

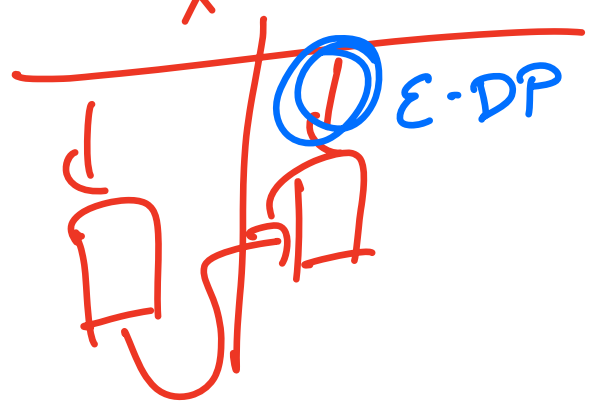
CSCI 699: Privacy Preserving Machine Learning - Week 5

Gaussian DP and Privacy Auditing

Sai Praneeth Karimireddy, Sep 27 2024



Recap



- Composition: simple $k\varepsilon$ -DP

ε -DP

Theorem. Advanced Composition

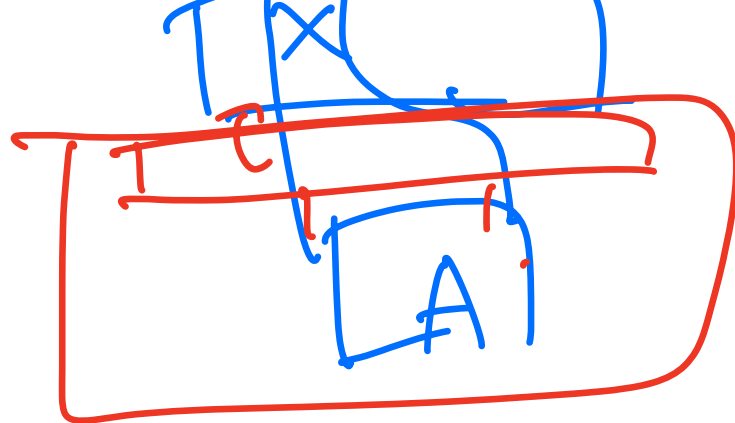
A combination of $A_1 \circ A_2 \circ \dots \circ A_k$, each of which is (ε, δ) -DP is $(\tilde{\varepsilon}, \tilde{\delta})$ -DP where

$$\tilde{\varepsilon} = \varepsilon \sqrt{2k \ln(1/\delta')} + k \frac{e^\varepsilon - 1}{e^\varepsilon + 1} \quad \text{and} \quad \tilde{\delta} = k\delta + \delta'$$

For any choice of δ' .



Recap



- Subsampling amplification

Theorem. Subsampling Amplification

Composing an (ϵ, δ) -DP A with a sampling rate of q results in an $(\tilde{\epsilon}, \tilde{\delta})$ -DP algorithm where

$$\tilde{\epsilon} = \log(1 - q + qe^\epsilon) = O(q\epsilon) \quad \text{and} \quad \tilde{\delta} = q\delta$$

Recap

- Private SGD with clipping L1 norm:

- $\theta_t = \theta_{t-1} - \gamma \text{Clip}_\tau (\nabla_\theta \ell(f(x_t; \theta), y_t)) + \text{Lap}(2\tau/\epsilon)$

- With $q = 1/n$, k rounds satisfies $(O(\epsilon/n\sqrt{k \ln(1/\delta)}), \delta)$ -DP for any $\delta > 0$.

- Can also clip L2 norm and use Gaussian mechanism.

- Q: what did you observe empirically L1 vs. L2?

$\forall \delta$

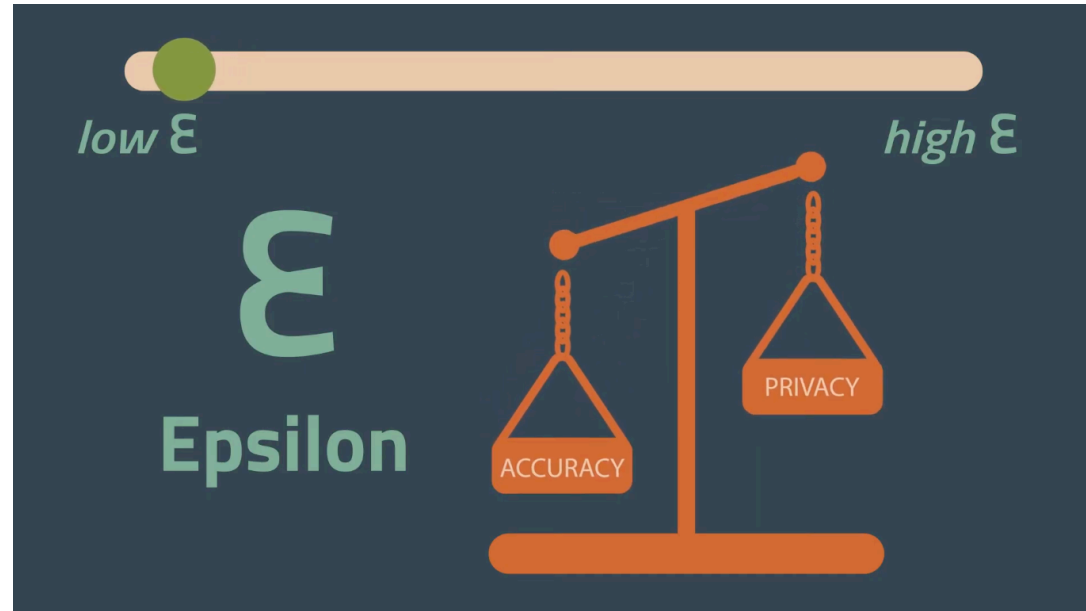
$$\mathcal{N}\left(0, \left(\frac{C}{\epsilon}\right) \log \frac{1}{\delta}\right)^2$$

Agenda for today

Analyzing privacy of ML training

- Gaussian DP
- Privacy Auditing
- Presentations + discussions
- Auditing Practical - next week (needs HW2 soln)

Gaussian Differential Privacy



Drawbacks of Approximate DP

- After k steps of Lap-SGD, we were able to show $(\epsilon\sqrt{2k \ln(1/\delta)}, \delta)$ -DP
- But advanced composition is too loose.

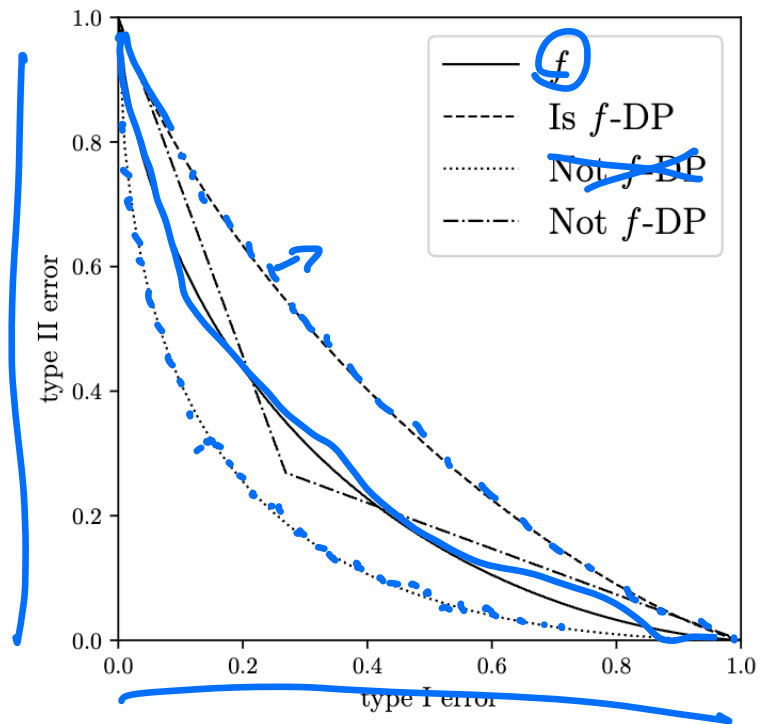
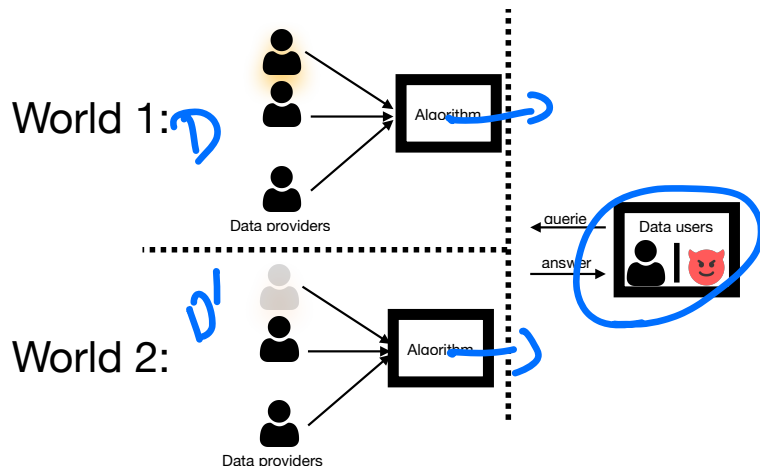
$$+ \frac{k(e^\epsilon - 1)}{e^\epsilon}$$

- Renyi DP
- Concentrated DP

f-DP

Most general privacy definition

- **Definition.** Given a function f , we say an algorithm is f -DP if the tradeoff curve of an optimal distinguisher is strictly above f .



f-DP

Generalization (ϵ, δ) -DP

- **Prop 2.5 [WZ10]**. A is (ϵ, δ) -DP iff it satisfies $f_{\epsilon, \delta}$ -DP for

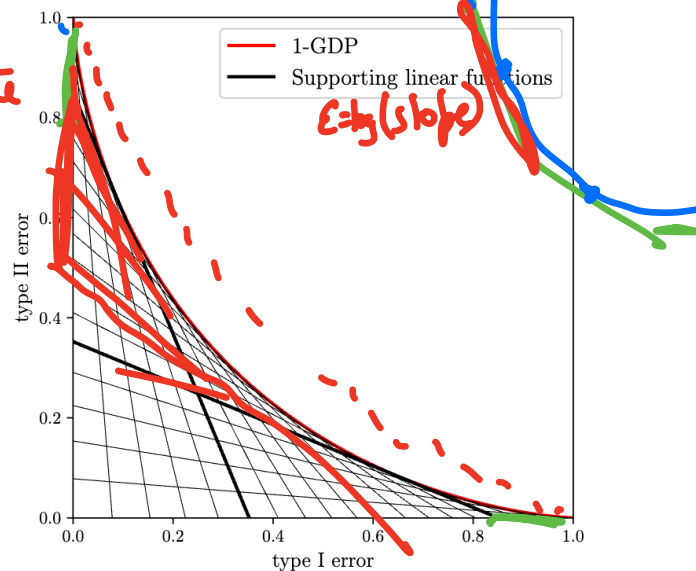
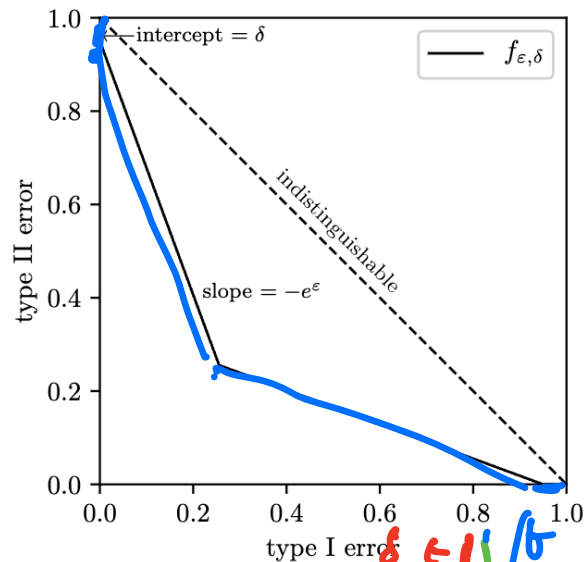
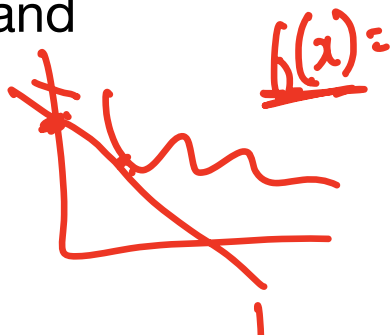
$$f_{\epsilon, \delta} = \max(1 - \delta - e^\epsilon x, (1 - \delta - x)/e^\epsilon)$$

Primal - Dual / convex conjugate

- **Prop 2.12 [DRS19]** A is f -DP iff it satisfies $(\epsilon, \delta_f(\epsilon))$ -DP for $\forall \epsilon \geq 0$ and

$$\delta_f(\epsilon) = 1 + f^*(-e^\epsilon).$$

→ dual



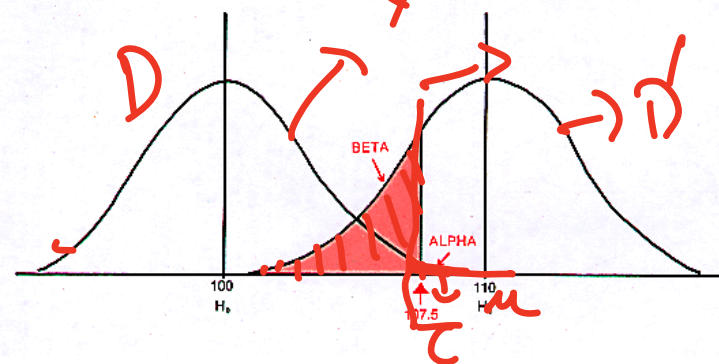
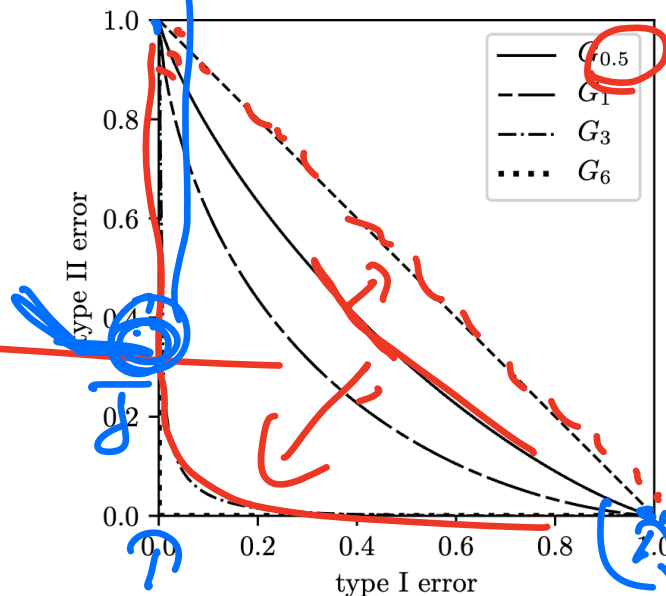
Gaussian-DP

- **Definition.** A is μ -GDP if it satisfies f_μ -DP for $f_\mu = T(\mathcal{N}(0,1), \mathcal{N}(\mu,1))$

- $$\frac{\Pr[A(D) = t]}{\Pr[A(D') = t]} \leq \frac{\Pr[\mathcal{N}(0,1) = t]}{\Pr[\mathcal{N}(\mu,1) = t]} = \exp\left(\frac{1}{2}(\mu^2 - 2\mu t)\right)$$

- $\alpha(\tau) = 1 - \Phi(\tau)$ and $\beta(\tau) = \Phi(\tau - \mu)$

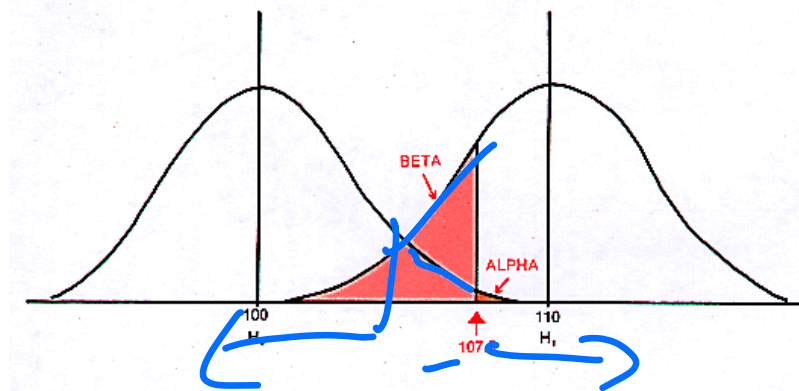
$$\beta(\alpha) = \Phi(\Phi^{-1}(1-\alpha) - \mu)$$



Gaussian-DP

Gaussian mechanism

- **Definition.** A is μ -GDP if it satisfies f_μ -DP for $f_\mu = T(\mathcal{N}(0,1), \mathcal{N}(\mu,1))$



Theorem. Gaussian mechanism

Given $f: \mathcal{X}^n \rightarrow \mathbb{R}^d$ with Δ bounded ℓ_2 -sensitivity, $f(D) + \mathcal{N}\left(0, \frac{\Delta^2}{\mu^2} I_d\right)$ is μ -GDP.

Gaussian Differential Privacy

Tight composition

Theorem. GDP Composition

Composition of $A_1 \circ A_2 \dots \circ A_k$, each of which is μ_i

-GDP is $\sqrt{\sum_{i=1}^k \mu_i^2}$ -GDP.

$\mu \sqrt{k}$ - GDP

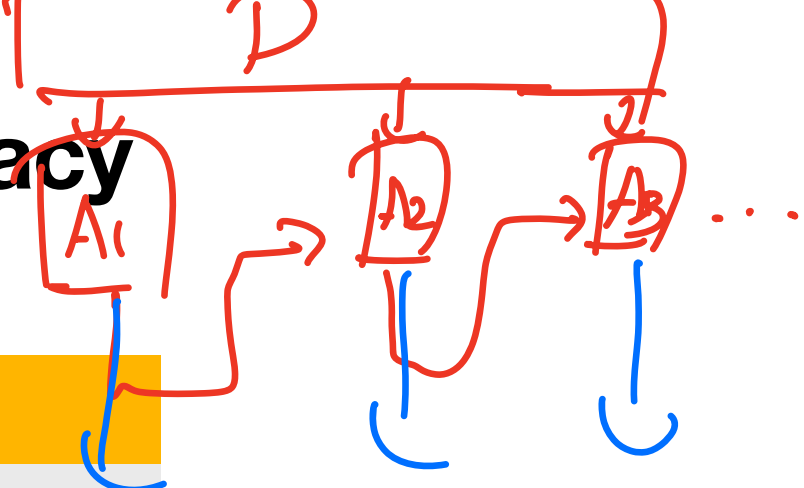
$A_1 \rightarrow \mu_1$ -GDP

$A_2 \rightarrow \mu_2$ -GDP

$(A_1(D), A_2(D))$ vs.

$(A_1(D'), A_2(D'))$

$P_2 \{A_1(D) = t_1, A_2(D) = t_2\}$



Gaussian Differential Privacy

Tight composition

$$\begin{aligned}
 & \overline{P_{\lambda}[A_1(D')=t_1, A_2(D)=t_2]} \\
 &= \overline{P_{\lambda}[A_1(D)=t_1]} \cdot \overline{P_{\lambda}[A_2(D)=t_2]} \\
 &= \overline{P_{\lambda}[A_1(D)=t_1]} \cdot \overline{P_{\lambda}[A_2(D')=t_2]}
 \end{aligned}$$

Theorem. GDP Composition

Composition of $A_1 \circ A_2 \dots \circ A_k$, each of which is μ_i

-GDP is $\sqrt{\sum_{i=1}^k \mu_i^2}$ -GDP.

$$\begin{aligned}
 & a_1 b_1 + a_2 b_2 \\
 & \leq \sqrt{a_1^2 + a_2^2} \cdot \sqrt{b_1^2 + b_2^2} \\
 & \exp\left(\sqrt{\mu_1^2 + \mu_2^2} - 2\sqrt{\mu_1^2 + \mu_2^2} \sqrt{t_1^2 + t_2^2}\right)
 \end{aligned}$$

$$\begin{aligned}
 & \leq \exp(\mu_1^2 - 2\mu_1 t_1) \cdot \exp(\mu_2^2 - 2\mu_2 t_2) \\
 & = \exp\left(\underbrace{(\mu_1^2 + \mu_2^2)}_{\sqrt{\mu_1^2 + \mu_2^2}^2} - 2(\mu_1 t_1 + \mu_2 t_2)\right) \\
 & \leq \exp\left(\sqrt{\mu_1^2 + \mu_2^2}^2 - 2\sqrt{\mu_1^2 + \mu_2^2} \sqrt{t_1^2 + t_2^2}\right)
 \end{aligned}$$

$$\exp\left(\sqrt{\mu_1^2 + \mu_2^2} - 2\sqrt{\mu_1^2 + \mu_2^2} \sqrt{t_1^2 + t_2^2}\right)$$

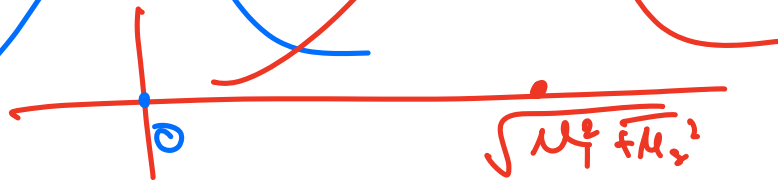
 X_2

$$vs. \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, I_2\right)$$

$$X_1 \sim \mathcal{N}(0, I_2)$$

$$P_{X_1}\left[X_1^T = \begin{bmatrix} t_1 \\ t_2 \end{bmatrix}\right]$$

$$P_{X_2}\left[X_2 = \begin{bmatrix} t_1 \\ t_2 \end{bmatrix}\right]$$



$$\Rightarrow \sqrt{\mu_1^2 + \mu_2^2} - \text{GDP} \quad \boxed{\text{E}}$$

Gaussian Differential Privacy

Canonical f

Theorem 3.4 [DRS19] Central limit theorem of composition

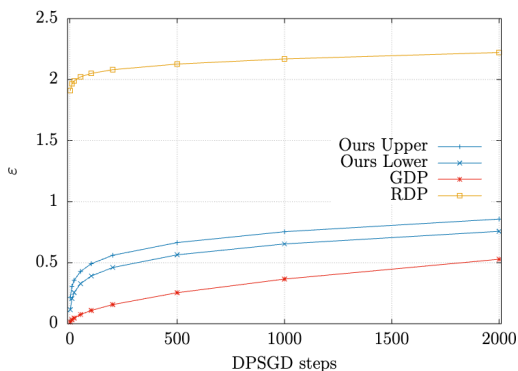
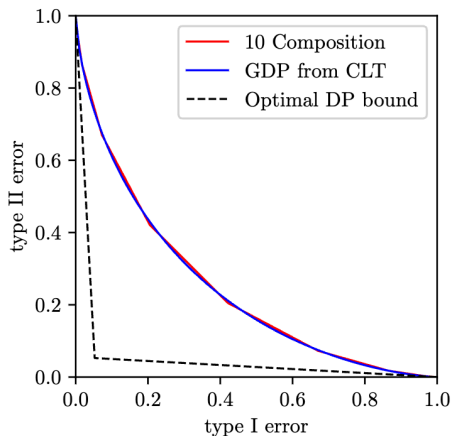
Given some regularity assumptions, composition of $A_1 \circ A_2 \dots \circ A_k$, each of which is f_i -DP is approximately μ -GDP for

$$\mu = \frac{2\sqrt{k}\kappa_1}{\kappa_1 - \kappa_2} \text{ for } \kappa_1 = -\int_0^1 \log |f'(x)| dx \text{ and } \kappa_2 = -\int_0^1 \log^2 |f'(x)| dx.$$

$$\sum_{i=1}^k R_i$$

Gaussian Differential Privacy

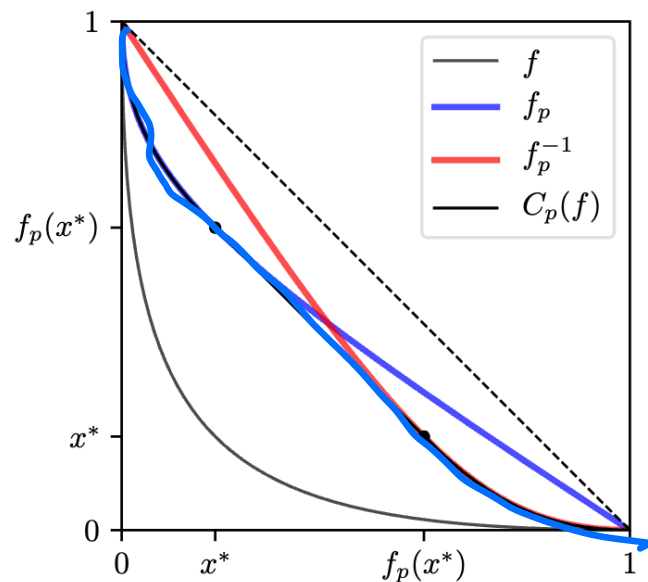
Canonical f



- In stats, combining many random variables \approx Gaussian by CLT. In DP, composing many DP steps \approx gDP.
- Caution: just like CLT sometimes fails, Thm 3.4 is sometimes fails and underestimates privacy [GLW21].

Gaussian Differential Privacy

Amplification by subsampling



- Define $f_q(x) = qf(x) + (1 - q)(1 - x)$ and f_q^{-1}
- **Theorem 4.2** [DRS19]
Composing q -sampling with f -DP, is $(\min(f_p, f_p^{-1}))^{**}$ -DP
- Unfortunately, no closed form for GDP, compute numerically.

Private SGD

Using Gaussian-DP

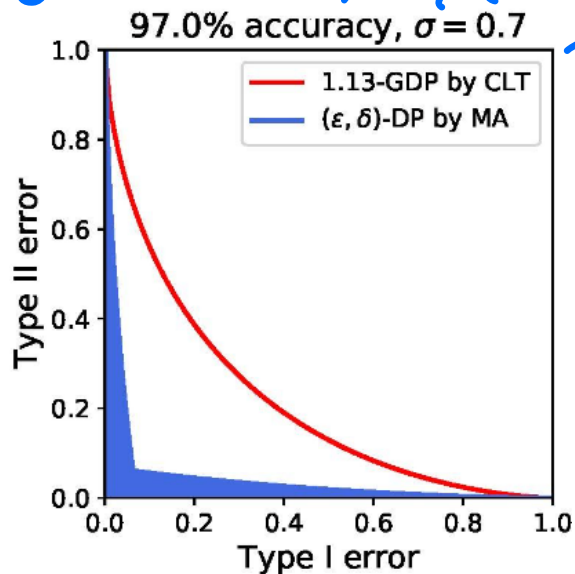
CLT $\| \cdot \|_2 - \bar{C}$

add $\mathcal{N}(0, \frac{\bar{C}^2}{\mu^2} I)$

Corollary 5.4 [DRS19] Subsampled Composition

Suppose each A_i is μ -GDP. Then, composing q -sampled A_i 's asymptotically

$$(q\sqrt{k}\sqrt{e^{\mu^2}\Phi(3\mu/2) + 3\Phi(-\mu/2) - 2})\text{-GDP.}$$



Tightest privacy bound [B+'20].
But, only asymptotically valid.

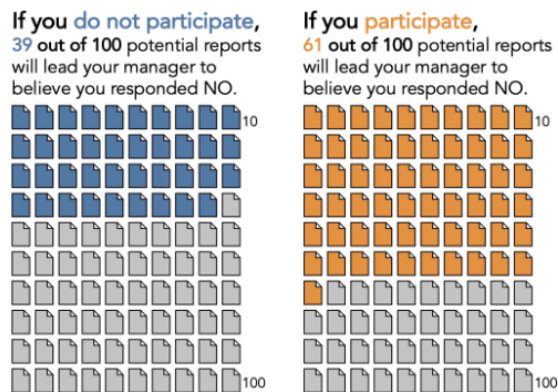
Aside: Communicating Privacy

Odds ratio

If you **do not participate**,
39 out of 100 potential reports will lead your manager to believe you responded NO.

If you **participate**,
61 out of 100 potential reports will lead your manager to believe you responded NO.

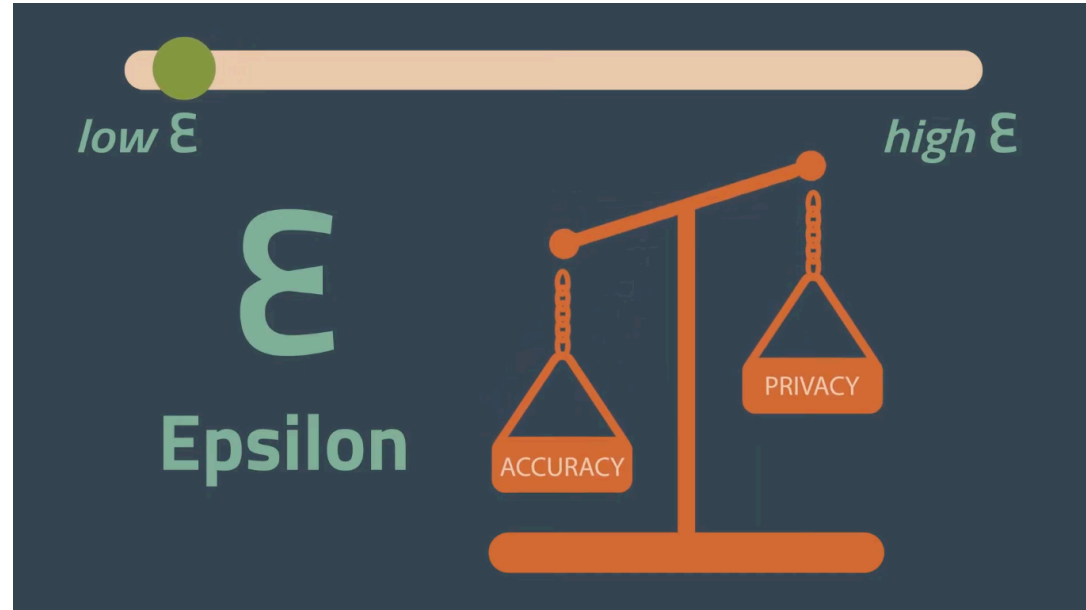
(a) ODDS-TEXT



(b) ODDS-VIS

- How do you communicate privacy risk to your friends?
- Excellent study: [[N+UseNIX'23](#)]
- Using odds ratio leads to increased understanding of risks and willingness to share data.
- How to explain ϵ -DP and μ -GDP? Need to incorporate prior knowledge of attacker.

Privacy Auditing



Drawbacks of pure theory

- Bounds always loose
 - people assume this and train models with high theoretical ϵ
- Maybe my implementation is incorrect
- Why should I trust your claim?

Backpropagation Clipping for Deep Learning with Differential Privacy

Timothy Stevens* Ivoline C. Ngong* David Darais Calvin Hirsch
University of Vermont *University of Vermont* *Galois, Inc.* *Two Six Technologies*

David Slater Joseph P. Near
Two Six Technologies *University of Vermont*

- In 2022, proposed to integrate clipping into forward/backward pass directly
- SOTA accuracy with 30x smaller ϵ

Privacy Auditing

Debugging Differential Privacy: A Case Study for Privacy Auditing

Florian Tramèr*, Andreas Terzis, Thomas Steinke, Shuang Song, Matthew Jagielski, Nicholas Carlini
Google Research

- Consider the following test:
 - D = MNIST dataset: 60k images
 - D' = Add (x', y') .
 - Train a CNN θ using [S+22] to get 0.98 acc and $(0.21, 10^{-5})$ -DP.
 - Check $\ell_{\theta}(x', y') \leq \tau$. If D' will be smaller.
 - Repeat 100k on D and 100k on D' .

