Lots of data, lots of "analysis" one can do.

Eg :

{ Understand causality b/w smoking and }
{                        lung disease. }

↓
Population wide studies.

Q :
{ How to conduct useful population-wide }
{ studies / "data analysis" without }
{ compromising individual privacy }
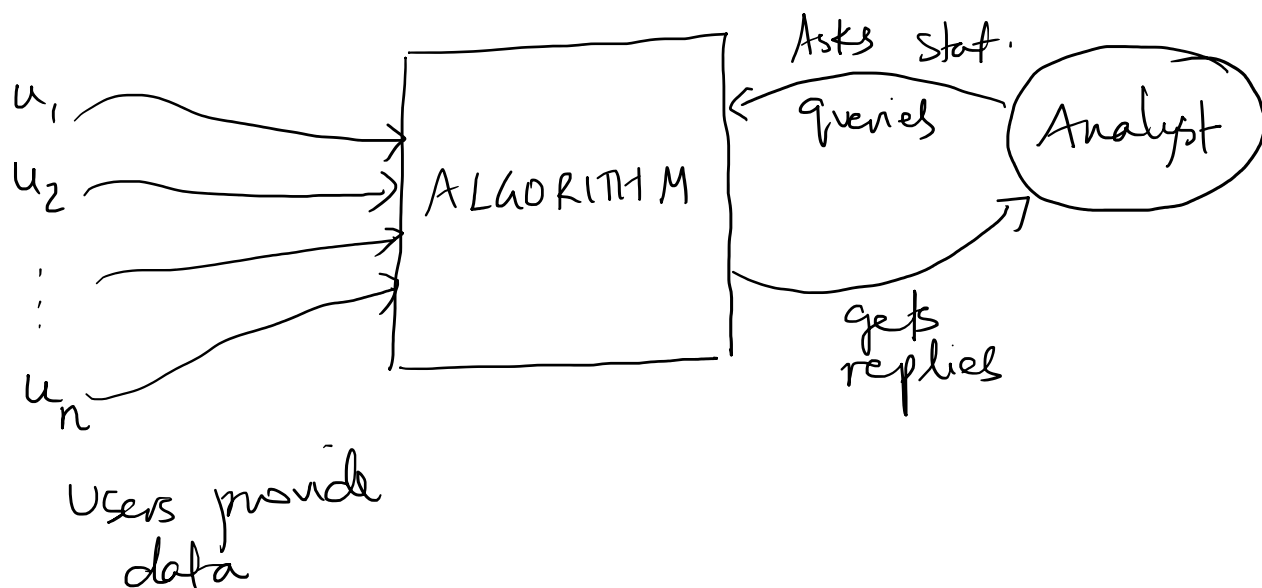
↓

Means users are incentivized to
join the study

[ Collectively we can learn something,
but individually don't lose anything]

{ different from "CRYPTOGRAPHY" }
Which is like a vault + key.

{ Here we want to contribut data }
{ but the analyst learns }

{ but the analyst learns
nothing "private"

---

## INITIAL ATTEMPT AT MODELING PRIVACY



$u_1$
$u_2$
$\vdots$
$u_n$

ALGORITHM

Asks stat. queries

Analyst

gets replies

Users provide data

## Attempt ①

{ Want the analyst to not learn
anything new about any
individual.

## Natural Issue

You will learn something new from
the answer of the query.

## Example

# Example

Imagine each $u_i = (\#\text{left feet}, \#\text{right feet})$

- Analyst asks $Avg(\#\text{left ft} - \#\text{right ft})$

- Reply (likely) $= 0$

$\Rightarrow$ Analyst "learns" that $u_1$ has equal
$\#$ of left & right ft.
(Maybe need few more queries to
get $Min(\#\text{ left ft})$
$Min(\#\text{ right ft})$, etc.

- Why is this definition vacuous?
A: what we learnt is something "global".
Nothing specific about the
particular individual

Q: How do we capture "Individual Privacy"?
[The earlier definitions of privacy is broken]

DMNS — Dwork, McSherry, Nissim, Smith

[Differential Privacy]

[What if the Algorithm guarantees that the answer to the query is "almost" the same regardless of whether $u_1$ (or any fixed individual) is present in the input or not ?]

The analyst can't even distinguish if $u_1$ was part of the study or not, so how can he/she learn anything about $u_1$'s data ?

## The Differential Privacy Model

Input : Database $X$ of 'n' rows, each corr. to a user; fn f.
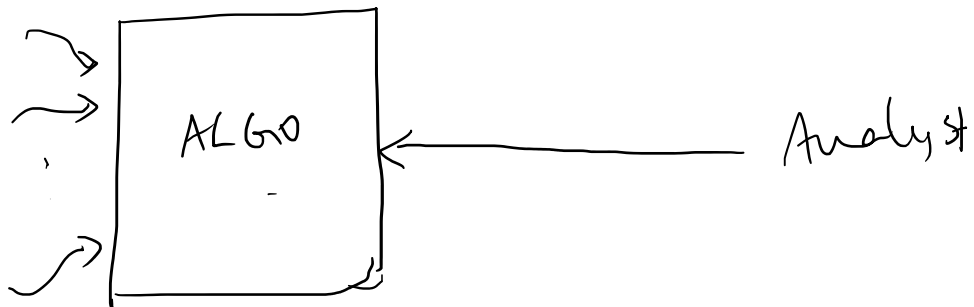
Algo takes $X$ and outputs $\tilde{f}(X)$

[It can be scalar output, vector output, etc]

$f$ is the statistical query

$\tilde{f}$ is the output.

$\tilde{f}$ is the output.

$\tilde{f}(x)$ can be some "noisy / approximate" response to $f(x)$.

If $x$ and $x'$ are two databases which differ in a single row, then

want

$$\boxed{|\tilde{f}(x) \;"\approx"\; \tilde{f}(x')|} \quad \textcircled{1}$$

$\uparrow$ guarantees privacy.

_and_

$\textcircled{2}$ $|f(x) - \tilde{f}(x)|$ is "small" for all $x$

$\uparrow$ guarantees utility of study.

Need to formalize "$\approx$" meaning and "small"

$\Big\{$ Just $\textcircled{1}$ is easy to satisfy :
       Output $\tilde{f}(x) = 0$ always.
       Full privacy, No utility $\Big\}$

$\Big\{$ How to get both together ? $\Big\}$

$\tilde{}$ ...... more "simple queries" for

There are many simple queries for which we need to introduce noise , else we break privacy.

Ex ①



All Microsoft
Employees
SALARY details

Analyst asks " count # people with
Salary $\geq$ 2M $ "

Sps Algo replies ①.

{ Then analyst can learn that CEO's
Salary $\geq$ 2M $. }

Output will be 0 if CEO doesn't
take part in survey.

' take part in survey.

$\Rightarrow$ CEO's privacy is compromised acc. our definition

---

Is this just a $\underline{1}$ v $0$ issue ?

Perhaps not.

Sps Answer $= 1000.$

and tomorrow, the answer is $1001.$

Maybe Microsoft hired a new employee,

$\Rightarrow \left\{ \text{likely that this person has salary} \geq 2M \$. \right\}$

we may learn something about the new Individual .

So our Model / definition captures these issues well
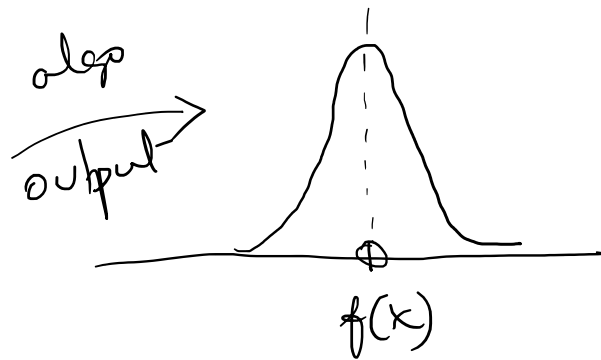
$\Rightarrow$ we must add some noise to $f(x)$

So that

$$\tilde{f}(x) \approx \tilde{f}(\dot{x}) \text{ for all}$$

$$\tilde{f}(x) \approx f(x) \text{ for all neighboring datasets.}$$

---

## Model Formally

$\boxed{\text{ALGO}}$ is randomized, adds noise acc. distribution



Distribution of $\tilde{f}(x)$.

### Error of Algo

$$\underset{\substack{\text{random} \\ \text{choices}}}{E} \left[ \left( \tilde{f}(x) - f(x) \right)^2 \right]$$

## Privacy Requirement :-

$\forall \ x, x'$ differing in one row, want distributions to be nearly identical.

nearly _identical_.



$\tilde{f}(x)$ & $\tilde{f}(x')$ distribution.

∀ $x, x'$ , and for all subsets $R$ of possible outputs of $\tilde{f}$,

$$\boxed{Pr\left[\tilde{f}(x) \in R\right] \le e^{\varepsilon} Pr\left[\tilde{f}(x') \in R\right]}$$

like $(1+\varepsilon)$

$\varepsilon$ − Differential Privacy.

__Goal__ :

$$\left[\begin{array}{c} \text{How to answer queries with} \\ \varepsilon - DP, \quad \text{with MIN. ERROR ?} \end{array}\right]$$

↓

TO MORROW

Such a scheme for "counting" queries

Database X

| Users | Salary |
|-------|--------|
| $u_1$ | $s_1$ |
| $u_2$ | $s_2$ |
| . | |
| | |
| $u_n$ | $s_n$ |

Algo

$f(x)$

Query = "Count # users with Salary $\geq V$"

What should a good $\tilde{f}(x)$ be ?

Want  ① $\tilde{f}(x)$ close to $f(x)$

② $\tilde{f}(x)$ close to $\tilde{f}(x')$ for all neighboring $x'$

## Idea
- Add noise to real answer.

Compute $f(x)$ , output $f(x) + R$

R is some suitable

where $R$ is some suitable noise.

eg $\left\{\begin{array}{l} R \sim N(0, \sigma^2) \text{ for suitable } \sigma \\ \text{will give privacy with} \\ \text{low error for suitable parameters} \end{array}\right\}$

## Our Algo  [ Technical Difference ]

Add noise $z$ w.p $\propto e^{-\frac{|z|}{\sigma}}$

$\nearrow$  ( in contrast gaussian noise has prob $\propto e^{-z^2/\sigma^2}$ )

"LAPLACIAN DISTRIBUTION"

$$f_R(z) = \frac{1}{2\sigma} e^{-\frac{|z|}{\sigma}}$$

## Check

① $\displaystyle\int_{-\infty}^{\infty} f_R(z) \, dz = 1$

② $\displaystyle\int_{-\infty}^{\infty} z \, f_R(z) \, dz = E[R] = 0$

③ $\displaystyle\int_{-\infty}^{\infty} z^2 \, f_R(z) \, dz = E[R^2] = 2\sigma^2$

$$\int_{-\infty}^{\infty} 2 \, f_R(z) dz - \cdots$$

eg, Integration by parts

---

If we add noise acc. Laplacian$(\sigma)$,
what is the squared error like?

$$\tilde{f}(x) = f(x) + R$$

$$\Rightarrow \text{error} = E\left[ \left( \tilde{f}(x) - f(x) \right)^2 \right]$$

$$= E\left[ R^2 \right] \qquad \text{where } R \sim \text{Lap}(\sigma)$$

$$= 2\sigma^2$$

Want to set $\sigma$ sufficiently large to get desired privacy.

$\forall x, x'$ differing in a row, and any subset $S$ of output values

$$\left\{ \overset{\text{want}}{\phantom{.}} \quad Pr\left[ \tilde{f}(x) \in S \right] \leq e^{\varepsilon} \, Pr\left[ \tilde{f}(x') \in S \right] \right.$$

Privacy
Requirement.

$\uparrow\uparrow$

we'll interpret set noise such that

We'll inferet set noise such that
the pdf of Algo output for
$X$ and $X'$ are <u>very similar</u>.

Fix an output value `t`.

Let $f_{Alg}(X, t) = $ PDF of Alg Outputting $t$ on input $X$

$$= \frac{1}{2\sigma} \exp\left( \frac{-|f(x) - t|}{\sigma} \right)$$

Similarly,

$$f_{Alg}(X', t) = \frac{1}{2\sigma} \exp\left( \frac{-|f(x') - t|}{\sigma} \right)$$

$$\Rightarrow \quad \frac{f(x, t)}{f(x', t)} \leq e^{-\frac{|f(x) - f(x')|}{\sigma}}$$

$$\leq e^{-\frac{1}{\sigma}}$$

So we can set $\sigma = \frac{1}{\varepsilon}$ ☺

$\forall x, x', t$
neighboring databases

$$\boxed{\left| \frac{f(x,t)}{f(x',t)} \right| \leq e^{\varepsilon}} \quad \leftarrow \text{Privacy}$$

$$\text{And} \quad E\left[ \left( \tilde{f}(x) - f(x) \right)^2 \right] \leq \frac{2}{\varepsilon^2}$$
$$\uparrow \text{Utility}$$

Only thing we used in proof is how much $f$ can change from $x \rightsquigarrow x'$.

SENSITIVITY of fn.

$$\Delta_f = \text{Max}_{\substack{x, x' \\ \text{differing} \\ \text{in one} \\ \text{row}}} \left| f(x) - f(x') \right|$$

Noise will simply depend on $\Delta_f$ by setting $\sigma$ appropriately to ensure $\varepsilon$ - Privacy.

SUMMARY

Simple scheme which works not just for counting queries, but for any low-sensitivity function

any low - sensitivity function

Algo is called
"Laplace Mechanism"

Similarly, Gaussian Mechanism also exists
(w/ gaussian noise).

In above example, query answer was a
numerical value. what if it's not?

Example

| classroom/days | M | T | W | Th | Fr |
|---|---|---|---|---|---|
| Students | | | | | |
| 1 | | x | | | |
| 2 | x | * | | | |
| | x | * | | x | |
| . | | | | * | |
| . | * | | | | |
| . | | | | | * |
| | x | * | | | |
| n | | x | | | * |

Each student has a preference for
when to conduct exam.

We want to select "Most popular day"
in a differentially private
manner.

Q1 - Cant simply add noise, (Meaningless)

Q1: - Cant simply add noise, (Meaningless)

Q2: - How to measure utility of a scheme.

Privacy is easy to extend

$$\frac{\text{Pr}[\text{Alg selects Monday for X}]}{\text{Pr}[\text{" " " ' X'}]} \leq e^{\varepsilon}$$

Similarly for each other day.

⎡ How to get utility ? ⎤

{ Want output to be a popular
  day if not the Most
  popular day . }

## Idea

For each possible output (days on our example)

| Day | 1 | 2 | 3 | 4 | 5 |
|-----|---|---|---|---|---|

Compute
# people
who prefer        $n_1$   $n_2$   $n_3$   $n_4$  $n_5$

Output day 'i' as answer
with prob $= e^{\varepsilon \cdot n_i}$

$$\sum e^{\varepsilon n_i}$$

$\Rightarrow$ Intuitively popular days are more likely to be output.

## Privacy + Error Analysis

For any day $i$ and inputs $x$ and $x'$

$$\frac{\Pr[\text{Alg selects } i \text{ for } x]}{\Pr[\text{Alg selects } i \text{ for } x']}$$

$$= \left( \frac{e^{\varepsilon n_i(x)}}{\sum_j e^{\varepsilon n_j(x)}} \right) \cdot \frac{\sum_j e^{\varepsilon n_j(x')}}{e^{\varepsilon n_i(x')}}$$

$$= e^{\varepsilon(n_i(x) - n_i(x'))} \cdot \frac{\sum e^{\varepsilon n_j(x')}}{\sum e^{\varepsilon n_j(x)}}$$

Since $x$ and $x'$ differ in a single row,

each $-1 \le n_i(x) - n_i(x') \le 1$

for all $i$.

Overall, $$\frac{\Pr[\text{Alg output } i \text{ on } X]}{\Pr[\text{Alg output } i \text{ on } X']} \leq e^{\varepsilon} \cdot e^{\varepsilon} = e^{2\varepsilon}.$$

Satisfies $2\varepsilon$- Differential Privacy.

what about error?

Let database has $n$ people.

and suppose $n_1 = OPT$ is the day with largest count.

Ideally : Want Alg to output a day with count close to $n_1$

Let's calculate
$$\Pr[\text{Alg outputs a day with count} \leq n_1 - t]$$

Let's fix a day $i$ with count $\leq n_1 - t$

$$\Pr[\text{Alg outputs this day}] = \frac{e^{\varepsilon n_i}}{\sum_j e^{\varepsilon n_j}}$$

No [ Alg outputs ... ] 

$$\leq \frac{e^{\varepsilon(v_1 - t)}}{e^{\varepsilon v_1}} \leq e^{-\varepsilon t}$$

with $\overline{\sum e^{\varepsilon v_j}}$ pointing to numerator.

So for $t = \dfrac{\log(n/\delta)}{\varepsilon}$

This probability is $\leq \dfrac{\delta}{n}$

$\Rightarrow$ Union bound over all bad days gives

$$\Pr[\text{outputting a bad day}] \leq \delta$$

$\Rightarrow$ ① with good probability, Algo chooses a day with score $\geq$ OPT $- \dfrac{\log(n/\delta)}{\varepsilon}$

$$\left[\text{very useful if OPT} \gg \dfrac{\log n}{\varepsilon}\right]$$

AND

② Preserves $2\varepsilon$- Privacy of users.

EXPONENTIAL MECHANISM.