

lecture 23

~~Guest lecture Brian Z.~~

Differential Privacy

Evaluations =

Do Well on the Final!

Sat 7, 8 Monday
HW7 - \leq Sunday
HW8 TB posted

Q3 almost ready

OH Friday to
Mon ~~Mon~~
5:30-6:20
+ extra OH

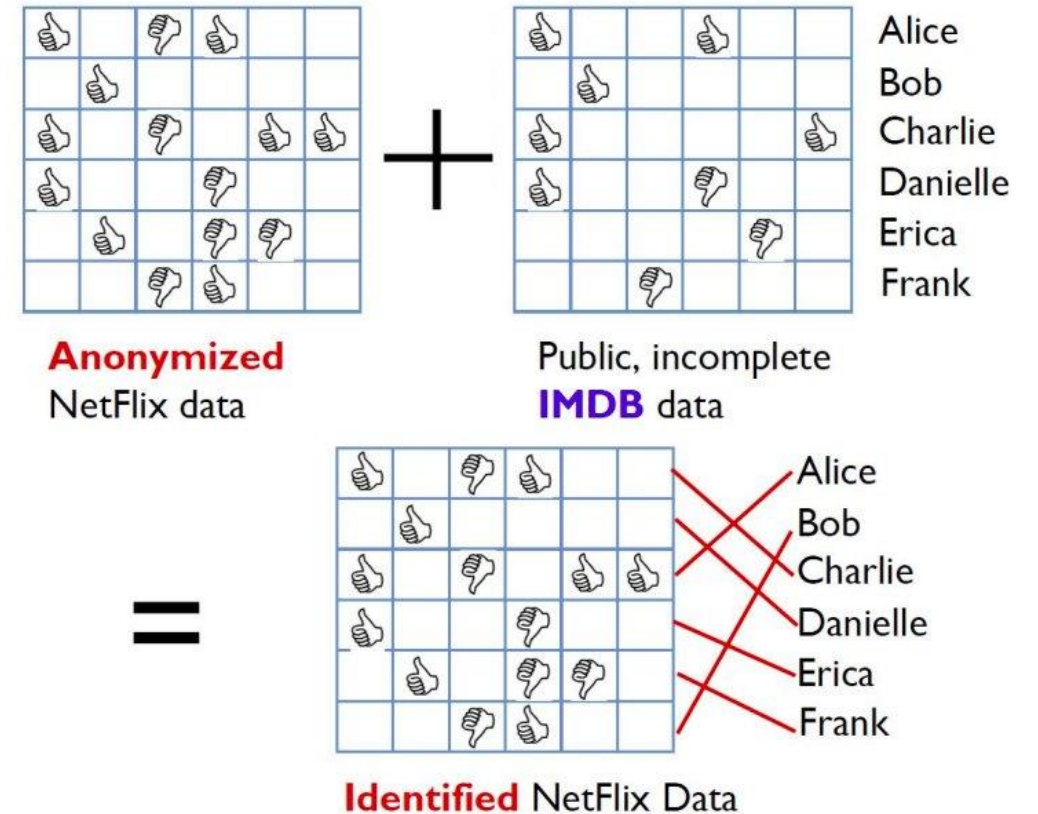
Differentially Private Machine Learning

Gautam Kamath



NetFlix Prize

- Recommendation engine competition (2006-2009)
- Training data: (anonymized) user ID, movie, rating, date
- Matched with public IMDb data: real name, movie, rating, date
- Class action lawsuit, cancellation of sequel



Privacy Concerns

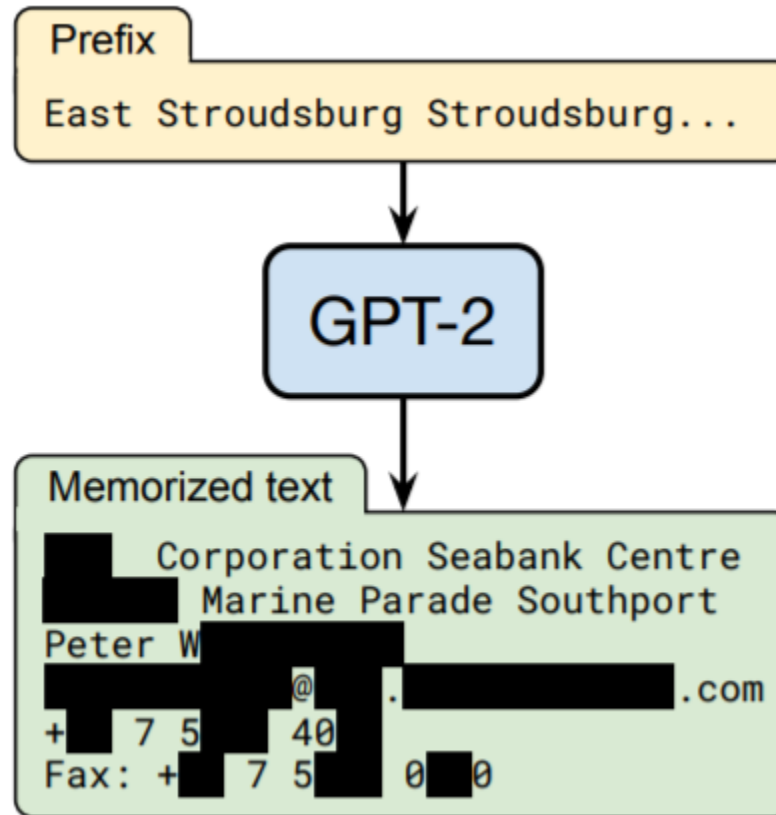
- Machine learning models are trained on very large datasets
- Can be coerced to reproduce training data verbatim!

In our [recent paper](#), we evaluate how large language models *memorize* and *regurgitate* such rare snippets of their training data. **We focus on GPT-2 and find that at least 0.1% of its text generations (a very conservative estimate) contain long verbatim strings that are “copy-pasted” from a document in its training set.**

Blog post: [Wallace, Tramer, Jagielski, Herbert-Voss], 2020

Paper: [Carlini, Tramer, Wallace, Jagielski, Herbert-Voss, Lee, Roberts, Brown, Song, Erlingsson, Oprea, Raffel], 2021

Personal Information



Copyrighted content

The escape of the Brazilian boa constrictor earned Harry his longest-ever punishment. By the time he was allowed out of his cupboard again, the summer holidays had started and Dudley had already broken his new video camera, crashed his remote-control aeroplane, and, first time out on his racing bike, knocked down old Mrs Figg as she crossed Privet Drive on her crutches.

Below, we prompt GPT-4o to generate a paragraph of text from the book *Harry Potter and the Philosopher's Stone*. **The model correctly generates the text before making its first**

Harry was glad school was over, but there was no escaping Dudley's gang, who visited the house every single day. Piers, Dennis, Malcolm, and Gordon were all big and stupid, but as Dudley was the biggest and stupidest of the lot, he was the leader. The rest of them were all quite happy to join in Dudley's favourite sport: Harry Hunting.

er and the Philosopher's Stone (about 240 words)

This was why Harry spent as much time as possible out of the house, wandering around and thinking about the end of the holidays, where he could see a tiny ray of hope. When September came he would be going off to secondary school and, for the first time in his life, he wouldn't be with Dudley. Dudley had been accepted at Uncle Vernon's old private school, Smeltings. Piers Polkiss was going there too. Harry, on the other hand, was going to Stonewall High, the local public school. Dudley thought this was very funny.

Also for Diffusion Models

Training Set



*Caption: Living in the light
with Ann Graham Lotz*

Generated Image



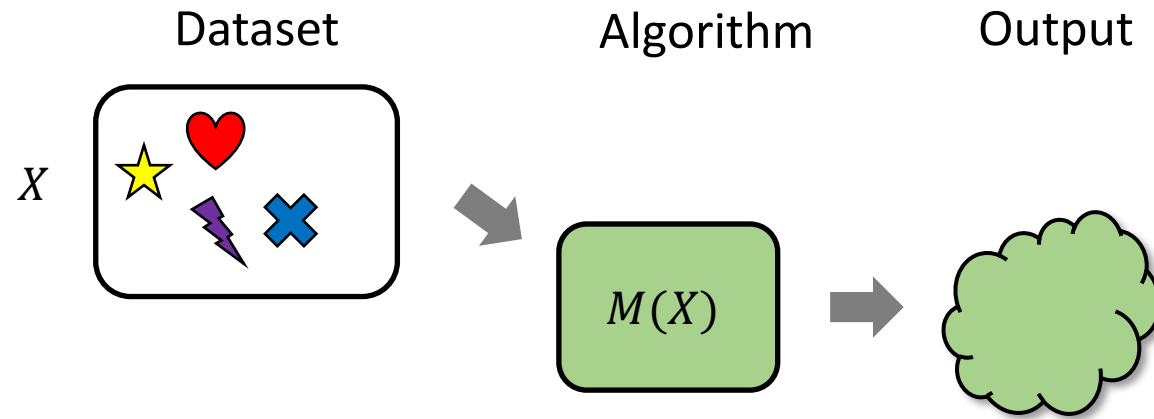
*Prompt:
Ann Graham Lotz*

Implications for copyright...?

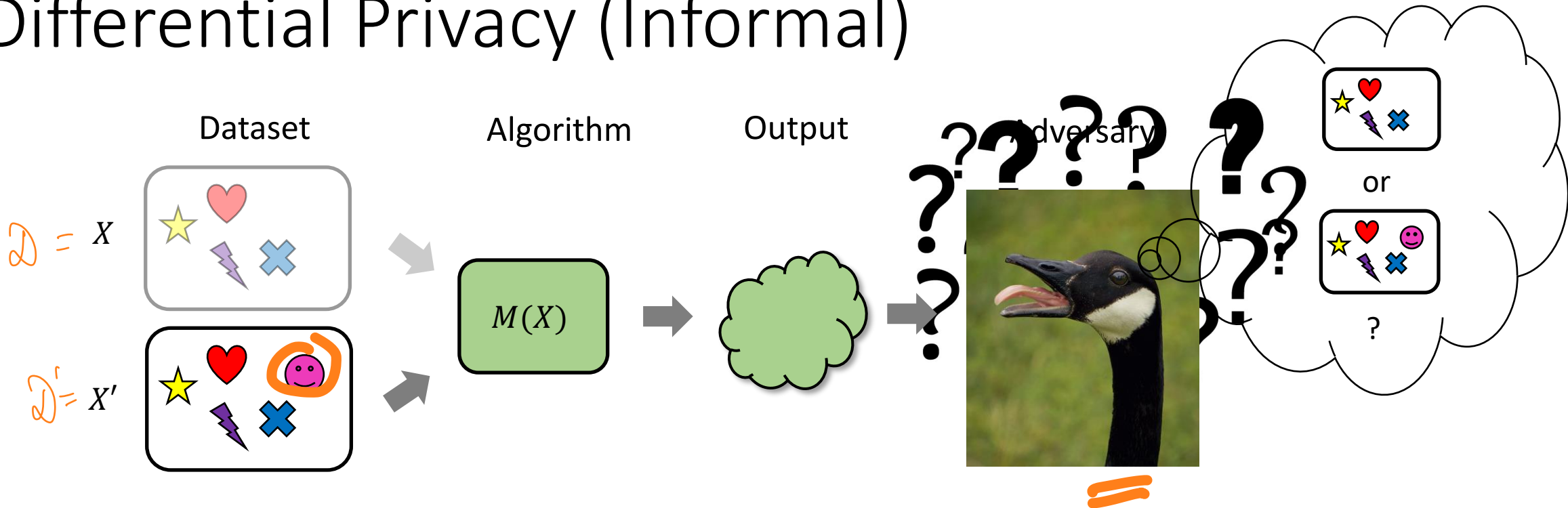
- Getty Images vs. Stability AI?



Differential Privacy (Informal)



Differential Privacy (Informal)

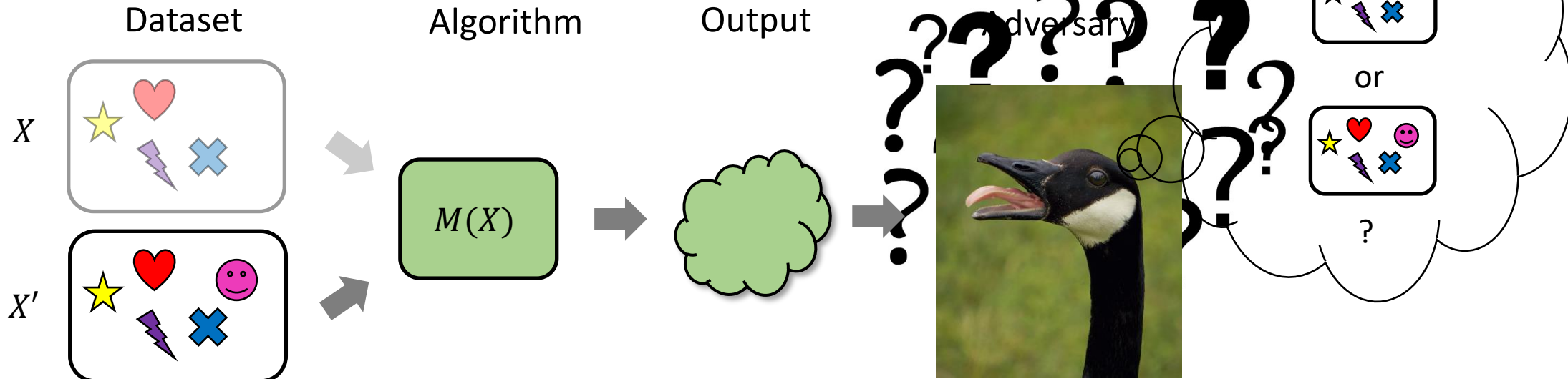


“An algorithm is differentially private if its distribution over outputs doesn’t change much after adding/removing one point.”

$M = g = g_{\text{reky}}$

$$-\epsilon \leq \ln \frac{\Pr[g(D) = y]}{\Pr[g(D') = y]} \leq \epsilon$$

Differential Privacy



Algorithm (D) \rightarrow y
Query

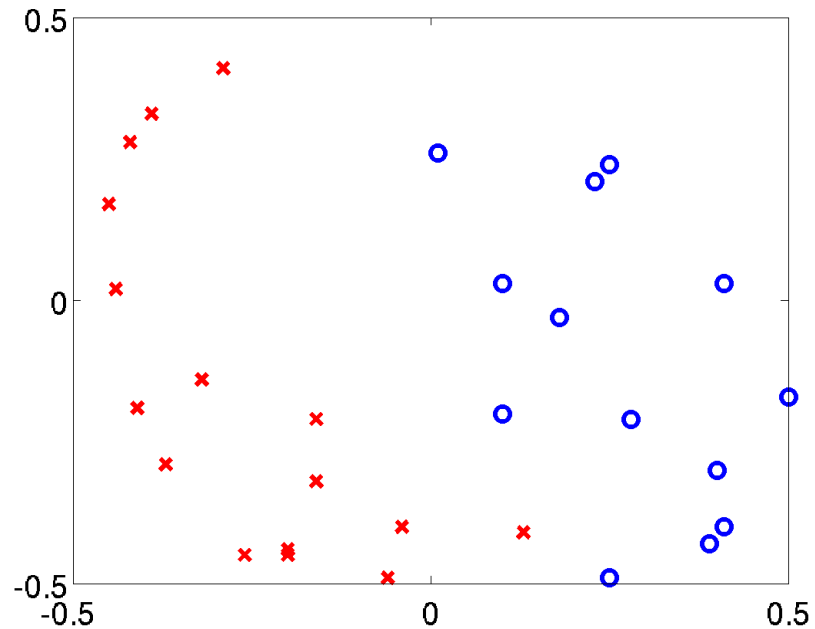
- $M: D^n \rightarrow R$ is (ϵ, δ) -DP if for all inputs X, X' which differ on one entry:

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \leq e^\epsilon \Pr[M(X') \in S] + \delta$$

e.g. {y} approx. DP

Widely Adopted

Accepted in theory...



And practice!




Endorsed by the White House



[Administration](#) [Priorities](#) [The Record](#)

OCTOBER 30, 2023

FACT SHEET: President Biden Issues Executive Order on Safe, Secure, and Trustworthy Artificial Intelligence

 [BRIEFING ROOM](#) [STATEMENTS AND RELEASES](#)

Today, President Biden is issuing a landmark Executive Order to ensure that America leads the way in seizing the promise and managing the risks of artificial intelligence (AI). The Executive Order establishes new standards for AI safety and security **protects Americans' privacy**, advances equity and civil rights, stands up for consumers and workers, promotes innovation and competition, advances American leadership around the world, and more.

(b) Within 365 days of the date of this order, to better enable agencies to use PETs to safeguard Americans' privacy from the potential threats exacerbated by AI, the Secretary of Commerce, acting through the Director of NIST, shall **create guidelines for agencies to evaluate the efficacy of differential-privacy-guarantee protections, including for AI.** The guidelines shall, at a minimum, describe the significant factors that bear on differential-privacy safeguards and **common risks to realizing differential privacy** in practice.

Differential Privacy (DMNS06)

learning alg (\mathcal{D}) = model
Query

$M: D^n \rightarrow R$ is (ϵ, δ) -DP if for all inputs X, X' which differ on one entry:

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \leq e^\epsilon \Pr[M(X') \in S] + \delta$$

• $\epsilon \approx 1$ and $\delta < 1/n$



• Worst-case guarantee ←

• $e^{\epsilon_1} e^{\epsilon_2} = e^{\epsilon_1 + \epsilon_2}$

• Symmetric definition

• M must be randomized

What DP does and does not mean

- Outcome is the same whether or not your data is in the dataset
- Protects against linkage and membership inference attacks
- Does *not* prevent statistics and machine learning 
 - “Smoking causes cancer”
- Not suitable when we need to identify a specific individual  *of course*
- Information-theoretic notion

Properties of Differential Privacy

- Post-processing

- If $M(X)$ is (ϵ, δ) -DP, then $f(M(X))$ is (ϵ, δ) -DP

- Group Privacy

- If M is (ϵ, δ) -DP, and X and X' differ in k entries,

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \leq e^{k\epsilon} \Pr[M(X') \in S] + \delta$$

- Composition

- If $M = (M_1, \dots, M_k)$ is a sequence of k (ϵ, δ) -DP algorithms

- M is $(k\epsilon, k\delta)$ -DP (Basic Composition)

- M is $(O(\sqrt{k}\epsilon \log(1/\delta')), k\delta + \delta')$ -DP (Advanced Composition)

Gaussian Mechanism

$X = \mathcal{D}$

- ℓ_2 -sensitivity of f

$$\Delta_2^{(f)} = \max_{X \sim X'} \|f(X) - f(X')\|_2$$

- If $\|f(X)\|_2 \leq C$, then $\Delta_2^{(f)} \leq 2C$

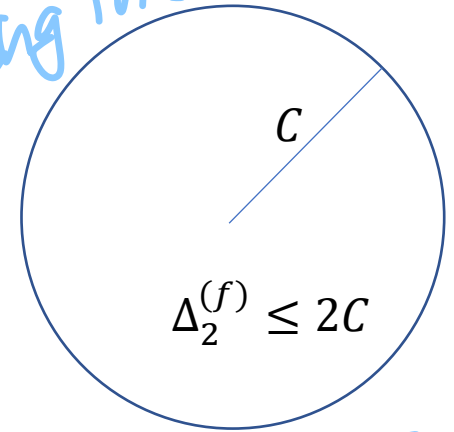
- Gaussian Mechanism

$$M(X) = f(X) + (Y_1, \dots, Y_k) \sim \text{noise}$$

Where $f(X) \in \mathbb{R}^k$, and the Y_i 's are $\approx N(0, \Delta^2 / \varepsilon^2)$

Thm (ε, δ) -DP

maxing inf



queky honest output

$N(0, \frac{\Delta^2}{\varepsilon^2} I_k)$

Stochastic Gradient Descent

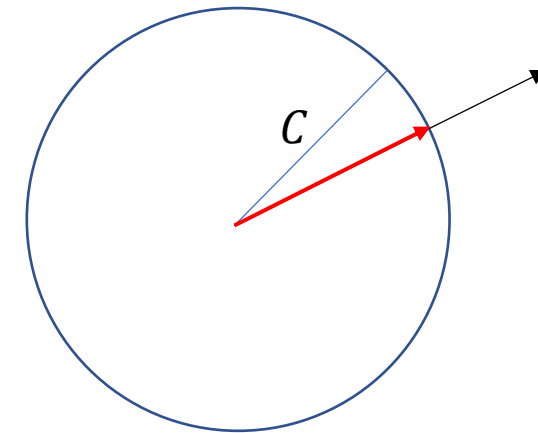
$$B \ll n$$

1. Choose a random minibatch B of points from the dataset
2. Compute the average gradient $\frac{1}{|B|} \sum_{(x,y) \in B} \nabla \ell(\theta_t, x, y)$
3. Take a step in the negative direction of the gradient
4. Repeat k times

parameters to learn

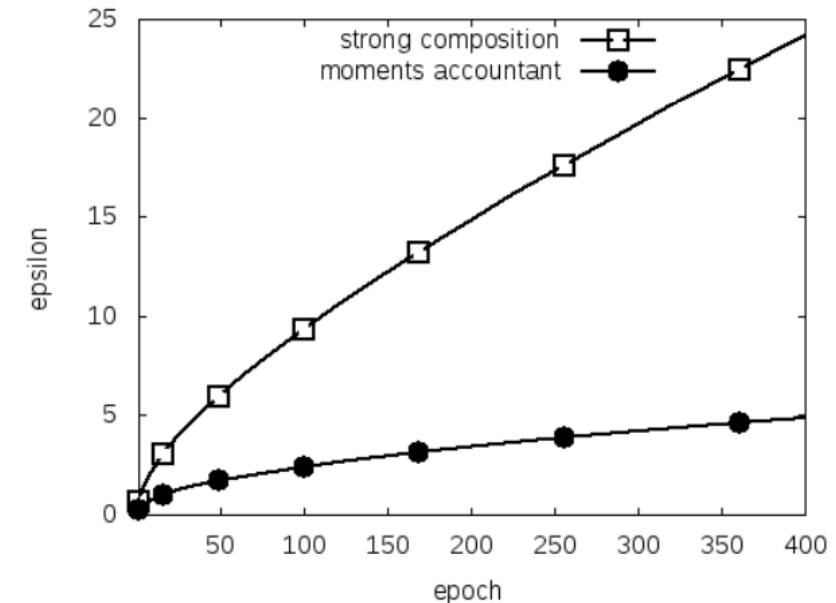
Differentially Private Stochastic Gradient Descent

1. Sample a “lot” of points of (expected) size L by selecting each point to be in the lot with probability L/n
2. For each point in the lot, compute the gradient $\nabla \ell(\theta_t, x, y)$ and “clip” it to have ℓ_2 norm at most C
3. Average the clipped gradients and add **Gaussian noise**
 - Apply the Gaussian Mechanism
4. Take a step in the negative direction of resulting vector
5. Repeat k times



Privacy of DPSGD (Informal)

- Suppose one step of DPSGD has privacy with parameter ϵ
- Since we subsample with probability L/n , each step is $\epsilon L/n$
 - “Privacy amplification by subsampling”
- k steps have privacy with parameter of $\epsilon \sqrt{k} L/n = \epsilon'$
 - Advanced composition
- Better analysis: “Moments accountant”



Does it work?

Data	ϵ -DP	Source	Test Accuracy (%)		
			CNN	ScatterNet+linear	ScatterNet+CNN
MNIST	1.2	Feldman & Zrnic (2020)	<u>96.6</u>	98.1 \pm 0.1	97.8 \pm 0.1
	2.0	Abadi et al. (2016)	95.0	98.5 \pm 0.0	98.4 \pm 0.1
	2.32	Bu et al. (2019)	96.6	98.6 \pm 0.0	98.5 \pm 0.0
	2.5	Chen & Lee (2020)	90.0	98.7 \pm 0.0	98.6 \pm 0.0
	2.93	Papernot et al. (2020a)	<u>98.1</u>	98.7 \pm 0.0	98.7 \pm 0.1
	3.2	Nasr et al. (2020)	96.1	–	–
	6.78	Yu et al. (2019b)	93.2	–	–
Fashion-MNIST	2.7	Papernot et al. (2020a)	<u>86.1</u>	89.5 \pm 0.0	88.7 \pm 0.1
	3.0	Chen & Lee (2020)	82.3	89.7 \pm 0.0	89.0 \pm 0.1
CIFAR-10	3.0	Nasr et al. (2020)	<u>55.0</u>	67.0 \pm 0.1	69.3 \pm 0.2
	6.78	Yu et al. (2019b)	44.3	–	–
	7.53	Papernot et al. (2020a)	<u>66.2</u>	–	–
	8.0	Chen & Lee (2020)	53.0	–	–

Public Data Helps for Private Vision

Dataset	Pre-Training Data	Top-1 Accuracy (%)				δ	Section
		$\epsilon = 1$	$\epsilon = 2$	$\epsilon = 4$	$\epsilon = 8$		
CIFAR-10	ImageNet	94.7	95.4	96.1	96.7	10^{-5}	4.1
CIFAR-100	ImageNet	70.3	74.7	79.2	81.8	10^{-5}	4.1
ImageNet	JFT-4B	84.4	85.6	86.0	86.7	$8 \cdot 10^{-7}$	4.2
Places-365	JFT-300M	-	-	-	55.1	$5 \cdot 10^{-7}$	4.3

