

# Differentially Private Data Release under Partial Information

David Zeber   Martin Lopatka

Systems Research Group, Mozilla

PEPR, 12 August 2019

# Data and privacy

User data!

Utility:

- Learning about users
  - build insights or models by computing summaries & aggregations
- Sharing
  - internally across different parts of the org
  - with business partners
  - publicly

What are the privacy implications for sensitive data?

# Differential Privacy as a solution

**Differential Privacy** (DP) is the state-of-the-art for releasing reports/views based on sensitive data

- avoids revealing “too much” about the users in the dataset

E.g. a DP view allows for further analysis without needing access to the raw data

- can be shared with a partner

# Differential Privacy

- A specific guarantee to individuals whose data is included in a dataset
- Limits how much can be learnt about an individual from a DP-protected dataset
- Limit is quantified by a number ( $\epsilon$ ): **privacy budget**

# DP mechanisms

To use DP, apply a randomized operation (**mechanism**) to the dataset.

- Original data is masked by injecting randomness in a strategic way
- Makes “influential” values in the dataset less influential

# Applying DP

In practice, DP is non-trivial to apply

- Real-world datasets are more complex than what is typically described in the literature
- Multiple mechanisms may be available

Selecting a “good” mechanism is important:

- a bad choice could end up giving low utility results while consuming privacy budget

# Selecting a DP mechanism

- Which mechanism will perform best depends on the properties of the particular dataset
- “Peeking” at the data violates DP

## Solutions:

- draw on literature set in the same domain
- allocate some of the privacy budget to learning characteristics in a DP way

# Existing partial information

In practice, we often do have prior partial knowledge about the dataset:

- may satisfy known constraints
  - e.g. limitations of an app or reasonable user behaviour
- other summary information about the users may already have been shared
  - e.g. `addons.mozilla.org` lists raw installation counts



# Partial information example






	<b>uBlock Origin</b> <span>Recommended</span> Finally, an efficient blocker. Easy on CPU and memory. ★★★★★ <small>Raymond Hill</small>	4,653,854 users
	<b>Video DownloadHelper</b> <span>Recommended</span> The easy way to download and convert Web videos from hundreds of YouTube-like sites. ★★★★★ <small>mig</small>	2,595,585 users
	<b>NoScript Security Suite</b> <span>Recommended</span> The best security you can get in a web browser! Allow active content to run only from sites you trust, and protect yourself against XSS other web security exploits. Disabled and can't reinstall? <a href="https://tinyurl.com/moz-ext-cert">https://tinyurl.com/moz-ext-cert</a> ★★★★★ <small>Giorgio Maone</small>	1,514,432 users
	<b>Facebook Container</b> <span>Recommended</span> Prevent Facebook from tracking you around the web. The Facebook Container extension for Firefox helps you take control and isolate your web activity from Facebook. ★★★★★ <small>Mozilla</small>	652,648 users
	<b>Privacy Badger</b> <span>Recommended</span> Automatically learns to block invisible trackers. ★★★★★ <small>EFF Technologists</small>	599,855 users

Figure: addons.mozilla.org (recommended add-ons)

# Proposed DP selection approach

**Data model:** describe a dataset probabilistically in terms of how it was produced

- i.e. as a generative distribution over possible datasets
- known constraints built in via conditioning

**Mechanism selection:** evaluate candidate DP mechanisms by simulating over a sample of possible datasets

# Proposed DP selection approach

## Benefits:

- Can incorporate partial information of differing specificities
- Does not consume any privacy budget
- Does not require access to the private dataset (e.g. could be performed by a third party)

# Problem setting

## Data model: user-item data

- $n \times m$  user-item matrix
- entry  $(i, j) = 1$  if user  $i$  interacted with item  $j$

	<i>(items)</i>			
	A	B	C	D
<i>(users)</i> 1	1	0	1	1
2	0	1	0	0
3	1	0	1	0
4	1	0	0	0
5	1	1	0	0

# User-item data: output

**Goal:** report the item frequency counts with privacy protection

- labeled column sums of the matrix

$$\begin{array}{c} \phantom{1} \\ \phantom{2} \\ \phantom{3} \\ \phantom{4} \\ \phantom{5} \end{array} \begin{array}{cccc} A & B & C & D \\ \left[ \begin{array}{cccc} 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{array} \right] \end{array} \xrightarrow{\text{colsums}} \begin{array}{cccc} A & B & C & D \\ \left[ \begin{array}{cccc} 4 & 2 & 2 & 1 \end{array} \right] \end{array}$$

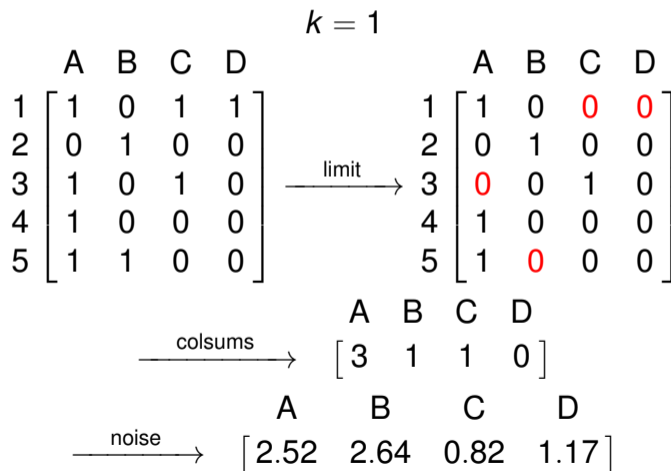
Influence of a single user: adds 1 to the count for each item they report

# User-item data: mechanism

- 1 Limit each user to  $k$  items selected at random
- 2 Compute item frequency counts over the limited items
- 3 Add  $Lap(k/\epsilon)$  random noise to each count

Satisfies  $\epsilon$ -DP

# User-item data: mechanism



# User-item mechanism tradeoff

This introduces a bias-variance tradeoff tuned by the parameter  $k$

- If many users have close to  $m$  items, would want  $k$  to be large (lower bias)
- If most users have few items, would want  $k$  to be small (lower noise)



# Case study

- A team at Mozilla curates a dataset containing a list of “favourite” sites (URLs) from a subset of Firefox users\*
- Access to raw data is highly restricted for privacy reasons
- Data curators are willing to release a DP version of the frequency distribution if they are provided a query to run

\* who have opted in to this data collection

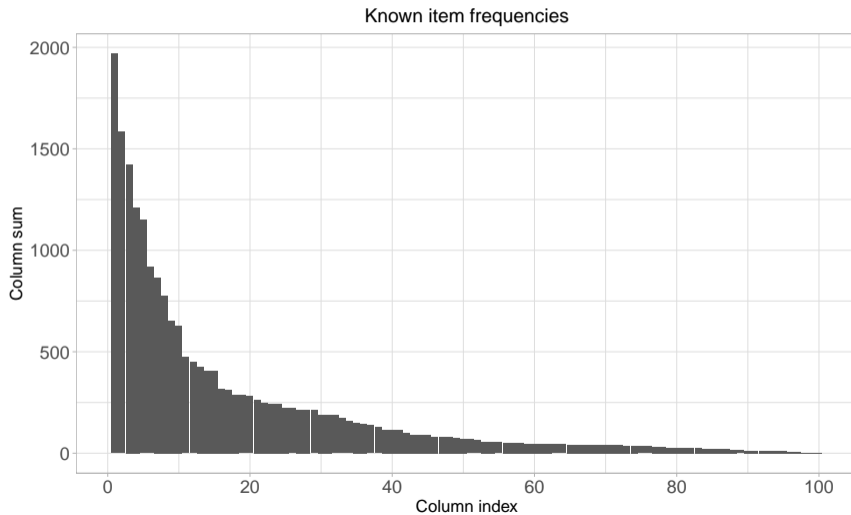
# Case study

What do we know about this data?

- Under the policies governing this dataset, the curators are able to share internally the site counts (i.e. frequency distribution) *with URLs scrubbed*

We can use this to build a generative model for this dataset with constraints.

# Item frequencies



# Generative model: setting

- Dataset is an  $n$  (users)  $\times$   $m$  (items) binary matrix
- Fix column ordering from highest to lowest column sum, same for rows
- Known:  $m$ , column sums  $n \geq c_1 \geq \dots \geq c_m \geq 1$

# Generative model: number of rows

In this case, number of users  $n$  is not known directly. However, column sums impose constraints:

- $n \geq c_1 = \max c_j$  (single copy of each item per user)
- $n \leq \sum_j c_j$  (at least 1 item per user)

# Examples: number of rows

$$\begin{array}{cccc} A & B & C & D \\ [4 & 2 & 2 & 1] \end{array}$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

$$n = 4$$

$$\begin{bmatrix} 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$n = 5$$

$$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$n = 8$$

# Generative model: row sums

Given  $n$ , we have constraints on the row sums (number of items reported by each user):

- $m \geq r_1 \geq \dots \geq r_n \geq 1$
- $\sum_{i=1}^n \min(r_i, p) \geq \sum_{j=1}^p c_j$  for  $p = 1, \dots, m$  (Gale-Ryser conditions for binary matrices)

These can be reformulated into sequential bounds  $L_i \leq r_i \leq U_i$  where  $L_i = L(r_1, \dots, r_{i-1})$  and  $U_i = U(r_1, \dots, r_{i-1})$

# Examples: row sums

$$\begin{array}{cccc} A & B & C & D \\ [4 & 2 & 2 & 1] \end{array}, \quad n = 5$$

$$\begin{array}{c} 4 \\ 2 \\ 1 \\ 1 \\ 1 \end{array} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{array}{c} 3 \\ 2 \\ 2 \\ 1 \\ 1 \end{array} \begin{bmatrix} 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$\begin{array}{c} 2 \\ 2 \\ 2 \\ 2 \\ 1 \end{array} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$



# Examples: matrices

$$\begin{array}{cccc} \text{A} & \text{B} & \text{C} & \text{D} \\ [4 & 2 & 2 & 1] \end{array}, \quad n = 5, \quad \mathbf{r} = [3, 2, 2, 1, 1]$$

$$\begin{array}{c} 3 \\ 2 \\ 2 \\ 1 \\ 1 \end{array} \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{array}{c} 3 \\ 2 \\ 2 \\ 1 \\ 1 \end{array} \begin{bmatrix} 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$\begin{array}{c} 3 \\ 2 \\ 2 \\ 1 \\ 1 \end{array} \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

# Generative model

- 1 Select  $n$  uniformly among possible choices
- 2 Given  $n$ , select  $r$  uniformly subject to constraints
- 3 Given marginal counts, select a binary matrix uniformly

Performed using MCMC sampling (cf. Verhelst, 2007)

# Mechanism selection

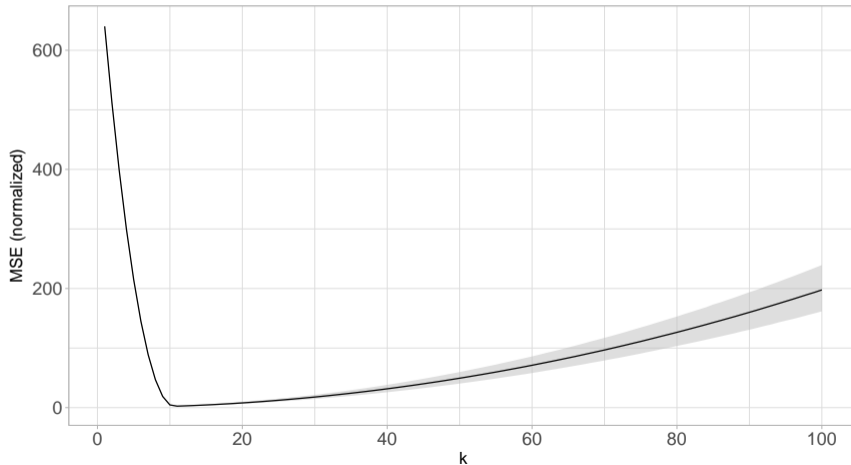
Given a feasible matrix:

- 1 Simulate the DP mechanism for different values of  $k$  with the desired  $\epsilon$
- 2 Compute squared-error loss relative to the true item counts
- 3 Estimate MSE by averaging over multiple replicates of (1) and (2)

Replicate the above over a sample of feasible matrices

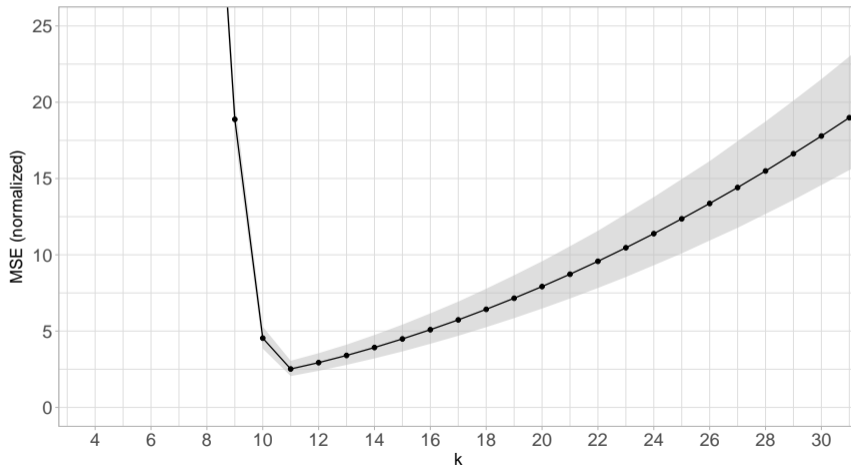
# Evaluation results

Average error of DP item frequencies for each  $k$   
 $n = 2,000$ ,  $m = 100$ , sample size = 100,000



# Evaluation results

Average error of DP item frequencies for each  $k$  (zoomed)  
 $n = 2,000$ ,  $m = 100$ , sample size = 100,000



# Conclusion

We present an alternative approach to mechanism selection for DP:

- does not consume privacy budget
- does not require access to the private dataset
- leverages any available ancillary information about the problem setting

# Conclusion

Vision: a software package applicable to a number of domains assisting in

- formulating constraints
- running simulations
- visualizing and interpreting results

Thank you!