Machine Learning with Differential Privacy

Elliot Creager May 2019 AMMI @ AIMS Kigali

Learning Objectives

We want you to understand the following:

- DP+ML: Why?
- Basic tradeoffs and design choices at play
- DP+ML: How? (applying basic mechanisms)
- The Gaussian Mechanism and its applications:
 - Stochastic Gradient Descent
- Composition: how to analyze a sequence of mechanisms
- Private neural nets:
 - Implementation difficulties
 - Recent advances

Motivating Differential Privacy

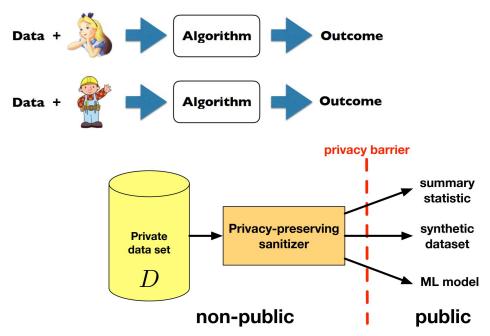
Anonymizing data is unsafe

- Netflix challenge
- Narayanan & Shmatikov 2008 showed a handful of Netflix recommendation plus IMDB info reveals dataset membership and other ratings.

Publishing summary stats is unsafe

• Given an individual's gene sequence, querying a genome-wide association study reveals whether the individual is in the "case group" of the study [Dwork & Smith 2014]

What is Differential Privacy?



An individual's participation in a dataset/study doesn't change the outcome

Privacy ensured through calibrated randomness

Adversary

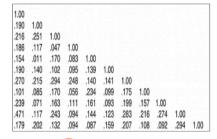


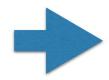
Prior Knowledge:

A's Genetic profile

A smokes

Case I: Study





A has cancer

Cancer

[Study violates A's privacy]

Case 2: Study



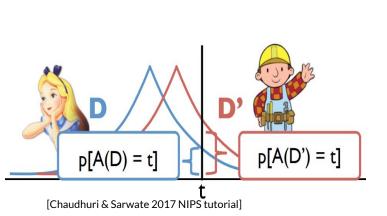


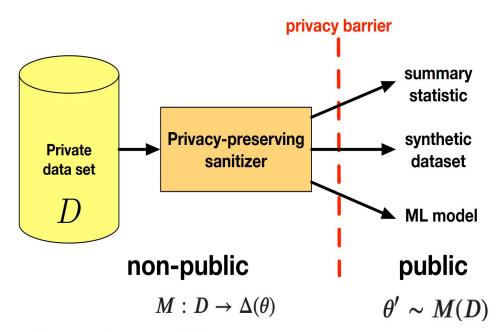
A probably has cancer

Smoking causes cancer

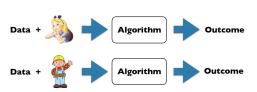
[Study does not violate privacy]

The Setting





note: θ' is a specific outcome; $\Delta(\theta)$ denotes space of distributions over possible outcomes

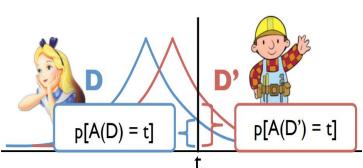


The Definition

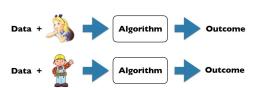
Algorithm A is (ϵ, δ) -differentially private iff \forall "neighboring" datasets D, D' and \forall outcomes $t \in \mathrm{Range}(A)$:

$$\mathbb{P}\left[A(D) \in t
ight] \leq e^{\epsilon} \mathbb{P}\left[A(D') \in t
ight] + \delta$$

In other words:



$$\left| \underbrace{\log rac{\mathbb{P}[A(D) \in t]}{\mathbb{P}[A(D') \in t]}}_{ ext{privacy loss}}
ight| \leq \epsilon \quad ext{with probability} \quad 1 - \delta$$



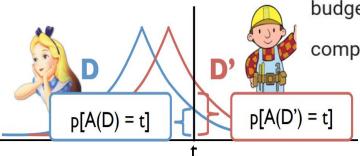
Properties of Differential Privacy

- Invariance to post-processing
 - Computation on algorithm outputs that don't touch data are "free"
- Composability

If algorithms $\{A_1,A_2,\ldots,A_N\}$ are DP with respective privacy

budgets $\{(\epsilon_1, \delta_1), (\epsilon_2, \delta_2), \dots, (\epsilon_N, \delta_N)\}$, then their

composition is $(\sum_{i=1}^N \epsilon_i, \sum_{i=1}^N \delta_i)$ -DP.



Recall Some Basic Mechanisms

Laplace Mechanism - add Laplacian noise to fn output

Gaussian Mechanism - add Gaussian noise to fn output

Exponential Mechanism - sample fn output prop. to score

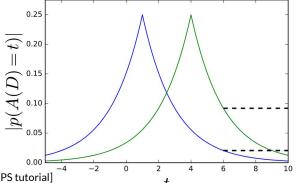
The Laplace Mechanism

Global Sensitivity of f is
$$S(f) = \max_{dist(D, D') = 1} |f(D) - f(D')|$$

Output f(D) + Z, where

$$Z \sim \frac{S(f)}{\epsilon} \text{Lap}(0,1)$$

 ϵ -differentially private



Laplace distribution:

$$p(z|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|z-\mu|}{b}\right)$$

The Laplace Mechanism

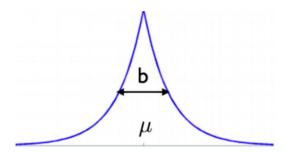
Example: compute private sample mean from n samples

f(D) = Mean(D), where each record is a scalar in [0,1]

Global Sensitivity of f = I/n

Laplace Mechanism:

Output
$$f(D) + Z$$
, where $Z \sim \frac{1}{n\epsilon} \mathrm{Lap}(0,1)$

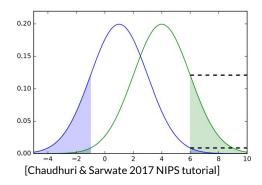


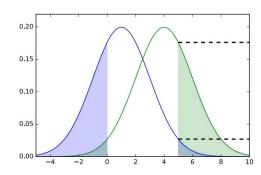
The Gaussian Mechanism

Global Sensitivity of f is S(f) =
$$\max_{D, D' \text{ "neighbouring"}} ||f(D) - f(D')||_2$$

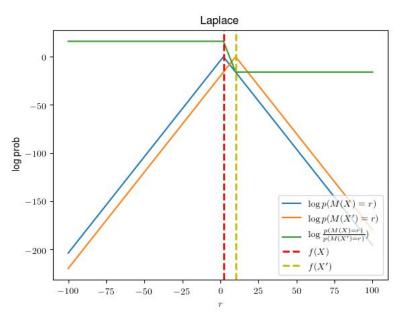
Output
$$f(D) + Z$$
, where

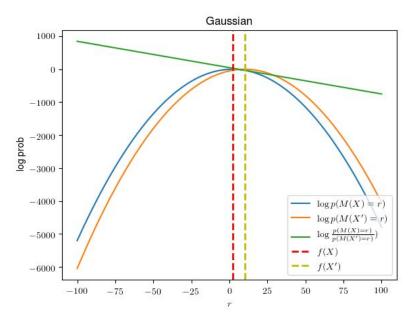
$$Z \sim \frac{S(f)}{\epsilon} \mathcal{N}(0, 2\ln(1.25/\delta)) \qquad (\epsilon, \delta) \text{-differentially}$$
 private





Why Approx. DP for GM?





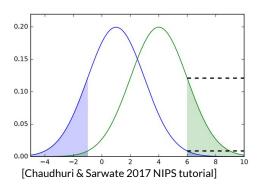
GM cannot achieve (eps, 0) DP...

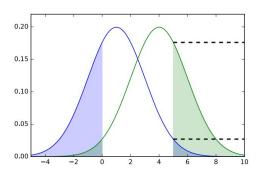
The Gaussian Mechanism

Global Sensitivity of f is S(f) =
$$\max_{D, D' \text{ "neighbouring"}} |f(D) - f(D')|_2$$

See Dwork and Roth Appendix A for proof that GM is private...

Output
$$f(D)+Z$$
, where
$$Z\sim \frac{S(f)}{\epsilon}\mathcal{N}(0,2\ln(1.25/\delta)) \qquad (\epsilon,\delta)\text{-differentially}$$
 private





The Exponential Mechanism

$$\Delta u = \max_{r \in R} \max_{x,y:||x-y| \le 1} |u(x,r) - u(y,r)|$$

assuming x and y have same normalization term, we have

$$\log \frac{\exp(\frac{\epsilon u(x,r)}{2\Delta u})}{\exp(\frac{\epsilon u(y,r)}{2\Delta u})} = \epsilon(\frac{u(x,r) - u(y,r)}{2\Delta u}) \le \frac{\epsilon}{2}$$

The actual mechanism saves have the privacy budget to account for changes in the normalization term

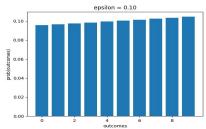
$$\mathbb{P}(M_{\epsilon}(x) = r) = \frac{\exp(\frac{\epsilon u(x,r)}{2\Delta u})}{\sum_{r' \in R} \exp(\frac{\epsilon u(x,r')}{2\Delta u})}$$

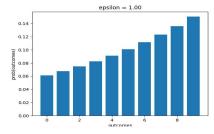
cf. Dwork and Roth section 3.4

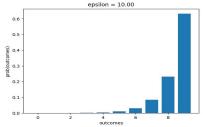
The Exponential Mechanism

$$\mathbb{P}(M_{\epsilon}(x) = r) \propto \exp(\underbrace{\epsilon}_{\text{"temperature"}} \cdot \underbrace{\frac{u(x, r)}{2\Delta u}}_{\text{"logits"}})$$

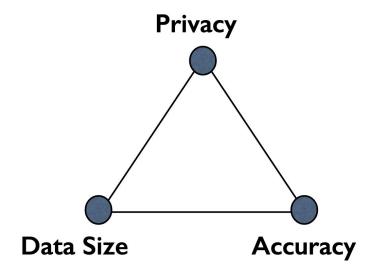
```
1 import matplotlib
 2 matplotlib.use('Agg')
 3 import matplotlib.pyplot as plt
 5 import torch
 6 from torch.distributions.categorical import Categorical
 8 scores = torch.arange(10.)
 9 sensitivity = torch.tensor([5.])
10
11 for eps in (0.1, 1., 10.):
12
      eps = torch.tensor([eps])
13
      mech_distn = Categorical(logits=(eps * scores / (2. * sensitivity)))
14
      f, a = plt.subplots()
15
      a.bar(torch.arange(10), mech_distn.probs)
16
      a.set_title('epsilon = {:.2f}'.format(eps.item()))
17
      a.set_xlabel('outcomes')
18
      a.set_ylabel('prob(outcomes)')
19
       f.savefig('./probs{:.2f}.png'.format(eps.item()))
```



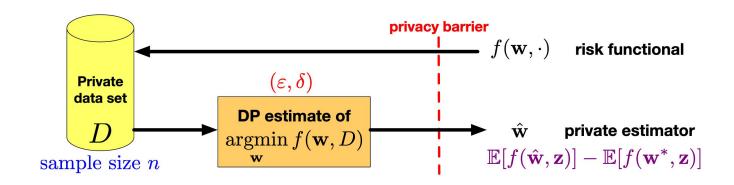




Tradeoffs in DP+ML



Tradeoffs in DP+ML



Statistical estimation: estimate a parameter or predictor using private data that has good expected performance on future data.

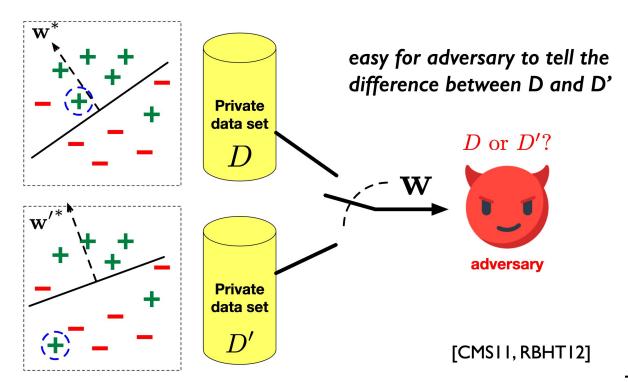
Goal: Good privacy-accuracy-sample size tradeoff

Private Empirical Risk Min.

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \ell(\mathbf{w}, (\mathbf{x}_i, y_i)) + \lambda R(\mathbf{w})$$

- Empirical Risk Minimization (ERM) is a common paradigm for prediction problems.
- Produces a predictor w for a label/response y given a vector of features/covariates x.
- Typically use a convex loss function and regularizer to "prevent overfitting."

Private ERM

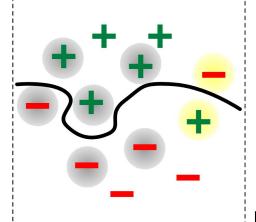


Kernel Approaches Even Worse

• Kernel-based methods produce a classifier that is a function of the data points.

• Even adversary with black-box access to w could potentially learn those points.

$$\mathbf{w}(\mathbf{x}) = \sum_{i=1}^{n} \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$



Privacy & Learning Compatible

- Good learning algorithms generalize to the population distribution, not individuals.
- Stable learning algorithms generalize [BE02].
- Differential privacy can be interpreted as a form of stability that also implies generalization [DFH+15,BNS+16].
- Two parts of the same story:
 Privacy implies generalization asymptotically.
 Tradeoffs between privacy-accuracy-sample size for finite n.

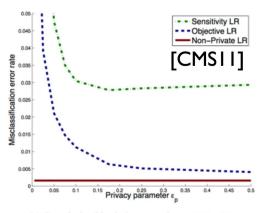
Revisiting ERM

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}, (\mathbf{x}_i, y_i)) + \lambda R(\mathbf{w})$$

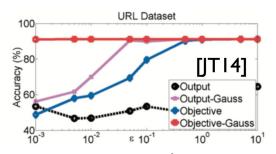
- Learning using (convex) optimization uses three steps:
 - I. read in the data

- input perturbation
- 2. form the objective function
- objective perturbation
- We will discuss this one --> 3. perform the minimization
- output perturbation
- We can try to introduce privacy in each step!

Typical Empirical Results



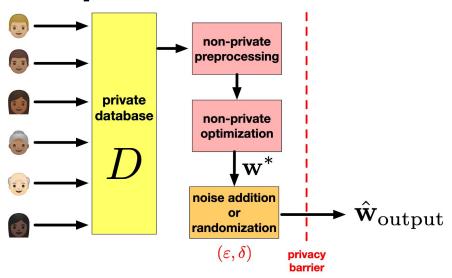
(a) Regularized logistic regression, KDDCup99



In general:

- Objective perturbation empirically outperforms output perturbation.
- Gaussian mechanism with (ε, δ)
 guarantees outperform Laplace like mechanisms with ε guarantees.
- Loss vs. non-private methods is very dataset-dependent.

Output Perturbation



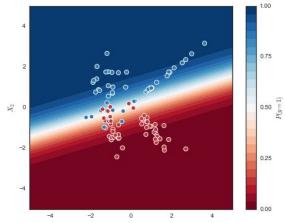
- Compute the minimizer and add noise.
- Does not require re-engineering baseline algorithms

Noise depends on the sensitivity of the argmin.

[CMS11, RBHT12]

$$L(w) = \sum_{i=1}^{N} \text{CrossEntropy}(\sigma(w^{T} x_{i}), y_{i}) + \frac{1}{2} \lambda w^{T} w$$

What is $\Delta L(w)$, the sensitivity of the loss?



https://stackoverflow.com/questions/28256058/plotting-decision-boundary-of-logistic-regression

$$L(w) = \sum_{i=1}^{N} \underbrace{\text{CrossEntropy}(\sigma(w^{T}x_{i}), y_{i})}_{\text{convex}} + \underbrace{\frac{1}{2}\lambda w^{T}w}_{\lambda-\text{strongly convex}}$$

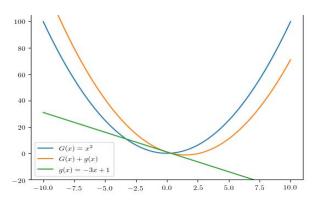
h(w) is convex iff $\nabla^2 h(w) \ge 0$ (its Hessian is positive semi-definite)

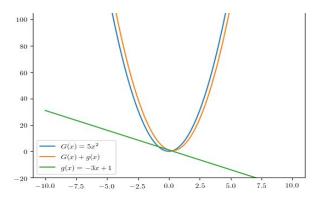
$$h(w)$$
 is λ -strongly convex iff $(\nabla h(x) - \nabla h(y))^T(x - y) \ge \lambda ||x - y||_2^2$

The sum of a convex function and a λ -strongly convex function is λ -strongly convex.

Lemma 7 Let $G(\mathbf{f})$ and $g(\mathbf{f})$ be two vector-valued functions, which are continuous, and differentiable at all points. Moreover, let $G(\mathbf{f})$ and $G(\mathbf{f}) + g(\mathbf{f})$ be λ -strongly convex. If $\mathbf{f}_1 = \operatorname{argmin}_{\mathbf{f}} G(\mathbf{f})$ and $\mathbf{f}_2 = \operatorname{argmin}_{\mathbf{f}} G(\mathbf{f}) + g(\mathbf{f})$, then

$$\|\mathbf{f}_1 - \mathbf{f}_2\| \leq \frac{1}{\lambda} \max_{\mathbf{f}} \|\nabla g(\mathbf{f})\|.$$





$$h(w)$$
 is λ -strongly convex iff $(\nabla h(x) - \nabla h(y))^T (x - y) \ge \lambda ||x - y||_2^2$

The sum of a convex function and a λ -strongly convex function is λ -strongly convex.

$$L(w) = \frac{1}{N} \sum_{i \in \{1..N\}} \text{CrossEntropy}(\sigma(w^T x_i), y_i) + \frac{1}{2} \lambda w^T w$$

$$L'(w) = \left(\frac{1}{N} \sum_{i \in \{1..N\} \setminus j} \text{CrossEntropy}(\sigma(w^T x_i), y_i)\right) + \text{CrossEntropy}(\sigma(w^T x_{j'}), y_{j'}) + \frac{1}{2} \lambda w^T w$$

$$L'(w) = L(w) + g(w)$$

$$g(w) = \frac{1}{N} \left(\text{CrossEntropy}(\sigma(w^T x_{j'}), y_{j'}) - \text{CrossEntropy}(\sigma(w^T x_{j}), y_{j}) \right)$$

Assuming $||x|| \le 1$, we have $\nabla g(w) \le \frac{2}{N}$

$$\rightarrow \Delta L = \frac{2}{N\lambda}$$

Algorithm:

sensitivity $\Delta L = \frac{2}{N}$

- 1) solve $w^* = \arg\min_w L(w)$
- 2) draw η with $p(\eta) \propto \exp(-\frac{\Delta L}{\epsilon}||\eta||)$ (vector analogue of Laplace draw)
- 3) return $w = w^* + \eta$