

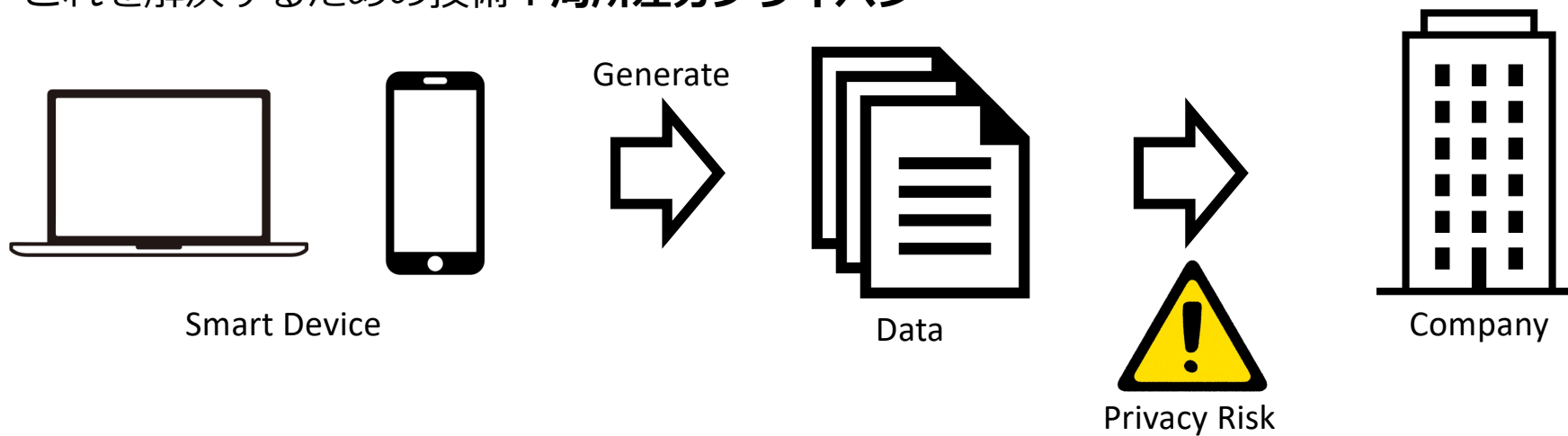
# Key-ValueデータのLDPプロトコルPCKVの 推定値操作攻撃の提案と評価

谷口輝海, 菊池浩明

明治大学

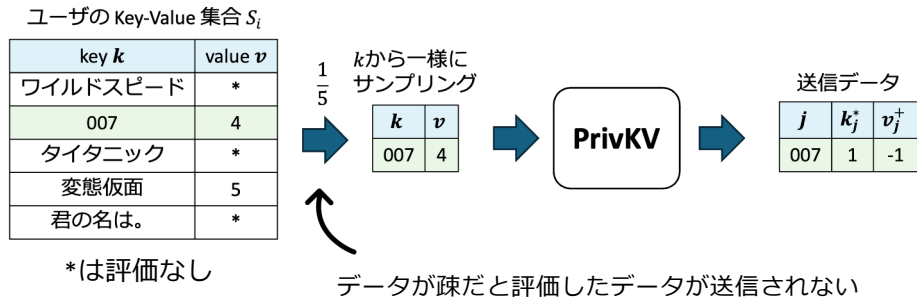
# 背景

- 企業は、ユーザデバイス上で生じるデータをサービス向上のために活用したい。
- データをそのまま収集するとプライバシーが侵害される恐れ。
- ユーザはプライバシーを守りたい。
- これを解決するための技術：**局所差分プライバシー**

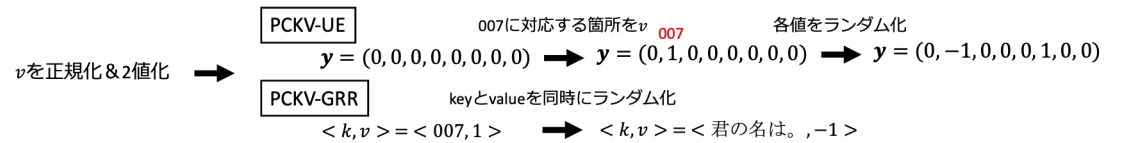
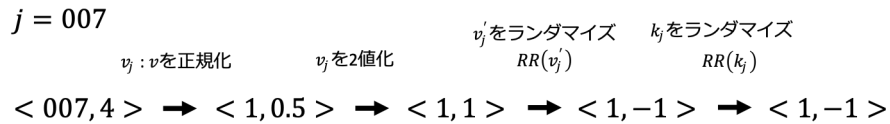
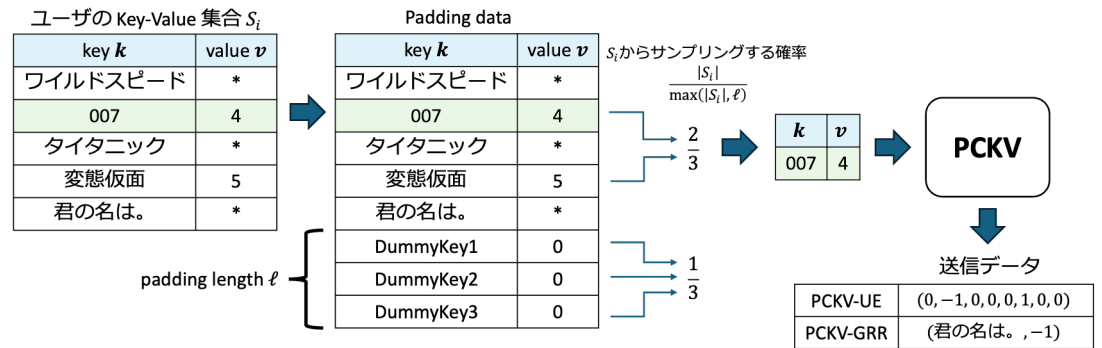


# Key-ValueデータのLDP方式

## PrivKV [Ye, et al. 2019]

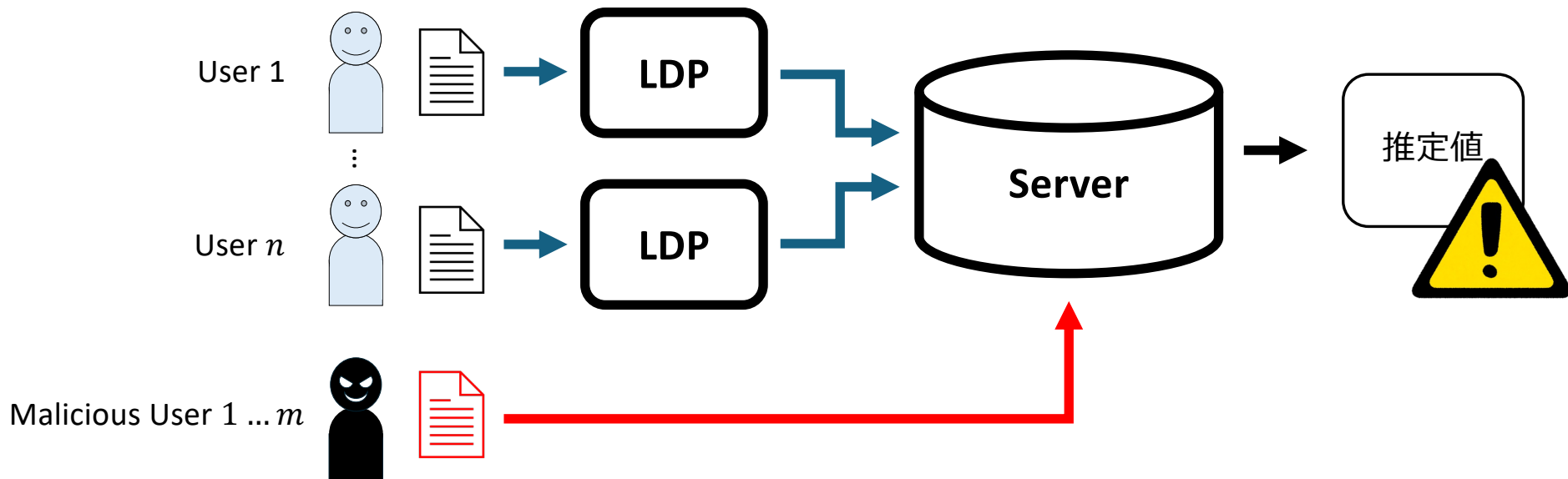


## PCKV [Gu, et al. 2020]



# ポイズニング攻撃

- 送信データを細工して推定値を操作する



代表的な攻撃

特定のカテゴリの頻度を増加させる。[Cao, et al. 2021]

特定のkeyの頻度と平均を増加させる。[Wu, et al. 2022]

平均値を攻撃者が意図した値に近づける。[Li, et al. 2023]

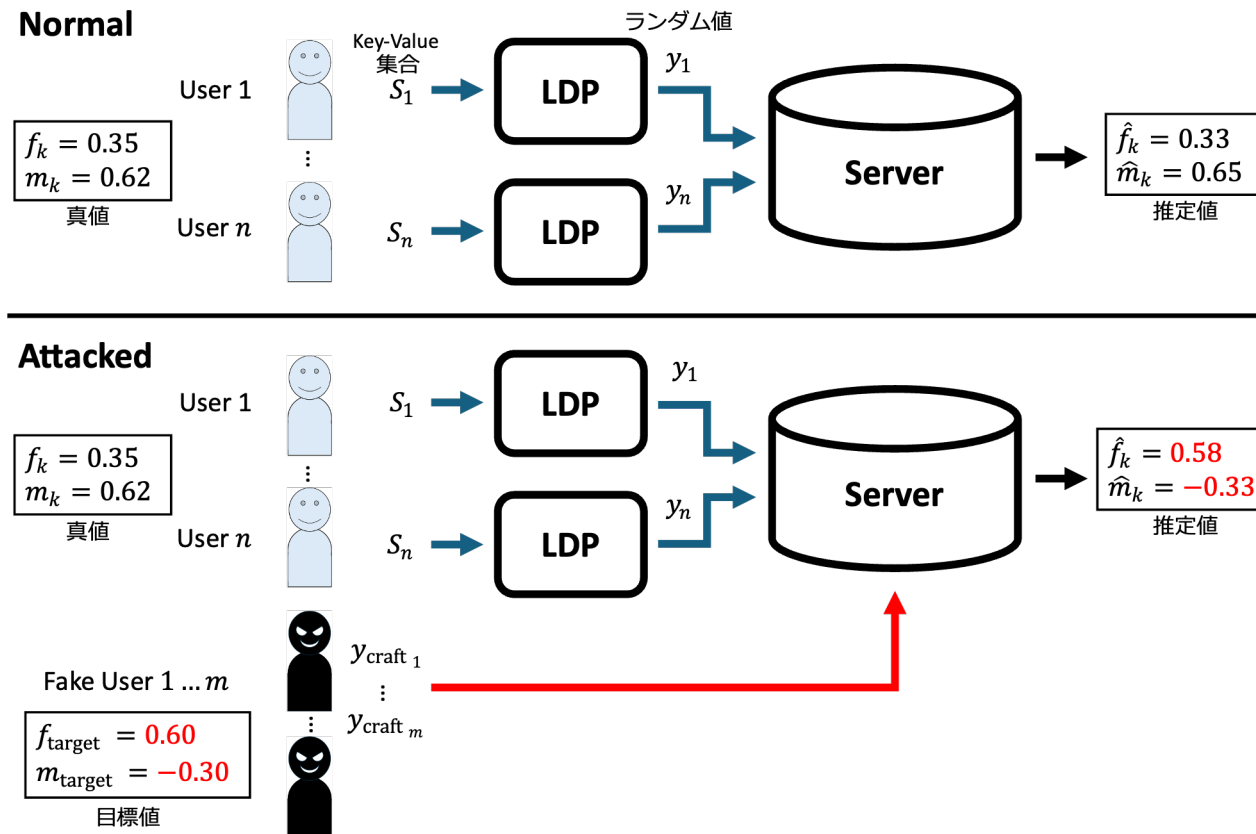
# 本研究の貢献

	攻撃目的\対象統計量	分散	平均値	頻度
利得最大化攻撃	真値から増加させる		Wu et al 2022	Cao et al 2021
推定値操作攻撃	攻撃者の目標値にする	Li et al 2023	本研究	

# 研究目的

- Key-ValueデータのためのLDP方式に対し、keyの推定頻度と推定平均を攻撃者の意図した値に近づけるポイズニング攻撃を提案する。
- 提案攻撃をPrivKV, PCKV-UE, PCKV-GRRに適用し、攻撃精度を比較評価する。

# 研究概要 推定値操作攻撃



# 問題設定

- 正規ユーザ $n$ 人の頻度 $f_k$ と平均 $\mu_k$ とする。ここで、標的となる頻度 $f_{k,t}$ と平均 $\mu_{k,t}$ にするために、 $m$ 人の偽ユーザを挿入する。このとき、 $m$ 人の偽ユーザの送信データを決定したい。
- 最終的に達成したい式

$$\begin{cases} E[\hat{f}_k] = f_{k,t} \\ E[\hat{\mu}_k] = \mu_{k,t} \end{cases}$$



# 提案方式

## PCKV-UEに対する攻撃の適用例

PCKV-UEにおける頻度推定

$$E[\hat{f}_k] = E\left[\frac{n_1^k + n_{-1}^k - b}{a - b} \ell\right]$$

PCKV-UEにおける平均値推定

$$E[\hat{\mu}_k] = E\left[\frac{\ell(n_1^k - n_{-1}^k)}{a(2p - 1)n\hat{f}_k}\right]$$

$m$ 人の偽ユーザを挿入( $m_1^k, m_{-1}^k, m_0^k$ )

$$E[\hat{f}_k] = E\left[\frac{n_1^k + n_{-1}^k + m_1^k + m_{-1}^k - b}{n + m} \ell\right] = f_{k,t}$$

攻撃者の目標値

$$E[\hat{\mu}_k] = E\left[\frac{\ell(n_1^k - n_{-1}^k + m_1^k - m_{-1}^k)}{a(2p - 1)(n + m)\hat{f}_k}\right] = \mu_{k,t}$$

攻撃者が推定する  
真の頻度と平均

$$\begin{cases} m_1^k + m_{-1}^k = \frac{a - b}{\ell} \left( (m + n)f_{k,t} - f_k \right) + mb \\ m_1^k - m_{-1}^k = \frac{a(2p - 1)}{\ell} \left( (m + n)f_{k,t}\mu_{k,t} - nf_k\mu_k \right) \\ m_1^k + m_{-1}^k + m_0^k = m \\ 0 \leq m_1^k, m_{-1}^k, m_0^k \leq m \end{cases}$$

# Research Question

RQ 1. 与えられた目標値を達成するための不正者数 $m$ はどれくらいか？

RQ 2. 3つのLDP方式で、推定値操作攻撃に対して最も安全な方式はどれか？

# 実験・評価指標

- 推定精度評価：推定値 $\hat{f}_k, \hat{m}_k$ と真値 $f_k, m_k$ の平均二乗誤差を計算

$$MSE_f = \frac{1}{|K|} \sum_{k \in K} (\hat{f}_k - f_k)^2, MSE_\mu = \frac{1}{|K|} \sum_{k \in K} (\hat{\mu}_k - \mu_k)^2$$

- 攻撃精度評価：推定値 $\hat{f}_k, \hat{m}_k$ と目標値 $f_{k,t}, m_{k,t}$ の平均二乗誤差を計算

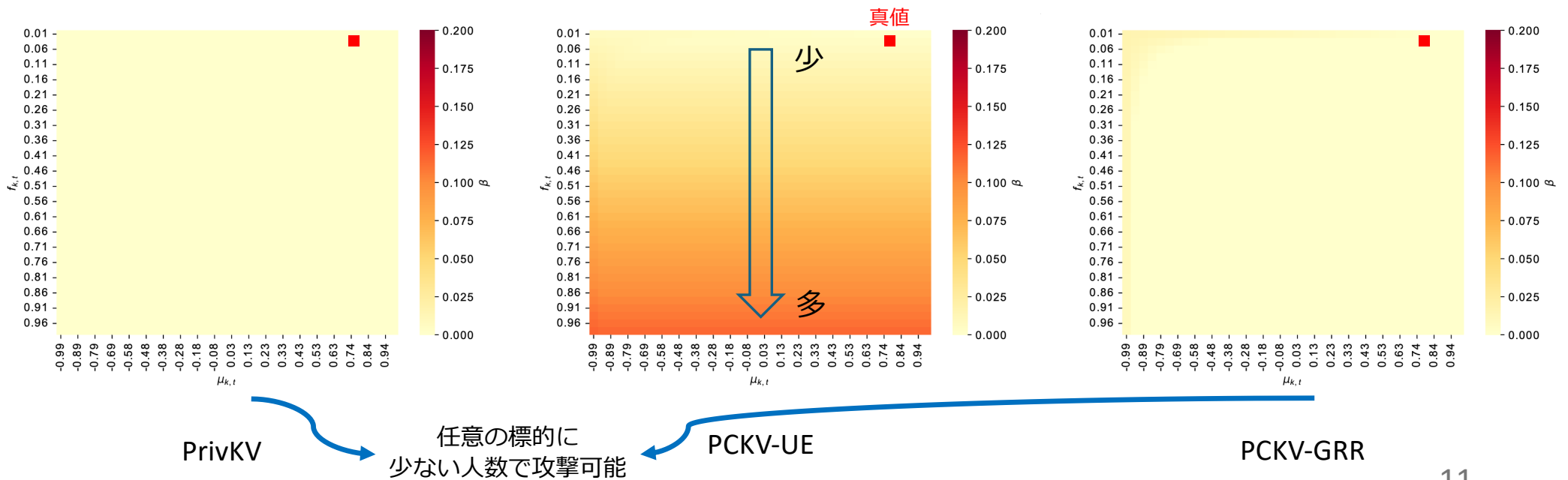
$$MSE_{f_k} = \frac{1}{N} \sum_{i \in [N]} (\hat{f}_k^{(i)} - f_{k,t})^2, MSE_{\mu_k} = \frac{1}{N} \sum_{i \in [N]} (\hat{\mu}_k^{(i)} - \mu_{k,t})^2$$

- 使用データセット：Clothing Dataset

User 数	key 数	record 数
105,508	5850	192,198

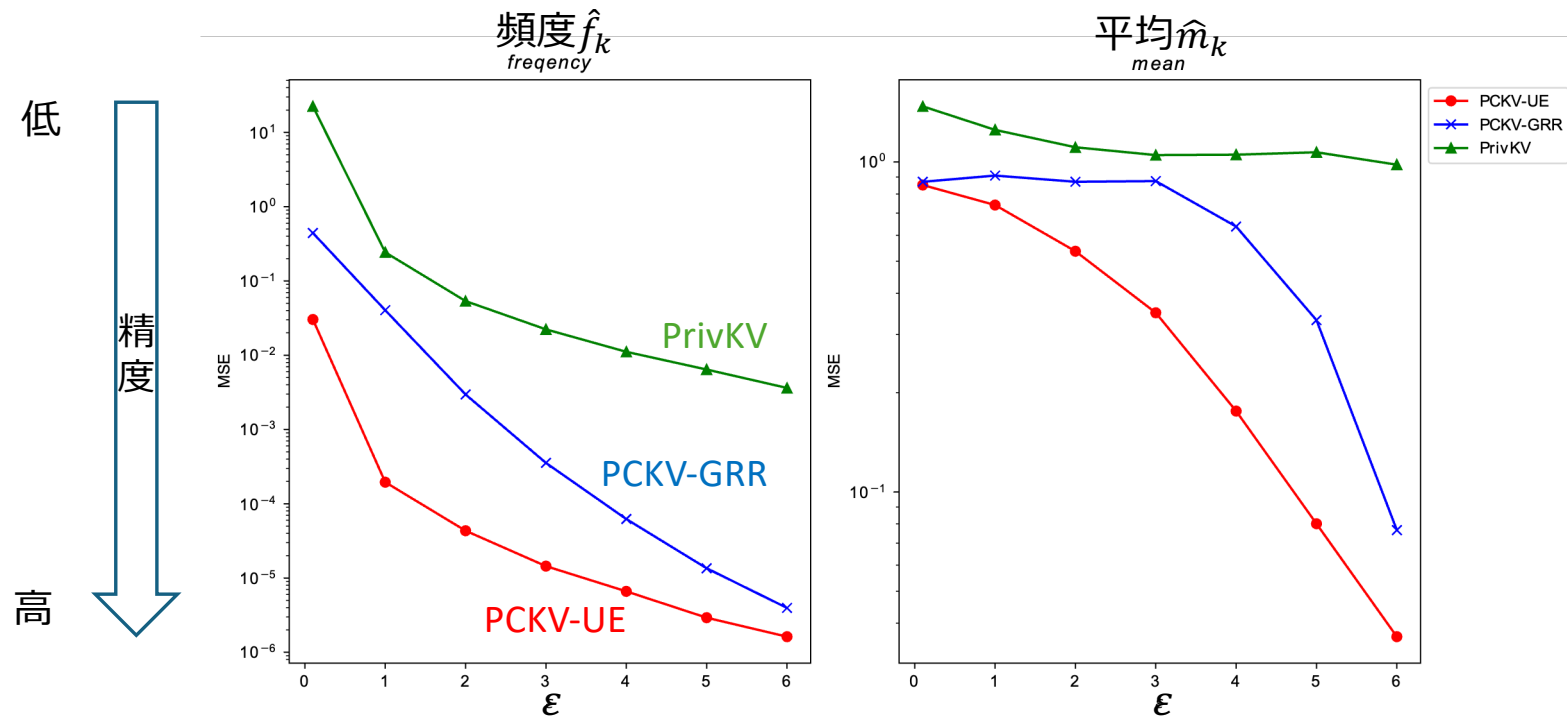
# 方程式が解を持つための $m$ の必要条件

- 前述の連立式を満たす必要十分な $m$ の範囲を、目標値 ( $f_{k,t}, m_{k,t}$ ) を変更して計算
- $m$ の下限を用いて $\beta$  (偽ユーザの割合) を計算したときのヒートマップ



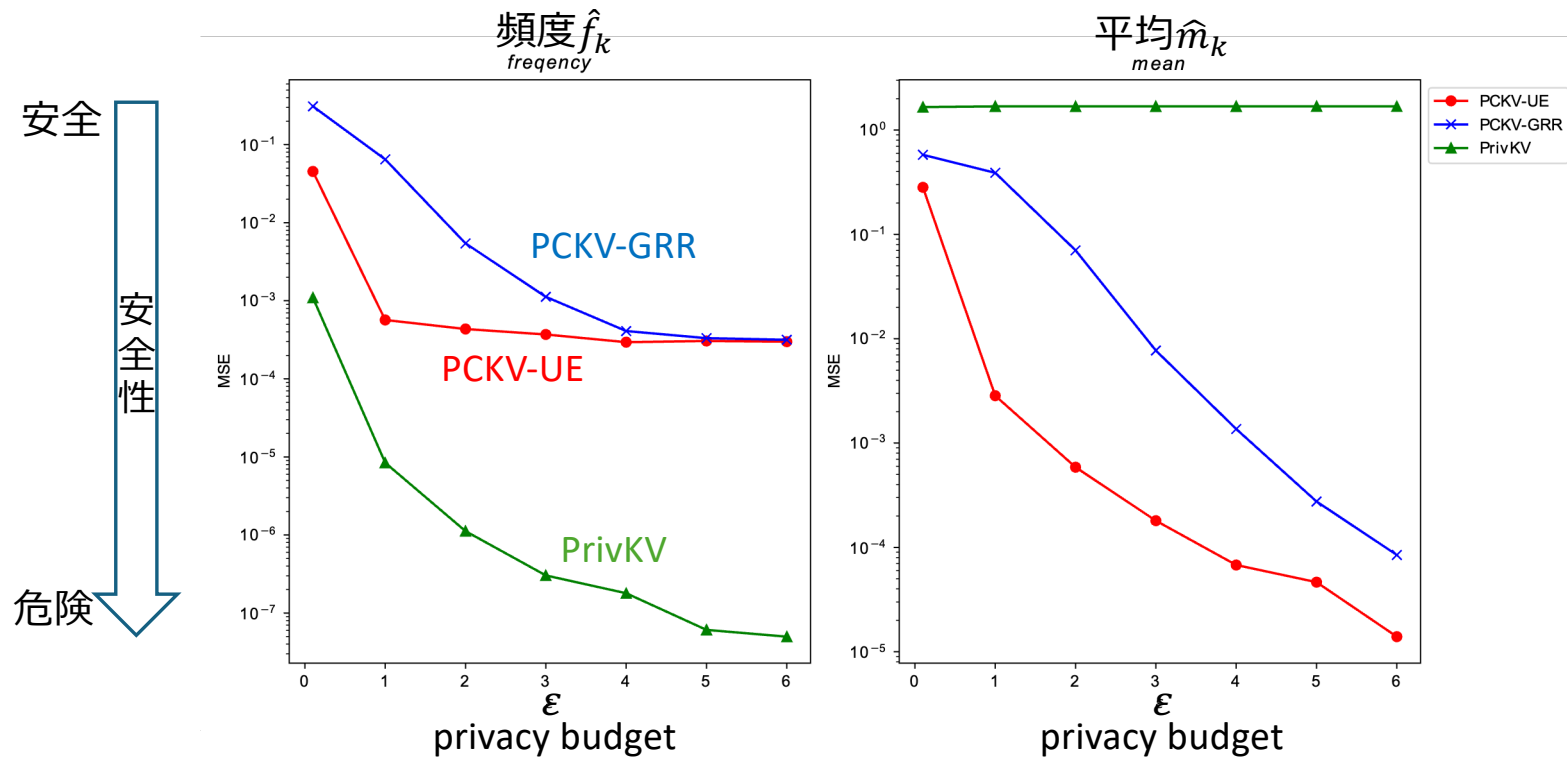
# 実験結果① 推定精度

- $\epsilon$ による3方式の推定精度推移



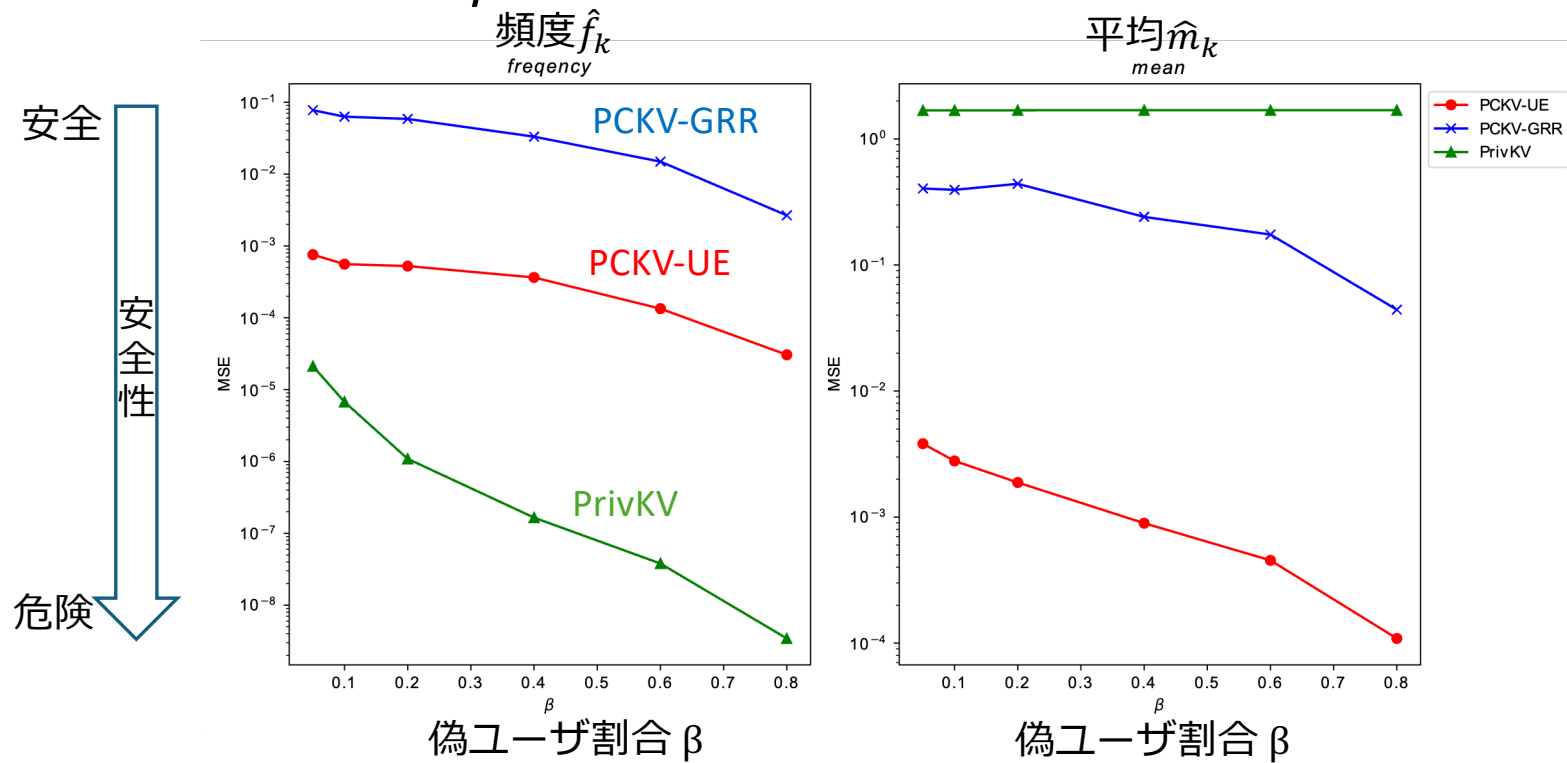
# 実験結果② 攻撃精度 ( $\epsilon$ )

- $\epsilon$  による3方式の攻撃精度推移



# 実験結果③ 攻撃精度 ( $\beta$ )

- 偽ユーザの割合  $\beta$  による攻撃精度推移



# 考察

- PrivKVが平均値への攻撃に対して最も安全な理由
  - $\varepsilon$ を大きくしても（攻撃をしなかった場合の）平均値の推定精度は大きく改善せず、ランダム化のノイズにより推定値の操作が困難になるため
- 頻度の操作結果
  - keyの摂動方式における値域の違いが影響している
  - PrivKV : 2値のRR
  - PCKV-GRR : K値のRR
  - PCKV-UE : K値のUE



# まとめ

- Key-Valueデータ収集のためのLDP方式PrivKV, PCKV-UE, PCKV-GRRに対して、keyの頻度とvalueの平均を攻撃者の目標値に操作する推定値操作攻撃を提案した
- RQ 1. 与えられた目標値を達成するための不正者数 $m$ はどれくらいか？
  - PrivKV, PCKV-GRRでは少数(全体の**1%**未満)で攻撃可、PCKV-UEでは目標値によって最大**10%**程度必要
- RQ 2. 3つのLDP方式で、推定値操作攻撃に対して最も安全な方式はどれか？
  - 総合的に見て**PCKV-GRR**が最も安全。
  - 推定精度ではPCKV-UEに劣り、平均値への攻撃ではPrivKVが最も安全であるが、PrivKVは平均値の推定精度が低い。
- 今後
  - 攻撃誤差の理論的分析
  - 防御手法についての検討