

CS208: Applied Privacy for Data Science

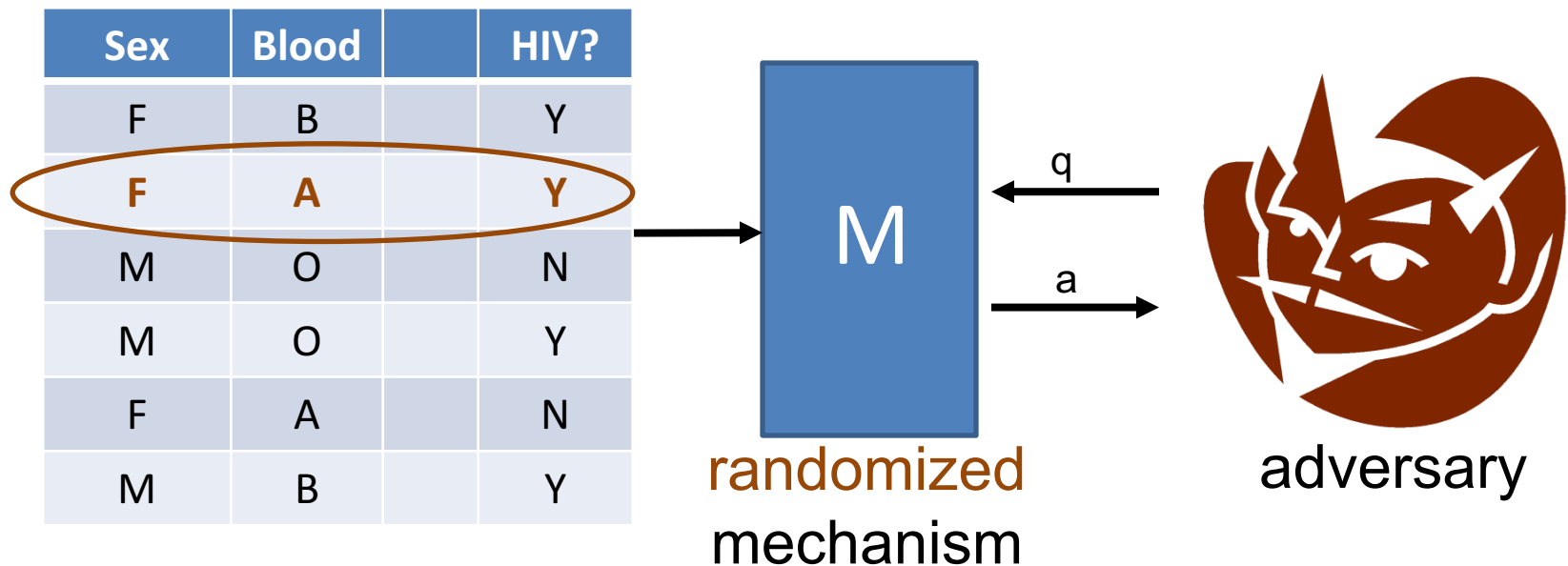
Introduction to Differential Privacy (cont.)

School of Engineering & Applied Sciences
Harvard University

February 15, 2022

DP for one query/release

[Dwork-McSherry-Nissim-Smith '06]



Def: M is ϵ -DP if for all D, D' differing on one row, and all q

$$\forall \text{ sets } T, \quad \Pr[M(D, q) \in T] \leq e^\epsilon \cdot \Pr[M(D', q) \in T]$$

(Probabilities are (only) over the randomness of M.)

The Laplace Mechanism

[Dwork-McSherry-Nissim-Smith '06]

- Let \mathcal{X} be a data universe, and \mathcal{X}^n a space of datasets. (For now, we are treating n as known and public.)
- For $x, x' \in \mathcal{X}^n$, write $x \sim x'$ if x and x' differ on at one row.
- For a query $q : \mathcal{X}^n \rightarrow \mathbb{R}$, the global sensitivity is
$$GS_q = \max_{x \sim x'} |q(x) - q(x')|.$$
- The Laplace distribution with scale s , $\text{Lap}(s)$:
 - Has density function $f(y) = e^{-|y|/s} / 2s$.
 - Mean 0, standard deviation $\sqrt{2} \cdot s$.

Theorem: $M(x, q) = q(x) + \text{Lap}(GS_q/\epsilon)$ is ϵ -DP.

Real Numbers Aren't

[Mironov '12]

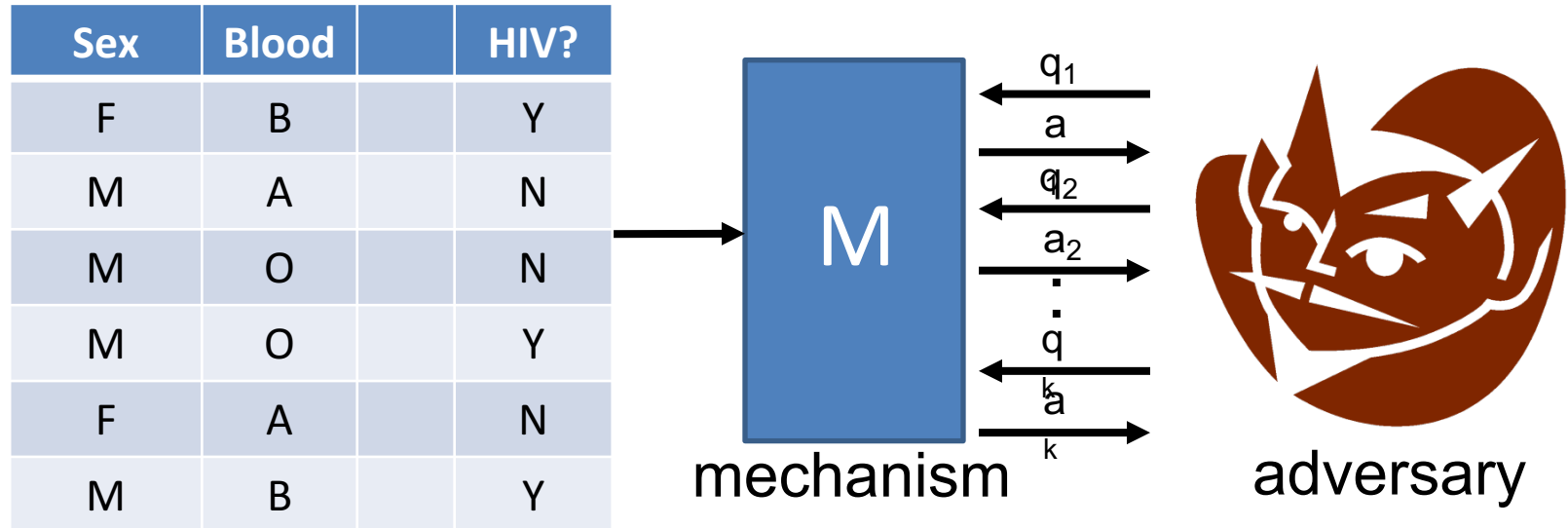
- Digital computers don't manipulate actual real numbers.
 - Floating-point implementations of the Laplace mechanism can have $M(x, q)$ and $M(x', q)$ disjoint \rightarrow privacy violation!
- **Solutions:**
 - Round outputs of M to a discrete value (with care).
 - Or use the **Geometric Mechanism**:
 - Ensure that $q(x)$ is always an integer multiple of g .
 - Define $M(x, q) = q(x) + g \cdot \text{Geo}(GS_q/g\varepsilon)$, where $\Pr[\text{Geo}(s) = k] \propto e^{-|k|/s}$ for $k \in \mathbb{Z}$.

Properties of the Definition

- **Suffices to check pointwise:** M is ϵ -DP if and only if
$$\forall x \sim x', \forall q, \forall t \Pr[M(x, q) = t] \leq e^\epsilon \cdot \Pr[M(x', q) = t]$$

← Replace with densities for continuous distributions →
- **Closed under post-processing:** if M is ϵ -DP and f is any function, then $M'(x, q) = f(M(x, q))$ is also ϵ -DP.
- **(Basic) composition:** If M_i is ϵ_i -DP for $i = 1, \dots, k$, then
$$M(x, (q_1, \dots, q_k)) = (M_1(x, q_1), \dots, M_k(x, q_k))$$
is $(\epsilon_1 + \dots + \epsilon_k)$ -DP.
 - Use independent randomness for k queries.
 - Holds even if q_i 's are adaptively chosen by an adversary.

Composition & Privacy Budgeting



Thm: If M is ϵ -DP if for one query, then it is $k\epsilon$ -DP for k queries.

- To maintain global privacy loss at most ϵ_g , can set $\epsilon = \epsilon_g/k$ and stop answering after k queries.
- More queries \Rightarrow Smaller $\epsilon \Rightarrow$ Less accuracy.
Some query-accuracy tradeoff is necessary! (why?)

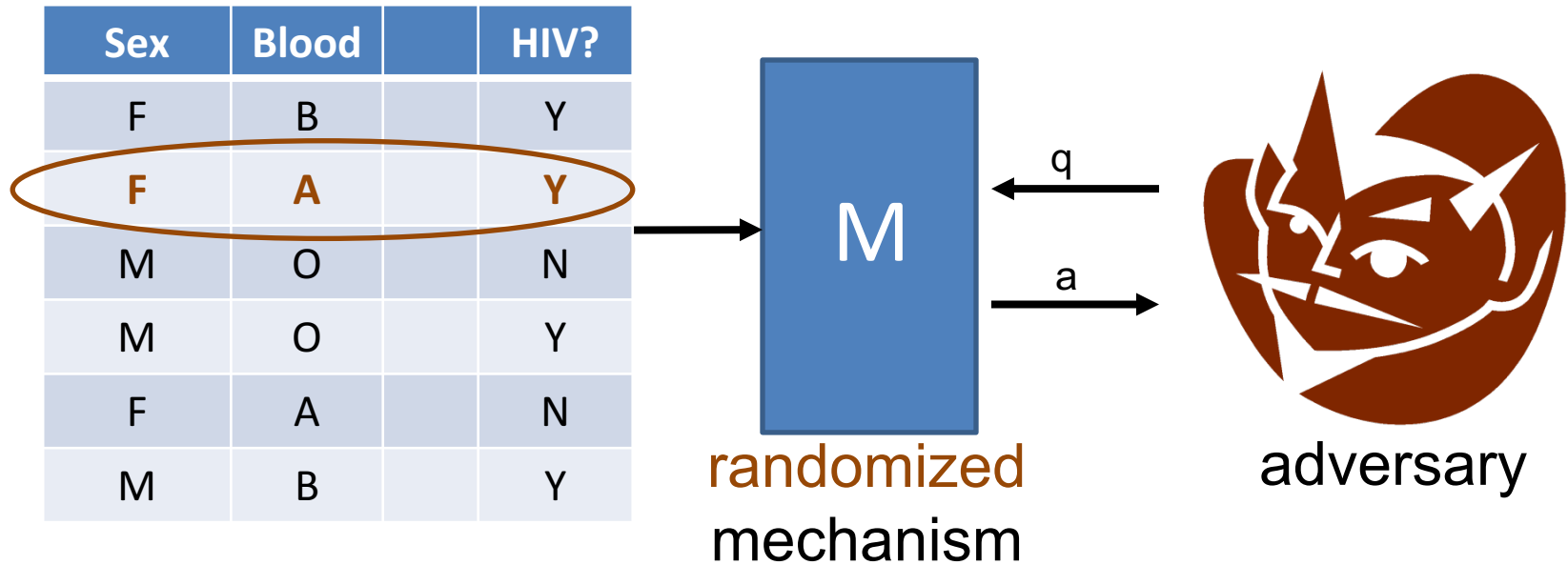
Composition for Algorithm Design

Composition and post-processing allow designing more complex differentially private algorithms from simpler ones.

Example:

- Many machine learning algorithms (e.g. stochastic gradient descent) can be described as sequence of low-sensitivity queries (e.g. averages) over the dataset, and can tolerate noisy answers to the queries. (The “Statistical Query Model.”)
- Can answer each query by adding Laplace noise.
- By composition and post-processing, trained model is DP and safe to output.

Interpreting the Definition



Def: M is ϵ -DP if for all D, D' differing on one row, and all q

$$\forall \text{ sets } T, \quad \Pr[M(D, q) \in T] \leq e^\epsilon \cdot \Pr[M(D', q) \in T]$$

(Probabilities are (only) over the randomness of M.)

Interpreting the Definition

- Whatever an adversary learns about me, it could have learned from everyone else's data.
- Mechanism cannot leak “individual-specific” information.
- Above interpretations hold regardless of adversary's auxiliary information or computational power.

But:

- No guarantee that adversary won't infer sensitive attributes.
- No guarantee that subjects won't be “harmed” by results of analysis.
- No protection for information that is not localized to a few rows.

A Bayesian Interpretation

- Let $X = (X_1, \dots, X_n) \in \mathcal{X}^n$ be a random variable distributed according to an adversary's "prior beliefs" about a dataset, and let $X_{-i} = (X_1, \dots, X_{i-1}, \perp, X_i, \dots, X_n)$ have person i 's data removed or replaced with a dummy value in \mathcal{X} .
- Suppose $M : \mathcal{X}^n \rightarrow \mathcal{Y}$ is ε -DP, and let $y \in \mathcal{Y}$ be any possible output. Then for every $x_i \in \mathcal{X}$,

$$\underbrace{\Pr[X_i = x_i | M(X) = y]}_{\text{Posterior belief about person } i \text{ after seeing output } y} \in e^{\pm \varepsilon} \cdot \underbrace{\Pr[X_i = x_i | M(X_{-i}) = y]}_{\text{Posterior belief about person } i \text{ after seeing output } y \text{ if person } i \text{'s data wasn't used}}$$

- Explains choice of multiplicative distance in def of DP.

Group Privacy & Setting ϵ

- **Thm:** If M is ϵ -DP for one query, then it is $k\epsilon$ -DP for k groups of size k : for all x, x' that differ on at most k rows,
$$\forall q \forall T \Pr[M(x, q) \in T] \leq e^{k\epsilon} \cdot \Pr[M(x', q) \in T]$$
 - Meaningful privacy for groups of size $O(1/\epsilon)$.
- **Cor:** Need $n \geq 1/\epsilon$ for any reasonable utility.
- Typical recommendation for “good” privacy guarantee:
 $.01 \leq \epsilon \leq 1$.

Variants of the Definition

- When n is not publicly known:
 - **Datasets:** multisets D of elements of \mathcal{X} , can represent as a histogram $D \in \mathbb{N}^{\mathcal{X}}$, where D_x = number of copies of x .
 - **Neighbors:** $D \sim D'$ iff $|D \Delta D'| = 1$ (add or remove an elt)
In histogram notation: $|D \Delta D'| = \sum_x |D_x - D'_x| \stackrel{\text{def}}{=} \|D - D'\|_1$
- Social network data:
 - **Datasets:** graphs G , possibly with labels on nodes and edges
 - **Neighbors v1:** $G \sim G'$ if differ by modifying one edge
 - **Neighbors v2:** $G \sim G'$ if differ by modifying one node & incident edges.
 - **Q:** which choice provides better privacy protection?

Approximate Differential Privacy

Def: M is (ϵ, δ) -DP if for all $D \sim D'$, and all q

\forall sets T , $\Pr[M(D, q) \in T] \leq e^\epsilon \cdot \Pr[M(D', q) \in T] + \delta$

- Intuitively: ϵ -DP with probability at least $1 - \delta$.
- Picking a random person from dataset and publishing their data is $(0, 1/n)$ -DP, so want $\delta \ll 1/n$.
- Ideally set δ to be cryptographically small (e.g. 2^{-50}).
- Satisfies postprocessing, basic composition (adding δ_i 's).
- Group privacy for groups of size up to $O(1/\epsilon)$.
- Does not suffice to check pointwise (need to consider sets T).

Benefits of Approximate DP

- More mechanisms, e.g. **Gaussian Mechanism:**

$$M(x, q) = q(x) + \mathcal{N}(0, \sigma^2),$$
$$\text{for } \sigma = \frac{\text{GS}_q}{\varepsilon} \cdot \sqrt{2 \ln(2/\delta)}$$

- **Advanced Composition Thm:** If M_i is (ε, δ) -DP for $i = 1, \dots, k$ and $k < 1/\varepsilon^2$, then $\forall \delta' > 0$

$$M(x, (q_1, \dots, q_k)) = (M_1(x, q_1), \dots, M_k(x, q_k))$$

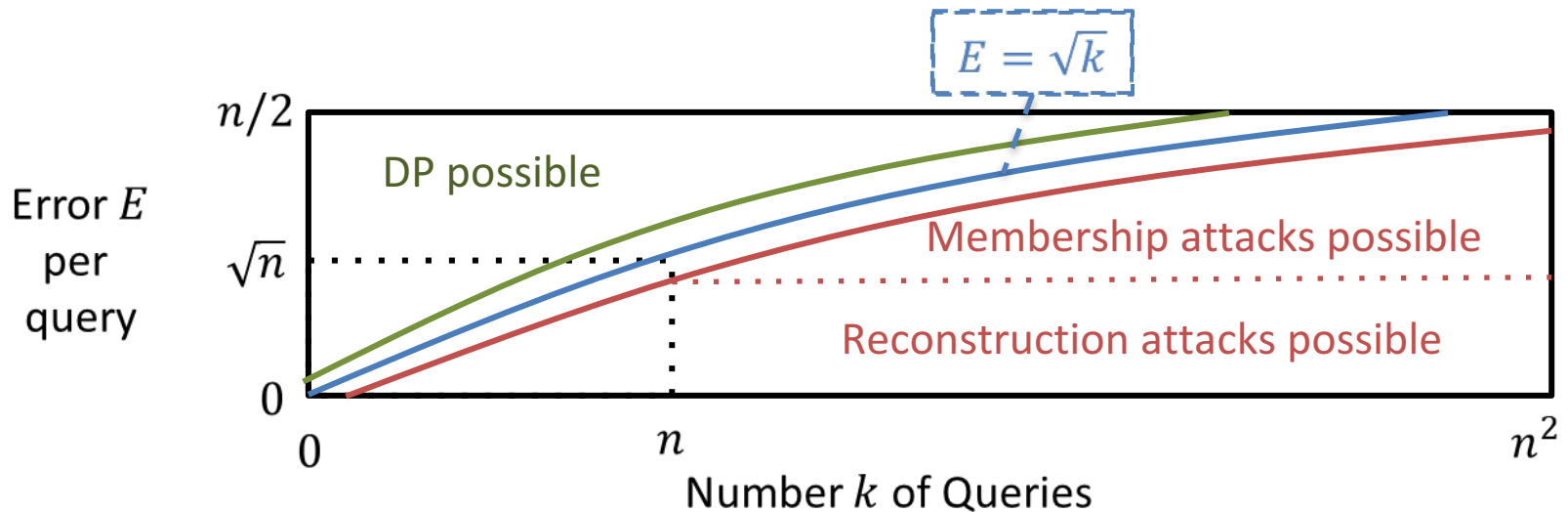
is $(\varepsilon', k \cdot \delta + \delta')$ -DP, for

$$\varepsilon' = O\left(\varepsilon \cdot \sqrt{k \cdot \log(1/\delta')}\right).$$

Queries vs. Accuracy Tradeoff

Using Laplace Mechanism to answer k queries, each with global sensitivity 1 (e.g. counts), under fixed privacy budget ϵ' :

- Set $\epsilon = 1/\tilde{O}(\sqrt{k})$ for each query (via Advanced Comp, hiding δ').
- Add noise of scale $E = 1/\epsilon \approx \tilde{O}(\sqrt{k})$ per query.



Note: DP prevents **all** membership & reconstruction attacks (not just those we've seen), e.g. $\Pr[\text{true pos}] \leq e^\epsilon \cdot \Pr[\text{false pos}] + \delta$

Doing Better than Composition

- Not all sequences of k queries require error growing as \sqrt{k} .
- **Example:** histograms
 - Let $B_1, \dots, B_k \subseteq \mathcal{X}$ be **disjoint** bins.
 - Define $q_j : \mathcal{X}^n \rightarrow \{0,1\}$ by $q_j(x) = \#\{i : x_i \in B_j\}$.
 - Define $M(x) = (q_1(x) + Z_1, q_2(x) + Z_2, \dots, q_k(x) + Z_k)$ where the Z_j 's are independent $\text{Lap}(2/\varepsilon)$ or $\text{Geo}(2/\varepsilon)$.
 - Then M is ε -DP.
- **Amazing result:** with **correlated** noise, can answer k **arbitrary** bounded averaging queries on a **finite** data universe \mathcal{X} w/error

$$\alpha = O \left(\frac{\sqrt{\log|\mathcal{X}| \cdot \log(1/\delta)} \cdot \log k}{\varepsilon n} \right)^{1/2}$$