

PRIVACY:

Privacy means "Proof of Privacy".

SENSITIVE DATASETS:

→ Medical Records

→ Genetic Data

→ Search Logs

SEARCH LOGS:

NAME	EMAIL	SEARCH SENTENCE	LOCATION	DATE
⋮				

HIDDEN

AOL: Released such a database.

Con: Can be reverse-engineered.

NY TAXICAB:

* FOIL Request to get taxi fare data.

> Medallion, ^{HASHED} License, Vendor_id, rate_type, pickup_time, drop-off time, pickup_location, drop-off_location, fare ...

Con: To calculate income earned by various cab drivers.

NETFLIX CHALLENGE DATASET:

USER IDS	MOVIES				
	1	2	3	4	5
1					
2					
...					
...					

Narayan and Smatikov (2008)

IMDB Ratings Database:

MOVIES		

Alice
Bob
Charlie

LINK

Netflix Challenge 2: 10 Million \$ - was cancelled because of a lawsuit for violating privacy.

MASSACHUSETTS : GROUP INSURANCE COMMISSION:

William Weld Governor

→ Every state employees hospital visit records are available (but anonymized to preserve privacy) -

NAME	SSN	ID	SEX	AGE	ZIP	HEIGHT	WEIGHT	...

⏟

are available by voter records database.

Sweeney:

Sent the governor his medical record information.

"Reconstruction Attacks".

How to get "GUARANTEED" PRIVACY?

→ First example:

Simplest data analysis tool

→ Person 1 : $x_1 \in \{0,1\}$ "you like Star Wars"

Person 2 : $x_2 \in \{0,1\}$

\vdots

Person n : $x_n \in \{0,1\}$

→ What we want is to estimate

the average $\frac{x_1 + x_2 + \dots + x_n}{n}$

\nearrow
 p

SILLIEST:

No Privacy:

$\left. \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \right\}$ Released to public
 \downarrow
Can compute p exactly.

Full Privacy:

y_1 : "Noisy version of x_1 ,"

y_2

\vdots

y_n : "Noisy version of x_n ."

"Low p accuracy"

WARNER 1965: "RANDOMIZED RESPONSE" (RA)

$$\text{Each user } y_i = \begin{cases} x_i & \text{with prob } \frac{1}{2} + \delta \\ 1 - x_i & \text{with prob } \frac{1}{2} - \delta \end{cases}$$
$$0 < \delta < \frac{1}{2}$$

$\delta = 0$: "Full Privacy"

$\delta = \frac{1}{2}$: "No privacy" at all.

\hat{p} is going to be a function of y_1, y_2, \dots, y_n .

$$E[y_i] = \left(\frac{1}{2} + \delta\right) x_i + \left(\frac{1}{2} - \delta\right) (1 - x_i)$$
$$= \frac{1}{2} - \delta + 2\delta \cdot x_i$$

$$x_i = \frac{E[y_i] - \left(\frac{1}{2} - \delta\right)}{2\delta}$$

Suggests: Given the noisy information y_1, \dots

Estimate

$$\hat{p} = \frac{\left(\frac{y_1 + \dots + y_n}{n}\right) - \left(\frac{1}{2} - \delta\right)}{2\delta}$$

CLAIM:

$$\Pr [|\bar{p} - p| > \delta] \leq \frac{\delta}{\delta \cdot \sqrt{n}}$$

(Comes from computing variance of $\bar{p} - p$)

CLAIM:

With 75% chance my estimate \bar{p}

$$(\bar{p} - p) \leq \frac{1}{\delta \sqrt{n}}$$

If we need α accuracy

$$\frac{1}{\delta \sqrt{n}} \geq \alpha$$

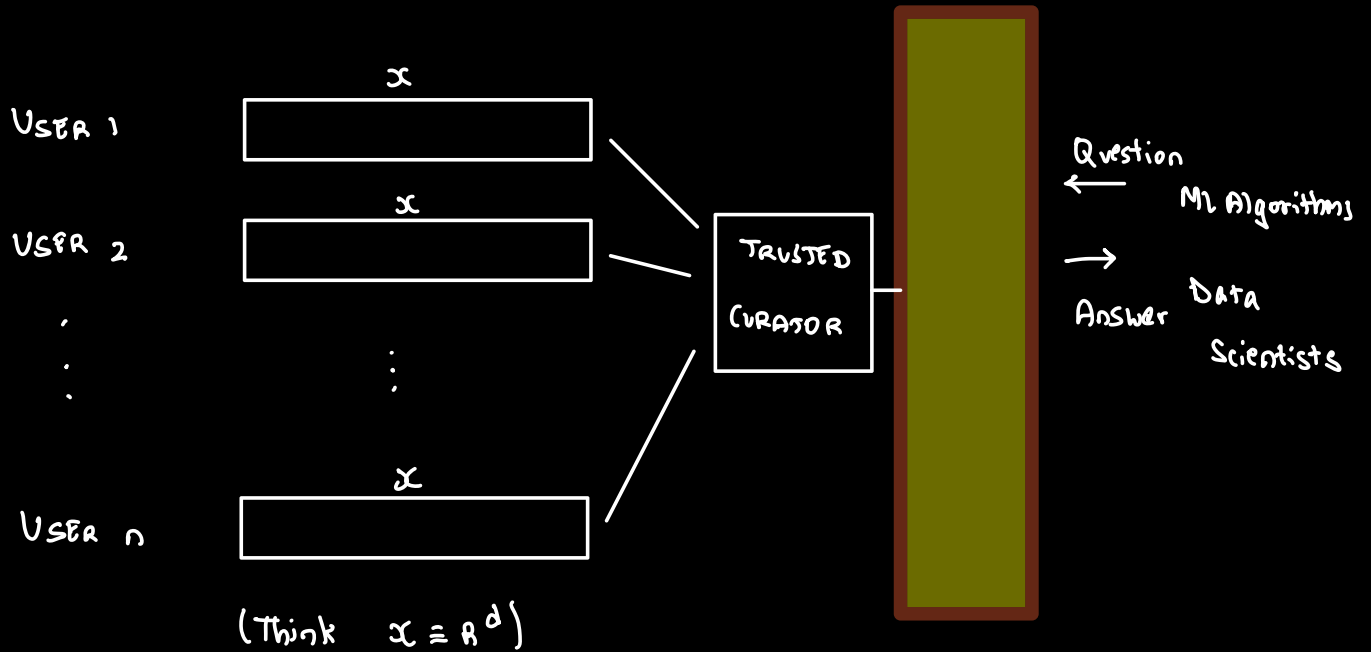
$$n \geq \frac{1}{\delta^2 \cdot \alpha^2}$$

You need at least $n \geq \frac{1}{\delta^2 \cdot \alpha^2}$ people.

DIFFERENTIAL PRIVACY:

"CENTRAL DIFFERENTIAL PRIVACY"

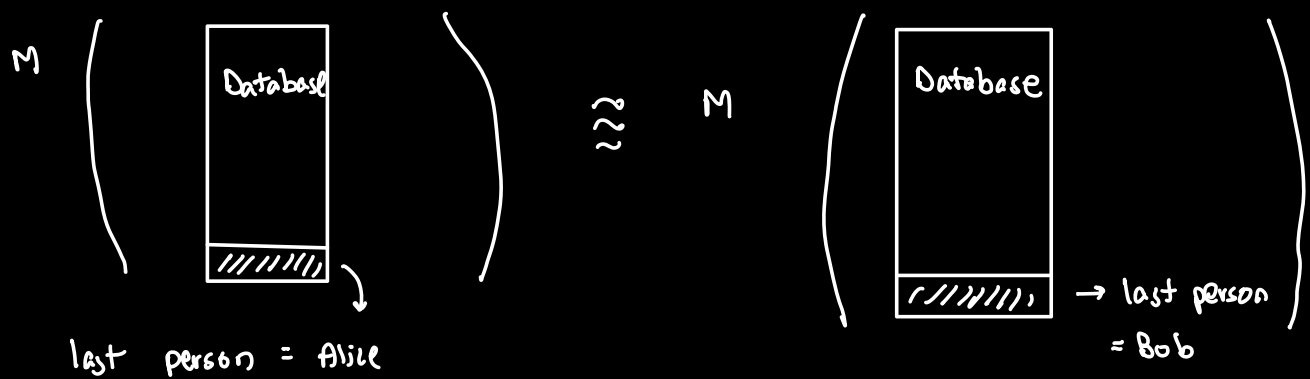
"TRUSTED CURATOR MODEL"



Curator adds noise based on sensitivity of question.

Curator $M: X^n \rightarrow Y \subseteq \mathbb{R}^k$
 ↓
 database

Idea: Removing one person's data should not change the answers much.
(or replacing)



DEFINITION:

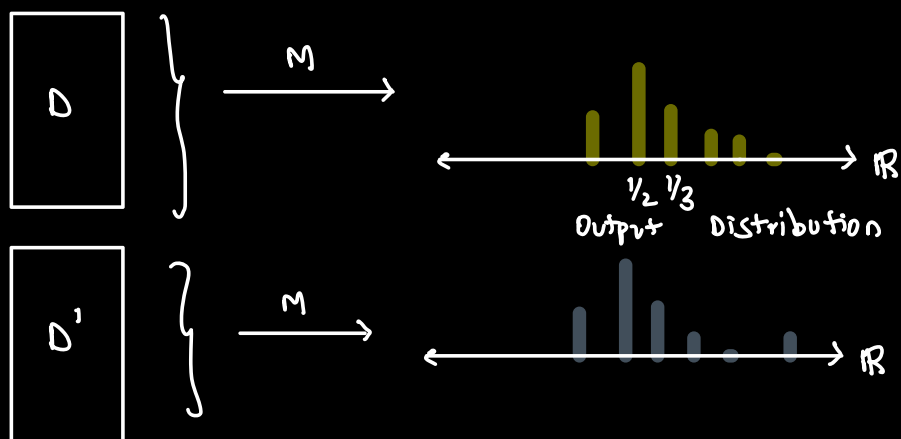
$D, D' \in \mathcal{X}^n$ are neighboring if they differ in exactly one row.

DIFFERENTIAL PRIVACY: $M : \mathcal{X}^n \rightarrow \mathcal{Y}$

is ϵ -Differentially private if \forall neighboring databases D, D' , $\forall y \in \mathcal{Y}$

$$\Pr[M(D) = y] \approx \Pr[M(D') = y]$$

$$\text{i.e. } e^{-\epsilon} \cdot \Pr[M(D') = y] \leq \Pr[M(D) = y] \leq e^{\epsilon} \Pr[M(D') = y]$$



It's not enough if outputs similar. Similar distribution is essential.

INTERPRETATION:

$$\begin{array}{ccc} e^{-\epsilon} \cdot \Pr[M(D') = y] & \leq & \Pr[M(D) = y] \leq e^{\epsilon} \Pr[M(D') = y] \\ \downarrow & & \downarrow \\ \approx (1 - \epsilon) & & \approx (1 + \epsilon) \end{array}$$

DMNSOR : Apple, Google, Microsoft, US Census Bureau 2020.

- * Differential Privacy is quantitative
- * Small ϵ corresponds to better privacy.
- * ϵ should be thought of as $\epsilon = 0.01$.
- * This is a worst-case guarantee on the data bases.
- * Probabilities are close multiplicatively.
- * e^{ϵ} vs $1 \pm \epsilon$ is just a convenience.