

DICOMO 2022

key-valueデータにおける局 所差分プライバシーアルゴリ ズムPrivKVの改良

堀込光, 菊池浩明 (明治大学)

Chia-Mu Yu (National Yang Ming Chiao Tung University)

背景

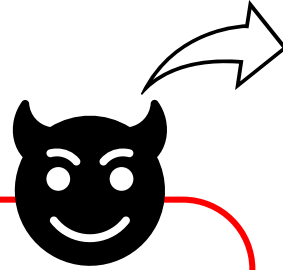
君の名は。
★★★★★★

A

もののけ姫
★★★★

B

“B”は「もののけ姫」と「タイタニック」を見ており、高く評価している



千と千尋の神隠し
★★

C

アナと雪の女王
★★★★

D

タイタニック
★★★★

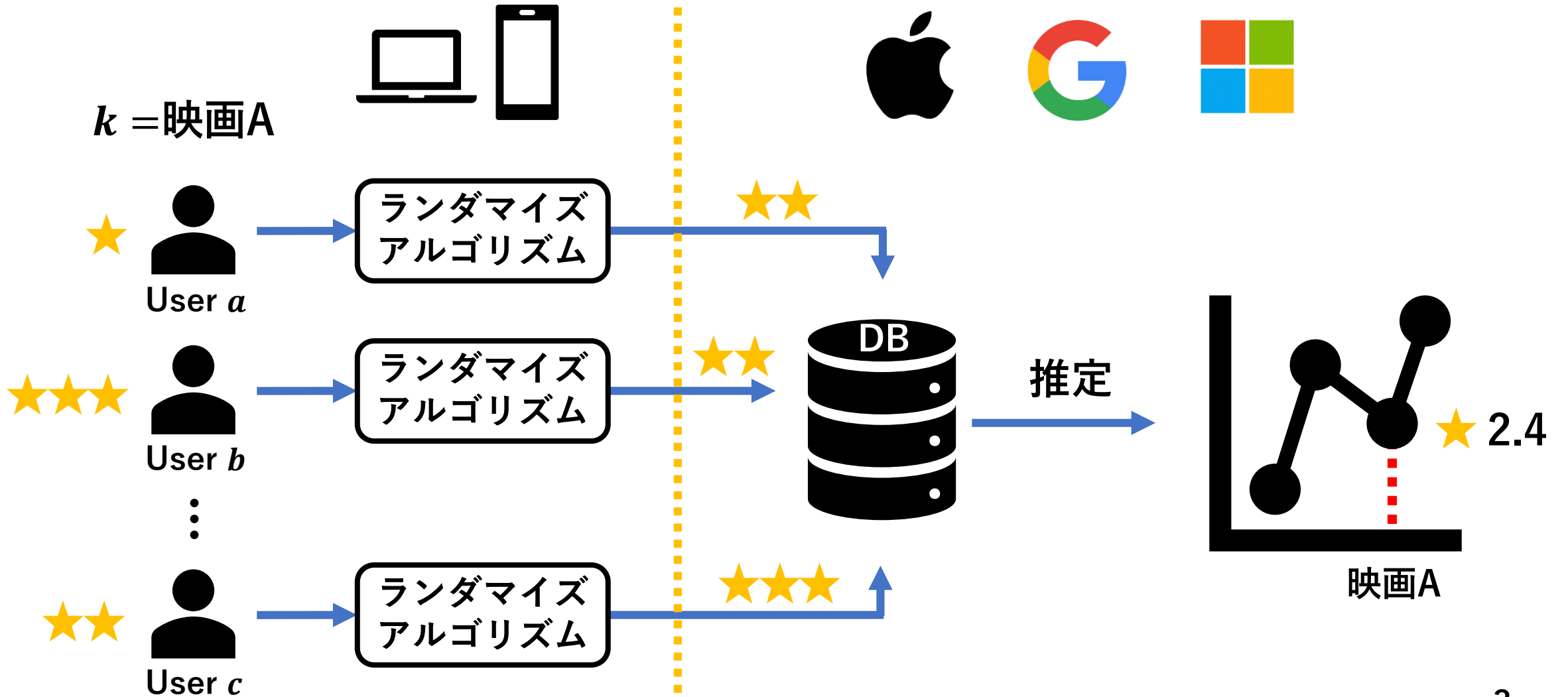
B

NETFLIX

- 1位
千と千尋の神隠し
★4.4
- 2位
タイタニック
★4.1
- 3位
アナと雪の女王
★3.8

⋮

局所差分プライバシー(LDP)



先行研究：PrivKV (1. 摂動)

key-valueデータ

k	v
千と千尋の神隠し	2
タイタニック	*
アナと雪の女王	*
君の名は。	4
もののけ姫	3

サンプリング $\langle k, v \rangle$
 $\langle \text{君の名は。}, 4 \rangle$

PrivKV

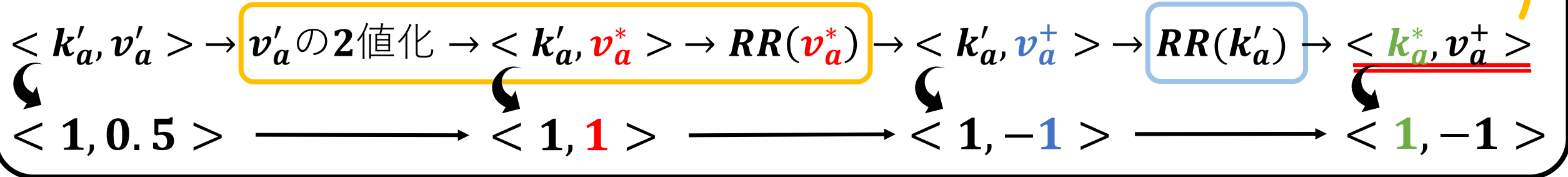
出力

a	k^* (フラグ)	v^+
君の名は。	1	-1

$a = \text{“君の名は。”}$

VPP (valueのランダムイズ)

RR(keyのランダムイズ)



先行研究：PrivKV (2. 推定)

- MLE (Maximum Likelihood Estimation)

集計(a = 君の名は。)

平均値推定

	k_a^*	v_a^+
u_1	1	-1
u_3	0	0
u_{20}	1	1
\vdots		
u_{51}	1	1
u_{77}	0	0

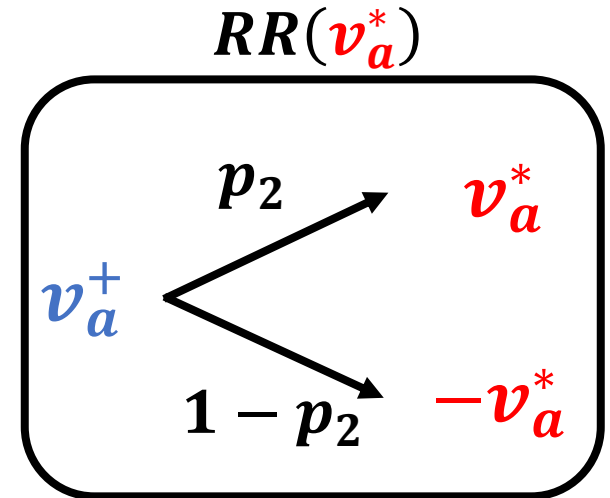
$$n'_1 = \text{count}(v_a^+ = 1)$$
$$n'_2 = \text{count}(v_a^+ = -1)$$
$$N = n'_1 + n'_2$$

$$L(\hat{n}_1) = \frac{N(p_2 - 1) + n'_1}{2p_2 - 1}$$

$$L(\hat{n}_2) = N - \hat{n}_1$$

key = a の推定平均値 \hat{m}_a は,

$$\hat{m}_a = \frac{n'_1 - n'_2}{n'_1 + n'_2}$$



$$E(n'_1) = n_1 p_2 + n_2 (1 - p_2)$$

MLEの問題点

k	v
君の名は。	*

$$\langle k'_a, v'_a \rangle = \langle 0, \text{random}[-1, 1] \rangle$$

VPP (valueのランダムイズ)

$$\langle k'_a, v'_a \rangle \rightarrow v'_a \text{の2値化} \rightarrow \langle k'_a, v_a^* \rangle \rightarrow RR(v_a^*) \rightarrow \langle k'_a, v_a^+ \rangle$$

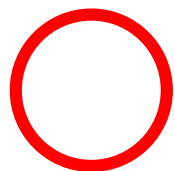
$$\langle 0, 0.4 \rangle \rightarrow v'_a \text{の2値化} \rightarrow \langle 0, -1 \rangle \rightarrow RR(-1) \rightarrow \langle 0, -1 \rangle$$

評価者の少ないkeyの遷移

$$\langle k, v \rangle \rightarrow \langle k_a^*, v_a^+ \rangle$$

$\langle \text{君の名は。}, * \rangle \rightarrow \langle 1, 1 \rangle \text{ or } \langle 1, -1 \rangle$

多



\hat{m}



\hat{m}

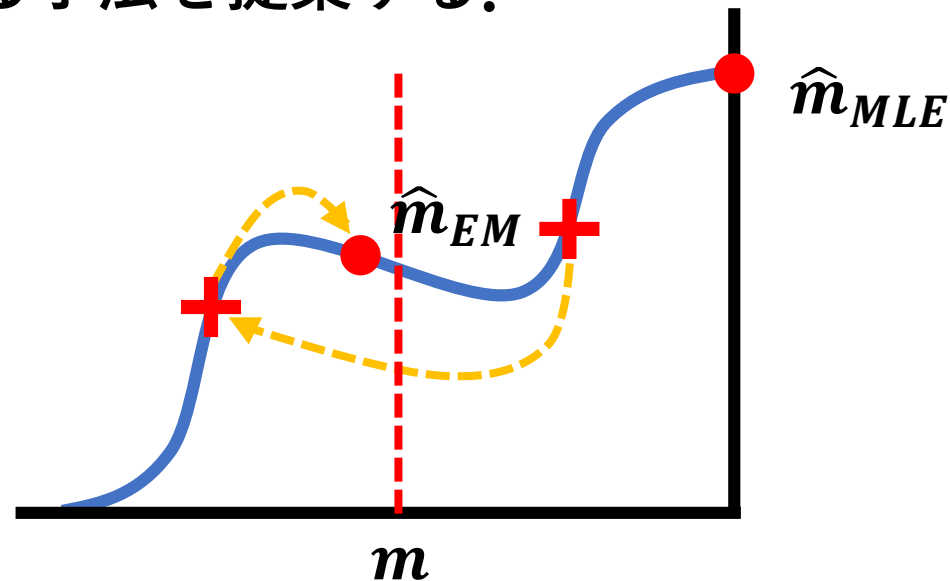
$v_a^+ = v_a^*$ であるか
 $v_a^+ = -v_a^*$ であるか
 の推定を考えており、
 $v = *$ を考慮して
 いない。

研究概要

	摂動化	推定
PrivKV	VPP+RR	MLE
本提案		EM

• 解決手法

PrivKVにEM(Expectation Maximization)アルゴリズム[宮川雅巳,1987]を適用し、推定する手法を提案する。



提案手法

EM(Expectation Maximization)アルゴリズムの適用

擾動の流れ

$a = \text{“君の名は。”}$

VPP (valueのランダムイズ)

keyのランダムイズ



事前確率を求める対象の設定

案1 $\langle k'_a, v'_a \rangle \rightarrow \langle k_a^*, v_a^+ \rangle \quad \times$

案2 $\langle k'_a, v_a^* \rangle \rightarrow \langle k_a^*, v_a^+ \rangle \quad \circ$

$$\langle k_a^*, v_a^+ \rangle = \{ \langle 1, 1 \rangle, \langle 1, -1 \rangle, \langle 0, 0 \rangle \}$$

$$\langle k'_a, v_a^* \rangle = \{ \langle 1, 1 \rangle, \langle 1, -1 \rangle, \langle 0, 1 \rangle, \langle 0, -1 \rangle \}$$

$$\langle k_a, v_a \rangle = \langle 0, _ * \rangle$$

・ 出力 $\langle k_a^*, v_a^+ \rangle$ を用いて $\langle k'_a, v_a^* \rangle$ の事前確率を推定する。

RQ.

- RQ1. 提案手法はPrivKVよりも高精度か？
- RQ2. 安全性 ϵ によって推定誤差に影響はあるのか？
- RQ3. データ規模 n によって推定誤差に影響はあるのか？

評価実験の概要

実験目的

- 3つのRQ.を調査する

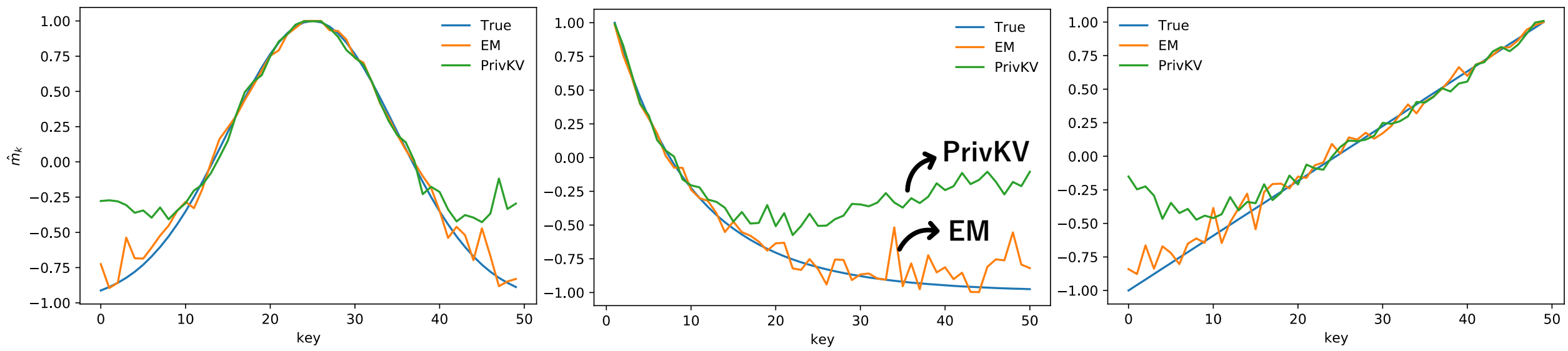
方法

- データ
 - keyとvalueがガウス分布，べき分布，線形分布に従う合成データ
- 評価実験
 - PrivKVと提案手法でkey-valueデータの度数と平均値を推定し，推定誤差MSEを算出する。
 - この試行を10回行いMSEの平均値をアルゴリズムの評価値とする。

実験結果1 推定平均値の分布

RQ1. 提案手法はPrivKVよりも高精度か？

平均値の推定分布 ($\epsilon = 4, n = 10^5$)



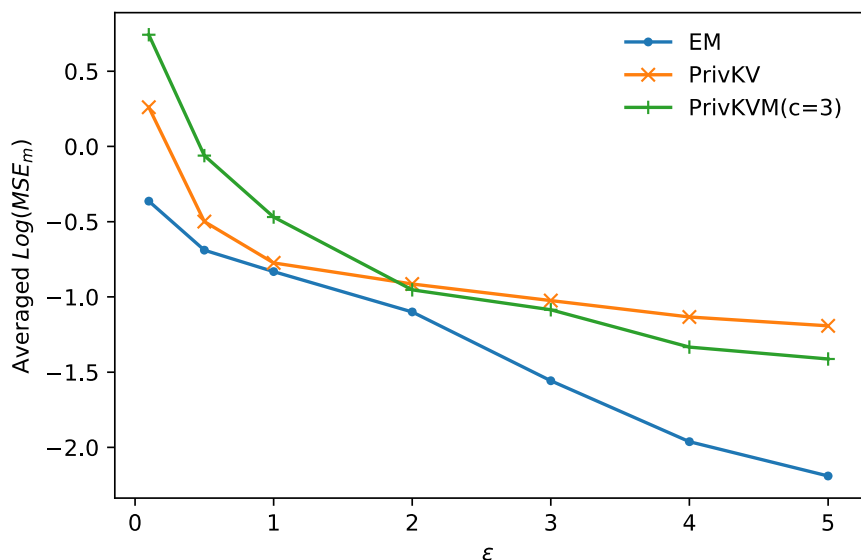
提案手法は低頻度のkeyに対しても高い推定精度

実験結果2

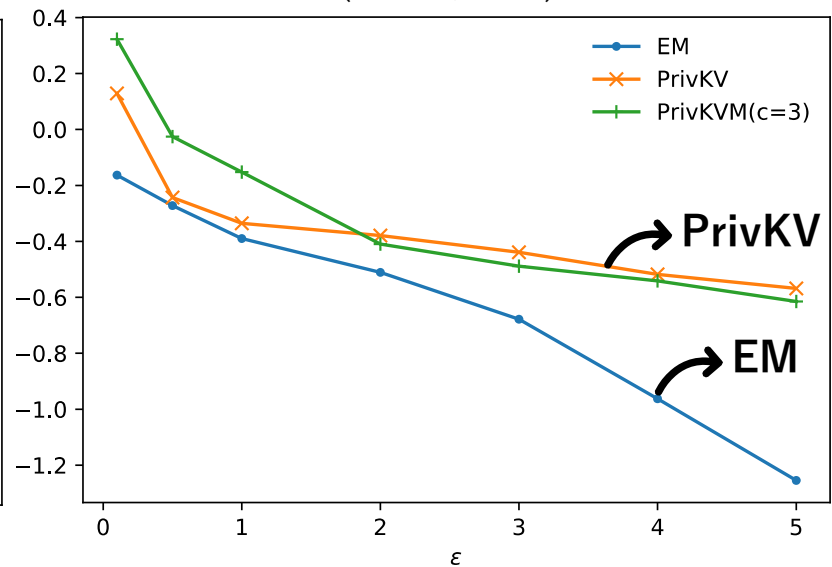
RQ2. 安全性 ϵ によって推定誤差に影響はあるのか？

key数=50, ユーザ数 $n = 10^5$

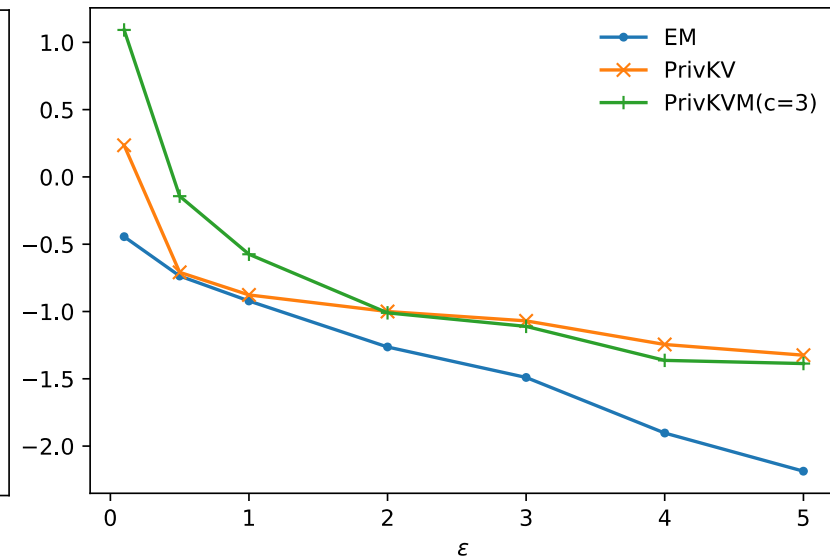
MSE_m(ガウス分布)



MSE_m(べき分布)



MSE_m(線形分布)

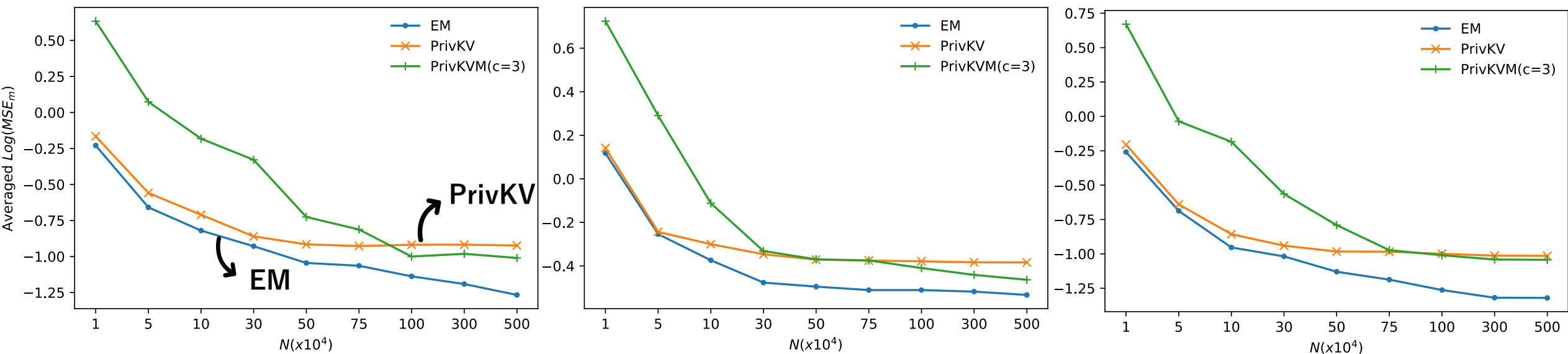


いかなる安全性 ϵ でも提案手法の誤差が小さい

実験結果3

RQ3. データ規模 n によって推定誤差に影響はあるのか？

key数=50, $\epsilon = 2$



データ規模に関わらず提案手法の誤差が小さい

まとめ

- PrivKVでは推定に最尤推定法が用いられており，度数が0.2以下の小さなkeyに関して，平均値の推定誤差が大きい。
- PrivKVで摂動化したkey-valueデータにEMアルゴリズムを適用して推定する手法を提案した。
- その結果，平均値推定では合成データの実験の平均で85.2%の改善が見られた。