## Lecture 3: Sparse Vector Technique and Linear Query Release (October 4)

*Lecturer: Yu-Xiang Wang*                                         *Scribes: Lawrence Lim*

**Note**: *LaTeX template courtesy of UC Berkeley EECS dept.*

**Disclaimer**: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

This lecture's notes illustrate some uses of various LaTeX macros. Take a look at this and imitate.

## 3.1 Mathematical notations representing a dataset or dataset space

- Standard data-table representation: For $n$ people and $d$ data points per person, we can represent a data-table as a $n \times d$ matrix. Each data point can be represented as $\mathcal{X}$, and the dataset is often represented as $\cup_{n=1}^{\infty} \mathcal{X}^n = \mathcal{X}^*$.

- Histogram representation: For each unique item in the dataset, counts the number of people with that item. Then we can represent our dataset as $x \in \mathbb{N}^{|\mathcal{X}|}$ for datapoints $\mathcal{X}$. A data-table representation can be converted into a histogram representation by summing indicator values.

- Indicator representation: a bit vector that means whether a person is inclusive with a feature.

## 3.2 Differential Privacy Definition

A randomized algorithm $\mathcal{M}$ with domain $\mathbb{N}^{|\mathcal{X}|}$ is $(\epsilon, \delta)$-differentially private if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$ and for all $x, y \in \mathbb{N}^{|\mathcal{X}|}$ such that $||x - y||_1 \leq 1$:

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

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

## 3.3 Utility of Laplace Mechanism

CDF of the Laplace distribution:

$$\begin{cases} \frac{1}{2} e^{\frac{x-\mu}{b}} & x \leq \mu \\ 1 - \frac{1}{2} e^{-\frac{x-\mu}{b}} & x \geq \mu \end{cases}$$

Let $f : \mathbb{N}^{|\mathcal{X}|} \to \mathbb{R}^k$, and let $y = \mathcal{M}_L(x, f(\cdot), \epsilon)$. Then for all $\delta \in (0, 1]$:

$$\Pr[||f(x) - y||_{\infty} = \max(f(x)_i - y_i) \geq \ln(\frac{k}{\delta}) \cdot \frac{\Delta f}{\epsilon}] \leq \delta$$

Intuitively, this means that for some probability $\delta$, we can bound the error. The chance that the error is greater than $\ln(\frac{k}{\delta}) \cdot \frac{\Delta f}{\epsilon}$ is less than $\delta$ probability for $\Delta f$ defined as L1 sensitivity for $f$, $k$ number of queries, and $\epsilon$-differential privacy.

### 3.3.1   L1 sensitivity

L1 sensitivity of a function $f$ represents how much the output of $f$ changes with the removal or addition of one datapoint cell.

A linear query, which is just counting the number of people that satisfy a property, has a sensitivity of one because the removal or addition of one datapoint results in at most one change in the output. For the same reasons, the L1 sensitivity of a histogram is 1.

### 3.3.2   Applying the Laplace mechanism

1. Set privacy budget and number of queries. For example, we can choose privacy budget $\epsilon_{\text{budget}} = 0.5$, number of queries $k = |Q| = 100$, so privacy budget per query $\epsilon_{\text{per query}} = \epsilon_{\text{budget}}/|Q| = 0.005$.

2. Decide how much noise to add. We use the Laplace mechanism to add Laplace noise $Lap(b)$ with $b = \frac{\Delta f}{\epsilon_{\text{per query}}} = \frac{1}{0.005} = 200$.

3. Compute the error bound. Using the error bound from before, $\max_{q \in Q}(|y_q - q^T x|) = \alpha \leq \frac{|Q|}{\epsilon} \ln \frac{|Q|}{\delta}$ with probability $1 - \delta$. Then the normalized error is $\tilde{O}(\frac{|Q|}{n\epsilon})$ and the statistical error is $\tilde{O}(\frac{1}{\sqrt{n}})$.

4. Finally, we can convert an error bound to the number of samples needed to reach a bounded error, called sample complexity. Error $\alpha = \frac{|Q|\log(\frac{|Q|}{\delta})}{n\epsilon}$, so rearranging terms, we get $n = \frac{|Q|\log(\frac{|Q|}{\delta})}{\epsilon\alpha}$.

## 3.4   Applying Randomized Response to Answer Linear Queries

Suppose we had a dataset of secret bits for each person. We can use randomized response applied on each bit to answer a linear query. We'll call our original dataset $x$ and our randomized response dataset $\hat{x} = 0.5 + \frac{Y-0.5}{2p-1}$.

Then $q^T \hat{x} = \sum_{i=1}^n q_i \hat{x}_i$ and the error $= q^T \hat{x} - q^T x$.

$0.5 + \frac{-0.5}{2p-1} \leq q_i \hat{x}_i = 0.5 + \frac{Y-0.5}{2p-1} \leq 0.5 + \frac{0.5}{2p-1}$

We then apply Hoeffding's inequality. $\Pr(\frac{1}{n}|q^T \hat{x} - q^T x| > t) < 2e^{1\frac{2nt^2}{O(\frac{1}{\epsilon})}}$

$\ln(\frac{2}{\delta}) = \epsilon 2nt^2/C$ for some constant $C$. This implies that $t = \frac{C\sqrt{ln(2/\delta)}}{\sqrt{n}}$.

Intuitively, this means that with probability $1 - \delta$ the error is $O(\frac{\sqrt{\ln(1/\delta)}}{\sqrt{n}\epsilon})$.

Answering many randomized response queries does not cost any additional privacy using RR for linear queries.

For many queries the Laplace Mechanism has error $O(|Q|\frac{\log(\frac{|Q|}{\delta})}{n\epsilon})$ while the Randomized Response Mechanism has error $O(\frac{\sqrt{\log(|Q|/\delta)}}{\sqrt{n}\epsilon})$.

## 3.5 Apply Laplace Mechanism to Release Histograms

The L1-sensitivity $\Delta f$ of a histogram representation release of data is 1. So for each $x \in \mathcal{X}$, we can use the Laplace Mechanism to release the histograms.

The error is $\frac{1}{n}|q^T \hat{\mathbf{Hist}}(\text{Data}) - q^T \mathbf{Hist}(\text{Data})| = O(\frac{1}{n\epsilon}\sqrt{|\mathcal{X}|}\log(\frac{1}{\delta}))$.

## 3.6 AboveThreshold Mechanism and Sparse Vector Mechanism

Suppose we want to answer queries that meet a certain threshold. That is, release the identity of the first query in a stream that passes this threshold. How can we accomplish this?

The AboveThreshold Mechanism enables this by applying Laplace noise to the threshold (at each check) and if value is above threshold, then add Laplace noise to the value and return that. Intuitively the AboveThreshold Mechanism guarantees privacy by sometimes returning the wrong index. AboveThreshold guarantees $\epsilon$-DP.

The SparseVector mechanism allows one to query for multiple values above a threshold. Due to the Composition Theorem and the fact that AboveThreshold guarantees $\epsilon$-DP, the SparseVector mechanism chains multiple AboveThreshold mechanisms to obtain multiple values that pass a test that one may be interested it. For $c$ values, and total $\epsilon$ privacy budget, the SparseVector mechanism uses the AboveThreshold mechanism with $\epsilon/c$ privacy budget $c$ times or until the stream of data is used up.