# On Privacy Risk of Releasing Data and Models

Ashish Dandekar
*Supervised by: A/P Stéphane Bressan*
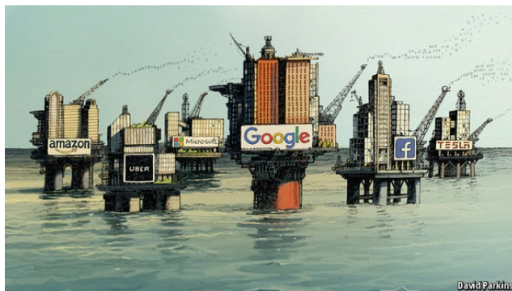
July 18, 2019

# Data is the new oil!



*(The Economist, 6 May 2017).*

# AI is the new electricity!



'AI IS THE NEW ELECTRICITY'

"Just as electricity transformed almost everything 100 years ago, today I actually have a hard time thinking of an industry that I don't think AI will transform in the next several years."

**Andrew Ng**
Former chief scientist at Baidu, Co-founder at Coursera

# Privacy risk: Publishing Data



Mr. Mark Zuckerberg

**Mea culpa, mea culpa, mea maxima culpa!**

'Facebooks failure to compel *Cambridge Analytica* to delete all traces of data from its servers  including any "derivatives"  enabled the company to retain predictive models derived from millions of social media profiles!' *(The Guardian, 6 May 2018)*.

## Privacy risk: Publishing Data

An arms race between anonymisation and re-identification!

- Re-identification of the governor of Massachusetts in 2000
- Re-identification of Thelma Arnold from AOL searches in 2006
- Re-identification of the users from Netflix dataset in 2006
- Re-identification of the cabs in New York City taxi dataset in 2014

# Privacy risk: Publishing Models

If machine learning models learn latent patterns in the dataset, what are the odds that they learn something that they are not supposed to learn?

## Attacks on machine learning models

- **Inference attack.** [Homer et al., 2008] infer presence of a certain genome in the dataset from the published statistics of genomic mixture dataset.

- **Model inversion attack.** [Fredrikson et al., 2014] infer genetic marker of patients given the access to machine learning model trained on the warfarin drug usage dataset.

- **Membership inference attack.** [Shokri et al., 2017] infer the presence of a data-point in the training dataset given the access to machine learning models hosted on cloud platforms.

## Our contributions

*"Synthetic datasets put a full stop on the arms race between anonymisation and re-identification."* [Bellovin et al., 2018]

### Publication of data

- We illustrate partially and fully synthetic dataset generation techniques using a selection of discriminative models.
- We adapt and extend Latent Dirichlet Allocation, a generative model, to work with spatiotemporal data.

## Our contributions

We use *differential privacy* [Dwork et al., 2014] to provide quantifiable privacy guarantee while releasing machine learning models.

### Publication of models

- We illustrate use of the functional mechanism to provide differential privacy guarantees for releasing regularised Linear regression.
- We illustrate use of perturbation of model functions to provide differential privacy guarantees for a selection of non-parametric models.
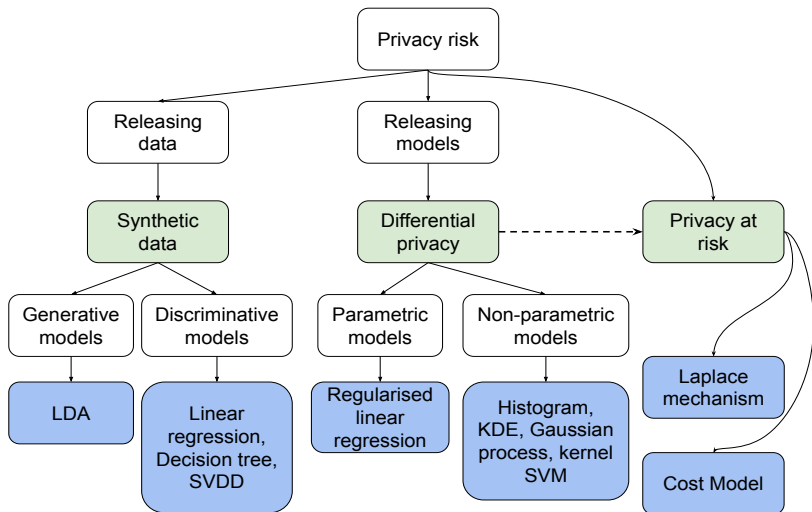
## Our contributions

In the spirit of making differential privacy amenable to business entities, we propose *privacy at risk*. It is a probabilistic relaxation of differential privacy.

### Privacy at risk

- We define *privacy at risk* that provides probabilistic bounds on the privacy guarantee of differential privacy by accounting for various sources of randomness.
- We illustrate privacy at risk for Laplace mechanism.
- We propose a cost model that bridges the gap between the abstract guarantee and compensation budget estimated by a GDPR compliant business entity.

# Summary

# Publication of data

*(Privacy risk of re-identification)*

## Synthetic Data

As authentic as these "Nike" shoes!

## Synthetic dataset generation techniques

With the help of a domain expert, a data scientist classifies features of any data-point into two categories.

- **Identifying features.** This is a set of attributes that are not typical to the dataset under study. These attributes can be publicly available as a part of other datasets.

- **Sensitive features.** This is a set of attributes that are typical to the dataset under study. These attributes contain data that is deemed to be sensitive.
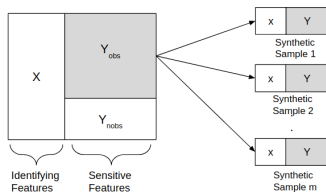
$$\underbrace{\text{DOB, Marital Status, Gender,}}_{\text{Identifying features}} Income \rightarrow \text{Census dataset}$$
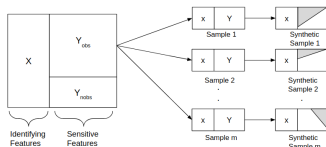$$\text{DOB, Marital Status, Gender,} HIVStatus \rightarrow \text{Health dataset}$$

# Synthetic dataset generation techniques

## Fully synthetic dataset generation [Rubin, 1993]



- Sample $m$ datasets from imputed population and release them publicly

## Partially synthetic dataset generation [Reiter, 2003]



- Instead of imputing all values of sensitive features, only impute those values that bear higher cost of disclosure

# Experimental evaluation

We extend comparative study of [Drechsler and Reiter, 2011] by using linear regression as well as neural networks as data synthesisers on the US Census dataset of 2003[1].

## Utility Evaluation

| Feature | Data Synthesisers | Original Sample Mean | Fully Synthetic Data | | |
|---------|-------------------|----------------------|----------------|---------|--------------|
| | | | Synthetic Mean | Overlap | Norm KL Div. |
| Income | Linear Regression | 27112.61 | 27074.80 | 0.52 | 0.55 |
| | Decision Tree | 27081.45 | 27091.02 | 0.55 | 0.58 |
| | Random Forest | 27107.04 | 28720.93 | 0.54 | 0.64 |
| | Neural Network | 27185.26 | 26694.54 | 0.54 | **0.99** |

| Feature | Data Synthesisers | Original Sample Mean | Partially Synthetic Data | | |
|---------|-------------------|----------------------|----------------|---------|--------------|
| | | | Synthetic Mean | Overlap | Norm KL Div. |
| Income | Linear Regression | 27112.61 | 27117.99 | 0.98 | 0.54 |
| | Decision Tree | 27081.45 | 27078.93 | 0.98 | **0.99** |
| | Random Forest | 27107.04 | 27254.38 | 0.95 | 0.58 |
| | Neural Network | 27185.26 | 27370.99 | 0.81 | **0.99** |

[1]https://usa.ipums.org/usa/

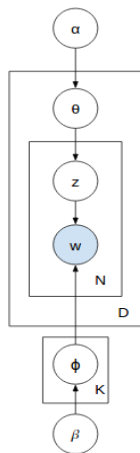# Experimental evaluation

## Disclosure risk evaluation scenario

Consider, an intruder who is interested in people who are born in US and earn more than $250,000. We consider a tolerance of 2 when matching on the age of a person. *We assume that the intruder knows that the target is present in the publicly released dataset.*

| Data Synthesisers | True match rate | False match rate |
| --- | --- | --- |
| Linear Regression | 0.06 | 0.82 |
| Decision Tree | 0.18 | 0.68 |
| Random Forest | 0.35 | 0.50 |
| **Neural Network** | **0.03** | **0.92** |

## Why generative models?

- Generative models learn $\mathbb{P}(Data|pattern)$ unlike discriminative models that learn $\mathbb{P}(pattern|Data)$

- Generative models do not tend to overfit the training data

- Generative models have a data-generative process at the heart of its inception

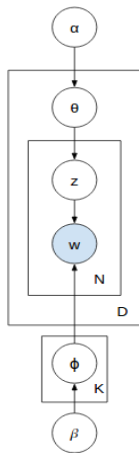# Latent Dirichlet Allocation [Blei et al., 2003] (LDA)



## Notation

- $N$ : Vocabulary size
- $D$ : Total number of Documents
- $K$ : Total number of Topics

## Intuition

- Bag of Words assumption
- A document is a distribution over topics
  - $\boldsymbol{\theta_m} \to K$-dim vector; $m \in [1...D]$
- A topic is a distribution over words
  - $\boldsymbol{\phi_k} \to N$-dim vector; $k \in [1...K]$

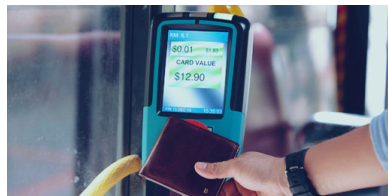# Latent Dirichlet Allocation [Blei et al., 2003] (LDA)



### Generative Process

1. Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$ for a document
2. For each word in the document:
   1. Draw a topic $z \sim \text{Mult}(\theta_d)$
   2. Draw a word $w_{d,z} \sim \text{DirMult}(\phi_z | \beta)$

# Generating travelling records of commuters

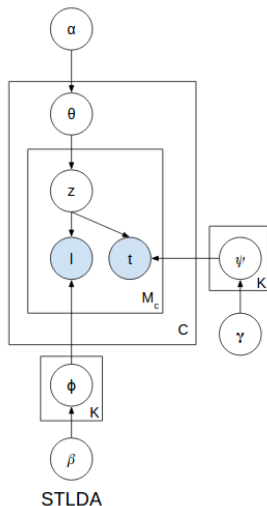| Card Number | In-Timestamp | Out-timestamp | In-ID | Out-ID |
|:---:|:---:|:---:|:---:|:---:|
| c530524 | 2012-02-12;07:22:49.0 | 2012-02-12;07:28:50.0 | 2383 | 1467 |
| c530545 | 2012-02-12;12:09:40.0 | 2012-02-12;12:29:40.0 | 1464 | 8 |
| c630568 | 2012-02-12;13:10:30.0 | 2012-02-12;13:40:50.0 | 2413 | 99 |
| c534554 | 2012-02-12;20:08:12.0 | 2012-02-12;20:28:07.0 | 2384 | 2 |
| c837483 | 2012-02-12;16:02:10.0 | 2012-02-12;16:34:33.0 | 1467 | 185 |



Credit: *home.ezlink.com.sg*
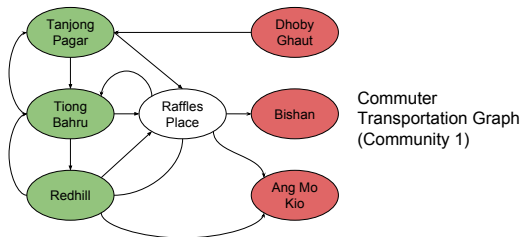


Credit: *mustsharenews.com*
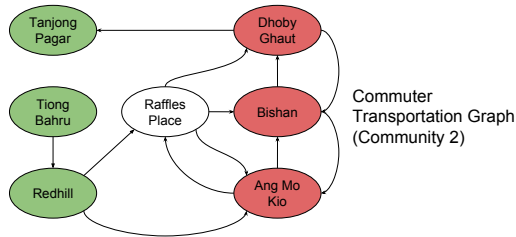
# Adapting and extending LDA



STLDA

We adapt LDA to SLDA and TLDA that work with spatial and temporal data respectively. We extend LDA to STLDA that works with spatiotemporal data. We call topics found by these models as *communities*.

| Model | Documents | Words | Topics |
|-------|-----------|-------|--------|
| SLDA | Commuters | Visits | Spatial mobility patterns |
| TLDA | Commuters | Timestamps | Temporal mobility patterns |
| STLDA | Commuters | Spatiotemporal events | Spatiotemporal mobility patterns |

# Adapting and extending LDA



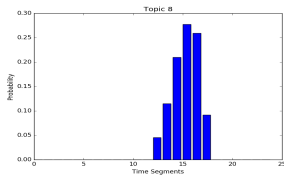Commuter
Transportation Graph
(Community 1)

Commuter
Transportation Graph
(Community 2)
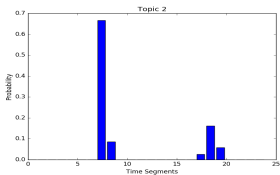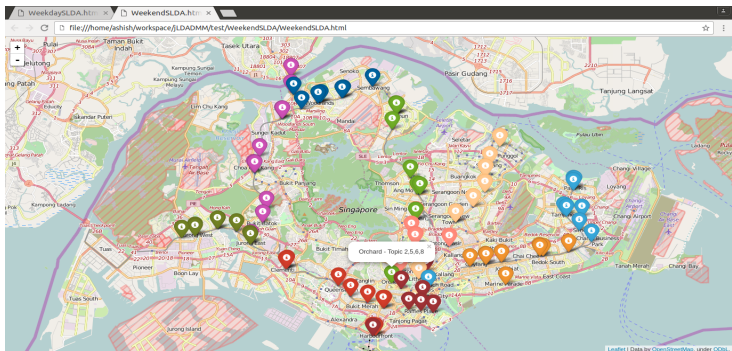
We synthesise commuter
records for a community by
performing a random walk on
the commuter transportation
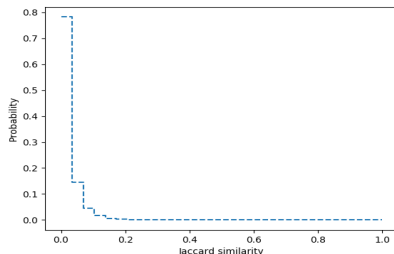graph for the specified
community.
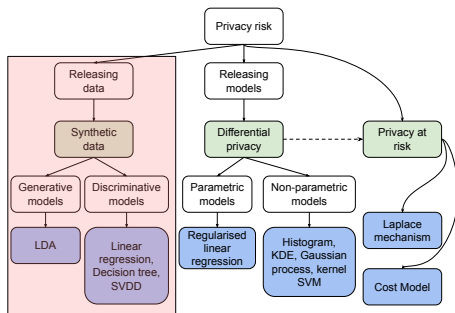
# Results

## Results

### Privacy risk evaluation

We evaluate privacy by using the metric of *Jaccard similarity*.

For every community, we generate travelling records of 1000 commuters. For every travelling record, we compute its Jaccard similarity with training documents.

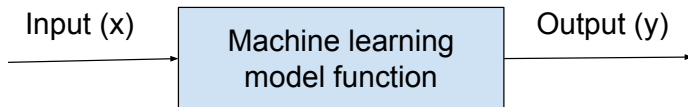## Summary

**Synthetic dataset generation.**



- Risk evaluation measures heavily rely on the disclosure scenario
- Risk evaluation is performed for a generated instance of dataset

Synthetic datasets may not save you from an attribute disclosure and inferential attacks!

# Publication of models

*(Privacy risk of leakage of information)*

# Why do we publish models?



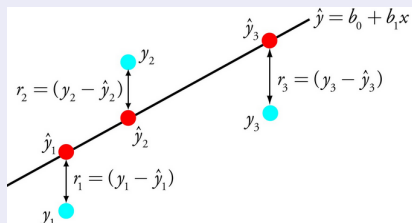Input (x) → | Machine learning model function | → Output (y)

$$y = f(\theta, x)$$

Each component of a machine learning model has become an asset for the organisations!

- **Features**
- **Hyper-parameters**
- ✓ **Parameters**
- ✓ **Predictions**

# Two kinds of machine learning models

**Parametric models.** Values of the parameters are sufficient to compute outputs for new data inputs. Parameters are often estimated by minimising a *loss function* on the training dataset.

## Example: Linear regression



Linear regression predicts $y \in \mathbb{R}$ for a specfied $\mathbf{x} \in \mathbb{R}^d$ as

$$y = \boldsymbol{\theta} \cdot \mathbf{x}$$
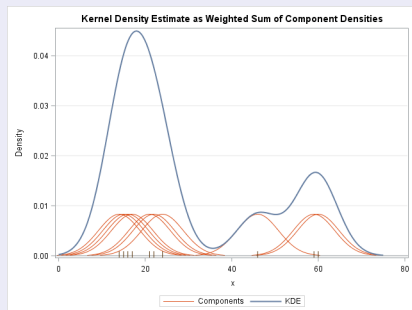
where $\boldsymbol{\theta} \in \mathbb{R}^d$ are parameters of the model.

$$\boldsymbol{\theta}^* = \arg\min \ell(\theta, D) = \arg\min \sum_{i=1}^{n} \frac{1}{n}(y_i - \boldsymbol{\theta} \cdot \mathbf{x_i})^2$$

# Two kinds of machine learning models

**Non-parametric models.** Values of the parameters, known as *hyperparameters*, along with the training dataset are required to compute outputs for new data input.

## Example: Kernel density estimation



Kernel Density Estimate as Weighted Sum of Component Densities

Kernel density estimation (KDE) estimates the probability of a data-point to be from a specified dataset.

$$f_D(\cdot) = \frac{1}{n} \sum_{x_i \in D} k(\cdot, x_i)$$

$$= \frac{1}{n} \sum_{x_i \in D} \frac{1}{(2\pi h_i^2)^d} \exp\left(-\frac{\langle \cdot, x_i \rangle}{2h_i^2}\right)$$

# Differential privacy [Dwork, 2006] (DP)

A randomised algorithm $\mathcal{M}$ with domain $\mathcal{D}$ is $(\epsilon, \delta)$-differentially private if for all $S \in \mathrm{Range}(\mathcal{M})$ and $D, D' \in \mathcal{D}$ such that $D$ and $D'$ are *neighbouring datasets*

$$Pr(\mathcal{M}(D) \in S) \leq e^{\epsilon} Pr(\mathcal{M}(D') \in S) + \delta$$

$(\epsilon, 0)$-differential privacy is often referred as $\epsilon$-differential privacy.

Differential privacy quantifies the degree of *indistinguishability* in the outputs when two *neighbouring* inputs are given to a randomised algorithm.

# DP for machine learning models



Privacy-preserving mechanisms that add random noise from probability distributions can be calibrated to satisfy $\epsilon$-differential privacy.

# DP for machine learning models

Once a machine learning model is released, either by releasing its parameters or by publishing it as a service, we cannot control the number of times an analyst uses the model.

**Sequential Composition [Dwork et al., 2014].** Privacy guarantee of differentially private privacy-preserving mechanism linearly degrades with the number of evaluations of the mechanism.

### Our approach while releasing the models

Privacy-preserving mechanisms that add noise in the model functions are well-suited for the scenario of releasing machine learning models.

## Releasing parametric models

Training of parametric models comprises of estimation of parameters $\theta$ by optimising an objective function $F(\theta, D)$ on a specified dataset $D$. Privacy-preserving mechanisms perturb this objective function to $F'(\theta, D)$.

### Functional mechanism [Zhang et al., 2012]

Suppose $F(\theta, D) = \ell(\theta, D)$ has an expansion in a functional basis, say Taylor Basis. Functional mechanism perturbs $F$ by appropriately scaled noise from the Laplace distribution.

$$F'(\theta, D) = A' + B'\theta + C'\theta^2 + ...$$

where $A', B', C', ...$ are perturbed Taylor coefficients.

# Releasing parametric models

Training of parametric models comprises of estimation of parameters $\theta$ by optimising an objective function $F(\theta, D)$ on a specified dataset $D$. Privacy-preserving mechanisms perturb this objective function to $F'(\theta, D)$.

## Objective perturbation [Chaudhuri et al., 2011, Kifer et al., 2012]

Objective perturbation mechanism perturb $F$ by explicitly adding noise terms in the objective function.

$$F'(\theta, D) = F(\theta, D) + \lambda \|\theta\|_2^2 + b^t \theta$$

where $b$ is sampled from $e^{-\frac{\|\theta\|_2}{\Delta}}$.

# Calibration of the privacy-preserving mechanisms

Appropriately calibrated privacy-preserving mechanisms satisfy $\epsilon$-differential privacy. The calibration involves computation of the *sensitivity* of a function and a privacy level $\epsilon > 0$.

### Sensitivity

Sensitivity of a function $f$, $\Delta_f$, with domain $\mathcal{D}$ is an upper bound on fluctuation in the value of function on any pair of neighbouring datasets $D, D' \in \mathcal{D}$.

$$\Delta_f \geq \max_{D,D'} \|f(D) - f(D')\|_1$$

# Calibration of the privacy-preserving mechanisms

### That's the hard part!

Computation of sensitivity is highly non-trivial and involves enforcement of the constraints on the boundedness and smoothness of the objective function!

Equally harder is the task to provide, a tight bound on the privacy level of differential privacy with the appropriate calibration of the noise!

# DP for regularised linear regression

Ridge ($L_2$ normed regularisation)

$$\theta^* = \arg\min_{\theta} \ell(D, \theta) + \lambda\|\theta\|_2^2$$

LASSO ($L_1$ normed regularisation)

$$\theta^* = \arg\min_{\theta} \ell(D, \theta) + \lambda\|\theta\|_1$$

Elastic net (Convex combination of $L_2$ and $L_1$ normed regularisation)

$$\theta^* = \arg\min_{\theta} \ell(D, \theta) + \lambda(\alpha\|\theta\|_2^2 + (1 - \alpha)\|\theta\|_1), \quad 0 \le \alpha \le 1$$

# DP for regularised linear regression

We extend [Zhang et al., 2012] by providing differential privacy guarantees for LASSO and elastic net regression using the functional mechanism.

## Contributions

- **Sensitivity calculation.** We adopt the sensitivity calculation from [Zhang et al., 2012] to compute sensitivity for the LASSO and elastic net loss function.
- **Optimising non-differentiable loss function.** Existence of $L_1$ norm in LASSO and Elastic net regression leads to non-differentiable loss function. We use Conic programming solvers for optimisation.
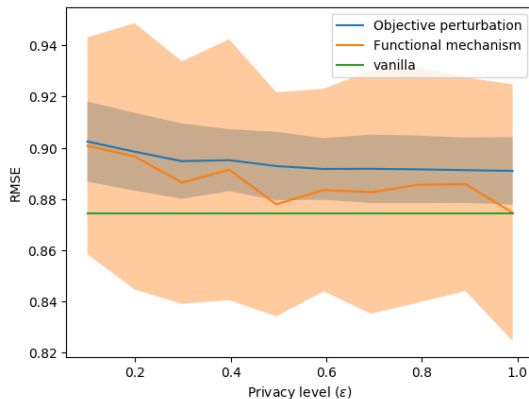
# Performance evaluation



Figure: Comparative evaluation of functional mechanism and objective perturbation mechanism on the wine quality testing dataset for *Ridge regression*

## Releasing non-parametric models

We focus on an important class of non-parametric models that make use of kernels. Model functions of machine learning models that use the "kernel trick" lie in an regenerating kernel Hilbert space (RKHS) spanned by the specified kernel.

| Model | Model function |
|---|---|
| Kernel density estimator | $f_D(\cdot) = \frac{1}{n} \sum_{i=1}^{n} k(\cdot, \mathbf{x_i})$ |
| Gaussian process regression | $\bar{f}_D(\cdot) = \sum_{d_i \in D} \left( \sum_{d_j \in D} (K_D + \sigma_n^2 \mathbb{I})_{ij}^{-1} y_j \right) k(\cdot, \mathbf{x_i})$ |
| Kernel SVM | $w_D = \sum_{i=1}^{n} \alpha_i^* y_i k(\cdot, \mathbf{x_i})$ |

# DP for kernel methods

## Functional perturbation [Hall et al., 2013]

Functional perturbation mechanism perturbs the model function by appropriately scaled noise from the Gaussian process with zero mean and the associated kernel $k$.
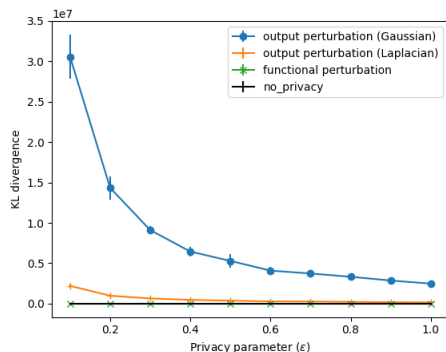
$$f'_D = f_D + \Delta \frac{c(\delta)}{\epsilon} G.$$

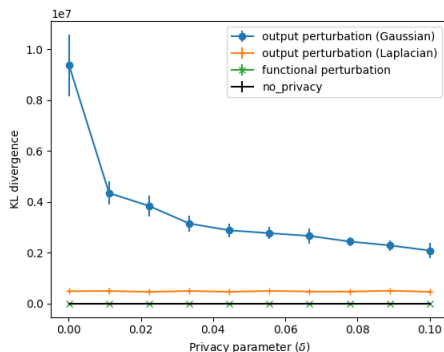With appropriate calibration, it satisfies $(\epsilon, \delta)$-differential privacy.

| Model | Model function | Implementation |
|---|---|---|
| Kernel density estimator | $f_D(\cdot) = \frac{1}{n} \sum_{i=1}^n k(\cdot, \mathbf{x_i})$ | [Hall et al., 2013] |
| Gaussian process regression | $\bar{f}_D(\cdot) = \sum_{d_i \in D} \left( \sum_{d_j \in D} (K_D + \sigma_n^2 \mathbb{I})_{ij}^{-1} y_j \right) k(\cdot, \mathbf{x_i})$ | [Smith et al., 2016] |
| Kernel SVM | $w_D = \sum_{i=1}^n \alpha_i^* y_i k(\cdot, \mathbf{x_i})$ | Partly by [Hall et al., 2013] |

# Performance evaluation

We conduct extensive empirical evaluation of these models on a real-world census dataset as well as benchmark datasets.
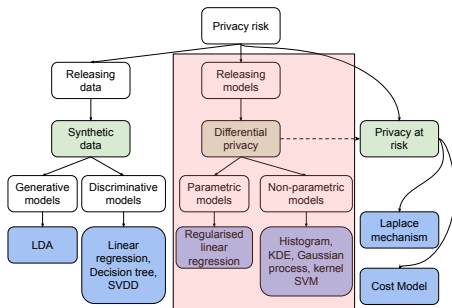


(a) Privacy level ($\epsilon$)                (b) Privacy level ($\delta$)

# Summary

**Differential privacy.**



## Pros

- Privacy guarantee is independent of any disclosure scenario
- Privacy guarantee is for a model that generates data rather than a generated dataset

## Cons

- Privacy guarantee is for the worst-case privacy loss
- Privacy guarantee is too abstract to be actionable

# Privacy at risk

# Privacy at risk

### Problem

Differential privacy accounts for the worst-case privacy loss!

### Motivation

Risk analysts use *Value at Risk* [Jorion, 2000] to quantify the loss in investments for a given portfolio and an acceptable confidence bound. Motivated by the formulation of *Value at Risk*, we define *privacy at risk*.

$(\epsilon, \gamma)$-**privacy at risk.**

$$\mathbb{P}\left[\underbrace{\log\left|\frac{\mathbb{P}(\mathcal{M}(f,\Theta)(x) \in Z)}{\mathbb{P}(\mathcal{M}(f,\Theta)(y) \in Z)}\right| < \epsilon}_{\mathcal{M} \text{ is DP}}\right] \geq \overbrace{\gamma}^{\text{confidence level}}$$

## Privacy at risk

|  | **Source** | **Privacy definition** |
|---|---|---|
| Implicit randomness | Data-generation distribution | Random differential privacy [Hall et al., 2012] |
| Explicit randomness | Noise distribution | Probabilistic differential privacy [Machanavajjhala et al., 2008] |

### Contribution

- We extend the existing works by accounting for the combined effect of implicit randomness and explicit randomness.
- We instantiate privacy at risk for Laplace mechanism.

# Privacy at risk for Laplace mechanism

Laplace mechanism [Dwork et al., 2014]

$$\text{noisy output} \leftarrow f(x) + Lap\left(0, 2\left(\frac{a}{b}\right)^2\right) \quad \{a \leftarrow \Delta_f, b \leftarrow \epsilon_0\}$$

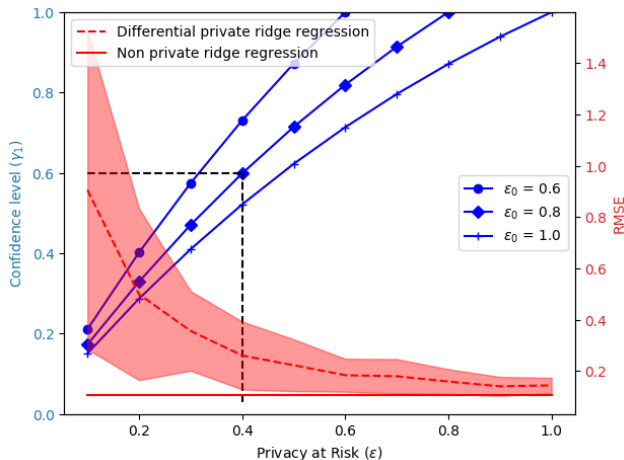| The source of randomness | Analytical result | Contribution |
|---|---|---|
| Laplace distribution | Closed form solution | Overlap computation under the sensitivity constraint. |
| Data-generation distribution | Upper bound on the confidence level | Sensitivity estimation using data-generation distribution. |
| Laplace distribution and data-generation distribution | Upper bound on the confidence level | Overlap computation under the estimated sensitivity. |

# Privacy-utility tradeoff



Figure: Utility, measured by RMSE (right y-axis), and privacy at risk for selected Laplace mechanism (left y-axis) for varying confidence levels

## Cost model

### Problem

Differential privacy guarantee is too abstract to be actionable.

### Art. 82 GDPR
# Right to compensation and liability

(1)   Any person who has suffered material or non-material damage as a result of an
      infringement of this Regulation shall have the right to receive compensation from the
      controller or processor for the damage suffered.

We assume that the compensation budget secured by a GDPR compliant
business entity is commensurate to the differential privacy guarantee
provided by the business entity while processing the data.

## Cost model

Let, $E$ and $E_\epsilon^{dp}$ be compensation budgets per stakeholder in absence of privacy measures and under $\epsilon$-differential privacy guarantee respectively.

### Properties of a cost model

- For all $\epsilon \in \mathbb{R}^{\geq 0}$, $E_\epsilon^{dp} \leq E$.
- As $\epsilon \to 0$, $E_\epsilon^{dp} \to 0$.
- As $\epsilon \to \infty$, $E_\epsilon^{dp} \to E$.
- $E_\epsilon^{dp}$ is a monotonically increasing function of $\epsilon$.

## Cost model

Cost model for $\epsilon$-differential privacy

$$E_\epsilon^{dp} \triangleq E e^{-\frac{c}{\epsilon}}$$

Let $E_{\epsilon_0}^{par}$ be compensation budget per stakeholder when a business entity uses $\epsilon_0$-differentially private privacy preserving mechanism that satisfies $(\epsilon, \gamma)$-privacy at risk.

Cost model for $(\epsilon, \gamma)$-privacy at risk

$$E_{\epsilon_0}^{par}(\epsilon, \gamma) \triangleq \gamma E_\epsilon^{dp} + (1 - \gamma) E_{\epsilon_0}^{dp}$$

# Cost model

Let $E_{\epsilon_0}^{par}$ be compensation budget per stakeholder when a business entity uses $\epsilon_0$-differentially private privacy preserving mechanism that satisfies $(\epsilon, \gamma)$-privacy at risk.

### Cost model for $(\epsilon, \gamma)$-privacy at risk
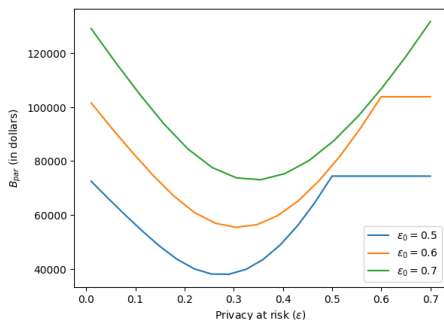
$$E_{\epsilon_0}^{par}(\epsilon, \gamma) \triangleq \gamma E_{\epsilon}^{dp} + (1-\gamma)E_{\epsilon_0}^{dp} \quad \leftarrow \text{Convex function!}$$

There exists a smallest value of $\epsilon$ for a specified $\epsilon_0$-differentially private mechanism that yields the smallest compensation budget!
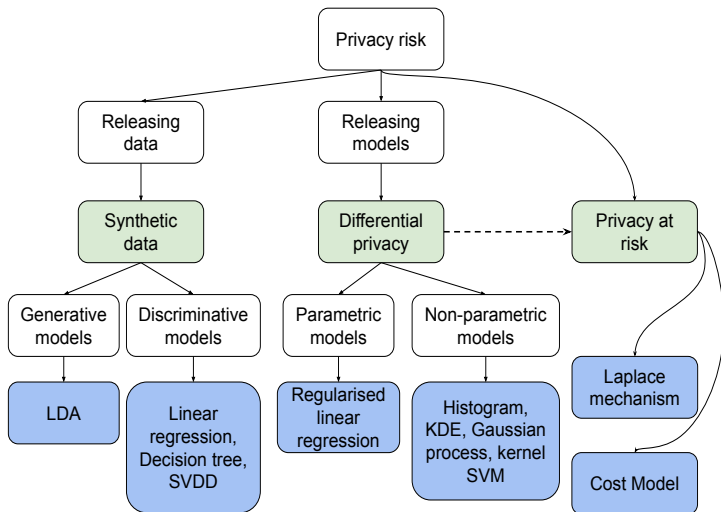
## Illustration

Consider, a case of obesity related data breach in an organisation. Research shows that the incremental cost in the premiums for health insurances with morbid obesity is \$5500 on an average [Moriarty et al., 2012].



$$E_{0.5}^{dp} = \$74434.40$$
$$E_{0.5}^{par}(0.29, 0.64) = \$37805.86$$
$$\text{Savings} = \$36628.54$$

# Conclusion

## Future directions

- Privacy in distributed machine learning
    - Differentially private federated learning (*SAP labs, 2017*)
    - Towards federated learning at scale: System design (*Google, 2018*)

- Privacy in law studies
    - Privacy and synthetic datasets (*Stanford Tech. Law Review, 2018*)
    - Synthetic data, privacy, and the law (Science, 2019)

## Publications

### Publication of data

- Ashish Dandekar, Remmy A. M. Zen, and Stéphane Bressan. A comparative study of synthetic dataset generation techniques. In *DEXA 2018, Proceedings, Part II*, pages 387-395

- Ashish Dandekar, Remmy A. M. Zen, and Stéphane Bressan. Comparative evaluation of data generation methods. In *Deep Learning Security Workshop, Singapore, December 2017. (Poster)*

- Ashish Dandekar, Stéphane Bressan, Talel Abdessalem, Huayu Wu, Wee Siong Ng. Detecting communities of commuters: graph based techniques versus generative models. In *CoopIS 2016, Proceedings*, pages 482-502

- Ashish Dandekar, Stéphane Bressan, Talel Abdessalem, Huayu Wu, Wee Siong Ng. Trajectory simulation in communities of commuters. *ICACSIS 2016, Proceedings*, pages 39-42 *(Invited paper)*

## Publications

### Publication of models

- Ashish Dandekar, Debabrota Basu, and Stéphane Bressan. Differential privacy for regularised linear regression. In *DEXA 2018, Proceedings, Part II, pages 483-491*
- Ashish Dandekar, Debabrota Basu, Thomas Kister, Geong Sen Poh, Jia Xu, and Stéphane Bressan. Privacy as a service. *DASFAA 2019 (Demo Paper)*
- Ashish Dandekar, Debabrota Basu, and Stéphane Bressan. Evaluation of differentially private non-parametric models as a service. *DEXA 2019 (Under review)*

### Privacy at risk

- Ashish Dandekar, Debabrota Basu, and Stéphane Bressan. Differential privacy at risk. Submitted in *Journal of Privacy and Confidentiality (Under review)*

### Miscellaneous

- Ashish Dandekar, Remmy A. M. Zen, and Stéphane Bressan. Generating fake but realistic headlines using deep neural networks. In *DEXA 2017,Proceedings, Part II, pages 427-440*

National University of Singapore

Thank you!

# References I

Bellovin, S. M., Dutta, P. K., and Reitinger, N. (2018).
Privacy and synthetic datasets.
*Stanford Technology Law Review, Forthcoming.*

Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003).
Latent dirichlet allocation.
*Journal of machine Learning research,* 3(Jan):993–1022.

Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. (2011).
Differentially private empirical risk minimization.
*Journal of Machine Learning Research,* 12(Mar):1069–1109.

Drechsler, J. and Reiter, J. P. (2011).
An empirical evaluation of easily implemented, nonparametric methods for generating synthetic datasets.
*Computational Statistics & Data Analysis,* 55(12):3232–3243.

Dwork, C. (2006).
Differential privacy.
In *33rd International Colloquium on Automata, Languages and Programming, part II (ICALP 2006),* volume 4052, pages 1–12, Venice, Italy. Springer Verlag.

Dwork, C., Roth, A., et al. (2014).
The algorithmic foundations of differential privacy.
*Foundations and Trends® in Theoretical Computer Science,* 9(3–4):211–407.

# References II

Fredrikson, M., Lantz, E., Jha, S., Lin, S., Page, D., and Ristenpart, T. (2014).
Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing.
In *Proceedings of the... USENIX Security Symposium. UNIX Security Symposium*, volume 2014, pages 17–32. NIH Public Access.

Hall, R., Rinaldo, A., and Wasserman, L. (2012).
Random differential privacy.
*Journal of Privacy and Confidentiality*, 4(2):43–59.

Hall, R., Rinaldo, A., and Wasserman, L. (2013).
Differential privacy for functions and functional data.
*Journal of Machine Learning Research*, 14(Feb):703–727.

Homer, N., Szelinger, S., Redman, M., Duggan, D., Tembe, W., Muehling, J., Pearson, J. V., Stephan, D. A., Nelson, S. F., and Craig, D. W. (2008).
Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays.
*PLoS genetics*, 4(8):e1000167.

Jorion, P. (2000).
Value at risk: The new benchmark for managing financial risk.

Kifer, D., Smith, A., and Thakurta, A. (2012).
Private convex empirical risk minimization and high-dimensional regression.
In *Conference on Learning Theory*, pages 25–1.

# References III

Machanavajjhala, A., Kifer, D., Abowd, J., Gehrke, J., and Vilhuber, L. (2008).
Privacy: Theory meets practice on the map.
In *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*, pages 277–286. IEEE.

Moriarty, J. P., Branda, M. E., Olsen, K. D., Shah, N. D., Borah, B. J., Wagie, A. E., Egginton, J. S., and Naessens, J. M. (2012).
The effects of incremental costs of smoking and obesity on health care costs among adults: a 7-year longitudinal study.
*Journal of Occupational and Environmental Medicine*, 54(3):286–291.

Reiter, J. P. (2003).
Inference for partially synthetic, public use microdata sets.
*Survey Methodology*, 29(2):181–188.

Rubin, D. B. (1993).
Discussion statistical disclosure limitation.
*Journal of official Statistics*, 9(2):461.

Shokri, R., Stronati, M., Song, C., and Shmatikov, V. (2017).
Membership inference attacks against machine learning models.
In *Security and Privacy (SP), 2017 IEEE Symposium on*, pages 3–18. IEEE.

Smith, M. T., Zwiessele, M., and Lawrence, N. D. (2016).
Differentially private gaussian processes.
*arXiv preprint arXiv:1606.00720*.

# References IV

Zhang, J., Zhang, Z., Xiao, X., Yang, Y., and Winslett, M. (2012).
Functional mechanism: regression analysis under differential privacy.
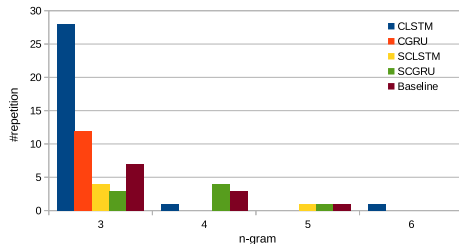*Proceedings of the VLDB Endowment*, 5(11):1364–1375.

## Datasets

**Census Dataset.**

| Attribute Name | Variable Type |
|---|---|
| House Type | Categorical |
| Family Size | Ordinal |
| Sex | Categorical |
| Age | Ordinal |
| Marital Status | Categorical |
| Race | Categorical |
| Educational Status | Categorical |
| Employment Status | Categorical |
| Income | Ordinal |
| Birth Place | Categorical |

- 1% random sample from US 2001 Census Data[1]

- Survey data of 316,277 heads of households

- We synthetically generate values for *Age* and *Income*

---

[1]https://usa.ipums.org/usa/
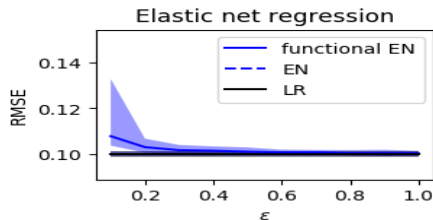
# News results



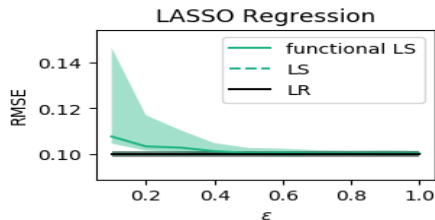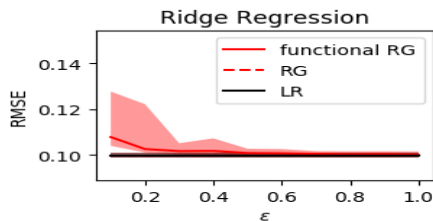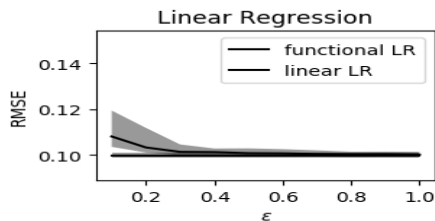*CLSTM tends to generate headlines with long repetitions*



*SCLSTM tends to generate novel headlines on an average*

# Experiments

# Privacy at risk

### Implicit randomness (Thm 3.2)

The confidence level $\gamma_1 \in [0, 1]$ of achieving a privacy at risk level $\epsilon \geq 0$ by a Laplace Mechanism $\mathcal{L}_{\epsilon_0}^{\Delta_f}$ for a query $f : \mathcal{D} \to \mathbb{R}^k$ is given by

$$\gamma_1 = \frac{\mathbb{P}(T \leq \epsilon)}{\mathbb{P}(T \leq \epsilon_0)},$$

where $T$ is a random variable dependent on the Laplace noise $\mathrm{Lap}(\frac{\Delta_f}{\epsilon_0})$, and follows the $\mathrm{BesselK}\left(k, \frac{\Delta_f}{\epsilon_0}\right)$ distribution.

## Privacy at risk

### Explicit randomness (Thm 3.13)

Analytical bound on the confidence level for empirical privacy at risk, $\hat{\gamma}_2$, for Laplace mechanism $\mathcal{L}_\epsilon^{\Delta_{S_f}}$ with privacy at risk level $\epsilon$ and sampled sensitivity $\Delta_{S_f}$ for a query $f : \mathcal{D} \to \mathbb{R}^k$ is

$$\hat{\gamma}_2 \geq \gamma_2(1 - 2e^{-2\rho^2 n})$$

where $n$ is the number of samples used for estimation of the sampled sensitivity and $\rho$ is the accuracy parameter. $\gamma_2$ denotes the confidence level for the privacy at risk.

## Privacy at risk

### Coupled effect (Lemma 3.19)

For Laplace Mechanism $\mathcal{L}_{\epsilon_0}^{\Delta_{S_f}}$ with sampled sensitivity $\Delta_{S_f}$ of a query $f : \mathcal{D} \to \mathbb{R}^k$ and for any $Z \subseteq Range(\mathcal{L}_{\epsilon}^{\Delta_{S_f}})$,

$$\hat{\gamma_3} \geq \frac{\mathbb{P}(T \leq \epsilon)}{\mathbb{P}(T \leq \eta\epsilon_0)} \gamma_2 (1 - 2e^{-2\rho^2 n})$$

where $n$ is the number of samples used to find sampled sensitivity, $\rho \in [0, 1]$ is a accuracy parameter and $\eta = \frac{\Delta_f}{\Delta_{S_f}}$.
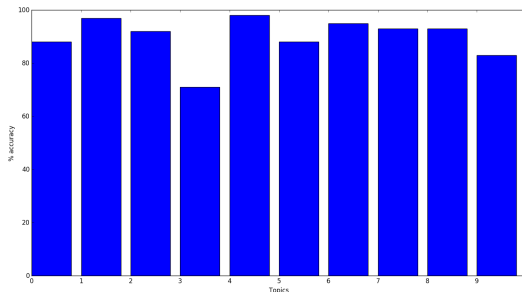
## Privacy-Utility tradeoff with cost model

**Mean absolute error**

$$\mathbb{E}\left[|\mathcal{L}_\epsilon^1(x) - f(x)|\right] = \frac{1}{\epsilon}$$

If we have a maximum permissible expected mean absolute error $T$, Equation 1 illustrates the upper and lower bounds that dictate the permissible range of $\epsilon$ that a data publisher can promise depending on the budget and the permissible error constraints.

$$\frac{1}{T} \le \epsilon \le \left[\ln\left(\frac{\gamma E}{B - (1-\gamma)E_{\epsilon_0}^{dp}}\right)\right]^{-1} \tag{1}$$

# Effectiveness of communities



### How effective the synthetic datasets are?

For every community, we generate travelling records of 1000 commuters. We use pre-trained generative models to classify these commuters into the communities. We use classification accuracy as the metric of effectiveness of the synthetically generated commuter.