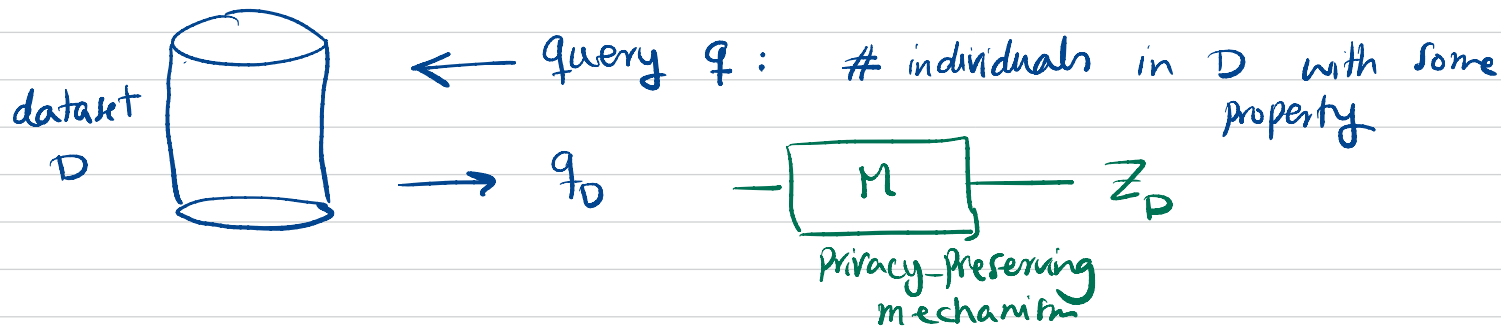# Chapter 2:   Privacy

**Privacy desideratum:** The output of the algorithm shouldn't change an adversary's knowledge about a dataset <u>at all</u>.

This statement ensures that <u>no</u> information should be revelead about the dataset; which is very unrealistic & goes against the point of <u>learning</u>. A more realistic goal is:

**Practical privacy desideratum:** The output of the algorithm shouldn't <span style="color:green">significantly</span> change an adversary's knowledge.
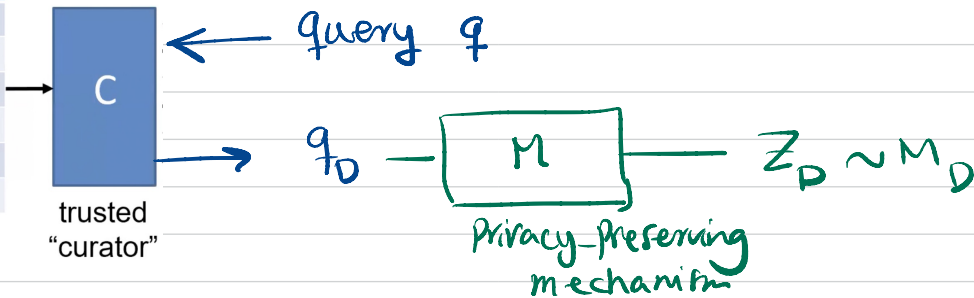
How to mathematically formulate this goal:

dataset D ⟵ query $q$: # individuals in D with some property

$\longrightarrow$ $q_D$ — $\boxed{M}$ — $Z_D$

Privacy-Preserving mechanism

* **Goal**: We wish to design M in such a way that an adversary observing $Z_D$ shouldn't be able to distinguish D from another dataset D' that differ in one entry!
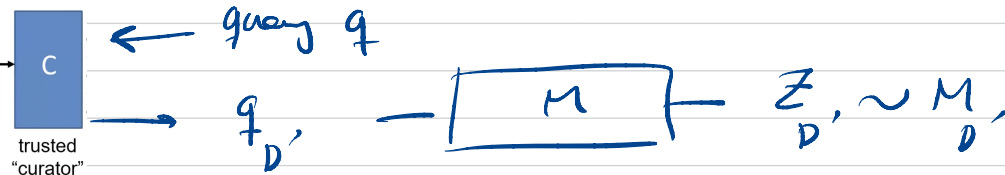
data set $D$

| Name | Sex | Blood | HIV? |
|------|-----|-------|------|
| Chen | F | B | Y |
| Jones | M | A | N |
| Smith | M | O | N |
| Ross | M | O | Y |
| Lu | F | A | N |
| Shah | M | B | Y |

C

trusted "curator"

← query $q$

$q_D$ — [ M ] — $Z_D \sim M_D$

Privacy-Preserving mechanism

data set $D'$

| Name | Sex | Blood | HIV? |
|------|-----|-------|------|
| Chen | F | B | Y |
| Jones | M | A | N |
| Smith | M | O | N |
| Bob ← | M | O | N |
| Lu | F | A | N |
| Shah | M | B | Y |

C

trusted "curator"

← query $q$

$q_{D'}$ — [ M ] — $Z_{D'} \sim M_{D'}$

If $M_D$ & $M_{D'}$ are 'close' $\Rightarrow$ adversary can't distinguish $D$ from $D'$

$D \sim D'$  neighboring dataset

Jones has a plausible deniability

**Def:** Randomized mechanism $M$ is said to be $\varepsilon$-DP for $\varepsilon \geq 0$

if

$$E_r(M_D \| M_{D'}) = 0 \quad \forall \; D \sim D'$$

$r = e^{\varepsilon}$

neighboring datasets

\* Since $M_D$ & $M_{D'}$ are close, no adversary can distinguish $M_D$ from $M_{D'}$

( & hence $D$ from $D'$ ).

## Remarks:

1. Why HS divergence? It's known that TV doesn't lead to

any meaningful privacy definition

(See lecture note of Sahil Vadhan & also

sec 1.6 "The complexity of

the 2nd assignment ) differential privacy"

The main reason for HS divergence is its connection with hypothesis testing, making the DP definition operational.

2- DP was originally proposed by Dwork, McSherry, Nissim, & Smith in the paper "Calibrating noise to Sensitivity in Private Data Analysis"

3- More Standard def:

$$M \text{ is } \varepsilon\text{-DP} \implies \boxed{M_D(A) \leq e^{\varepsilon} M_{D'}(A)} \quad \forall \text{ interval } A$$
$$\& \quad D \sim D'$$

* The equivalence between this formula & the above

definition comes from the Variational expression of $H_S$

divergence $\qquad E_\gamma(p\|q) = \sup_A \left[ P(A) - \gamma Q(A) \right]$

3. You may even see some definition like this:

M is $\varepsilon$-Dp $\quad \Longleftrightarrow \quad$ $e^{-\varepsilon} M_{D'}(A) \leq M_D(A) \leq e^{\varepsilon} M_{D'}(A)$ $\quad \forall$ interval A

$\underbrace{\phantom{e^{-\varepsilon} M_{D'}(A) \leq}}$ $\qquad\qquad$ & $D \sim D'$

This part is redundant! $\longleftarrow$ Prove this!

$*$ To prove a mechanism M is $\varepsilon$-Dp, it suffices to show

$E_{e^\varepsilon}(M_D \| M_{D'}) = 0$ or $\dfrac{M_D(A)}{M_{D'}(A)} \leq e^\varepsilon$ $\qquad$ for any event A

$\qquad\qquad\qquad\qquad$ & any $D \sim D'$

4- the definition **doesn't** depend on the dataset that you have hand.

* Note that the above definition should hold for any arbitrary pair of $D$ & $D'$

* **Perfect Privacy:** $\underline{\underline{\varepsilon = 0}}$

Let $M$ be an $\varepsilon$-Dp with $\varepsilon = 0$. Then since we have

$$M \text{ is } \varepsilon\text{-Dp} \iff E_{\underset{=1}{e^\varepsilon}}(M_D \| M_{D'}) = 0 \qquad \forall D \sim D'$$

$$\text{so} \quad TV(M_D, M_{D'}) = 0 \iff M_D = M_{D'}$$

So the distributions of $\underset{D}{Z}$ & $\underset{D'}{Z}$ are <u>exactly</u> the same.

Great in terms of privacy, but leads to very poor utility.

$\varepsilon > 0$, but sufficiently small:

$M$ is $\varepsilon$-DP $\Rightarrow$ $\xi_{e^\varepsilon}(M_D \| M_{D'}) = 0$

$\Longleftrightarrow$ $M_D(A) - e^2 M_{D'}(A) \leq 0$ $\forall$ event $A$

Since $\varepsilon$ is small then $e^\varepsilon \approx 1 + \varepsilon$ So

$M_D(A) - (1 + \varepsilon) \cdot M_{D'}(A) \leq 0$

$$\Rightarrow \qquad M_D(A) - M_{D'}(A) \leq \varepsilon \, M_{D'}(A) \leq \varepsilon$$

$$\Rightarrow \qquad \left| M_D(A) - M_{D'}(A) \right| \leq \varepsilon \qquad \text{for any event } A.$$

Thus, for small $\varepsilon$, DP means no matter what value output takes, $D$ & $D'$ can't be distinguished.

$\varepsilon = \infty$.

$$\frac{M_D(A)}{M_{D'}(A)} \leq e^{\varepsilon} \qquad ; \qquad \text{this requirement for}$$

large $\varepsilon$ is always

satisfied

$\varepsilon$ characterizes the ==privay gaurantee== of mechanism M

==the smaller $\varepsilon$ is the better privay is==

==Operational== Privay guarantee

we have a mechanism M that is $\varepsilon$-DP:

* Suppose M generates an output Z.

Null hypo: $H_0$ : Dataset = D ← Alice is in D

Alternative hyp: $H_1$ : Dataset = D' ← but not in D'

the goal is for an adversary to ==reliably== test $H_0$ against $H_1$.

* Assumption: Adversary knows everyone else in real dataset.

* what does it mean to test $H_0$ against $H_1$ reliably?

Let's $\varphi: \overset{\rightarrow}{Z} \longrightarrow \{0,1\}$
test function

$\varphi(Z)=0 \longrightarrow H_0$ is accepted

$\varphi(Z)=1 \longrightarrow H_1$ is accepted.

For any test function $\varphi$, we associate two errors :

1- Adversary rejects $H_0$ when $H_0$ was correct
« False positive »    « Type I error »
$F_p$

2_ Adversary accepts $H_0$ when $H_1$ was correct

"False negative"    "Type II error"

FN

when there exist a test function $\varphi$ such that both

FP & FN are small, we say that $H_0$ can be tested

against $H_1$ reliably.              FP, FN $\approx$ 0

* claim: If $Z$ is the output of an $\varepsilon$Dp mechanism

with small $\varepsilon$, then $H_0$ can't be tested

against $H_1$.

proof : Given $\varphi$, we can characterize FP & FN using the following set :

$$A = \{ z \in \mathbb{Z} : \varphi(z) = 0 \}$$

$z \in A \Rightarrow H_0$ is accepted

$z \in A^c \Rightarrow H_1$ is accepted

Think about it!

$FP = M_D(A^c) = 1 - M_D(A)$

$FN = M_{D'}(A)$

$FP + e^{\varepsilon} FN = 1 - M_D(A) + e^{\varepsilon} M_{D'}(A)$

$$= 1 - \left[ M_D(A) - e^{\varepsilon} M_{D'}(A) \right]$$

To find the best test function $\varphi$, we need to find $\varphi$ resulting in smallest FP & FN, or equivalently the best set A :

$$\inf_A \; FP + e^{\varepsilon} FN = \inf_A \left[ 1 - \left[ M_D(A) - e^{\varepsilon} M_{D'}(A) \right] \right]$$

$$\leq 1 - \sup_A \left[ M_D(A) - e^{\varepsilon} M_{D'}(A) \right]$$

$$= 1 - E_{e^{\varepsilon}} (M_D \| M_{D'})$$

Since $M$ is $\varepsilon$-DP, we have $\bar{E}_{e^\varepsilon}(M_D \| M_{D'}) = 0$ & hence

$$\boxed{\inf \; FP + e^\varepsilon FN \;\; = 1}$$

theorem: If $Z$ is the output of an $\varepsilon$-DP mechanism, then for any test function:

$$FP + e^\varepsilon FN \geq 1$$

$$FN + e^\varepsilon FP \geq 1$$

This theorem demonstrates that no test function can be found with $FP$ & $FN \approx 0 \implies H_0$ can't be tested against $H_1$ reliably!