

紛失通信とアダマール行列を用いてポイズニング安全性を強化した LDP 方式の提案

清水 正浩 †

明治大学総合数理学部 先端メディアサイエンス学科 菊池研究室 †

1 はじめに

近年, スマートデバイスの普及により, サービス事業者はユーザの使用履歴を盛んに収集し, 利活用している. しかし, サービス事業者はユーザの多くのデータを収集するため, ユーザのプライバシーを守ることができないという課題があった. その課題を解決するために Duchi らによって局所差分プライバシー (Local Differential Privacy Protocol, LDP) が提案された [6]. LDP はユーザが自身のデータにノイズを加えた後に, サービス事業者に送信する. サービス事業者はユーザから送信されたデータに脱ノイズ処理を施し, 集計を行う.

しかし, 局所差分プライバシーはユーザが局所的にノイズ処理を行うために, 悪意のあるユーザが意図的なデータをサーバに送信して, 集計結果を操作するポイズニング攻撃に対して脆弱であることが Cao らによって指摘されている [1].

そこで, 本研究では, 2017 年に Apple から提案された局所差分プライバシー方式 Count Mean Sketch(CMS) に対する, 3 つのポイズニング攻撃をする. ポイズニング攻撃に対するロバスト性を向上させるために, CMS に紛失通信プロトコルを適用した OT-CMS を提案する. しかし, OT-CMS はハッシュ関数の値域に比例して送信量と処理コストが増加してしまうという問題点がある. そこで, その問題点を解決するために, Hadamard Count Mean Sketch(HCMS)[3] を導入し, 紛失通信プロトコルの通信コストを削減した OT-HCMS を提案する. 本研究の概要を図を 1 に示す.

2 準備

2.1 基本定義

使用する記号を表 1 に整理する. 各ユーザーは自身のプライベートなデータ d にノイズを付与し, 摂動化されたデータ \tilde{d} をサービス事業者に送信する. サービス事業

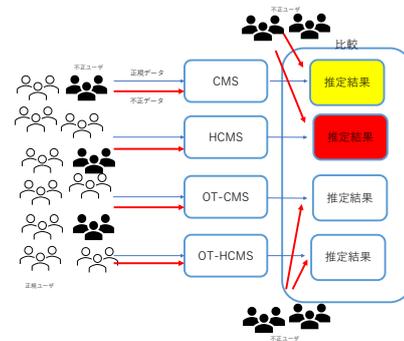


図 1 研究概要図

表 1 記号一覧

記号	意味
D	アイテムの集合
ϵ	プライバシー予算
H	ハッシュ関数の集合
k	ハッシュ関数の集合のサイズ
m	ハッシュ関数の値域
n'	不正ユーザーの数
β	不正ユーザーの割合
T	ターゲットアイテムの集合
t	ターゲットアイテム
r	ターゲットアイテムの数

者は各ユーザーから収集したデータを集計し, 頻度などを推定する. 各ユーザーはデータにノイズを付与する際, ランダム化アルゴリズム $M(d, \epsilon)$ を用いることで, プライバシー予算を ϵ と d を入力とする.

局所差分プライバシーは, 任意の異なる 2 つの入力に対して, ランダム化アルゴリズム M の出力確率に区別がつかないことを保証する. これにより, サービス事業者は送信されたデータを用いてユーザの真の入力を知ることができず, ユーザのプライバシーを保証する. ランダム化アルゴリズム M に対して局所差分プライバシーは以下のように定義される.

定義 2.1. D を入力の集合, Z を出力の集合と定義する. ランダム化アルゴリズム M は入力 $d \in D$ を受け取り, $z \in Z$ を出力とする. この時, 任意の異なる 2 つの入力 $d_1, d_2 \in D$ の出力 $z \in Z$ に対して,

†Kikuchi Laboratory, Department of Frontier Media Science, School of Interdisciplinary Mathematical Science, Meiji University.

$$\Pr(M(d_1, \varepsilon) = z) \leq e^\varepsilon \Pr(M(d_2, \varepsilon) = z)$$

が成立するとき、ランダムイズアルゴリズム ε -局所差分プライバシーを満たすという。

2.2 Count Mean Sketch

Count Mean Sketch(CMS) は 2017 年に Apple が提案した局所差分プライバシー方式の一つである [3]。CMS はユーザの使用履歴を収集し、その頻度を推定する。ユーザとサーバでハッシュ関数を共有する。

入力: ハッシュ関数の集合を $H = \{h_j | h_j : D \rightarrow [m], j \in [k]\}$ とする。各ユーザは自身のデータ $d \in D$ を H から一様ランダムに取得したハッシュ関数を用い、次のようにして、 m 次元ベクトル v に変換する。

(例 1) あるサービスを使用しているユーザの性別の頻度を推定する場合を考える。 $D = \{\text{“男”}, \text{“女”}\}, m = 2, k = 4$ とする。この時、ユーザが $d = \text{“男”} \in D$ というデータを持っている場合、 j を $\{1, 2, 3, 4\}$ から一様ランダムにサンプリングをして $j = 3 \in \{1, 2, 3, 4\}$ とする。 $h_3(\text{“男”}) = 2$ を計算する。2次元ベクトル v の2番目の要素を1に設定し、そのほかの要素を -1 とする。得られる2次元ベクトルは $v = (-1, 1)$ である。 $v = (-1, 1)$ に摂動化を施した \tilde{v} をサーバに送信する。

摂動: m 次元ベクトルを (v_1, \dots, v_m) とする。 $i = 1, \dots, m$ について確率 p で真の値 v_i を出力し、確率 $q = 1 - p$ で $-v_i$ を出力する。すなわち、

$$\tilde{v}_i = \begin{cases} v_i & w./p. \quad p, \\ -v_i & w./p. \quad q. \end{cases}$$

$$p = \frac{e^{\frac{\varepsilon}{2}}}{1 + e^{\frac{\varepsilon}{2}}}, \quad q = \frac{1}{1 + e^{\frac{\varepsilon}{2}}}$$

のとき、 ε -局所差分プライバシーを満たす。

集計: n 人のユーザからの出力を収集し、各 $d_i \in D$ の頻度を推定する。CMS は収集したデータを用いて、Sketch Matrix と呼ばれる $k \times m$ の行列を作成する。ユーザから収集したデータの集合を $S = \{(\tilde{v}^{(1)}, j^{(1)}), \dots, (\tilde{v}^{(n)}, j^{(n)})\}$, $c_\varepsilon = \frac{e^{\frac{\varepsilon}{2}} + 1}{e^{\frac{\varepsilon}{2}} - 1}$ と定義する。この時、 $\tilde{v}^{(i)}, k, m$ を用いて、 $\tilde{x}^{(i)} = k(\frac{c_\varepsilon}{2} \tilde{v}^{(i)} + \frac{1}{2} \mathbf{1})$ を計算する。 $\tilde{x}^{(i)}$ を累積して、Sketch Matrix M を構築する。 $i \in [n], \ell \in [k]$ とすると、 M は $j^{(i)}$ 行 ℓ 列の要素 $M_{j^{(i)}, \ell}$ に $\tilde{x}_{\ell^{(i)}}$ の累積する。Sketch Matrix M から、

$$\tilde{f}(d) = \left(\frac{m}{m-1} \right) \left(\frac{1}{k} \sum_{\ell=1}^k M_{\ell, h_\ell(d)} - \frac{n}{m} \right)$$

として、アイテム d のハッシュエントリを平均化することによって、頻度推定を行い $\tilde{f}(d)$ を求める。

(例 2) あるサービスを使用しているユーザの性別の頻度を推定する場合を考える。 $n = 4, D = \{\text{“男”}, \text{“女”}\}, m = 2, k = 2, \varepsilon = \infty, S = \{((1, -1), 1), ((1, -1), 2), ((-1, 1), 1), ((-1, 1), 2)\}$ とする。この時、 $x^{(1)} = \text{“男”}, x^{(2)} = \text{“男”}, x^{(3)} = \text{“女”}, x^{(4)} = \text{“女”}$ とし、 $\tilde{x}^{(1)}, \tilde{x}^{(2)}, \tilde{x}^{(3)}, \tilde{x}^{(4)}$ は以下のように計算される。

$$\tilde{x}^{(1)} = 2 \left(\frac{1}{2}(1, -1) + \frac{1}{2}(1, 1) \right) = (2, 0)$$

$$\tilde{x}^{(2)} = 2 \left(\frac{1}{2}(1, -1) + \frac{1}{2}(1, 1) \right) = (2, 0)$$

$$\tilde{x}^{(3)} = 2 \left(\frac{1}{2}(-1, 1) + \frac{1}{2}(1, 1) \right) = (2, 0)$$

$$\tilde{x}^{(4)} = 2 \left(\frac{1}{2}(-1, 1) + \frac{1}{2}(1, 1) \right) = (0, 2)$$

この時、Sketch Matrix は

$$M = \begin{pmatrix} 4 & 0 \\ 2 & 2 \end{pmatrix}$$

と得られる。従って、

$$\tilde{f}(\text{“男”}) = \frac{2}{2-1} \left(\frac{1}{2} \sum_{\ell=1}^2 M_{\ell, h_\ell(\text{“男”})} - \frac{4}{2} \right) = 2$$

$$\tilde{f}(\text{“女”}) = \left(\frac{2}{2-1} \right) \left(\frac{1}{2} \sum_{\ell=1}^2 M_{\ell, h_\ell(\text{“女”})} - \frac{4}{2} \right) = 2$$

3 Hadamard Count Mean Sketch

HCMS は Apple に提案された CMS の亜種である [3]。CMS は、ユーザは m 次元ベクトルを送信するため、ハッシュ関数の値域の大きさに比例して送信量が大きくなってしまふ。その問題点を解決するために、CMS に改良を施したのが HCMS である。HCMS はアダマール行列を適用することによって、送信量を減らすことが可能になる。また、アダマール行列はユーザとサーバで共通しているものとする。

入力: ハッシュ関数の集合を $H = \{h_j | h_j : D \rightarrow [m] : j \in [k]\}$ とおく。各ユーザは自身のデータ $d \in D$ を H から一様ランダムに取得したハッシュ関数を用いて m 次元ベクトル v に変換する。アダマール行列と積をとり m 次元ベクトル w に変換する。そのベクトル w の中から一様ランダムに1ビットを取得し、摂動化する。その1ビットをサーバに送信する。

ここで、アダマール行列は、以下のように再帰的に定義される。

$$H_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$

$$H_m = \begin{pmatrix} H_{m/2} & H_{m/2} \\ H_{m/2} & -H_{m/2} \end{pmatrix}$$

(例 3)(例 1) と同様にユーザの性別の頻度を推定する場合を考える。 $D = \{\text{“男”}, \text{“女”}\}, m = 2, k = 4$ とする。この時、あるユーザが $d = \text{“男”} \in D$ というデータを持っている場合、次のような処理をする。 j を $\{1, 2, 3, 4\}$ から一様ランダムにサンプリングをして $j = 3 \in \{1, 2, 3, 4\}$ とする。 $h_3(\text{“男”}) = 1$ とする。2次元ベクトル v の1番目の要素を1に、そのほかの要素は例1と異なり0とする。結果として得られる2次元ベクトルは $v = (0, 1)$ と表される。

$$w = H_2 v = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = (1, -1)$$

となり、2次元の要素の中から一様に $\ell = 1 \in \{1, 2\}$ をサンプリングし、 $\tilde{w}_\ell = 1, j = 3, \ell = 1$ をサーバに送信する。

摂動: w の中から一様ランダムに取得された1ビットを w とおく。CMSと同様に、確率 p で真の値 v を出力し、確率 $q = 1 - p$ で $-v$ を出力する。

$$\tilde{w} = \begin{cases} w & w./p. \quad p, \\ -w & w./p. \quad q. \end{cases}$$

$$p = \frac{e^\varepsilon}{1 + e^\varepsilon}, \quad q = \frac{1}{1 + e^\varepsilon}$$

のとき、 ε -局所差分プライバシーを満たす。

集計: n 人のユーザからの出力を収集し、各 $d_i \in D$ の頻度を推定する。ユーザから収集した摂動化データを $S = ((\tilde{w}^{(1)}, j^{(1)}, \ell^{(1)}), \dots, (\tilde{w}^{(n)}, j^{(n)}, \ell^{(n)}))$ 、プライバシー予算 $\varepsilon, c_\varepsilon = \frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1}$ と定義する。この時、 $w^{(i)}, k, m$ を用いて、 $\tilde{x}^{(i)} = k c_\varepsilon w^{(i)}$ を求める。 $\tilde{x}^{(i)}$ を用いて、Sketch Matrix M を構築する。 $i \in [n], \ell \in [m]$ とすると、 M の $j^{(i)}$ 行 $\ell^{(i)}$ 列の要素 $M_{j^{(i)}, \ell^{(i)}}$ に $\tilde{x}_{\ell^{(i)}}$ を累積する。Sketch Matrix に、アダマール行列を適用し、正規化する。

$$\tilde{f}(d) = \left(\frac{m}{m-1} \right) \left(\frac{1}{k} \sum_{\ell=1}^k M_{\ell, h_\ell(d)} - \frac{n}{m} \right)$$

各アイテムのハッシュエントリを平均化することによって、頻度推定を行う。

(例 4) $n = 4, D = \{\text{“男”}, \text{“女”}\}, m = 2, k = 2, \varepsilon = \infty, S = ((1, 1, 1), (1, 1, 2), (1, 2, 1), (1, 2, 1))$ とする。この時、 $\tilde{x}^{(1)}, \tilde{x}^{(2)}, \tilde{x}^{(3)}, \tilde{x}^{(4)}$ は以下のように計算される。

$$\tilde{x}^{(1)} = 2 \cdot 1 \cdot 1 = 2$$

$$\tilde{x}^{(2)} = 2 \cdot 1 \cdot 1 = 2$$

$$\tilde{x}^{(3)} = 2 \cdot 1 \cdot 1 = 2$$

$$\tilde{x}^{(4)} = 2 \cdot 1 \cdot 1 = 2$$

この時、Sketch Matrix は以下のように表される。

$$M = \begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 2 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 2 & 0 \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 4 & 0 \end{pmatrix}$$

M に、正規化を施すと、

$$M H_2^{-1} = \begin{pmatrix} 2 & 2 \\ 4 & 0 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} = \begin{pmatrix} 4 & 0 \\ 4 & 4 \end{pmatrix},$$

$$\tilde{f}(\text{“男”}) = \left(\frac{2}{2-1} \right) \left(\frac{1}{2} \sum_{\ell=1}^2 M_{\ell, h_\ell(\text{“男”})} - \frac{4}{2} \right) = 4$$

$$\tilde{f}(\text{“女”}) = \left(\frac{2}{2-1} \right) \left(\frac{1}{2} \sum_{\ell=1}^2 M_{\ell, h_\ell(\text{“女”})} - \frac{4}{2} \right) = 2$$

このように頻度を推定する。

3.1 1-out-of-2 Oblivious Transfer

1-out-of-2 Oblivious Transfer[4] は、受信者は送信者から送られた2つの情報のうち片方しか知ることができず、送信者は受信者に送信した2つの情報のうちの情報を得られたか知ることができないことを保証する2パーティの暗号プロトコルである。

3.2 ポイズニング攻撃

ポイズニング攻撃は、悪意のあるユーザが意図的なデータをサーバに送信して推定結果を操作する不正行為である。例えば、オンラインショッピングサービスの場合を考える。オンラインショッピングサービスの運営者はどの商品がよく売れているかを参考にして仕入れる商品を選択する。そのため、商品を製造しているメーカは自社商品がよく売れていると偽装させることによって利益の向上を狙う動機がある。不正なデータを送信することによって自社製品の売れ行きを操作する。

各不正ユーザは、局所差分プライバシー方式を改ざんすることができ、任意のデータをサーバに送信することができる想定する。

攻撃者は、システム上で n' 人の不正ユーザを操作することができる。不正ユーザの数を n' とする。攻撃者が、操作する r 個のアイテムをターゲットアイテムとし、その集合を $T = \{t_1, t_2, \dots, t_r\}$ とする。サーバは n 人の真のユーザと m 人の不正ユーザの出力から統計量を推定する。

n 人の真のユーザのアイテム t に対する頻度の推定値を \hat{f}_t , 不正ユーザを含めた $n + n'$ 人のアイテム t に対する頻度の推定値を \tilde{f}_t とする。ポイズニング攻撃による頻度の推定値の変化量を $\Delta\tilde{f}_t = \tilde{f}_t - \hat{f}_t$ とする。

[1] によるポイズニング攻撃には Random Perturb Attack(RPA), Random Item Attack(RIA), Maximal Gain Attack(MGA) がある。RPA は各不正ユーザが摂動されたデータ集合からランダムに一つ選びサーバに送信する攻撃である。RIA は各不正ユーザがターゲットアイテムの中からランダムに一つ選択し、そのデータを定められた方法で正しく摂動し、サーバに送信する。MGA は摂動されたデータを不正ユーザの意図したデータに置換してサーバに送信する。

4 提案手法

4.1 CMS, HCMS に対するポイズニング攻撃の評価

4.1.1 Random Perturb Attack

RPA は、出力をランダムに選択する攻撃である。CMS では、 2^m の数だけ選択肢があり、各選択肢を $\frac{1}{2^m}$ の確率一つ選びでサーバに送信する。例えば、 $m = 2$ の場合、不正ユーザは $(1, 1), (-1, 1), (1, -1), (-1, -1)$ から一様ランダムに取得したデータをサーバに送信する。HCMS の場合では、1 または -1 を $\frac{1}{2}$ の確率でサーバに送信する。

4.1.2 Random Item Attack

RIA は、摂動対象を操作する攻撃である。不正ユーザはランダムに $t \in T$ を選択する。また、不正ユーザは t を CMS または HCMS を適用し、得られた出力をサーバに送信する。

4.1.3 Maximal Gain Attack

MGA は、出力を操作する攻撃である。アイテム t に対するポイズニング後の推定値を \tilde{f}_t , ポイズニング前の推定値を \hat{f}_t とする。このとき、不正ユーザは Frequency Gain(FG) を

$$FG = \sum_{t \in T} E[\tilde{f}_t - \hat{f}_t]$$

と定める。FG が最大になるように出力を生成する。

CMS に対する FG の期待値は次のように計算される。

$$\begin{aligned} FG &= \sum_{t \in T} E[\hat{f}_t - \tilde{f}_t] \\ &= \sum_{t \in T} E \left[\left(\frac{m}{m-1} \right) \left(\frac{1}{k} \sum_{\ell=1}^k \tilde{M}_{\ell, h_t(d)} - \frac{n}{m} \right) \right. \\ &\quad \left. - \left(\frac{m}{m-1} \right) \left(\frac{1}{k} \sum_{\ell=1}^k M_{\ell, h_t(d)} - \frac{n}{m} \right) \right] \\ &= \frac{1}{k} \left(\frac{m}{m-1} \right) \sum_{t \in T} E \left[\sum_{\ell=1}^k (\tilde{M}_{\ell, h_t(d)} - M_{\ell, h_t(d)}) \right] \end{aligned} \quad (1)$$

このとき、 $d \in D$ として、 i 番目のユーザによる、Sketch Matrix の ℓ 行 $h_\ell(d)$ 列のエントリ $M_{\ell, h_\ell(d)}$ を $Y_\ell^{(i)}(d)$ とおく。また、 i 番目の不正ユーザによる、Sketch Matrix の ℓ 行 $h_\ell(d)$ 列のエントリ $M_{\ell, h_\ell(d)}$ を $X_\ell^{(i)}(d)$ とおくと、

$$\begin{aligned} \tilde{M}_{\ell, h_\ell(d)} &= \sum_{i=1}^n Y_\ell^{(i)}(d) + \sum_{i=1}^{n'} X_\ell^{(i)}(d) \\ M_{\ell, h_\ell(d)} &= \sum_{i=1}^n Y_\ell^{(i)}(d) \end{aligned} \quad (2)$$

より、(2) を用いて (1) を変形すると、

$$FG = \frac{1}{k} \left(\frac{m}{m-1} \right) \sum_{t \in T} E \left[\sum_{\ell=1}^k \sum_{i=1}^{n'} X_\ell^{(i)}(d) \right] \quad (3)$$

これより、 $X_\ell^{(i)}(d)$ を最大すればよい。つまり、攻撃者は H から任意のハッシュ関数 h_j を選択し、 $T = \{t_1, t_2, \dots, t_r\}$ に対応するハッシュ関数の出力 $h_j(t_1), h_j(t_2), \dots, h_j(t_r)$ を調べ、得られた出力の位置を 1 に設定し、その他は -1 を設定して、 (v, j) を送信する。例えば、 $D = \{\text{“男”}, \text{“女”}\}$, $m = 2, k = 2, H = \{h_j : D \rightarrow [2], j \in [2]\}$ のとき、不正ユーザが男の集計結果を増加させるシナリオを考える。攻撃者は任意のハッシュ関数 h_j を選択し、 $h_j(\text{“男”}) = 0$ であれば、 $(1, -1)$ とハッシュ関数番号 j を送信する。

一方、HCMS に対する MGA は複数の戦略が考えられる。2.4 節で述べたように、あるアイテムの入力が他のアイテムの推定結果に影響を及ぼすことに起因している。ここでは最も簡単な戦略を述べる。

$d \in D$ として、 i 番目のユーザによる、Sketch Matrix の ℓ 行 $h_\ell(d)$ 列のエントリ $M_{\ell, h_\ell(d)}$ を $Z_j^{(i)}(d)$ とおく。また、 i 番目の不正ユーザによる、Sketch Matrix の ℓ 行 $h_\ell(d)$ 列のエントリ $M_{\ell, h_\ell(d)}$ を $\tilde{Z}_j^{(i)}(d)$ とおくと、(1) は次のように変形できる。

$$FG = \frac{1}{k} \left(\frac{m}{m-1} \right) \sum_{t \in T} E \left[\sum_{\ell=1}^k \sum_{i=1}^{n'} \tilde{Z}_\ell^{(i)}(d) \right] \quad (4)$$

この時,

$$\tilde{Z}_l^{(i)}(d) = k \cdot c_\varepsilon \cdot 1 \quad (5)$$

とすれば, FG が最大となる. つまり, 攻撃者は $w = 1, l = 0, j$ は任意 を送信すればよい.

4.2 OT-CMS, OT-HCMS

MGA は, 正規の摂動化プロセスの後, 不正ユーザが送信するデータを意図的なものに変更することによって実現される. つまり, 送信されたデータは摂動化が行われていないまま送信される.

そこで, Oblivious Transfer を用いて, ユーザに摂動化を強制させる OT-CMS と OT-HCMS を提案する. 本来, CMS は摂動化した m 次元ベクトル \tilde{v} とハッシュ関数の番号 j を送信するが, OT-CMS においては, ユーザは m 次元ベクトル v の摂動化を単独に行わず, 1-out-of-2 OT を用いてベクトルの要素をサーバの協力により選択する. 従って, 不正な摂動化が防止される. ハッシュ関数の番号は従来通り送信する.

例えば, ユーザが $v = (1, -1)$ をサーバに送信する場合を考える. ユーザが 1 番目の要素 1 を送信するとき, ユーザは 1, -1 のどちらも OT の送信候補とする. サーバはその内, 真のデータを確率 p , 偽データを確率 $q = 1 - p$ で取得する. このとき,

$$p = \frac{e^{\frac{\varepsilon}{2}}}{1 + e^{\frac{\varepsilon}{2}}}, \quad q = \frac{1}{1 + e^{\frac{\varepsilon}{2}}}$$

この試行を m 回繰り返す. この際, ユーザとサーバは 1-out-of-1/p OT を用いて通信をしているため, サーバはどちらか片方の情報のみ得ることができ, ユーザはサーバがどちらの情報を取得したか知ることができない.

OT-HCMS では,

$$p = \frac{e^\varepsilon}{1 + e^\varepsilon}, \quad q = \frac{1}{1 + e^\varepsilon}$$

に変更し, 試行は 1 ビット分行えば十分である.

5 CMS と HCMS の送信量

CMS は m に比例して, 送信量が大きくなってしまふ. その課題を解決するために考案された局所差分プライバシー方式が Hadamard Count Mean Sketch である. 通信速度 B を $1 \times 10^9 \text{bps}$ であると仮定し, 図 3 に CMS と HCMS の m に対する送信時間の推移を示す. ここで, 送信時間は $\frac{m}{B}$ と表せる.

表 2 評価に用いるデフォルトパラメータ

パラメータ	
プライバシー予算 ε	1.0
ベクトル長 m	2^7
ハッシュ関数の数 k	2^{10}
不正ユーザの割合 β	0.01
ターゲットアイテムの数 r	1

6 実験

6.1 実験目的

アダマール行列を用いて削減する送信量と安全性を評価する. 有用性と安全性を用いて評価する. また, 提案手法の効果を評価する.

6.2 データセット

オンラインショッピングサービスの購入履歴のデータセット [7] を使用する. 図 2 にアイテムの購入頻度の分布を示す. 例えば, 商品 A2 は 3013 件ある. アイテム数は $|D| = 43$ である.

6.3 評価方法

実験で使用するデフォルトパラメータを表 2 に整理する. 真の分布と推定分布の平均二乗誤差 MSE を用いて有用性を評価する. FG を用いて安全性を評価する. 試行を 50 回繰り返しその平均値を評価値とする. ただし, OT-CMS, OT-HCMS の安全性評価では, 試行は 10 回行い, アイテムは “A18”, “A34” について評価する. MGA ポイズニング攻撃を評価する.

6.4 実験結果

有用性: 図 4, 図 5, 図 6, 表 3 にそれぞれプライバシー予算 ε , CMS の符号化ベクトル長 m , ハッシュ関数の数 k に対する CMS と HCMS の MSE を示す.

プライバシー予算については CMS, HCMS どちらも単調に MSE が減少している. CMS は HCMS に比べ $\varepsilon = 0.1$ の場合を除き, MSE が小さくなっている. 特に $\varepsilon = 1.0$ のとき, 89.7% だけ CMS の方が精度が良い. ベクトル長については全ての場合で CMS の方が HCMS よりも MSE が小さい. 特に $m = 1024$ のとき, 52.8% だけ CMS の方が精度が良い. ハッシュ関数の数についても CMS, HCMS どちらも単調に MSE が減少している.

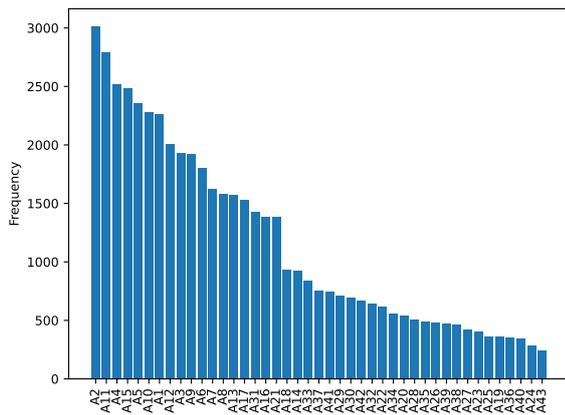


図2 データセットのアイテムの購入頻度の分布

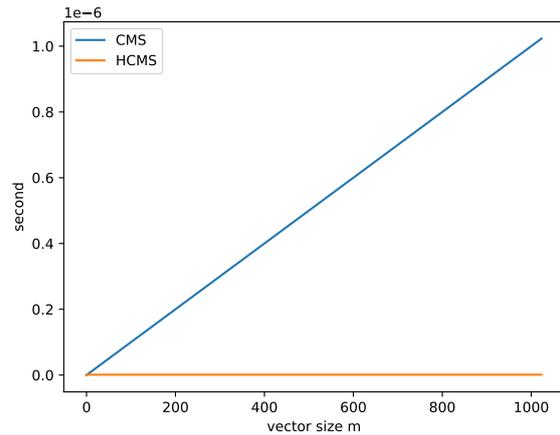


図3 CMS と HCMS の送信時間

安全性: 図7, 図8, 図9, 表4, 表5, 表6にそれぞれプライバシー予算 ϵ , 不正ユーザの割合 β , ターゲットアイテムの数 r に対する CMS と HCMS の FG を示す。

MGA については, ϵ, β, r の全ての条件で HCMS の方が CMS より FG が平均で 16.0% 小さい。RPA についても同様に, ϵ, β, r の全ての条件で HCMS の方が CMS より FG が小さい。特に $\beta = 0.1$ の時, 最大で 40.35 小さい。一方で, RIA では, ϵ, β, r における FG は CMS と HCMS の差は小さい。

OT-CMS と OT-HCMS の安全性: 図10 と図11, 表7に不正ユーザの割合に対する OT-CMS と OT-HCMS の FG を示す。

OT-CMS, OT-HCMS のどちらの場合でも, CMS と HCMS よりも FG が小さくなった。特に, $\beta = 0.10$ のとき, OT-CMS は CMS より 267.83 安全であり, OT-HCMS は HCMS にくらべ 211.83 安全であった。

6.5 考察

3 より, HCMS は CMS に比べて誤差が大きく, 有用性が低い。これは, ユーザが行う HCMS の処理に m 次元ベクトルから 1 ビットをランダムにサンプリングすることが原因だと考えられる。HCMS はサンプリングが一様ランダムになるとき CMS と同じ性能となることが Apple により示されている [3]。そのため, 実験的に行うとサンプリングにある程度の偏りが生じ, HCMS は CMS に比べ有用性が低下する。

MGA に対して, CMS に比べ HCMS の方が安全であった。その原因はノイズ除去の方法によって引き起こされている。

CMS はユーザから送信されたデータのノイズを除去

するときに, $\frac{e^{\frac{\epsilon}{2}+1}}{e^{\frac{\epsilon}{2}-1}}$ と積をとる。一方, HCMS は $\frac{e^{\epsilon}+1}{e^{\epsilon}-1}$ と積をとる [3]。 ϵ が同じ値であれば, $\frac{e^{\frac{\epsilon}{2}+1}}{e^{\frac{\epsilon}{2}-1}}$ の方が大きな値をとる。そのため, 不正ユーザから MGA を用いて攻撃された際に, 攻撃の効果がより増幅される。RIA についてはほぼ同様の安全性を示した。HCMS の方が CMS よりも安全と言える。

OT-CMS と OT-HCMS では, 1-out-of-2 OT を適用し, 摂動を強制的に行うため不正ユーザは意図的なデータを送信することができない。FG が小さくなり, 安全性が向上した。

6.6 OT-CMS と OT-HCMS の限界

局所差分プライバシー方式は, ユーザのプライバシーを保護するために, サーバに信頼しないモデルとして提案された。本研究で提案した OT-CMS と OT-HCMS はサーバがユーザのデータを摂動することを補助しているため, サーバをある程度信頼しているモデルである。そのため, 本来の局所差分プライバシー方式の考え方と異なる。私は実際の利用環境によって使用するモデルを変更するべきだと考える。

7 おわりに

本稿では, 局所差分プライバシープロトコル Count Mean Sketch(CMS), Hadamard Count Mean Sketch(HCMS) の有用性を比較し, ポイズニング攻撃に対する安全性を評価した。ポイズニング攻撃に対するロバスト性を向上させるために, Oblivious Transfer を適用する手法を提案した。実験に基づき, 有用性では, CMS の方が高く, 安全性では, HCMS の方が高いことが分かった。

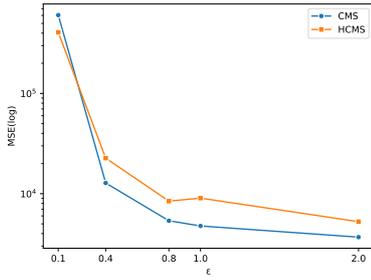


図4 プライバシー予算 ϵ

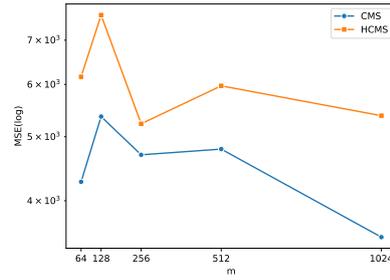


図5 ベクトル長 m

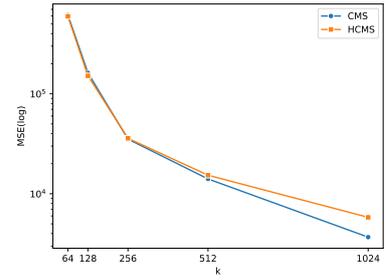


図6 ハッシュ関数の数 k

表3 各パラメータによる MSE($\times 10^3$) の変化

ϵ	CMS	HCMS
0.1	604.66	407.07
0.4	12.78	22.63
0.8	5.37	8.43
1.0	4.76	9.03
2.0	3.69	5.26

m	CMS	HCMS
64	4.27	6.16
128	5.36	7.64
256	4.70	5.23
512	4.79	5.97
1024	3.52	5.38

k	CMS	HCMS
64	613.03	596.29
128	162.13	151.32
256	35.22	35.89
512	14.07	15.34
1024	3.68	5.81

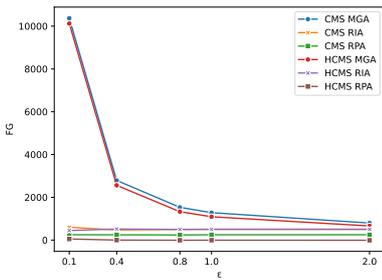


図7 プライバシー予算 ϵ についての安全性 FG

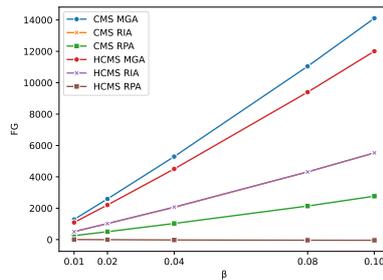


図8 不正ユーザの割合 β についての安全性 FG

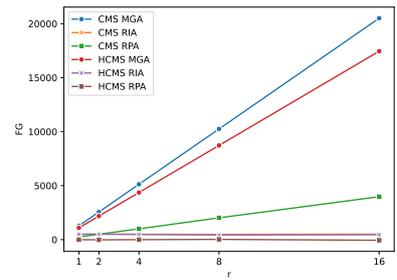


図9 ターゲットアイテムの数 r についての安全性 FG

参考文献

- [1] X. Cao, J. Jia, N. Z. Gong, "Data poisoning attacks to local differential privacy protocols", USENIX Security Symposium, pp. 947-964, 2021.
- [2] Hikaru Horigome, Hiroaki Kikuchi and Chia-Mu, Yu, "Local Differential Privacy Protocol for Key-Value Data Robust against Poisoning Attacks", Modeling Decisions for Artificial Intelligence, 2023, Volume 13890
- [3] Differential Privacy Team, Learning with Privacy at Scale, <https://machinelearning.apple.com/research/learning-with-privacy-at-scale>.
- [4] Even, Shimo Goldreich, Oded Lempel, Abraham. (1982). A Randomized Protocol for Signing Con- tracts.. Communications of the ACM. 28. 205-210. 10.1145/3812.3818.
- [5] Gadotti, Andrea Houssiau, Florimond Annamalai, Meenatchi Montjoye, Yves-Alexandre. (2023). Pool Inference Attacks on Local Differential Privacy: Quantifying the Privacy Guarantees of Apple's Count Mean Sketch in Practice, 31st USENIX Security Symposium (USENIX Security 22), 2022.
- [6] J. C. Duchi, M. I. Jordan and M. J. Wainwright, "Local Privacy and Statistical Minimax Rates," 2013 IEEE 54th Annual Symposium on Foundations of Computer Science, Berkeley, CA, USA, 2013, pp. 429-438, doi: 10.1109/FOCS.2013.53.
- [7] clickstream data for online shopping, (2019), UCI Machine Learning Repository.

表4 プライバシー予算 ε による FG($\times 10^2$) の変化

ε	CMS RPA	HCMS RPA	CMS RIA	HCMS RIA	CMS MGA	HCMS MGA
0.1	2.51	0.53	6.01	4.46	103.70	101.24
0.4	2.47	0	4.27	5.17	27.87	25.58
0.8	2.40	-0.08	4.99	4.98	15.31	13.28
1.0	2.50	-0.05	5.08	5.00	12.82	10.91
2.0	2.50	-0.08	5.03	5.04	7.96	6.60

表5 不正ユーザの割合 β による FG($\times 10^2$) の変化

β	CMS RPA	HCMS RPA	CMS RIA	HCMS RIA	CMS MGA	HCMS MGA
0.01	2.51	0.04	5.18	5.01	12.82	10.91
0.02	5.03	-0.05	10.19	10.21	25.92	22.06
0.04	10.28	-0.24	20.52	20.73	52.91	45.03
0.08	21.39	-0.40	43.15	43.18	110.44	93.99
0.10	27.67	-0.43	55.19	55.30	141.11	120.09

表6 ターゲットアイテムの数 r による FG($\times 10^2$) の変化

r	CMS RPA	HCMS RPA	CMS RIA	HCMS RIA	CMS MGA	HCMS MGA
1	2.50	-0.01	4.93	5.04	12.82	10.91
2	4.99	-0.17	5.09	5.07	25.64	21.82
4	10.02	-0.02	4.95	4.84	51.28	43.64
8	20.20	0.25	4.89	4.29	102.55	87.27
16	39.77	-0.58	5.10	4.56	205.11	174.54

表7 OT-CMS と OT-HCMS の不正ユーザの割合 β による FG の変化

β	CMS	OT-CMS	HCMS	OT-HCMS
0.01	38.30	18.56	32.6	17.15
0.02	76.61	30.59	65.19	38.75
0.04	158.33	59.70	134.73	65.34
0.08	331.97	129.44	282.51	135.75
0.10	423.90	156.07	360.74	148.91

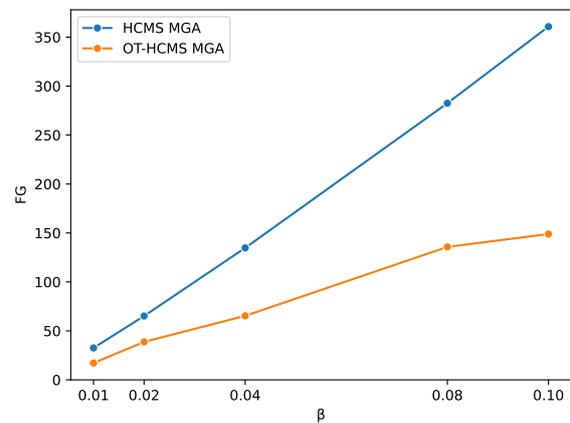
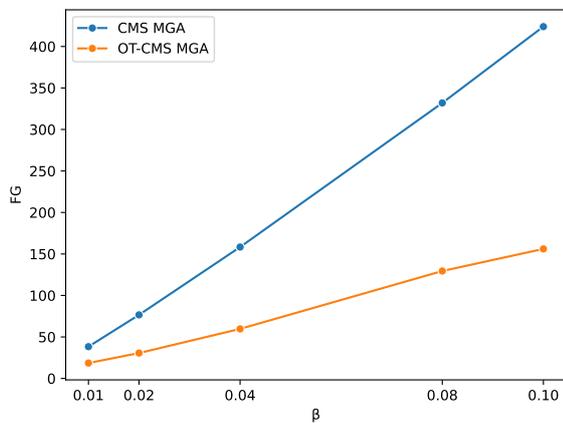


図10 不正ユーザの割合 β についての OT-CMS の安全性 FG

図11 不正ユーザの割合 β についての OT-HCMS の安全性 FG