

Paillier 暗号を用いたデータベース演算実装方式における 性能解析に関する検討

内藤 華†

中野 美由紀‡

小口 正人†

†お茶の水女子大学

‡津田塾大学

1 はじめに

第三者のサーバ上で機密情報を含むデータを扱う際、情報漏洩や改ざんを防ぐために暗号化が必要である。準同型暗号という方式を用いると、処理のために復号する必要がなく、データを暗号化したまま演算結果を取り出すことができる。加算または乗算のどちらかのみ計算可能な部分準同型暗号や演算回数に制限のない完全準同型暗号などが存在し、行える演算の種類や回数の制限が少ないほど、有用である反面、処理負荷が大きくなることが知られている。本研究で用いる Paillier 暗号 [1] は、Paillier によって提案された公開鍵暗号方式で、加法準同型性を持つ。したがって、他の方式と比べてデータサイズが小さく低いコストでの演算が可能である。さらに、Chowdhury らの Crypte [2] では、Paillier 暗号を用いて暗号文同士の加算に加えて乗算も行う手法が提案されている。本稿では、Crypte を基に 7 種類のデータベース演算プリミティブを実装し、演算の手法や処理速度について考察する。さらに、Crypte とは異なる手法での演算方法の提案を行う。

2 関連研究

Crypte は、暗号化されたデータに対して 2 種類のサーバを用いてプログラムを実行することで、差分プライバシーを満たした演算結果を出力するシステムである。

演算を実行する Analytics Server (以下 AS) と、

鍵やプライバシーコストの管理を行う Cryptographic Service Provider (以下 CSP) の 2 つのサーバを中心に動作する。CSP はデータ所有者のプライバシーコストの見積もり (以下 ϵ^B) を記録し、Paillier 暗号で秘密鍵と公開鍵の鍵ペアを生成する。各データ所有者は、CSP が生成した公開鍵で自分のデータを暗号化し AS に送信する。暗号化に Labeled Homomorphic Encryption (以下 labHE) [3] を用いることで暗号文同士の乗算が可能となる。AS はそれらのデータを集めてプログラムを実行、ノイズを付与した演算結果を CSP に送信する。CSP は、演算によって発生したプライバシーコストが ϵ^B を超えていないことを確認できると、保管しておいた秘密鍵で復号し結果を出力する。

3 実験

3.1 処理の手順

図 1 は、Gender について GroupByCount を実行する際の暗号化から演算、復号までの流れを示している。赤い矢印は、暗号化、演算、復号の実行時間として計測した箇所を表している。2 台の PC 上に作成した AS と CSP が通信を行い、演算処理前と処理後の暗号データは Google Cloud Storage を介して共有を行う。CSP はデータを one-hot-encoding に変換後 Paillier 暗号で暗号化し、クラウドに送信する。AS は、取り出した暗号データに対して演算を実行した結果を暗号文のままクラウド上に保存、CSP が読み込んで復号する。Adult データセット [4] から属性 (Age, Gender, NativeCountry, Race) を抽出し、レコード数 10000 件のテーブルを作成した。暗号化は、Paillier 暗号を用いて labHE で行なった。なお、本稿では差分プライバシーのためのノイズ付

Considerations on a Method of Paillier-Encrypted Database Primitives

†Hana Naito

‡Miyuki Nakano

†Masato Oguchi

†Ochanomizu University

‡Tsuda University

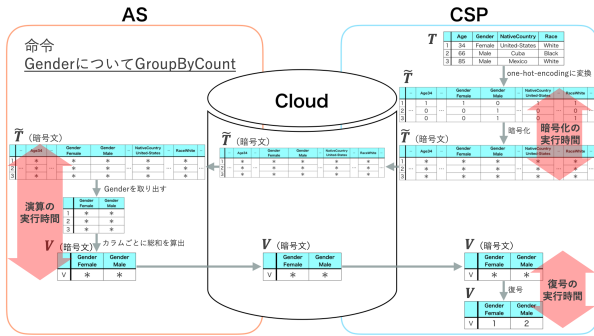


図 1: システム構成

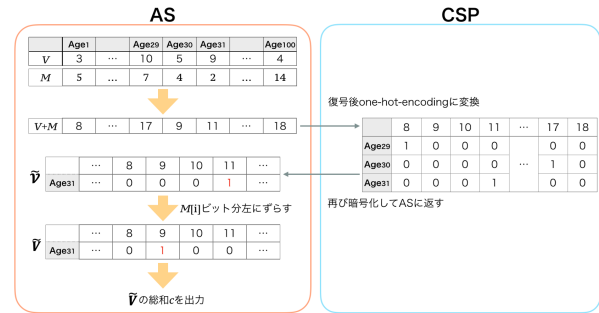
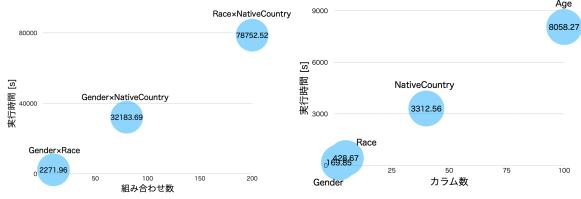


図 3: CountDistinct 演算方法



(a) CrossProduct (b) GroupByCount

図 2: カラム数と実行時間の関係

の個数を求める. Cryptε と異なる手法で演算を行うことにより簡略化を図った. 図 3 に属性 Age について CountDistinct を行う方法を示す. 他の演算と同様, one-hot-encoding で表現したときのカラム数と実行時間はほぼ比例した.

4 まとめと今後の課題

labHE を用いることで, 比較的小さい処理コストで加算と乗算を行うことができた. CrossProduct のように複数の属性の直積を用いる演算では処理時間が増大した. 一方で, Count では B を効率的に利用することで短時間での処理が可能であった. 今後の課題として, 直積演算やデータの暗号化の効率的な手法を検討したい.

参考文献

- [1] P. Paillier, "Public-key cryptosystems based on composite degree residuosity classes," Proceedings of the 17th International Conference on Theory and Application of Cryptographic Techniques, EUROCRYPT' 99, pp. 223–238, Berlin, Heidelberg, 1999. Springer-Verlag.
- [2] A. Roy Chowdhury, C. Wang, X. He, A. Machanavajjhala, and S. Jha, "Cryptε: Crypto-Assisted Differential Privacy on Untrusted Servers," Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data, pp. 603–619, 2020
- [3] M. Barbosa, D. Catalano, and D. Fiore. Labeled homomorphic encryption - scalable and privacy-preserving processing of outsourced data. In ESORICS, 2017.
- [4] Dua, D. and Graff, C. (2019). UCI Machine Learning Repository. Irvine, CA: University of California, School of Information and Computer Science.

与とプライバシーコストの計算は行わない.

3.2 基本性能評価結果

本節では, CrossProduct, GroupByCount, CountDistinct の 3 種類の演算の実行時間について考察する.

CrossProduct: テーブル \tilde{T} に含まれる属性 A_i と A_j の直積をとり新たな属性 A' で置き換える. 2つの属性のカラム数の積 s が小さい方から 3通りの組み合わせについて実行時間を計測した. 2つの属性の直積を求めるためにレコード 1 件につき乗算を s 回行う必要があり, ほかの演算と比べて処理時間が長くなった. 2つの属性に含まれる値の組み合わせ総数と実行時間は, 図 2(a) の通りほぼ比例している. このデータセットで最も時間のかかる組み合わせは NativeCountry と Age で, 400 時間以上かかると予想できる.

GroupByCount: \tilde{T} のカラムごとにすべての要素を足し合わせることでそれぞれの値の個数を求め, ベクトル V を出力する. それぞれの属性を one-hot-encoding で表現したときのカラム数と GroupByCount の処理時間の関係は図 2(b) の通りである. カラム数に伴って処理時間も増加していることが読み取れる.

CountDistinct: 属性 A について重複しない値