# LECTURE 15

$$F(x) = f(x) + Lap\left(\frac{\Delta}{e}\right) \quad \left[\text{where } \Delta = S = \text{sensitivity}\right]$$

The above eqn = laplace distribution centered at '0'.

→ SENSITIVITY IS THE PROPERTY OF THE
QUERY [THAT is for COUNT QUERY THE sensitivity was 1 sensitivity), ~~Noedf~~

~~we too~~

Q) Now let's imagine we have various queries and the privacy budget (E) is fixed, so which query will involve more noise? Will It be the one with higher Δ or the one with lower Δ?

Ans FOR HIGH VALUE Δ MORE NOISE
FOR LOW VALUE Δ LESS NOISE

So the conclusion is if we have query with High sensitivity (Δ), the Laplace privacy scheme will inject more noise, and give us privacy at value of E we have set.

[This is logical, because see if a ~~data~~ query means it has a high sensitivity to mark/hide the records of the individual user and inorder to hide the record we need more noise, because high sensitivity also means that the presence or absence of that user in the database will make a huge difference to the database, so we have to ensure high privacy of that user, Hence we have to ~~hide~~ add more noise to hide the record]

Q) Now if for a query $\Delta$ is fixed and $\varepsilon$ is varied, what will happen?

HIGH VALUE OF $\varepsilon$ = LESS IS THE VALUE OF $b$ = LESS IS NOISE

LOW VALUE OF $\varepsilon$ = ~~⬤⬤⬤~~ IS THE VALUE OF $b$ = MORE IS THE MORE NOISE

[If we have to understand this logical if $\varepsilon$ / privacy budget is high, that

means we Have a Higher tolerance for privacy loss hence we can make_do with a little bit of noise

Q EXAMPLE WAY OF UNDERSTANDING THE EQN

Suppose you do a count query on a Hospital database for certain database

$$F(x) = f(x) + lap\left(\Delta/\varepsilon\right)$$

↑ so[orpq] count

↑ s[the noise added]

$f(x) = 55$

[the this noise is properly generated and not randomly generated to disclose the DP bound]

Q IN A CODING FORMAT. HOW CAN LAPLACE MECHANISM BE WRITTEN

Ans
sensitivity = 1
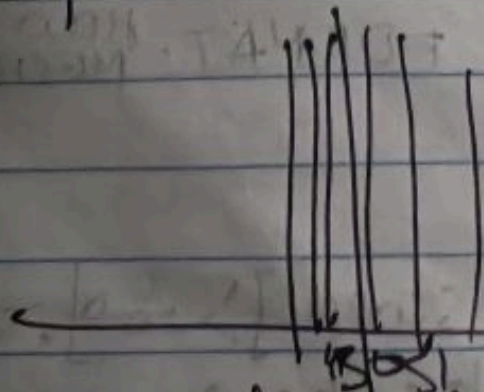epsilon = 0.1
CS5340[CS5340['score'] >= 80].shape[0] + np.random.laplace (loc=0, scale = sensitivity/ε)

In the code sensitivity is set to 1 as
it is a columnry query and 2 to set
to 0.1, ___ ___ which set it to 1
of ___ → which is set randomly here
we can set different values gs depending
on our privacy tolerance

Q     In the previous examples way of
understanding the eqn → we got the
ans as 55 (were the actual ans to0l
so), so, How will the user know whats
is the correct ans for the query)
He will simply run the query
$F(x)$ multiple time, on running
it multiple times, he might get
the answers as 48, 49, 50, 51, 52...55
, so no ao if we plot this in the
graph

we know where the values are centered
around.

On simply taking the average of the query results [that is 48, 49, 50-53] we get the ans = 50

Q. BUT WHAT PROBLEM DOES THE ABOVE SOLN CAUSE?
So suppose we query 'n' times in the above scenario, the privacy budget ($\varepsilon$) also increases by $n\varepsilon$, this means we are compromizing on our privacy

Q. SO HOW TO FIX THIS PROBLEM? [i.e the threat of the person doing repeated queries and finding out the ans]

Ans. we will simply put a limit on How much a person can query. i.e if we Have a privacy budget $\varepsilon$ we can design a system that such that if a person does a query that the privacy budget as $\varepsilon_1$ and He does query 'h' times

$$n\varepsilon_1 = \varepsilon$$

the person cannot do

[The limit here is that the person cannot do more than n queries]

⊗ Another way of saying is that each query should be $\boxed{\dfrac{\varepsilon}{n}}$

## PROPERTIES OF DP

① Sequential Composition :
$F_1(x)$ satisfies $\varepsilon_1 - DP$
$F_2(x)$ satisfies $\varepsilon_2 - DP$
The mechanism $G = (F_1(x), F_2(x))$ which releases both the results satisfied $(\varepsilon_1 + \varepsilon_2) DP$

[Above eg can be related to query multiple times. i-e $F_1(x) = 50$ $F_2(x) = 52$
so using the ABOVE SCENARIO IT HELP US IN KNOWING HOW TO DESIGN OUR SYSTEM [Eg in the previous case we place a limit on the no of queries]

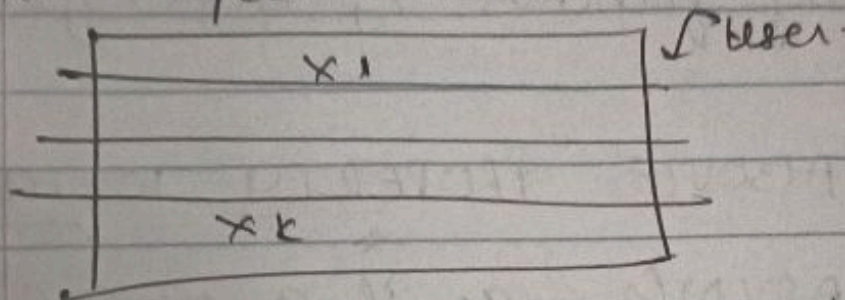② Parallel Composition
If $f(x)$ satisfies $\varepsilon - DP$ and we split a dataset $x$ into $k$ disjoint chunks $x_1 \cup x_2 --- \cup x_2 = x$.

→ The mechanism which releases $F(x_1), F(x_2) \cdots F(x_k)$ satisfies $\varepsilon - DP$

[Eg: There is a dataset, and its rows are split into chunks

$$\boxed{\begin{array}{c} x_1 \\ \\ \\ x_k \end{array}} \int user.$$

In ABOVE CASE, THE user is not present in each chunk but is present in one of the chunk maybe $x_1/x_2 \cdots /x_k$, so if the user does query on $x_1$ then query on $x_2$ -- till $x_k$ then only a particular user is affected, as if only one query core run on the dataset, this eliminates the [Because user's data is only in one chunk] [Hence $\varepsilon$ does not change here]

$\underline{\varepsilon}$ adding up

THUS PARALLEL COMPOSITION is BETTER THAN SEQUENTIAL COMPOSITION

## (3) POST PROCESSING

→ We cant reverse DP through post processing.

→ If $\{$ (x) satisfies $\frac{\varepsilon - D}{h}$, h(F(x)) satisfied

→ For any function $h$, $\varepsilon$-DP

$$\varepsilon - DP$$

[THE ABOVE PROPERTY IS HIGHLY ENCOURAGING, as if a person wants to release his database, he can rest assured that his database will still have $\varepsilon$-DP privacy guarante despite of any post processing done his data or any kind of manipula done to his data such as h]