child}(n_p)} \text{Var}(\hat{f}_c)} \\ &= \text{Var}(\tilde{f}_p) - \frac{\text{Var}^2(\tilde{f}_p)}{\text{Var}(\tilde{f}_p) + \sum_{n_c \in \text{child}(n_p)} \text{Var}(\hat{f}_c)} \\ &\leq \text{Var}(\tilde{f}_p) - \frac{\text{Var}^2(\tilde{f}_p)}{\text{Var}(\tilde{f}_p) + b \cdot \text{Var}(\tilde{f}_p)} \end{aligned}$$

$$\begin{aligned}
&= \frac{b_k}{b_k + 1} \cdot \text{Var}(\tilde{f}_p) \\
&\leq \frac{K}{K + 1} \cdot \text{Var}(\tilde{f}_p) \\
&= O\left(\frac{K}{K + 1} \cdot \frac{\log K}{(1 - \alpha)N\epsilon^2}\right).
\end{aligned}$$

Next, we analyze the error variance following the frequency consistency step. Let n_c represent the child node pending update, and n_p its parent. At this stage, the nodes $n_c \in \text{child}(n_p)$ are divided into two groups, D_0 and D_+ . Nodes in D_0 typically display frequencies \hat{f}_c close to or below zero, which are then adjusted to zero. This modification might slightly increase or decrease their variances but does not alter their variance bound. For nodes in D_+ , frequency updates are based on both parental and sibling frequencies, leading to a variance as follows.

$$\begin{aligned}
\text{Var}(\tilde{f}_c) &= \text{Var}\left(\frac{|D_+| - 1}{|D_+|} \cdot \hat{f}_c + \frac{1}{|D_+|} \cdot \tilde{f}_p - \frac{1}{|D_+|} \sum_{n_s \in \text{child}(n_p)/n_c} \hat{f}_s\right) \\
&= \frac{(|D_+| - 1)^2}{|D_+|^2} \cdot \text{Var}(\hat{f}_c) + \frac{1}{|D_+|^2} \cdot \text{Var}(\tilde{f}_p) \\
&\quad + \frac{1}{|D_+|^2} \cdot \sum_{n_s \in \text{child}(n_p) \setminus n_c} \text{Var}(\hat{f}_s) + \frac{2(|D_+| - 1)}{|D_+|^2} \cdot \text{Cov}(\tilde{f}_p, \hat{f}_c) \\
&\quad - \frac{2}{|D_+|^2} \cdot \text{Cov}\left(\tilde{f}_p, \sum_{n_s \in \text{child}(n_p) \setminus n_c} \hat{f}_s\right) \\
&\leq \frac{(|D_+| - 1)^2}{|D_+|^2} \cdot \text{Var}(\hat{f}_c) + \frac{1}{|D_+|^2} \cdot \text{Var}(\tilde{f}_p) \\
&\quad + \frac{1}{|D_+|^2} \cdot \sum_{n_s \in \text{child}(n_p) \setminus n_c} \text{Var}(\hat{f}_s) \\
&\quad + \frac{2(|D_+| - 1)}{|D_+|^2} \cdot \sqrt{\text{Var}(\tilde{f}_p) \cdot \text{Var}(\hat{f}_c)} \\
&\leq \frac{|D_+|^2 + |D_+| - 1}{|D_+|^2} \max\left(\{\text{Var}(\tilde{f}_p)\} \cup \{\text{Var}(\hat{f}_j) | j \in \text{child}(n_p)\}\right) \\
&\leq \frac{|D_+| + 1}{|D_+|} \max\left(\{\text{Var}(\tilde{f}_p)\} \cup \{\text{Var}(\hat{f}_j) | j \in \text{child}(n_p)\}\right) \\
&\leq \left(\frac{|D_+| + 1}{|D_+|}\right)^{h_{\max}} \max\left(\{\text{Var}(\tilde{f}_p)\} \cup \{\text{Var}(\hat{f}_j) | j \in \text{child}(n_p)\}\right) \\
&\leq 2^{\log_2 K} \cdot O\left(\frac{\log K}{(1 - \alpha)N\epsilon^2}\right) \\
&= O\left(\frac{K \log K}{(1 - \alpha)N\epsilon^2}\right).
\end{aligned}$$

□

C.2 The Numerical Method for Calculating Noise and Sampling Error

The variances of node n_k 's frequency \tilde{f}_k after PriPL-Tree Refinement can be viewed as a weighted average of frequency estimates of all nodes from the first two phases [26], as we exemplified in Fig. 16. Here, the vector \mathbf{V} stores the original independent variances of each node, and each node n_k maintains a weight vector \mathbf{W}_k corresponding to each value in \mathbf{V} . Let $w_{k,j}$ denote the j -th value in

\mathbf{W}_k and v_j denote the j -th value in \mathbf{V} . For node n_k , its final variance can be computed as $\sum_{j \leq |\mathbf{V}|} w_{k,j}^2 \cdot v_j$.

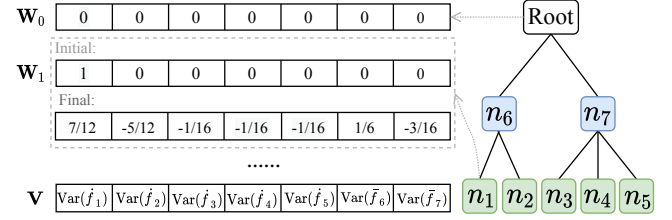


Figure 16: An Example of Variance Computation.

(For convenience, assuming all values in \mathbf{V} are equal to v^* and all frequencies remain positive during the frequency consistency step, we can derive weight \mathbf{W}_1 as illustrated in the figure. Consequently, the final updated variance for \tilde{f}_1 is $\text{Var}(\tilde{f}_1) = \frac{113}{192} v^*$.)

Following this idea, the main challenges in calculating the variance of each node's frequency involve computing the basic variance vector \mathbf{V} and each node's weight vector \mathbf{W}_k . We provide an algorithm in Algorithm 4 and present it as follows.

For the vector \mathbf{V} , each non-leaf node n_k records the variance $\text{Var}(\tilde{f}_k)$ (i.e., OUE's variance) within it, as shown in lines 6~7 in Algorithm 4. For each leaf node n_k , it stores the variance of the weighted updated frequency \hat{f}_k , which is computed by combining $\text{Var}(\tilde{f}_k)$ derived according to OUE and $\text{Var}(\hat{f}_k)$ derived according to SW, as shown in line 5 in Algorithm 4. Unlike OUE, which yields an unbiased estimate with constant variance, SW leads to estimates that are biased and affected by the unknown original data distribution, making it difficult to express the variance directly [23]. As an alternative, we can compute its empirical error $(\hat{f}_k - \tilde{f}_k)^2$ if we approximate \hat{f}_k by \tilde{f}_k . Considering $\mathbb{E}((\hat{f}_k - \tilde{f}_k)^2) = \mathbb{E}(((\hat{f}_k - \tilde{f}_k) + (\tilde{f}_k - \hat{f}_k))^2) = \mathbb{E}(\hat{f}_k - \tilde{f}_k)^2 + \text{Var}(\tilde{f}_k)$ and that SW typically derives a better estimate with smaller variance than OUE [23], we can use $(\hat{f}_k - \tilde{f}_k)^2 - \text{Var}(\tilde{f}_k)$ to approximate its square error further when $(\hat{f}_k - \tilde{f}_k)^2 > \text{Var}(\tilde{f}_k)$. For convenience, we use this approximate square error to approximate the variance, as shown in lines 3~4 in Algorithm 4.

For the weight vector \mathbf{W}_k of node n_k , after *weighted averaging*, it is updated as shown in lines 10~11 in Algorithm 4, with the optimal variance presented in Section 3.5. During *frequency consistency*, node $n_k \in D_0$, which has a frequency close to or below zero, is set to zero, and its variance changes minimally. For node $n_k \in D_+$, its updated frequency is the weighted average of itself, its parent's frequency, and its sibling $n_s \in D_+$'s frequencies. As such, its weight can be updated as indicated in lines 14~15 in Alg. 4.

D ERROR ANALYSIS ON MULTI-DIMENSIONAL QUERIES

Multi-dimensional queries primarily depend on the estimated frequency of 2-D grids, which are adaptively constructed and refined based on the PriPL-Tree for each attribute. As such, we analyze the error in queries on these 2-D grids. As discussed in Section 4.2, for a range query Q that selects a portion r of the area of a grid of size $g_i \times g_j$, the squared error is $2rg_i g_j \sigma_c^2 \binom{m}{2} + r\eta \sum_{c \in G} o_c^2 \tilde{f}_c^2$.

Algorithm 4: Numerical Method for Variance Estimation

Input: Frequencies \hat{f}_k or \tilde{f}_k for each node and the node number $|\mathcal{T}|$

Output: The error variance of each node's refined frequency.

// Initialize $V_{|\mathcal{T}| \times 1}$

```

1 Initialize vector  $V$ ;
2 for leaf node  $n_k$  do
3    $\text{Var}(\hat{f}_k) = \frac{4e^\epsilon}{\alpha_k \cdot N \cdot (e^\epsilon - 1)^2}$ ;
4    $\text{Var}(\hat{f}_k) \approx (\hat{f}_k - \tilde{f}_k)^2 > \text{Var}(\tilde{f}_k) ? (\hat{f}_k - \tilde{f}_k)^2 - \text{Var}(\tilde{f}_k) : (\hat{f}_k - \tilde{f}_k)^2$ ;
5    $v_k = \text{Var}(\hat{f}_k) = \frac{\text{Var}(\hat{f}_k) \text{Var}(\tilde{f}_k)}{\text{Var}(\hat{f}_k) + \text{Var}(\tilde{f}_k)}$ ;
6 for non-leaf node  $n_k$  do
7    $v_k = \text{Var}(\hat{f}_k) = \frac{4e^\epsilon}{\alpha_k \cdot N \cdot (e^\epsilon - 1)^2}$ ;
// Calculate  $W$ 
8 Initialize vector  $W_0 = 0_{|\mathcal{T}| \times 1}$  for the root and  $W_k$  with only the
    $k$ -th value being one and others being zero for node  $n_k$ ;
9 for non-leaf node  $n_k$  visited in postorder traversal do
10   $\theta = \frac{\sum_{n_c \in \text{child}(n_k)} \text{Var}(\hat{f}_c)}{\text{Var}(\hat{f}_k) + \sum_{n_c \in \text{child}(n_k)} \text{Var}(\hat{f}_c)}$ ;
11   $W_k = \theta \cdot W_k + (1 - \theta) \cdot (\sum_{n_c \in \text{child}(n_k)} W_c)$ ;
12 for node  $n_k$  visited in level order traversal do
13  if  $n_k \in D_+$  then
14    Let  $n_p$  represent  $n_k$ 's parent;
15     $W_k = \frac{|D_+|-1}{|D_+|} \cdot W_k + \frac{1}{|D_+|} \cdot W_p - \sum_{n_s \in D_+} \frac{1}{|D_+|} \cdot W_s$ ;
// Calculate variances
16 for each node  $n_k$  do
17   $\text{Var}(\hat{f}_k) = \sum_{1 \leq j \leq |\mathcal{T}|} w_{k,j}^2 \cdot v_j$ ;
18 return  $\text{Var}(\hat{f}_k)$  of each node  $n_k$ ;

```

where $\sigma^2 = \frac{4e^\epsilon}{N \cdot (e^\epsilon - 1)^2}$. For simplicity, let g represent the maximum grid partition size on each attribute; the asymptotic bound of the squared error is $O(\frac{r \cdot g^2 \cdot m^2}{N \cdot \epsilon^2})$. Furthermore, Theorem D.1 provides the maximum absolute error, derived from Lemma D.2 and the analysis in Section 4.2.

THEOREM D.1. For any range query Q selecting a portion r of the area of a grid with size g^2 , with at least $1 - \beta$ probability,

$$\max_c |\hat{Q}(r) - Q(r)| = O\left(\frac{rg^2 m \sqrt{\log(g^2/\beta)}}{\epsilon \sqrt{N}}\right) \quad (11)$$

LEMMA D.2. For any cell c in 2-D grids whose frequency is estimated by OUE with $\frac{N}{2 \cdot \binom{m}{2}}$ users, let g denotes the maximum grid partition size on each dimension, \hat{f}_c denotes the estimated frequency and f_c denotes the actual frequency, with at least $1 - \beta$ probability,

$$\max_c |\hat{f}_c - f_c| = O\left(\frac{m \sqrt{\log(g^2/\beta)}}{\epsilon \sqrt{N}}\right) \quad (12)$$

PROOF. For a frequency \hat{f}_c estimated by OUE using $n = \frac{N}{2 \cdot \binom{m}{2}}$ users, it is unbiased and has a variance given by $\text{Var}(\hat{f}_c) = \frac{4e^\epsilon \cdot \binom{m}{2}}{N \cdot (e^\epsilon - 1)^2}$.

By Bernstein's inequality, when ϵ is small,

$$\Pr[|\hat{f}_c - f_c| \geq \delta] \leq 2 \cdot \exp\left(-\frac{n^2 \delta^2}{2 \cdot n \cdot \frac{4e^\epsilon}{(e^\epsilon - 1)^2} + \frac{1}{3} \cdot n \cdot \delta \cdot \frac{2 \cdot e^\epsilon + 1}{e^\epsilon - 1}}\right)$$

$$= 2 \cdot \exp\left(-\frac{n \delta^2}{2 \cdot O(\frac{1}{\epsilon^2}) + \frac{1}{3} \cdot \delta \cdot O(\frac{1}{\epsilon})}\right).$$

By the union bound, there exist $\delta = O\left(\frac{\sqrt{\log(g^2/\beta)}}{\epsilon \sqrt{n}}\right) = O\left(\frac{m \sqrt{\log(g^2/\beta)}}{\epsilon \sqrt{N}}\right)$ such that $\max_c |\hat{f}_c - f_c| \leq \delta$ holds with at least $1 - \beta$ probability. \square

E TIME COMPLEXITY ANALYSIS

In this section, we theoretically and experimentally compare the time complexity between PriPL-Tree and its competitors. For convenience, we assume all attributes have the same domain size d .

E.1 1-D Time Complexity

Theoretical Analysis of PriPL-Tree: The construction time follows three phases. Phase 1 is relatively complex, involving the distribution estimation and private PL fitting two steps. The first step mainly includes the users' values aggregation time $O(N \cdot \alpha)$ and the distribution estimation time $O(d \cdot T)$ based on the EM or EMS algorithm within the SW protocol, where T denotes the number of iterations in EM and EMS. The second step includes matrix-based segment fitting, taking time $O(K^2 \cdot d)$, where calculating the inverse of matrices by the Gaussian elimination method, and interval partitioning invoking at most $O(K \cdot d)$ segment fitting. When we apply the search acceleration strategy in interval partitioning with a granularity factor ϕ (refer to Section 3.3.2), $O(K \cdot \phi \cdot \log_\phi d)$ times segment fitting is required. Next, in phase 2, involving the construction of the PriPL-Tree and the estimation of each node's frequency, the time complexity is $O(K + N \cdot (1 - \alpha) \cdot K)$. Phase 3, refining the PriPL-Tree through two traversals, leads to a complexity of $O(K)$. Overall, the construction time complexity primarily stems from private user data aggregation, frequency estimation, and private PL fitting during phases 1 and 2, totaling $O(N \cdot K + d \cdot T + d \cdot \log d \cdot K^3)$. This has been confirmed across various datasets and domain sizes d , as illustrated in Fig. 17. The legend labeled "LDP frequency estimation" represents user data aggregation and frequency estimation using existing LDP mechanisms.

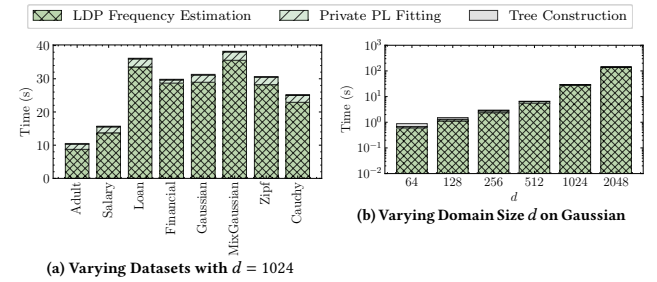


Figure 17: Construction Time of PriPL-Tree

Theoretical Analysis of Competitors: In Tab. 3, we compare the construction and query time complexities of different methods in 1-D scenarios. Because all compared methods utilize hierarchical tree structures, the construction time complexity of them primarily stems from processing users' reports, while the query time complexity depends on the tree height. In this table, the construction time complexities for AHEAD and DHT are sourced from [7], and that

for PrivNUD is analyzed by us using the same methodology. Since that the maximum segment number K in PriPL-Tree is typically small, our construction time complexity is competitive with others, and our query time is significantly lower than our competitors.

Table 3: Time Complexity Comparison for 1-D Range Queries

Methods	Construction Time	Query Time
PriPL-Tree	$O(N \cdot K + d \cdot T + d \cdot \log d \cdot K^3)$	$O(\log_2 K)$
PrivNUD	$O(\log_d 2 \cdot N \cdot d)$	$O(\log_2 d)$
AHEAD	$O(\log_d 2 \cdot N \cdot d)$	$O(\log_2 d)$
DHT	$O(N + d^3)$	$O(\log_2 d)$

Experimental Comparison: To ensure a fair comparison of method runtimes, we reimplemented DHT, originally in C++, in Python to match the programming language of the other methods. We have analyzed the construction and query times across different datasets in Section 5.2. Additionally, we evaluate the construction times of different methods on varying domain sizes d and the query times on varying query volumes.

For construction time, all methods are mainly influenced by the time of aggregating users' perturbed messages, a step inherent in employed LDP frequency estimation mechanisms, and thus increase with domain size d . In our comparisons, AHEAD [7] and PrivNUD [33] share the same asymptotic time complexity in Tab. 3, but show significant discrepancies in experimental performance. This variation stems from their different implementations, where AHEAD utilizes the Treelib package for tree structures, similar to ours, while PrivNUD uses a nested array.

For query time, DHT shows a rising trend with increasing query volume $Vol(Q)$, whereas others, including PriPL-Tree, AHEAD, and PrivNUD, initially rise then decline as query volume increases. This trend is significantly pronounced for PrivNUD but slight for PriPL-Tree and AHEAD. This observation aligns with the theoretical analysis that in the hierarchical trees, the query time correlates with the differing numbers of nodes accessed at various query volumes, typically peaking when the query volume is around 50%.

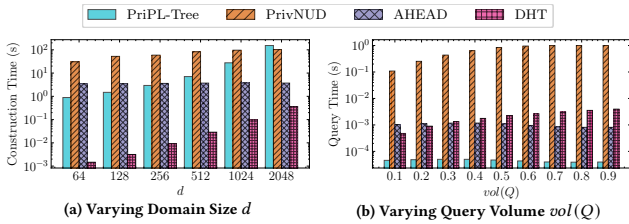


Figure 18: Runtime Evaluation on Gaussian

E.2 Multi-D Time Complexity

Theoretical Analysis: In multi-dimensional scenarios, leveraging the 1-D and 2-D estimations is a common strategy across all methods. The construction time complexity of these hybrid low-dimensional data structures is primarily from processing N user's reports, where each user reports multiple perturbed counts for nodes in the tree or cells in the grid. The query time complexity arises from computing the frequencies of the corresponding $\binom{\lambda}{2}$

2-D grids and generating the response matrix with 2^λ values. Because generating the response matrix with given $\binom{\lambda}{2}$ 2-D grids is uniform for all methods, we focus on the time complexity of 2-D grids estimation for comparison, omitting the time consumption of the response matrix step. Following this idea, we analyze the time complexity of our PriPL-Tree and other competitors in Tab. 4. During our analysis, we note that the granularity of the grids in the PriPL-Tree method is typically smaller than the maximum number of segments per attribute, denoted as K ; therefore, we set the granularity to $O(K)$. For other methods, we use the optimized granularity reported in their respective papers.

Table 4: Time Complexity Comparison for λ -d Range Queries on m -d Datasets

Methods	Construction Time	Query Time
PriPL-Tree	$O(N \cdot K^2 + m \cdot d \cdot T)$	$O(\lambda^2 \cdot K^2)$
PrivNUD	$O(\log_d 2 \cdot N \cdot d/m + N\sqrt{N}/m)$	$O(\lambda^2 \cdot \sqrt{N}/m)$
AHEAD	$O(\log_d 4 \cdot N \cdot d^2)$	$O(\lambda^2 \cdot \log_d d^2)$
HDG	$O(N\sqrt{N}/m)$	$O(\lambda^2 \cdot \sqrt{N}/m)$
PRISM	$O(N\sqrt{N}/m)$	$O(\lambda^2 \cdot \sqrt{N}/m)$

*: K is the maximum segment number in PriPL-Tree, and T is the number of iterations in the EM algorithm within the SM protocol [23].

Experimental Comparison: In Fig. 19, we assess the actual runtime of these methods. For construction time across eight different datasets, i.e., Fig. 19 (a), our PriPL-Tree performs comparably to PrivNUD and is slightly longer than HDG, which is the fastest and only uses grid structures without trees. In Fig. 19 (c), the construction times for all methods increase with data dimensionality. PRISM is an exception, with significantly larger time consumption in 2-D cases. This is because it generates significantly finer grids in the 2-D case than in other dimensions with its granularity optimization strategy [37]. For query time, i.e., Fig. 19 (b) and (d), our PriPL-Tree method consistently achieves the fastest responses, with times increasing as query dimension λ increases.

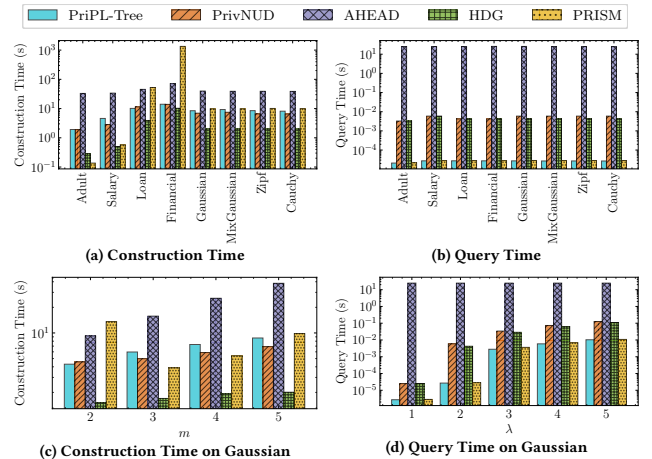


Figure 19: Runtime Evaluations (Defaultly, $m = 5$ and $\lambda = 2$)

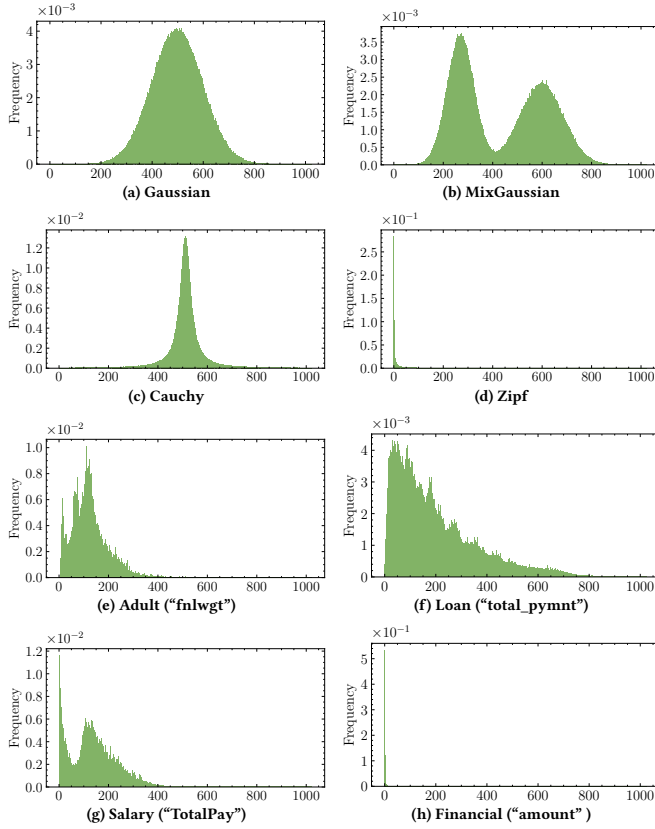


Figure 20: Data Distributions of 1-D Datasets

F DATASETS DESCRIPTION

In this section, we provide additional descriptions of both synthetic and real-world datasets. For our main 1-D scenario, we illustrate the distribution for the default attribute of each dataset in Fig. 20. For multi-dimensional scenarios, we present the mean, variance, skewness, and kurtosis for each attribute across all datasets. Skewness indicates data distribution asymmetry, with positive values suggesting a long right tail and negative values a long left tail. Kurtosis measures the peak sharpness and tail thickness of the distribution. We bold the Kurtosis values exceeding 3 in the following tables, indicating an excessively peaked distribution. Given the similar statistical characteristics across dimensions, we summarize the average indicators for synthetic datasets in Tab. 5. For real-world datasets, we detail these characteristics per attribute in the following: Adult dataset in Tab. 6, Loan dataset in Tab. 7, Salary dataset in Tab. 8, and Financial dataset in Tab. 9.

Table 5: Summary of Synthetic Datasets ($d = 256$)

Dataset	Mean	Variance	Skewness	Kurtosis
Gaussian	155.0	775.34	1.44	-0.28
MixGaussian	135.12	2516.62	1.15	-0.28
Cauchy	159.37	609.50	4.42	15.46
Zipf	24.40	2517.47	18.58	284.47

Table 6: Summary of Adult Dataset ($d = 256$)

Attribute	Mean	Variance	Skewness	Kurtosis
fnlwgt	30.36	336.88	2.19	4.03
age	21.58	186.06	-0.03	-1.64
capital-gain	2.71	353.74	15.90	250.75
capital-loss	5.11	555.48	15.90	250.93
hours-per-week	39.44	152.45	8.82	80.72

Table 7: Summary of Loan Dataset ($d = 256$)

Attribute	Mean	Variance	Skewness	Kurtosis
total_pymnt	48.67	1620.19	1.27	0.35
total_rec_int	21.69	596.95	2.98	8.38
installment	66.01	1595.08	1.12	0.37
total_il_high	5.02	29.44	5.61	32.28
_credit_limit	93.03	3661.59	4.75	26.27

Table 8: Summary of Salary Dataset ($d = 256$)

Attribute	Mean	Variance	Skewness	Kurtosis
TotalPay	33.59	515.88	1.80	2.56
TotalPayBenefits	42.15	797.60	1.64	2.64
BasePay	52.80	1174.51	1.46	2.33
OvertimePay	5.08	140.77	15.43	240.57
OtherPay	6.22	25.47	13.47	194.04

Table 9: Summary of Financial Dataset ($d = 256$)

Attribute	Mean	Variance	Skewness	Kurtosis
amount	0.23	2.65	15.76	247.68
oldbalanceOrg	3.40	151.81	15.82	249.22
newbalanceOrg	4.24	225.10	15.72	246.97
oldbalanceDest	0.57	5.61	15.57	243.49
newbalanceDest	0.64	6.60	15.47	241.30

G ADDITIONAL EVALUATION ON HIGH-DIMENSIONAL DATASET

In this section, we evaluate performance on the high-dimensional IPUMS dataset, as shown in Fig. 21. This dataset, sourced from the 2022 IPUMS repository¹, contains 50 dimensions and 15,721,123 samples. Due to the high runtime complexity of AHEAD, as detailed in Appendix E.2, it could not produce effective results on 50-D datasets within the limited time in our experimental setup, so we excluded its results from the figure. This omission does not impact our analysis, as previous experiments have already demonstrated its inferior performance compared to PrivNUD and our PriPL-Tree.

Overall, PriPL-Tree consistently achieves the lowest MSE among all competitors. For various privacy budgets (Fig. 21 (a)), PriPL-Tree shows a reduction in MSE ranging from 46.34% to 72.22%, with an average reduction of 60.67%. Across different query dimensions (Fig. 21 (b)), it reduces MSE by 56.80% to 91.72%, averaging 69.90%.

¹<https://usa.ipums.org/usa/>

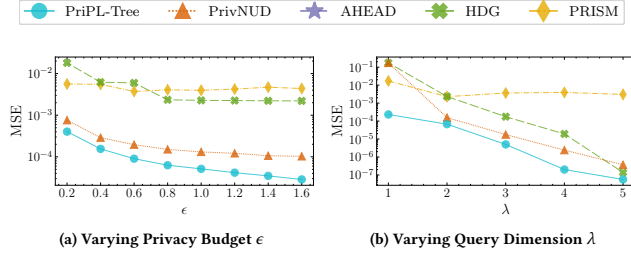


Figure 21: Evaluation of Range Queries on 50-D IPUMS

The superiority of PriPL-Tree here is more pronounced than in experiments on only 5-D datasets, as shown in Figs. 9 and 10.

Furthermore, Fig. 21 (b) shows a decrease in MSE with increasing query dimension λ , in contrast to the results from 5-D synthetic Gaussian datasets depicted in Fig.10 (b), where MSE increases with λ . This discrepancy arises from the different characteristics of the datasets. Methods estimating high-dimensional range queries using low-dimensional responses are more accurate when attribute correlations are minor. The Gaussian dataset, with high inter-attribute correlation (covariance of 0.6), shows decreasing accuracy as λ increases. In contrast, IPUMS, with generally lower correlations among attributes, allows for relatively precise estimates across varying dimensions. The observed decrease in MSE is primarily due to the lower actual frequency of higher-dimensional queries, where the error is proportional to this frequency.