

论文详解

Almost Instance-optimal Clipping for Summation Problems in the Shuffle Model of Differential Privacy

这篇论文提出了一种在 shuffle 模型下求和问题的全新方法，其核心目标是在保证差分隐私 (DP) 的前提下，不仅能够实现实例自适应 (instance-optimal) 的误差界，也能在通信成本上达到 $1 + O(1)$ 条消息/用户的高效标准。下面我将详细介绍论文中的主要内容与技术思想。

1. 研究背景与动机

• 问题背景

在差分隐私领域，传统的机制（如 Laplace 机制）通常针对最坏情况设计，因而其误差与数据取值上界 U 成正比，即 $O(U/\epsilon)$ 。然而，在许多实际数据中，大部分数据远小于 U ，因此更希望设计一种能根据数据的真实“最大值” $\text{Max}(D)$ 来自适应调整误差的机制，使得误差可以降低到 $O(\text{Max}(D) \cdot \log \log U/\epsilon)$ 的水平。

• Shuffle 模型的挑战

与中心化差分隐私不同，shuffle 模型通过引入一个可信的洗牌器，在客户端发送消息到分析器前随机打乱消息顺序，从而在保护隐私的同时获得比局部模型更优的隐私-精度权衡。过去在该模型下实现求和问题的的工作往往只能保证最坏情况的误差，或需要多轮交互（例如先求出合适的剪裁阈值 τ ，再用该 τ 来计算剪裁和），这会增加通信开销和延迟，还可能泄露额外信息。

2. 论文的主要贡献

论文的贡献主要体现在三个方面：

1. 单轮协议设计

作者设计了一种单轮 (one-round) 的协议，它能够同时确定一个合适的剪裁阈值 τ 并计算剪裁后的求和。核心思想是将数据域进行分区，然后在不同的分区上并行运行基本求和子协议 (BaseSumDP)，利用并行组合的特性避免了隐私预算的额外分摊。

2. 实例自适应误差

协议的误差界不仅依赖于最坏情况的 U ，而是依赖于数据中实际最大的元素 $\text{Max}(D)$ ，实现了 $O(\text{Max}(D) \cdot \log(\log U/\beta)/\epsilon)$ 的误差，这在实际应用中往往远优于 $O(U/\epsilon)$ 的最坏情况误差。

3. 扩展到高维与稀疏向量聚合

除了标量求和问题，论文还将方法推广到了高维求和和稀疏向量聚合问题。对于高维求和，论文通过随机旋转 (random rotation) 技巧，将 d 维数据“均匀分散”，然后在每个维度上独立使用 1D 协议；而对于稀疏向量聚合，则利用“剪裁在稀疏度上的”思路，将每个向量按照非零元素的个数进行域划分，从而实现对频率估计的高效处理。

3. 技术细节

3.1 单轮协议设计 (SumDP)

• 域分割 (Domain Partitioning)

传统的“试遍所有可能的 τ ”方法需要将隐私预算分割成 $O(\log U)$ 份，从而导致误差膨胀。论文的创新在于将数值域分成若干个不相交的区间，如 $[1,1]$ 、 $[2,2]$ 、 $[3,4]$ 、 $[5,8]$。在每个子域上，各个用户只有在数据落入该区间时才参与消息的发送，这样不仅利用了并行组合 (parallel composition) 的优势，还使得每个用户平均发送 $1+o(1)$ 条消息。

- **无额外成本的阈值选择**

在理想的非隐私情形下，我们可以通过观察最后一个非空子域来确定一个合适的 τ ，使得 $\text{Max}(D) \leq \tau \leq 2 \cdot \text{Max}(D)$ 。为了在有噪声的私有场景中避免因噪声“虚假激活”而导致过大的 τ ，论文提出了一个基于噪声估计的门槛条件（例如：只有当估计和超过 $1.3 \cdot 2^j \cdot \ln(2(\log U + 1)/\beta)/\epsilon$ 时，才将 τ 设为 2^j ）。这种方法保证了在大概率下不会过冲或严重低估，从而保持误差在实例最优的范围内。

3.2 高维求和 (HighDimSumDP)

- **随机旋转 (Random Rotation)**

为了处理 d 维向量求和问题，论文采用随机旋转矩阵 W ，将数据旋转到一个新的坐标系中。这一步骤利用了 Hadamard 矩阵与随机符号的乘积，使得经过旋转后每个维度的贡献近似均等，从而可以在每个维度上独立地应用 1D 求和协议，并最终通过逆旋转恢复原始坐标系的求和结果。理论证明表明，这样处理后的误差界为

$$O(\text{Max} \ell_2(D) \cdot \sqrt{d} \cdot \log(nd/\beta) \cdot \log(1/\delta) \cdot \log(d \cdot \log(U \ell_2)/\beta)/\epsilon)$$

3.3 稀疏向量聚合 (Sparse Vector Aggregation)

- **剪裁在稀疏度上的应用**

对于输入为二元向量且通常非常稀疏的情况（即每个向量中的 1 的个数远小于 d ），论文提出先对每个向量的稀疏度进行划分。将可能的稀疏度区间分为 $[1,1]$ 、 $[2,2]$ 、 $[3,4]$ 等，然后分别对每个区间内的向量进行求和和计数估计。通过一个额外的计数器，协议能够判断每个区间中是否有足够的非零元素，从而确定最终聚合时应该采用哪些区间的结果。最终，每个维度上的 ℓ_∞ 误差可以控制在

$O((\text{Max} \ell_2(D) \cdot \sqrt{(\log(1/\delta) + \log \log d)) \cdot \log(d/\beta)/\epsilon})$ 的水平，而通信成本则依赖于每个向量的非零个数。

4. 实验与优化

论文不仅在理论上证明了新协议的优越性，还通过实验验证了其实用性。实验部分主要涵盖：

- **标量求和实验**

在合成数据（如 Zipf 分布、Gauss 分布）以及真实世界数据（如薪资数据、贸易数据）上，SumDP 显著降低了相对误差（有时比最坏情况方法低 3000 倍以上），而且通信成本远低于之前的多轮协议。

- **高维求和与稀疏向量聚合实验**

对于 MNIST 数据集和 AOL 用户点击数据，新的高维求和和稀疏聚合方法均在误差与通信量上表现出色，与中心化 DP 的最优机制相当甚至更好。

此外，论文还讨论了如何在实际实现中根据参数 n 、 ϵ 、 δ 等选择 [GKM+21] 与 [BBGN20] 的子协议，以在不同场景下达到最佳的通信效率和误差平衡。

5. 结论与未来工作

论文总结了主要成果：设计了一种单轮、实例自适应且通信高效的差分隐私求和协议，并将其扩展到高维和稀疏数据聚合问题。未来研究方向包括：

- 将这种域分割技术扩展到其他模型（如多方安全计算）；
- 探索在更复杂的机器学习任务中，如何利用这种私有求和机制来进一步提升整体的精度与效率