



Not All Attributes are Created Equal: dx -Private Mechanisms for Linear Queries

IE-506 Course Project



Team - ML Isolation
Roll No.- 23N0453

Outline:



- ❑ Problem Description
- ❑ Summary of Mid-Term Work-Done
- ❑ Major Comments
- ❑ Addressing of Comment
- ❑ Description of Work Done
- ❑ Experiments
- ❑ Contribution
- ❑ Conclusion
- ❑ Future Work
- ❑ References

Problem Description:

❑ Privacy Concerns:

- Talks about keeping things private, especially if they're about personal stuff like race, gender, health, or money.
- Stresses how important it is to keep things private when sharing data, especially when it's sensitive.
- It's about respecting people's privacy and rights

❑ dX-Privacy:

- Introduces a way to measure how well we're keeping things private when we share data.
- Like making sure no one can figure out personal details from what we share the data..

Problem Description:

❑ Mechanisms for Privacy:

- Suggests ways to share data safely, like adding extra noise or using special techniques.
- Introduces diverse algorithms which gives to different privacy and accuracy requirements.
- Addresses the challenge of balancing privacy and accuracy in data release scenarios.
- It's about finding algorithm which provide a balance between sharing useful data and keeping it private.

Summary of Mid-Term Work Done:

❖ Introduction and Motivation:

- Introduction to the concept of differential privacy (DP) and its importance in balancing privacy and data analysis.
- Motivation behind the development of dX-privacy, emphasizing the need for flexible privacy protection for different attributes.
- Examples of sensitive data scenarios in healthcare, ride-sharing platforms, and recommendation systems.
- Overview of past works in the field of DP and the evolution of privacy models.

❖ **Solution Approach and Algorithms:**

- Description of privacy concerns related to releasing data with sensitive attributes and the focus on linear queries.
- Introduction of the dX-privacy concept and its application to linear queries over data sets.
- Overview of mechanisms proposed for achieving dX-privacy, including Laplace and Exponential mechanisms, and the MWEM algorithm.
- Discussion on the mathematical definitions and computational frameworks involved in the solution approach.

Summary of Mid Term Work Done



❖ **Experimental Evaluation and Conclusion:**

- Overview of experiments conducted to evaluate the proposed mechanisms' performance in terms of privacy, utility, and computational efficiency.
- Comparison of results between the Laplace mechanism, enhanced Laplace mechanism, and the MWEM algorithm.
- Conclusion highlighting the superior performance of MWEM in preserving data utility while providing dX -privacy guarantees.
- Future directions for research, including exploring other frameworks and addressing challenges in maintaining privacy while minimizing the impact on query results.

Major Comments:

- ❖ Slides are verbose and Definition of differential privacy .
- ❖ For final review the team must do the following: slides -- can be made better comparison with paper result.
- ❖ Should be provided dx blowfish, dx smooth to be implemented .
- ❖ For large datasets -- plot rmse against iterations or other relevant quantity blowfish, smooth -- has to be implemented .
- ❖ Report is missing in submission.
- ❖ Results on real datasets should be showcased in the final review.



Addressing the Comment:

- Definition of Dp:
 - Learning more about DP by some other Research paper and uses some open sources.
- Slides :
 - By using bullet points, simple content and limited text .
- dx blowfish and dx- smooth implementation:
 - Implement these two in colab and getting the results
- For Large Data Sets:
 - Using RMSE to get results and plotting to compare the algorithms .
- Apply on Real Data Sets:
 - Implementing the algorithm over Real data set of US cities in which location is privatised .

Description of Work Done:

- Algorithms on Real Data set:
 - ◆ Applying all the algorithm on real data set and getting plots for comparison.
 - ◆ Finding RMSE values for all algorithm
- Understanding the dx Blowfish and dx Smooth :
 - ◆ Developing understanding over both the algorithms and understanding what is the difference between them.
- Implementation of code:
 - ◆ Applying on both real and synthetic data set
 - ◆ Comparing using RMSE and plotting the graph

★ Requirement:

- Why we not apply MWEM Algorithm all the time?
 - As MWEM add the noise for all the location as we have to privatised only the neighbouring location or selected location.
 - Due to this RMSE increases and we are not able to get better results.

★ Terminology:

- Threshold “T”:
 - Defining distance “T” as threshold , as if we are within this T then algorithm uses some different notion of adding noise.
 - Otherwise, it behaves different and use some other metrics to add noise .

dx- Blowfish :

- ★ Using laplace noise method we developing this algorithm adding some updates in choosing epsilon that is privacy budget.

- ★ dx-Blowfish :

$$d_{\mathcal{X}}^{\text{Blow}} \text{ s.t. } d_{\mathcal{X}}^{\text{Blow}}(i, j) = \epsilon \text{ if } d_{\mathcal{X}}^{\text{Euc}}(i, j) \leq T,$$
$$d_{\mathcal{X}}^{\text{Blow}}(i, j) = \infty \text{ otherwise}$$

Where ϵ is privacy budget and T is threshold

dx- Smooth:

- ★ If we add some new updates on dx-blowfish then we get dx smooth .
- ★ dx-smooth increases the privacy budget proportional to the distance between the pair of points.
- ★ dx-smooth:

$$d_{\mathcal{X}}^{\text{Smooth}} \text{ s.t. } d_{\mathcal{X}}^{\text{Smooth}}(i, j) = \epsilon \text{ if } d_{\mathcal{X}}^{\text{Euc}}(i, j) \leq T$$

$$d_{\mathcal{X}}^{\text{Smooth}}(i, j) = \frac{\epsilon d_{\mathcal{X}}^{\text{Euc}}(i, j)}{T} \text{ otherwise}$$

where $d_{\mathcal{X}}^{\text{Euc}}(i, j) := \sqrt{(u_i - u_j)^2 + (v_i - v_j)^2}$ which is euclidean distance between two points.

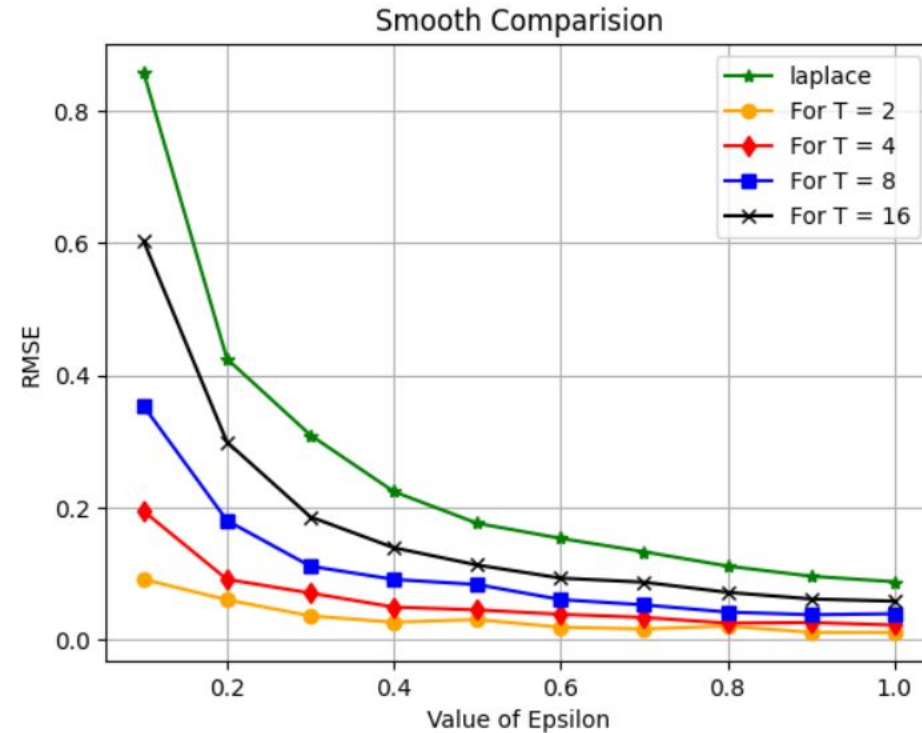
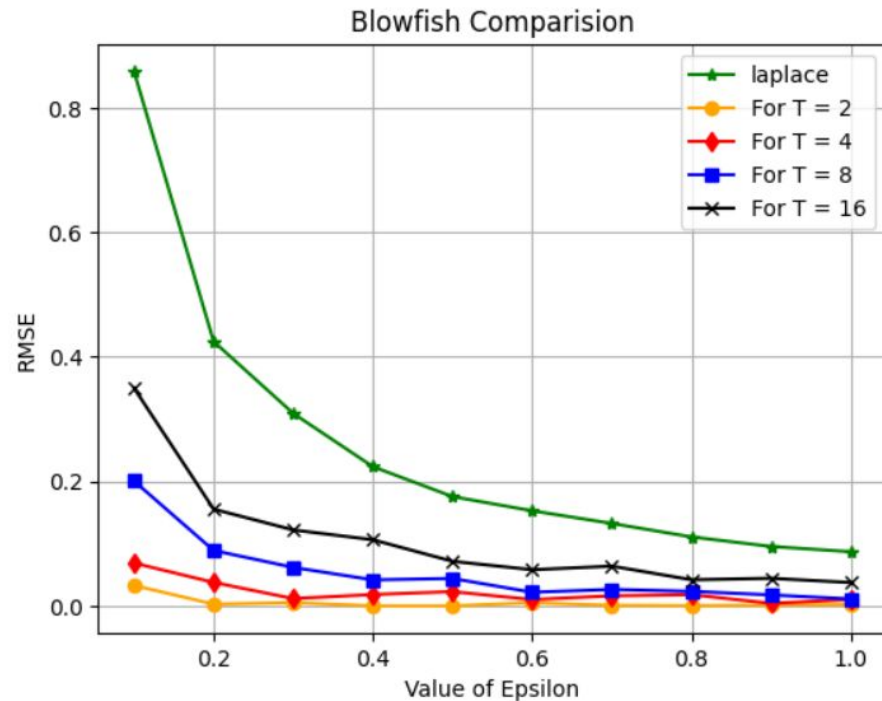
Implementation of Code:



- RMSE:
 - First of all applying code for finding RMSE of old laplace ,new laplace and MWEM on real Data set and synthetic data to compare the algorithm.
- **dx-Blowfish and dx-Smooth:**
 - Using Colab for writing the code of both algorithms in which uses laplace noise according to definition of updated privacy metrics
 - Plotting the RMSE for different values of “T” and Epsilon uses in paper.
 - Comparing the plots for both algorithms and visualise the difference between them.
 - Fixing the threshold value and comparing both algorithms.
 - Finally, doing the experiments and getting result similar to research paper.

Experiments:

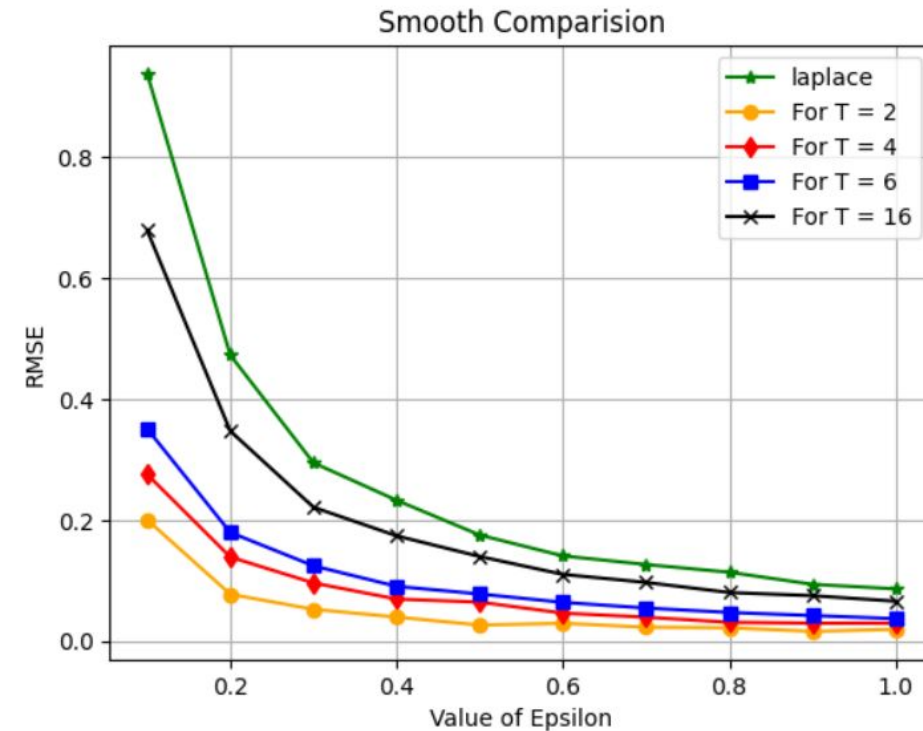
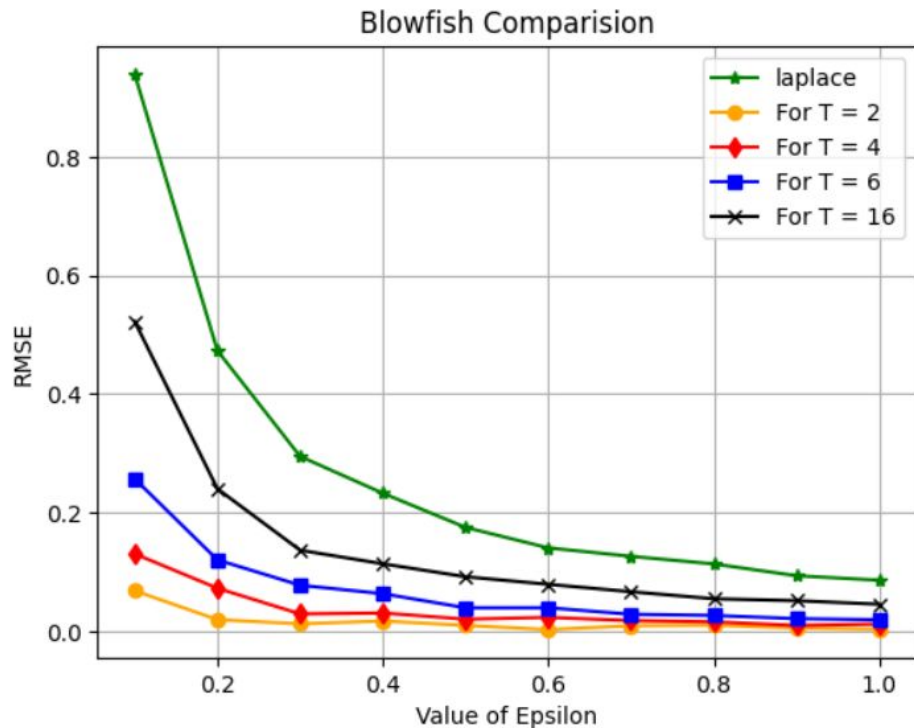
Experiment over synthetic data:



- ❖ Clearly from graphs for similar values of threshold and epsilon the RMSE of Blowfish is less than the Smooth It means the accurate result we get from Blowfish while maintaining the privacy.

Experiments:

❖ Experiments over Real Data set:



- As for different values of epsilon RMSE decreases it means that on Increasing the epsilon(privacy budget) our privacy decreases.
- As values of T increases the RMSE increases means more neighbours are protected

Experiments:

★ RMSE Values Comparison of dx-Blowfish and dx-Smooth :

Sr No	Value_of_T	dx-BlowFish	dx-Smooth
1	2	[0.012259353227633352]	[0.014318483398459581]
2	4	[0.014895104689777432]	[0.02606796302988975]
3	8	[0.024952217017176984]	[0.046566696417345184]
4	16	[0.0455358606632592]	[0.07261079237843399]

- Seeing the tabular value for fixing the Privacy budget=1 , We observe that for different value of “T” dx-BlowFish have less RMSE and Dx-Smooth have higher in every case , It means for better accuracy dx-BLowfish is better instead of dx-smooth

Contribution:

- ❑ As a single member , all work done is completely by self
- ❑ I have written all the code by myself
- ❑ There is no code available in research paper
- ❑ Data set is given in Research paper

Conclusions:

- ❖ For High Privacy and if we don't want to any how attackers can regenerate the data set then, from results we see MWEM is best.
- ❖ For just we have to add some noise in particular data points(less sensitive) we use Laplace mechanism.
- ❖ For preserving data from attackers and preventing regeneration and as well want to preserve accuracy of results on data set then our newly generated dx-BLowfish algorithm is extremely good.
- ❖ For preserving data from regeneration attacks and want some more privacy , We use dx-smooth .

Possible Future Work :

- ❖ Can we come up with new privacy metric for preserving the data sets and follow all triangular inequalities.
- ❖ Trying to find out best possible epsilon for any data attribute by which its privacy and accuracy can be determined together.
- ❖ Further , If we got some factor by which privacy and accuracy both measured then will try to improve existing algorithms or come up with some new one.

.

References:

Paper allotted: 1806.02389.pdf (arxiv.org) released - [1806.02389.pdf \(arxiv.org\)](https://arxiv.org/pdf/1806.02389.pdf)

- ❑ Not All Attributes are Created Equal: dx-Private Mechanisms for Linear Queries
 - ❑ Authors: Parameswaran Kamalaruban, Victor Perrier, Hassan Jameel Asghar, and Mohamed Ali Kaafar
 - ❑ Date of Publish: *28 Aug 2019*

- ❑ Differentially Private Fine-tuning of Language Models
 - ❑ Authors: Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, Huishuai Zhang
 - ❑ Link: <https://arxiv.org/pdf/2110.06500>
 - ❑ Date of Publish: 2021

References:

- ❑ Optimizing Linear Counting Queries Under Differential Privacy:
 - ❑ Authors: Chao Li[†] , Michael Hay[†] , Vibhor Rastogi[‡] , Gerome Miklau[†] , Andrew McGregor[†]
 - ❑ Link: <https://dl.acm.org/doi/abs/10.1145/1807085.1807104>
 - ❑ Date of Publish: June 2010

- ❑ Calibrating Noise to Sensitivity in Private Data Analysis:
 - ❑ Authors: Dwork, F. McSherry, K. Nissim, and A. Smith
 - ❑ Link: <https://people.csail.mit.edu/asmith/PS/sensitivity-tcc-final.pdf>
 - ❑ Date of Publish: 2006

- ❑ For basic Understanding Used Wikipedia and Youtube:
 - ❑ https://en.wikipedia.org/wiki/Differential_privacy