



Indian Institute of Technology
Hyderabad

Differentially Private Small Dataset Release Using Random Projections

Scalable Algorithms for Data Analysis

Team members:

KETHAVATH PRANEETH NAYAK
JEEVAN SAMMESWAR

Group number:
G12

Agenda

25.04.2023

- 01** motivation/importance
- 02** naive/existing solutions and their limitation
- 03** Algorithm details
- 04** Experimental results
- 05** Summary/Conclusion



MOTIVATION

Why we need data analysis?

Data aids reproducibility and promotes new discoveries

Problem?

- "as-is" sharing of data leads to privacy breach
- Direct measures of data sanitization, such as removing primary identifiers (like name.), and/or rounding of variables. Such data sanitization practices are ineffective and combining multiple such releases, an adversary can accumulate information about an individual, leading to uncontrolled privacy leakage or worse, a complete disclosure

Naive/existing solutions

- **Direct measures of data sanitization**, such as removing primary identifiers (like name.), and/or rounding of variables. Such data sanitization practices are ineffective and combining multiple such releases, an adversary can accumulate information about an individual, leading to uncontrolled privacy leakage or worse, a complete disclosure

Differentially Private GANs - (Generative Adversal Networks)

- Trained on real data.
- This guarantees the sampled data is “**synthetic**”, but still follows the distribution similar to real data. But as GANs are trained **unconstrained on real data**, they can implicitly or explicitly disclose sensitive information contained in the training set.

Naive/existing solutions [cont.]

Differentially Private GANs limitations

- Diff. priv. GANs can still fall short on the utility front for small-sized datasets especially (where a significant portion of releasable data is in high sensitive domains).

Algorithm details

Differential Privacy :

Definition 2.4 (Differential Privacy). A randomized algorithm \mathcal{M} with domain $\mathcal{N}^{|\mathcal{X}|}$ is (ϵ, δ) -differentially private if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$ and for all $x, y \in \mathcal{N}^{|\mathcal{X}|}$ such that $\|x - y\|_1 \leq 1$:

$$\Pr[\mathcal{M}(x) \in \mathcal{S}] \leq \exp(\epsilon) \Pr[\mathcal{M}(y) \in \mathcal{S}] + \delta,$$

where the probability space is over the coin flips of the mechanism \mathcal{M} . If $\delta = 0$, we say that \mathcal{M} is ϵ -differentially private.



Algorithm details

Random Projection :

Random projection is a method to project original d -dimensional data to a k -dimensional subspace ($k \neq d$ usually) through the origin, using a random $k \times d$ matrix.

JL lemma :

Lemma 1. (Johnson-Lindenstrauss Lemma [15]) *Let $\nu \in (0, 1/2)$. Let $Q \subset \mathbb{R}^d$ be a set of n points and $k = \frac{20 \log n}{\nu^2}$. There exist a Lipschitz mapping $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$, such that for all $u, v \in Q$, we have*

$$(1 - \nu) \|u - v\|_2^2 \leq \|f(u) - f(v)\|_2^2 \leq (1 + \nu) \|u - v\|_2^2$$

- The approximate distance between them is preserved.
- $X(n \times d)$ to a k -dimensional subspace.
- Create a random matrix $R(d \times k)$ and take the product, XR .
- Here, we use simple Gaussian R , where the entries are drawn from $N(0, 1/\sqrt{k})$.

Algorithm details

DPRP (Differentially Private Data Release via Random Projections) :

DPRP takes the original dataset (X) as input, computes the Singular Value Decomposition (SVD) of the covariance matrix X_C of X , and then uses the right singular vector (\hat{V}^T) in conjunction with a random projection P across the columns to reconstruct $X' \approx X$. For preserving differential privacy, we ensure that \hat{V}^T and P are differentially private (the only instances of real data needed for reconstruction).

Algorithm details

DPRP (Differentially Private Data Release via Random Projections) :

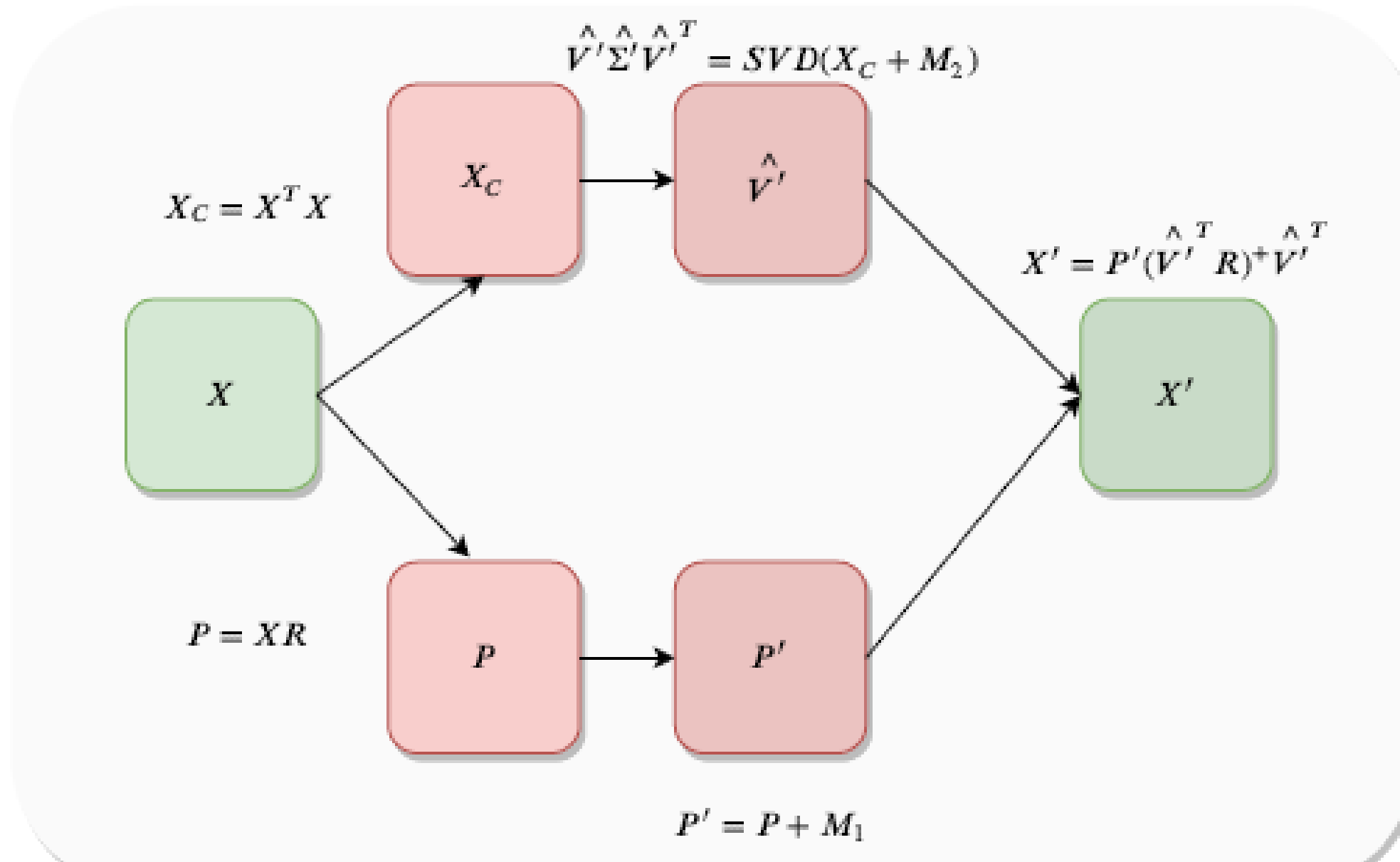


Figure 1: DPRP schema: Using X , we create a random projection P across the columns, and the covariance matrix X_C . We decompose the covariance matrix using SVD. Using noisy right singular vector (\hat{V}') from the decomposition along with noisy P' , we reconstruct $X' \approx X$. Details on noise addition (M_1, M_2) for differential privacy and the reconstruction are explained in detail in Section 3, Algorithm 1.

Algorithm details

Algorithm 1: DPRP: Differentially Private Reconstruction of Input Data

Input: Dataset: X ; Privacy parameters: ϵ, δ ;
Privacy budget allocation: $b_1\%$ for
random projection P , $1 - b_1\%$ for
 $SVD(X_C)$; Number of dimensions for
random projection P : k_1 ; Number of
values from right singular vector to keep
from $SVD(X_C)$: k_2

Algorithm details

Output: Differentially private dataset: X'

- 1 $R \sim \mathcal{N}(0, 1/\sqrt{k_1})^{d \times k_1}$
 - 2 $P = XR$
 - 3 $P' = P + M_1; M_1 \sim \mathcal{N}(0, \sigma_1^2)^{n \times k_1}$ // With budget $b_1\%$
 - 4 $X_C = X^T X$
 - 5 $\hat{V}' \hat{\Sigma}' \hat{V}'^T = \text{SVD}(X_C + M_2); M_2 \sim \mathcal{N}(0, \sigma_2^2)^{d \times d}$
// With budget $1 - b_1\%$
 - 6 $V'_{k_2} = \hat{V}'[1, \dots, k_2]$ // First k_2 columns
 - 7 $X' = P'(V'_{k_2}{}^T R)^+ V'_{k_2}{}^T$
-

Variance (σ_1, σ_2)

σ_1, σ_2 :

Lemma 2. [18] For two neighbouring datasets X and X' that only differ in one observation, i , with $\|X_i - X'_i\| \leq Z$, and a random Gaussian matrix P with entries drawn from $\mathcal{N}(0, \sigma_p^2)$, where $\sigma_p = 1/\sqrt{k_1}$. With probability at least $1 - \delta$, we have

$$\|XP - X'P\|_F \leq Z\sigma_p \sqrt{k_1 + 2\sqrt{k_1 \log(1/\delta)} + 2\log(1/\delta)}$$

Lemma 3. [20] The mechanism $M(D) = f(D) + G$, where G is a random Gaussian matrix with entries drawn from $\mathcal{N}(0, \sigma_1^2)$, satisfies (ϵ, δ) -differential privacy, if $\delta < \frac{1}{2}$, where $\sigma_1^2 = 2\Delta_2(f)^2(\log(1/2\delta) + \epsilon)/\epsilon^2$ and $\Delta_2(f)$ is the sensitivity

Theorem A P' is (ϵ_1, δ_1) -differentially private if we add noise from $\mathcal{N}(0, \sigma_1^2)$; where

$$\sigma_1 = Z\sigma_p \sqrt{k_1 + 2\sqrt{k_1 \log(2/\delta_1)} + 2\log(2/\delta_1)} \\ \sqrt{2(\log(1/2\delta_1) + \epsilon_1)}/\epsilon_1$$

Theorem B \hat{V}' is (ϵ_2, δ_2) -differentially private if we add noise to X_C from $\mathcal{N}(0, Z^2 \sqrt{2 \ln 1.25/\delta_2}/\epsilon_2)$

Using sequential composition[14], we get the Algorithm 1 as (ϵ, δ) -differentially private, where $\epsilon = \epsilon_1 + \epsilon_2$ and $\delta = \delta_1 + \delta_2$. \square

Contribution of Algorithm (Pros)

- A model-free, reconstruction based approach for diff. priv. datasets, utility bottle neck for generative models.
- DPRP is easy to implement, computationally cheap, & offers *one-shot* reconstruction.
- Avoids hyperparameter optimization (which is required in deep generative models).
- Extensive evaluation on *seven* diverse real-life datasets.

Experimental results on Indian Liver disease

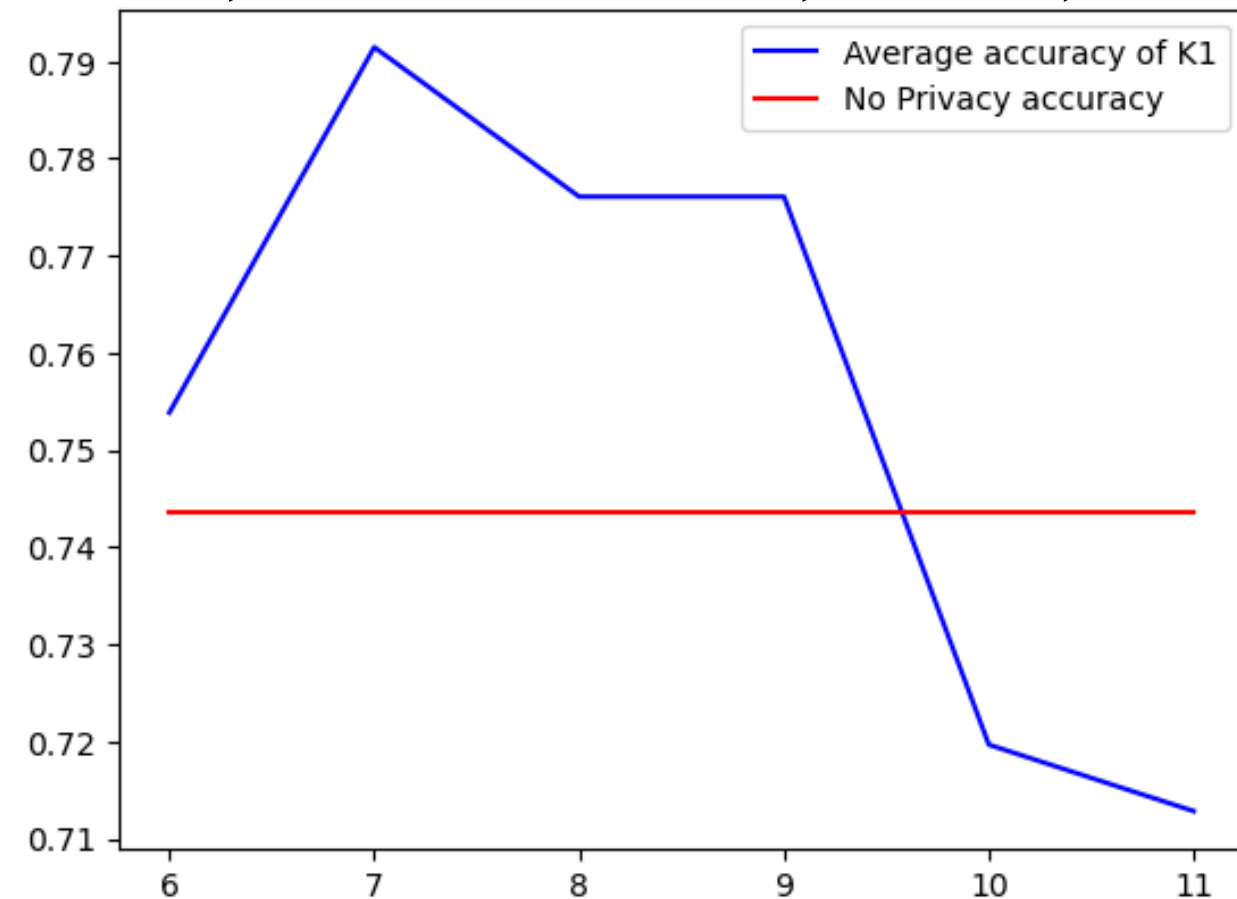
Used Classifier Random Forest Classifier



Indian Institute of Technology
Hyderabad

K1 Tuning

- K1 is the parameter describing final dimensions after random projection
- We got better result for $k1 = 10$
- Even after removing one element we got accuracy ~ 0.720 and ~ 0.725 for both
- This result satisfying $\Pr[M(x) \in S] \leq e^\epsilon \Pr[M(y) \in S] + \delta$
- $\epsilon = 4$, $\delta = 0.0001$, $k2 = 7$, $b = 0.8$



K1 vs accuracy for liver disease data



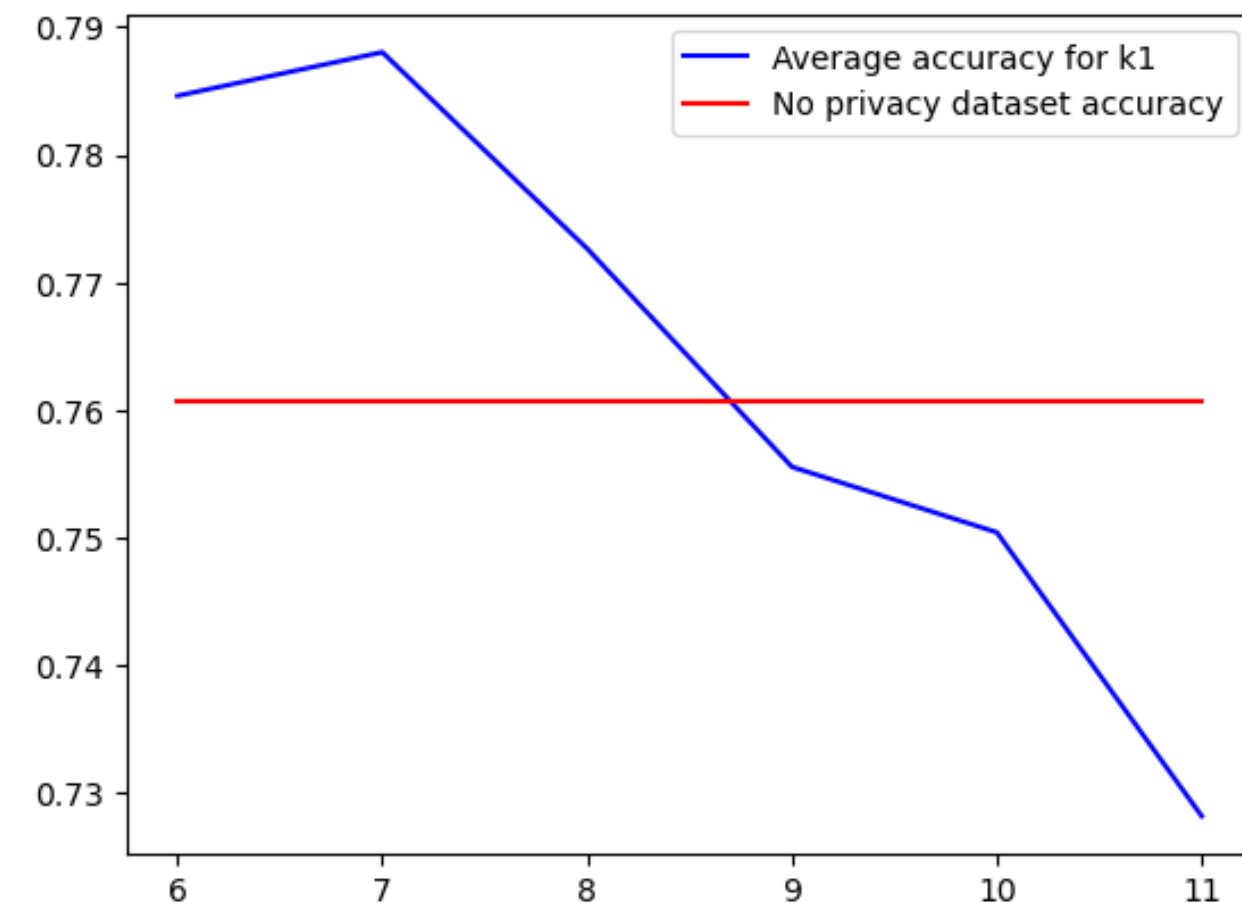
K1 vs accuracy for liver disease data
after removing one entry

K1 Tuning[cont.]

- After modifying one entry , the algorithm is applied on data.
- K1 is performing better at 10

Indian Liver	DPRP	0.79, 0.70	0.79, 0.66	0.77, 0.66	0.75, 0.65	0.72, 0.66
	DPGAN	0.55, 0.59	0.53, 0.60	0.51, 0.59	0.49, 0.57	0.46, 0.54
	DP-CGAN	0.54, 0.60	0.52, 0.59	0.50, 0.54	0.47, 0.52	0.45, 0.51
	No Privacy	0.83, 0.74	0.83, 0.74	0.83, 0.74	0.83, 0.74	0.83, 0.74

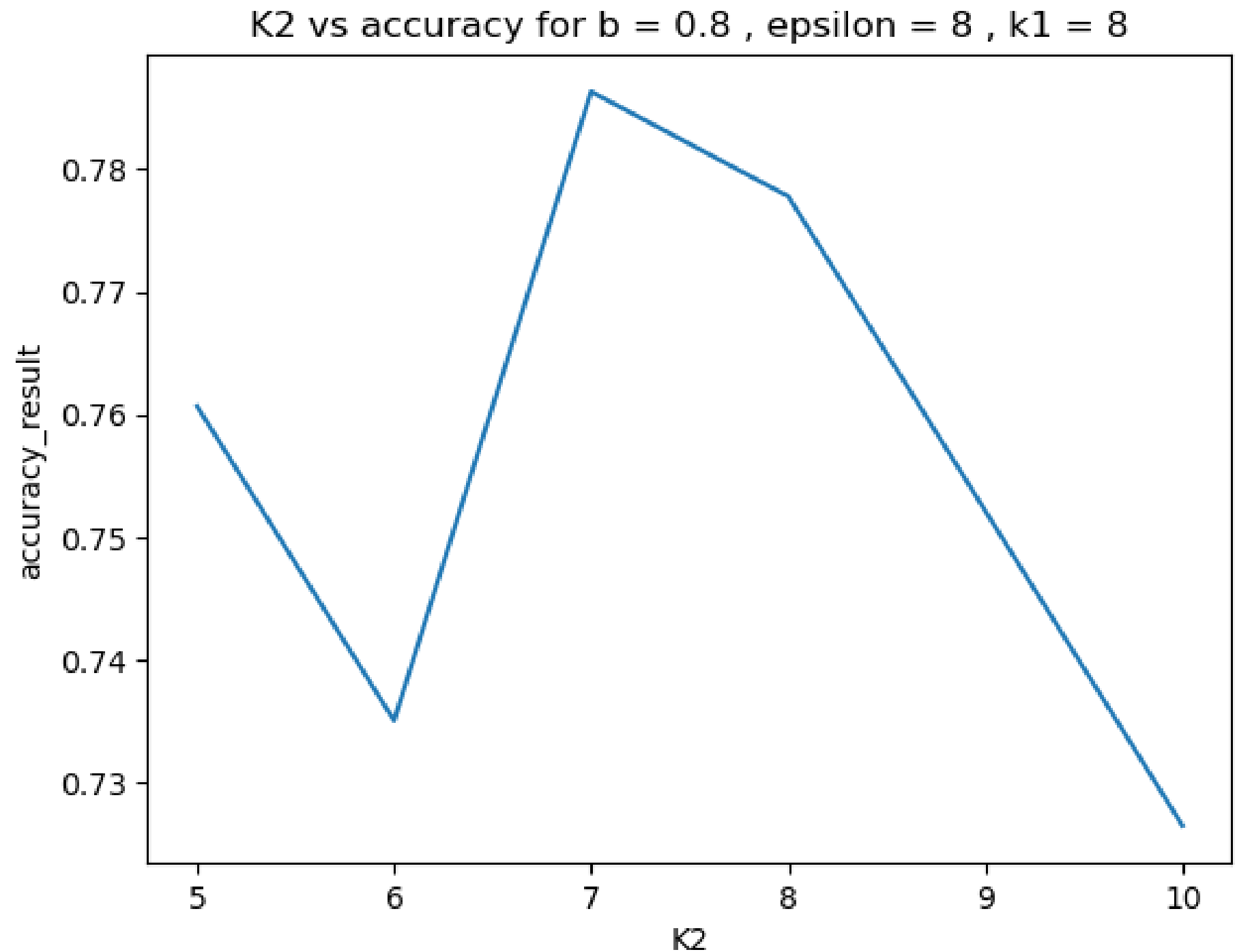
Comparisions



K1 vs accuracy for liver disease data after modifying one entry

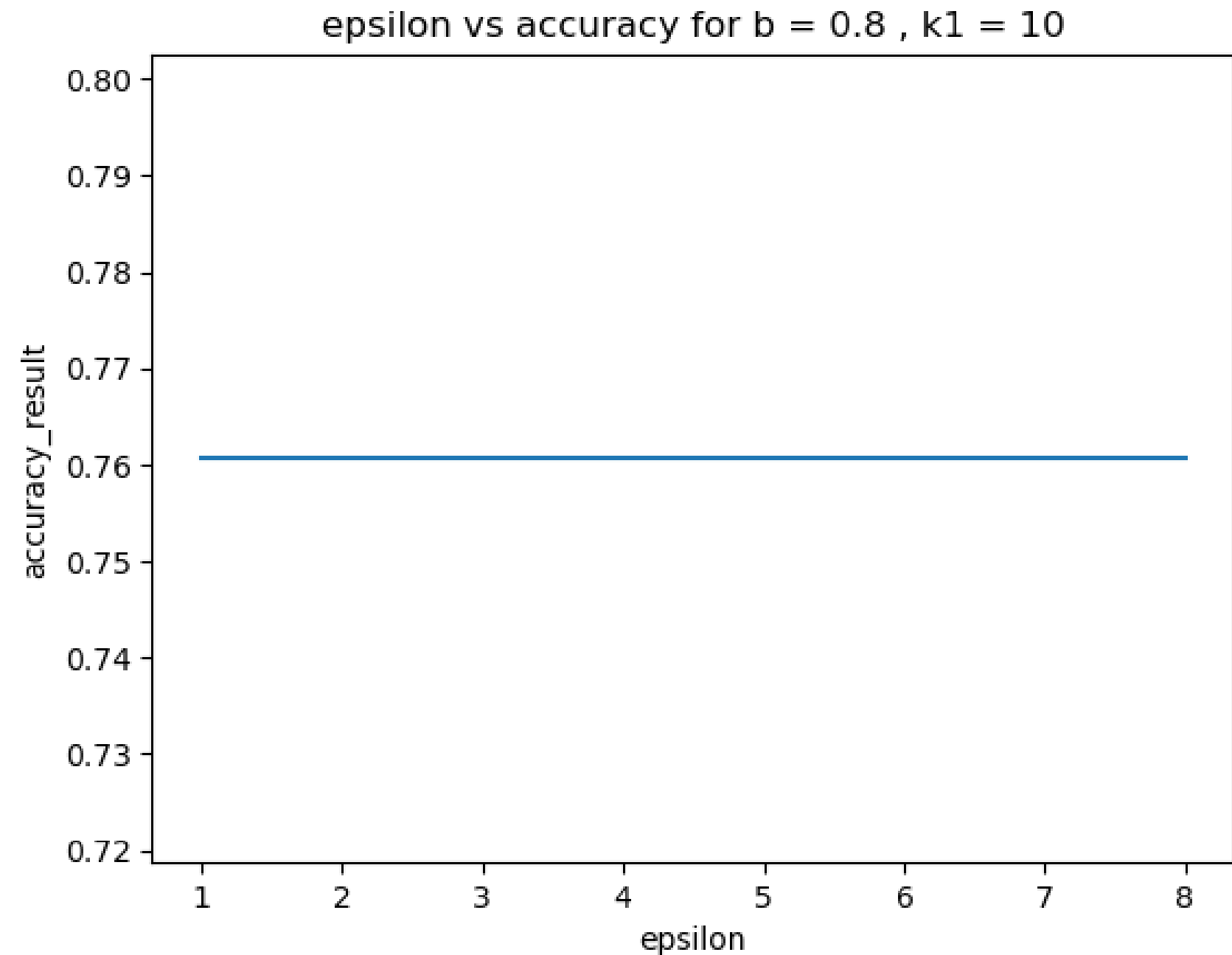
Effect of K2

- K2 is the parameter describing how many columns we need to pickup from right singular component after performing svd
- As per the theorem they achieved optimal value $0.6 \cdot d$. Here d is dimension of input vector.
- From this experimental results it is also giving better results at $0.6 \cdot 11 \sim 7$



Budget effect on epsilon, ϵ and delta, δ .

- Even with this large privacy budget the algorithm is performing better.
- Mostly Budget >40% is performing better
- epsilon can be [8,6,4,2,1] and delta = 0.0001



Reference



Lovedeep Gondara
Ke Wang

Department of Computing Science
Simon Fraser University



<http://proceedings.mlr.press/v124/gondara20a/gondara20a.pdf>

Differentially Private Small Dataset Release Using Random Projections

Lovedeep Gondara Ke Wang
Department of Computing Science
Simon Fraser University

Abstract

Small datasets form a significant portion of re-

1.2 CURRENT APPROACH

Differential privacy [3] offers a solution. Formalizing the notion of privacy as a mathematical definition, differential privacy promises any released data will not unduly

Thank you!



Indian Institute of Technology
Hyderabad