

## Group Privacy:

Thm: Let  $D, D' \in \mathcal{D}$  s.t.  $|D \Delta D'| = k$ .

Let  $A: \mathcal{D} \rightarrow \mathcal{R}$  be  $\epsilon$ -DP then  $\forall S \subseteq \mathcal{R}$ ,

$$\Pr[A(D) \in S] \leq e^{k\epsilon} \Pr[A(D') \in S]$$

Pf: Induction. Let  $D = D_0 \sim D_1 \sim D_2 \sim \dots \sim D_k = D'$ .

Claim:  $\forall i, \Pr[A(D_i) \in S] \leq e^{i\epsilon} \Pr[A(D_{i-1}) \in S]$

Pf: Base case:  $i=1$   
Def of  $\epsilon$ -DP ✓

inductive step: Is true for  $i-1$ .

$$\Pr[A(D_i) \in S] \leq e^{\epsilon} \Pr[A(D_{i-1}) \in S] \quad (\text{Def of } \epsilon\text{-DP})$$

$$\leq e^{\epsilon} \cdot e^{(i-1)\epsilon} \Pr[A(D) \in S] \quad (\text{induction})$$

$$\leq e^{i\epsilon} \Pr[A(D) \in S] \quad \checkmark$$

## Composition Theorems:

### Basic (adaptive) composition:

Q: If we run 2 (or more) DP algorithms,  
is the combination DP?

Setup: Let  $A_1: \mathcal{D} \rightarrow R_1$ ,  $\epsilon_1$ -DP

$A_2: \mathcal{D} \times R_1 \rightarrow R_2$ ,  $\epsilon_2$ -DP for every  $r \in R_1$ !  
 $\uparrow$   
adaptive!  $(A_2(\cdot, r): \mathcal{D} \rightarrow R_2 \text{ is } \epsilon_2\text{-DP } \forall r \in R_1)$

Let  $A: \mathcal{D} \rightarrow R_1 \times R_2$  be the algorithm that  
outputs  $A(D) = (r_1, r_2)$ , where  
 $r_1 = A_1(D)$ ,  $r_2 = A_2(D, r_1)$

Thm:  $A$  is  $(\epsilon_1 + \epsilon_2)$ -DP

Pf: Let  $D \sim D' \in \mathcal{D}$

Let  $(r_1, r_2) \in R_1 \times R_2$

$$\Pr[A(D) = (r_1, r_2)] = \Pr[A_1(D) = r_1] \cdot \Pr[A_2(D, r_1) = r_2]$$

$$\leq e^{\epsilon_1} \Pr[A_1(D') = r_1] \cdot e^{\epsilon_2} \Pr[A_2(D', r_1) = r_2]$$

$$\approx e^{\epsilon_1 + \epsilon_2} \Pr[A(D') = (r_1, r_2)] \quad \checkmark$$

Cor: Let  $A_1, \dots, A_k$  algorithms s.t.  $A_1: \mathcal{D} \rightarrow R_1$   
 is  $\epsilon_1$ -DP, and  $A_i: \mathcal{D} \times R_1 \times R_2 \times \dots \times R_{i-1} \rightarrow R_i$   
 is  $\epsilon_i$ -DP. Then  $A: \mathcal{D} \rightarrow R_1 \times R_2 \times \dots \times R_k$  which  
 runs the algorithms in sequence is  $\sum_{i=1}^k \epsilon_i$ -DP

Pf: induction, previous thm.

Tight in general!

Interpretation: Can build up larger algorithms from

DP steps!  $\epsilon$  is "privacy budget", so algorithm  
 runs  $k$   $\frac{\epsilon}{k}$ -DP algs, interspersed with computation  
 that does not look at input

$\Rightarrow$  Alg is  $\epsilon$ -DP!

Note: means that we'll often want to think of  
 $\epsilon$  as very small, not just small constant.  
 In particular,  $< \frac{1}{k}$

Ex (Adam Smith / Jon Allman): Lloyd's algorithm for  $k$ -means clustering.

Prob: Given  $x_1, x_2, \dots, x_n \in \mathbb{R}^d$ , find  $k$  points  $y_1, \dots, y_k$  minimizing  $\sum_{i=1}^n \min_{j=1}^k \|x_i - y_j\|_2^2$   
(min sum of squared distances to nearest center)

Lloyd's algorithm: (classical approx alg.)

First resolve so which all points in

$$\mathcal{U} = \{x \in \mathbb{R}^d : \|x\|_1 \leq 1\}$$

- init  $c_1^{(0)}, c_2^{(0)}, \dots, c_k^{(0)} \in \mathcal{U}$ .

- for  $i = 1$  to  $T$

- for  $j = 1$  to  $k$

-  $S_j \leftarrow \{i : c_j^{(i-1)} \text{ is closest pt-center to } x_i\}$

$$- c_j^{(i)} \leftarrow \frac{1}{|S_j|} \sum_{i \in S_j} x_i$$

return  $c_1^{(T)}, \dots, c_k^{(T)}$

What if we want to make it private?

Neighboring database: one datapoint added/removed.

Vectors in  $\mathbb{R}^d$ : could use Laplace Mechanism on final output.

But what is sensitivity of full algorithm?

Idea: use Laplace + composition!

Steps that use private info:

$$\sum_{i \in S_j} x_i, \quad |S_j|$$

$$- S_u + \epsilon' \sim \frac{\epsilon}{2T}$$

- for  $i = 1$  to  $T$

- for  $j = 1$  to  $k$

-  $S_j = \{i : c_i^{(t-1)} \text{ closest to } i\}$

$$- n_j = |S_j|$$

$$- a_j = \sum_{i \in S_j} x_i$$

- Release  $\hat{n}_j = n_j + Y$ , where  $Y \sim \text{Lap}(\frac{1}{\epsilon'})$

- Release  $\hat{a}_j = a_j + (Z_1, \dots, Z_d)$ , where  $Z_d \sim \text{Lap}(\frac{1}{\epsilon'})$  iid.

$$- c_j^{(t)} = \frac{\hat{a}_j}{\hat{n}_j} \quad \text{if } \hat{n}_j \geq 1$$

- Return  $c_1^{(T)}, \dots, c_k^{(T)}$  uniform in  $\mathcal{U}$  ok

Thm: This alg is  $\epsilon$ -DP.

Pf: By basic composition, if each iteration is  $2\epsilon'$ -DP, then overall alg. is  $2\epsilon' T = \epsilon$ -DP.

So consider some iteration, and  $x \sim x'$  differing points

$h_1, \dots, h_k$  form histogram!

$\Rightarrow$  adding  $\text{Lap}(\frac{1}{\epsilon'})$  noise to each is  $\epsilon'$ -DP

$\Rightarrow \hat{u}_i$ 's are  $\epsilon'$ -DP.

View  $a_1, \dots, a_k$  as one long vector of length  $k d$ .  
Call this vector  $A_x$  for  $x$ ,  $A_{x'}$  for  $x'$

In neighboring datasets  $x, x'$ , are different at most

$\Rightarrow \|A_x - A_{x'}\|_1 \leq 1$  (since each dataset has  $k$   $d$ -norm  
bounded by 1)

$\Rightarrow$  by Laplace mechanism, add noise  $\text{Lap}(\frac{1}{\epsilon'})$  to every  
entry is  $\epsilon'$ -DP

$\Rightarrow$  by composition, each iteration is  $2\epsilon'$ -DP

$\Rightarrow$  Alg is  $\epsilon$ -DP

## Advanced Composition:

what if allow  $(\epsilon, \delta)$ -DP?

- Basic composition still holds: if  $A_i$  is  $(\epsilon_i, \delta_i)$ -DP,  
set total alg is  $(\sum_{i=1}^k \epsilon_i, \sum_{i=1}^k \delta_i)$ -DP.

(exercise).  $(k\epsilon, k\delta)$  if all are  $(\epsilon, \delta)$ -DP

- (can do better: basically  $\sqrt{k}\epsilon$  !

Thm: If each  $A_i$  is  $(\epsilon, \delta)$ -DP, then

$\forall \epsilon, \delta \geq 0, \forall \delta' > 0$ , the composed alg.  $A$

is  $(\tilde{\epsilon}, \tilde{\delta})$ -DP, where

$$\tilde{\epsilon} = \epsilon \sqrt{2k \ln(\frac{1}{\delta})} + k\epsilon \frac{e^{\epsilon}-1}{e^{\epsilon}+1}, \quad \tilde{\delta} = k\delta + \delta'$$

get ahead for it before starting proof...

$\mathcal{F}_p$  have an algorithm which runs  $k$   $(\epsilon, \delta)$ -DP subalgorithms.

want overall alg to be  $(\tilde{\epsilon}, \tilde{\delta})$ -DP. How should we

set  $\epsilon, \delta$  as fn of  $\tilde{\epsilon}, \tilde{\delta}$ ?

Basic composition:  $\epsilon = \frac{\tilde{\epsilon}}{k}, \delta = \frac{\tilde{\delta}}{k}$

Advanced composition:  $\delta' = \tilde{\delta}/2 \Rightarrow \delta = \frac{\tilde{\delta}}{2k}$

For  $\varepsilon$ : For  $\varepsilon < 1$ ,  $\frac{e^\varepsilon - 1}{\varepsilon^2} \approx \frac{\varepsilon}{2}$

$$\Rightarrow \tilde{\varepsilon} = \Theta\left(\varepsilon \sqrt{k \ln\left(\frac{2}{\delta}\right)} + k \varepsilon^2\right)$$

$$\varepsilon^2 k \leq \tilde{\varepsilon} \Rightarrow \varepsilon \leq \sqrt{\frac{\tilde{\varepsilon}}{k}}$$

$$\varepsilon \sqrt{k \ln\left(\frac{2}{\delta}\right)} \leq \tilde{\varepsilon} \Rightarrow \varepsilon \leq \frac{\tilde{\varepsilon}}{\sqrt{k \ln\left(\frac{2}{\delta}\right)}}$$

$\sqrt{k}$  cost  
to compositions!

Note:  $\delta$  could equal 0  $\Rightarrow$  set  $\delta' = \tilde{\delta}$

still applies for pure DP mechanisms, but  
turns them into approximate DP!