

# Data Privacy

\* Simplest privacy method: Anonymization.

1. Remove "identifying" bits - eg, names, addresses, etc.
2. Publish data.

\* This has serious problems.

eg. AOL 2006, Netflix 2008.

\* Why? People's data tends to be very unique. For example:

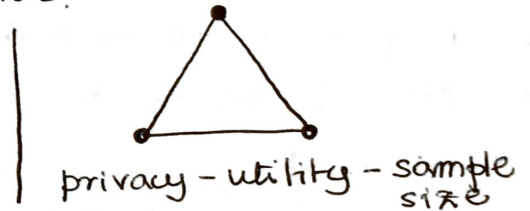
Gender	Position	Dept	Ethnicity
F	Faculty	CSE	S Asian

Only one person, Kamalika, fits description. "Linkage" information on CSE website.

\* Statistics from data also problematic.

eg, Wang et al study

Histogram with outliers.

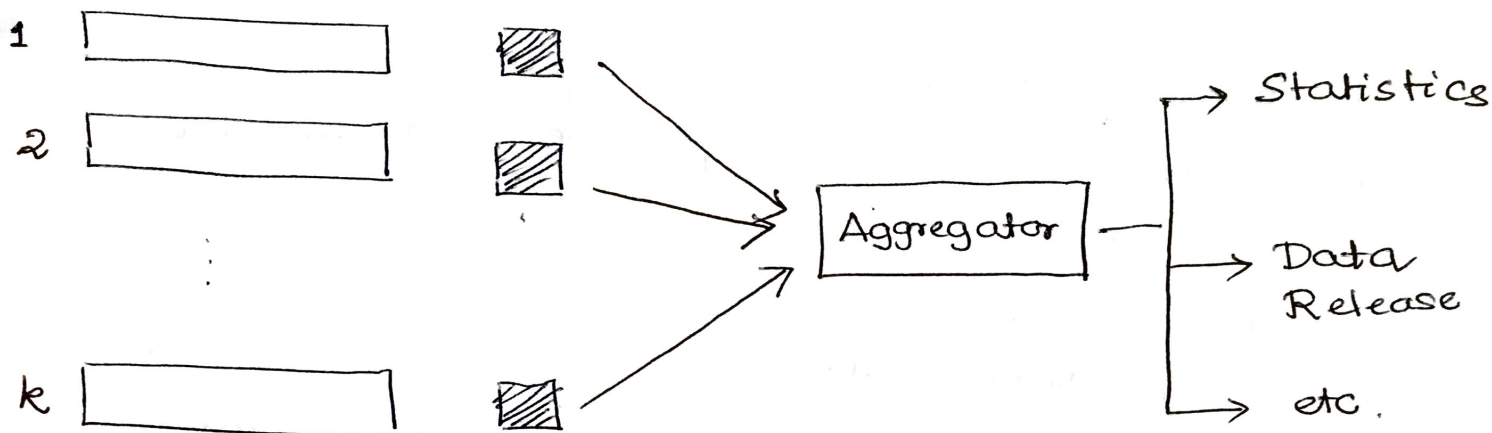


\* Need robust, rigorous measures to preserve privacy. tradeoffs.

## Three Privacy Settings:

1. Local sanitizer
2. Centralized sanitizer
3. Consent management.

# 1 Local Somitizer



1. Aggregator is untrusted.
2. Somitization happens locally before passing on to Aggregator.

Example: Randomized Response. [Werner 1965]

Each person is asked if they use an illegal drug (Yes or No). Everyone takes their answer, flips w.p  $p$  and returns it.

- Privacy offered : deniability.

$$\Pr(\text{Output} = Y \mid \text{True value} = Y) = 1 - p$$

$$\Pr(\text{Output} = N \mid \text{True value} = Y) = p$$

Using Bayes Rule,

$$\Pr(\text{True} = Y \mid \text{Output} = Y) = \frac{\Pr(\text{Output} = Y \mid \text{True} = Y) \Pr(\text{True} = Y)}{\Pr(\text{Output} = Y)}$$

assume  $= \frac{1}{2}$   
↑

$$\Pr(\text{Output} = Y \mid \text{True} = Y) \Pr(\text{True} = Y) + \Pr(\text{Output} = Y \mid \text{True} = N) \Pr(\text{True} = N)$$

$$= \frac{(1 - p) \frac{1}{2}}{\frac{1}{2} \cdot (1 - p) + \frac{1}{2} \cdot p} = 1 - p.$$

What is the utility offered?

Suppose we are aggregating  $N$  responses, and  $f$  is the fraction of drug users.

$$\mathbb{E}[\# \text{ Yes}] = N \Pr[\text{Output} = Y]$$

$$\Pr(\text{Output} = Y) = \Pr(\text{Output} = Y \mid \text{True} = Y) \Pr(\text{True} = Y) + \Pr(\text{Output} = Y \mid \text{True} = N) \Pr(\text{True} = N)$$

$$= (1-p)f + p(1-f) = p+f-2pf$$

$$\mathbb{E}[\# \text{ Yes}] = N(p+f-2pf) = Np + N(1-2p)f$$

$$\text{Var}(\# \text{ Yes}) \leq \sqrt{N}$$

$$\text{Let } T = \frac{\# \text{ Yes} - Np}{N(1-2p)}$$

$$\text{Then } \mathbb{E}[T] = f$$

$$\text{Var}(T) \leq \frac{1}{\sqrt{N}(1-2p)}$$

} Can estimate the fraction of drug users if  $p$  is not too close to  $\frac{1}{2}$ .

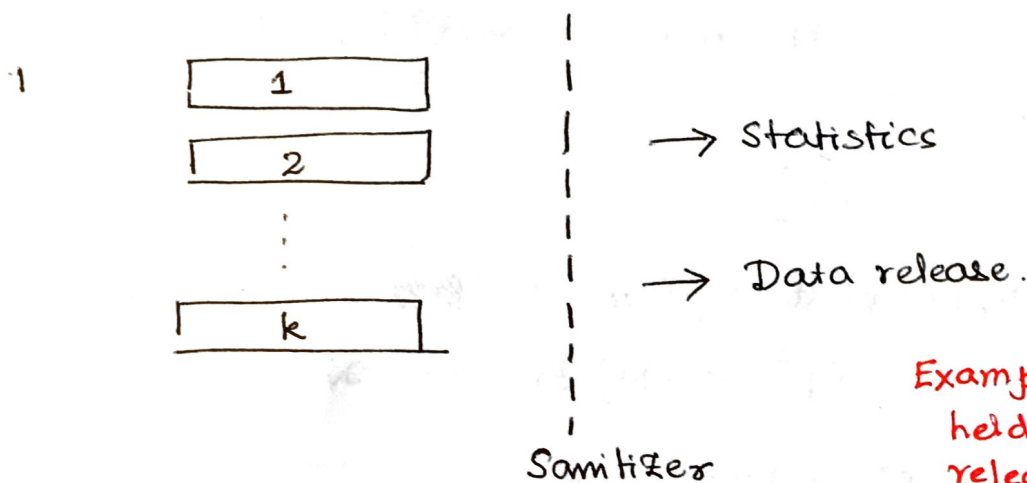
Note: Privacy - Utility - Datasize trade off

$p \approx \frac{1}{2}$  : High privacy, but high variance of  $T$ , so low utility.

Tradeoff is better when  $N$  is high

Applications: Data collection systems in companies, eg. Google, Apple, etc.

## [2] Centralized Somitizer



Example: Medical data, held at a hospital, summaries released.

Example 2: Census (2020 Census)

1. Somitizer is trusted.
2. Somitizer sees / collects raw sensitive data and "privatizes" it and passes the privatized version to the public sphere.

What is a good notion of privacy for such a setting?  
Differential privacy.

Main idea: Participation of a single person does not make a difference.

Alice + data  $\rightarrow$  <sup>Randomized</sup> Algorithm  $\rightarrow$  Output

Bob + data  $\rightarrow$  <sup>Randomized</sup> Algorithm  $\rightarrow$  Output

} similar.  
[~~output~~ output distributions are similar]

Adversary: whatever she can learn about Alice from algorithm's output, she can learn even if Alice is not in data.

Example 1: Study shows smoking causes cancer, Adversary knows Alice smokes  $\rightarrow$  infers Alice may have cancer.

NOT a privacy violation.



Example 2: Study in Wang et al  
with Alice in data.

Adversary knows about

Alice's genome  $\rightarrow$  solves equations and  
finds Alice in Cancer group

Privacy  
Violation

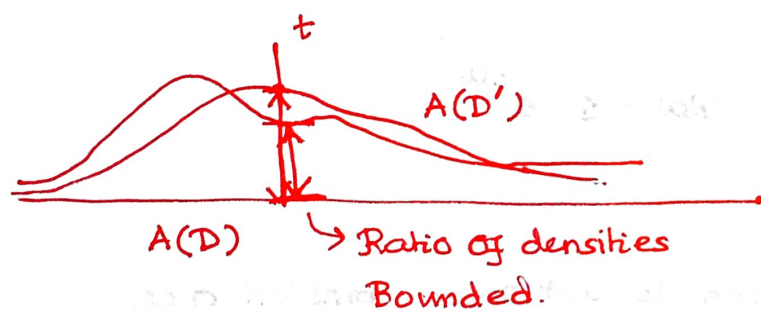
Formally, Differential privacy definition:

A mechanism  $A$  is  $\epsilon$ -DP if for all datasets  $D$  and  $D'$  that  
differ in the private value of a single person and for all  $t$ ,

$$\Pr(A(D) = t) \leq e^\epsilon \Pr(A(D') = t).$$

$\epsilon$  = privacy budget

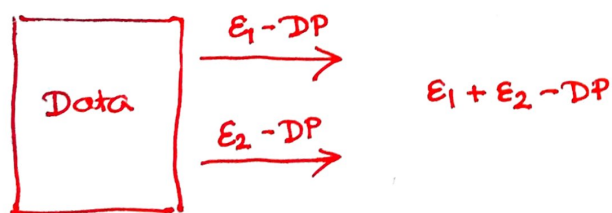
$\Pr$  = over randomization of the algorithm  $A$ .



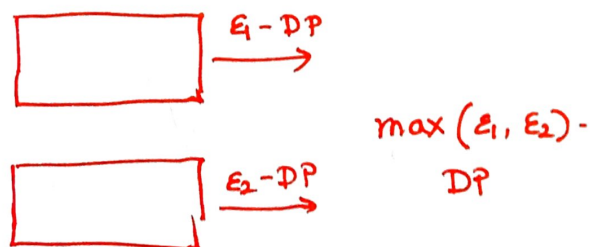
Properties: [1] Post processing Invariance.

Data  $\Rightarrow$   $\epsilon$ -DP answer  $\Rightarrow$   $\epsilon$ -DP (for any transformation)

[2] (Graceful) Composition



Sequential



Parallel

## How to get DP?

• The Global Sensitivity Mechanism.

• Global sensitivity of a function  $f$ :

$$GS(f) = \max_{D, D' \mid |D \setminus D'| = 1} |f(D) - f(D')|$$

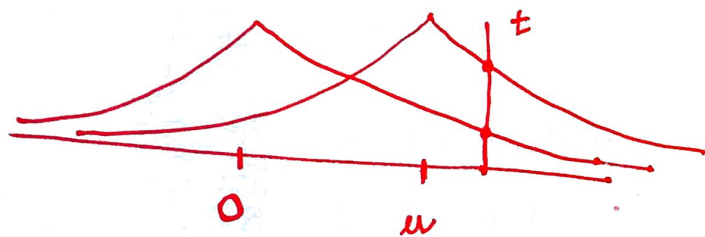
Mechanism:

Output  $f(D) + \frac{GS(f)}{\epsilon} Z$  where  $Z \sim \text{Lap}(\frac{1}{\epsilon})$ .

Laplace distribution:

$$f(z) = \frac{1}{2} e^{-|z|/b}$$

[mean 0, stdev:  $\sqrt{2}b$ ]



$$\text{Ratio} \leq e^{\epsilon u}$$

Use this mechanism + composition to get more complex mes.

Utility: Back to the drug user problem, if we use DP, then

GS of mean of  $n$  bits =  $\frac{1}{n}$ .

So stdev of noise added =  $\frac{\sqrt{2}}{n\epsilon}$  (as opposed to  $\frac{1}{\sqrt{n}(1-2p)}$  for RR)

So better privacy - utility - data size tradeoff