

# Lecture 5 Private selection (Part II), Privacy Loss RV and Advanced Composition

Yu-Xiang Wang



**COMPUTER SCIENCE**

UC SANTA BARBARA

*Computing. ReInvented.*

# A few logistic notes

- Time to think about your course project
  - I shared a list of ideas on Piazza
  - Form your own team (up to 3 people), or do an independent project.
  - Discuss your idea with me (piazza / email / in-person)!
- Scribing:
  - Scribes for lecture 1 and 2 now available.
  - Send me the latex file when you are done.

# Recap: last lecture

- Sparse Vector
  - Everything in the proof works for general low-sensitivity queries (possibly non-linear queries)
- The problem of linear query release
- Private multiplicative weights
  - Use sparse vector (“NumericSparse”)
  - A cute application of a no-regret online learning method

# Recap: Private Multiplicative Weight

- Learn a synthetic dataset while answering queries.
- Use SVT to check if the error of the synthetic data is large.

$$O\left(\frac{L}{\alpha^2}\right) \approx 5/Q$$

- Show that the number of rounds is  $1/\alpha^2$  by the regret bound, thus after that many rounds of queries, the synthetic dataset will take over.

# Summary of the problem of private query release

$$|\mathcal{X}| = 2^d$$

	Laplace (release query)	Laplace (release data / contingency table)	Private Multiplicative Weights
Error (normalized query)	$\frac{k \log(k/\delta)}{n\epsilon}$	$\frac{\sqrt{ \mathcal{X} } \log(k/\delta)}{n\epsilon}$	$\left( \frac{\log^{d}( \mathcal{X} ) \log(k/\delta)}{n\epsilon} \right)^{1/3}$
Computational complexity (per query)	$O(n)$	$O( \mathcal{X} )$	$O(\max\{ \mathcal{X} , n\})$



# Recap: Exponential mechanism

- Global sensitivity of the utility function

$$\underline{\Delta u} \equiv \max_{r \in \mathcal{R}} \max_{x, y: \|x-y\|_1 \leq 1} |u(x, r) - u(y, r)|.$$

- The exponential mechanism samples an output from a “Gibbs distribution”:

$$\mathcal{M}(x) \sim p(r|x) \propto \underline{\exp\left(\frac{\varepsilon u(x, r)}{2\Delta u}\right)}$$

$$\frac{e^{\frac{\varepsilon u(x, r)}{2\Delta u}}}{\sum_{r \in \mathcal{R}} e^{\frac{\varepsilon u(x, r)}{2\Delta u}}}$$

- Proof:** 1. Bound the ratio of the above exponentiated utility function (up to scale).  
2. Bound the ratio of the normalization constant.

# This lecture

- Privacy selection (Part II)
  - Utility of Exponential mechanism.
  - Application: SmallDB
  - ReportNoisyMax
- Advanced Composition
  - Privacy loss random variable
  - Advanced composition for pure-DP
  - Linear Query Release under Approximate DP

# Readings:

- Report Noisy Max / Exponential mechanism
  - Dwork and Roth 3.3 – 3.4
- SmallDB
  - Dwork and Roth 4.1
- Advanced Composition for pure-DP
  - Vadhan 2.2. (Specifically, Lemma 2.4)



# Randomized response and Laplace mechanism are instances of exponential mechanisms!

- Randomized Response

RR outputs  $\left\{ \begin{array}{l} x \text{ w.p. } \frac{e^\epsilon}{1+e^\epsilon} \\ 1-x \text{ w.p. } \frac{1}{1+e^\epsilon} \end{array} \right.$

$$u(x,y) = \mathbb{1}(y=x)$$
  

$$P(y|x) \propto e^{\epsilon \mathbb{1}(y=x)}$$
  

$$\Delta u = 1$$

- Laplace mechanism

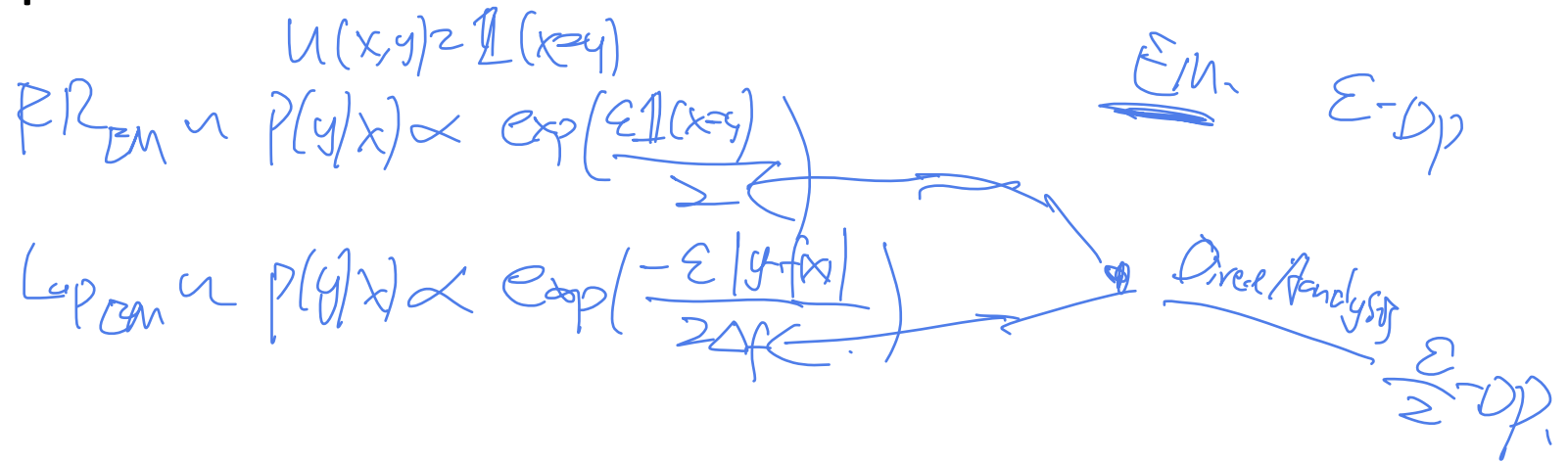
Lap Mech  $f(x) \Delta f$   
 outputs  $y := f(x) + \text{Lap}(\frac{\Delta f}{\epsilon})$

$$P(y|x) = \frac{\epsilon}{2\Delta f} \exp\left(-\frac{\epsilon |f(x)-y|}{\Delta f}\right)$$
  

$$u(x,y) = -|f(x)-y|$$
  

$$\Delta u \leq \Delta f$$

Detour: Applying the results of exponential mechanism to these two.



# Detour: What is really happening?

There are two different types of  
“exponential mechanism”

- Type I (“exponential mechanism” with insensitive utility function)

$$\mathcal{M}(x) \sim p(r|x) \propto \exp\left(\frac{\varepsilon u(x, r)}{2\Delta u}\right)$$

- Type II (“exponential mechanism” with insensitive log-probabilities)

$$\mathcal{M}(x) \sim p(r|x) = \exp\left(\frac{\varepsilon \tilde{u}(x, r)}{2\Delta \tilde{u}}\right)$$

$\forall u \text{ st } \Delta u$   
 $\exists \tilde{u} \text{ st } \Delta \tilde{u} \leq 2\Delta u$

Recommended further reading: Dong’s blog on exponential mechanism

<https://dongjs.github.io/2020/02/10/ExpMech.html>

Also, Durfee and Rogers (2019): <https://arxiv.org/abs/1905.04273>

# Utility of exponential mechanism

**Theorem 3.11.** Fixing a database  $x$ , let  $\mathcal{R}_{\text{OPT}} = \{r \in \mathcal{R} : u(x, r) = \text{OPT}_u(x)\}$  denote the set of elements in  $\mathcal{R}$  which attain utility score  $\text{OPT}_u(x)$ . Then:

$$\Pr \left[ u(\mathcal{M}_E(x, u, \mathcal{R})) \leq \text{OPT}_u(x) - \frac{2\Delta u}{\varepsilon} \left( \ln \left( \frac{|\mathcal{R}|}{|\mathcal{R}_{\text{OPT}}|} \right) + t \right) \right] \leq \frac{e^{-t}}{\delta}$$

**Proof:**

$$P[M(x) \notin \mathcal{R}_\alpha] = \sum_{r \in \mathcal{R}_\alpha} P[M(x) = r] = \frac{\sum_{r \in \mathcal{R}_\alpha} \exp\left(\frac{\varepsilon u(x, r)}{2\Delta u}\right)}{\sum_{r \in \mathcal{R}} \exp\left(\frac{\varepsilon u(x, r)}{2\Delta u}\right)}$$

$$\leq \frac{\sum_{r \in \mathcal{R}_\alpha} \exp\left(\frac{\varepsilon(\text{OPT} - \alpha)}{2\Delta u}\right)}{\sum_{r \in \mathcal{R}_0} \exp\left(\frac{\varepsilon \text{OPT}}{2\Delta u}\right)} \leq \frac{|\mathcal{R}|}{|\mathcal{R}_0|} \exp\left(\frac{-\varepsilon \alpha}{2\Delta u}\right)$$

$\frac{|\mathcal{R}|}{|\mathcal{R}_0|} = \frac{|\mathcal{R}|}{|\mathcal{R}_{\text{OPT}}|}$

$\alpha = \frac{2\Delta u}{\varepsilon} \left( \ln \left( \frac{|\mathcal{R}|}{|\mathcal{R}_{\text{OPT}}|} \right) + t \right)$

# Applying Exponential Mechanism to (offline) Linear Query Release

- Given a fixed set of queries of size  $k = |Q|$
- Run exponential mechanism to select a dataset most consistent with the answers to these queries

---

**Algorithm 4** The Small Database Mechanism

---

**SmallDB**( $x, Q, \varepsilon, \alpha$ )

Let  $\mathcal{R} \leftarrow \{y \in \mathbb{N}^{|\mathcal{X}|} : \|y\|_1 = \frac{\log |Q|}{\alpha^2}\}$

Let  $\tilde{u} : \mathbb{N}^{|\mathcal{X}|} \times \mathcal{R} \rightarrow \mathbb{R}$  be defined to be:

$$u(x, y) = - \max_{f \in Q} |f(x) - f(y)|$$

**Sample And Output**  $y \in \mathcal{R}$  with the exponential mechanism  $\mathcal{M}_E(x, u, \mathcal{R})$

---

# Analyzing the smallDB mechanism

- Privacy guarantee follows from that of the exponential mechanism

*f: normalized linear query*

$$f(x) = \frac{1}{n} \sum_i x_i$$

$$Q \in [0, 1]^n$$

- Analyzing the sensitivity  $u(x, y) = - \max_{f \in Q} |f(x) - f(y)|$

$$\Delta u = \max_{y, x, x'} |u(x, y) - u(x', y)|$$

$$= \max_y \max_{x, x'} \left| - \max_{f \in Q} |f(x) - f(y)| + \max_{f \in Q} |f(x') - f(y)| \right|$$

$$\leq \max_y \max_{x, x'} \left[ \max_{f \in Q} |f(x) - f(x')| \right]$$

$$= \max_y \max_{f \in Q} \max_{x, x'} |f(x) - f(x')| \leq \max_{f \in Q} \Delta f = \frac{1}{n}$$

*EM: output y  
P(y|x)*

$$P(y|x) \propto \exp\left(-\frac{\max_{f \in Q} |f(x) - f(y)|}{\epsilon/n}\right)$$

~~EM: output y~~

# Utility of the SmallDB in answering linear queries

$$\text{OPT} \approx \frac{\log |Q|}{2\alpha^2}$$

- Notice that we are restricting the sample size
  - If  $\frac{\log |Q|}{\alpha^2} < n$ , we need to work out the optimal solution (even if we output argmin in the clear)

• Claim: There always exists a smallDB that is  $\alpha$  accurate.

*iid sample from the  $\Delta(\mathcal{X})$ ,  $m = \frac{\log |Q|}{\alpha^2}$*

- Idea: randomly sample the dataset (with replacement).

*for a fixed  $f \in Q$ , w.p.  $1-\delta$ ,  $\frac{1}{n} \sum_{i=1}^n f(\phi_i) = \frac{1}{n} \sum_{i=1}^m f(\phi_i)$*

*Union over all  $f \in Q$ , w.p.  $1-\delta$*

*Population mean*  $\left| \frac{1}{n} \sum_{i=1}^n f(\phi_i) - \frac{1}{m} \sum_{i=1}^m f(\phi_i) \right| \leq \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$  *Hoeffding's inequality*

*Empirical mean*  $\max_{f \in Q} \left| \frac{1}{n} \sum_{i=1}^n f(\phi_i) - \frac{1}{m} \sum_{i=1}^m f(\phi_i) \right| \leq \sqrt{\frac{\log |Q|}{\delta}} = \alpha$

*Existence*  $\frac{\log |Q|}{2\alpha^2} = m$

# Apply the utility theorem of exponential mechanism

$$\Delta u = \frac{1}{n}$$

w.p.  $1-\delta$

$$u(x, M(x)) \leq \frac{2\Delta u}{\epsilon} \log\left(\frac{|R|}{\delta}\right)$$

$$= \frac{2}{n\epsilon} \log\left(\frac{|R|}{\delta}\right)$$

$$u(x, M(x)) \geq \text{OPT} - \text{OPT} + u(x, r^*) - \frac{2}{n\epsilon} \log\left(\frac{|R|}{\delta}\right)$$

$$|R| = \frac{1}{\alpha^2} |X|^m = |X|^{\frac{m \log |R|}{\alpha^2}}$$

$$= \frac{2}{n\epsilon} \frac{\log |R|}{\alpha^2} \log |X|$$

$$\max_{f \in \mathcal{F}} |f(M(x)) - f(x)| \leq \alpha + \frac{2 \log |R| \log |X|}{n\epsilon \alpha^2}$$

$$1 - \frac{4}{n\epsilon \alpha^3} \log |R| \log |X| = 0$$

$$\alpha = \sqrt{\frac{\log |R| \log |X|}{n\epsilon}}$$



# Checkpoint: SmallDB vs Private Multiplicative Weights

Q

- Both achieve the same asymptotic error

$$O\left(\frac{\log(|Q|) \log(X) \log(B)}{n \epsilon}\right)$$

- MW also works for an online sequence of adaptively chosen query

$|Q|$  is replaced by  $K$

~~$|Q|$~~   $\Rightarrow K$

- Neither is computationally efficient

# Alternative algorithm for private selection: ReportNoisyMax

- For each  $r \in \mathcal{R}$

$$\hat{u}(r, x) = u(r, x) + \text{Lap}(2\Delta u/\epsilon)$$

- Output  $\arg \max_{r \in \mathcal{R}} \hat{u}(r, x)$

# Privacy analysis is similar to that of SVT

- What is the output space?  
 $r \in \mathcal{R}$
- When does the algorithm output  $r$ ?

$V_r \sim \text{Lap}(\frac{2\epsilon}{\Delta u})$  is added to  $u(x,r)$

$$P(M(x)=r) = P(u(x,r) > \max_{r' \neq r} u(x,r'))$$

- The same trick of change of variable applies.

$$\begin{aligned}
 P(M(x)=r) &= E \left[ \prod_{r' \neq r} P(u(x,r) + V_r > u(x,r') + V_{r'}) \right] \\
 &= \int P(V_r) \left[ \prod_{r' \neq r} \int_{-\infty}^{u(x,r) - u(x,r') - V_r} P(V_{r'}) dV_{r'} \right] dV_r \\
 &= \int P(V-2\Delta) \left[ \prod_{r' \neq r} \int_{-\infty}^{u(x,r) - u(x,r') - V - 2\Delta} P(V_{r'}) dV_{r'} \right] dV_r \leq \int \frac{P(V-2\Delta)}{K e^{\epsilon P(V)}} \left[ \prod_{r' \neq r} \int_{-\infty}^{u(x,r) - u(x,r') - V - 2\Delta} P(V_{r'}) dV_{r'} \right] dV_r
 \end{aligned}$$

$e^{\epsilon P(M(x)=r)}$

# Remarks about RNM

1. If  $u(x,r)$  is a count for each  $r$ , then you can improve RNM by a factor of 2

- More generally, it applies when  $u(x,r)$  is **monotonic**.
- Similar proof.

$x = x' \cup \{ \text{independent} \}$

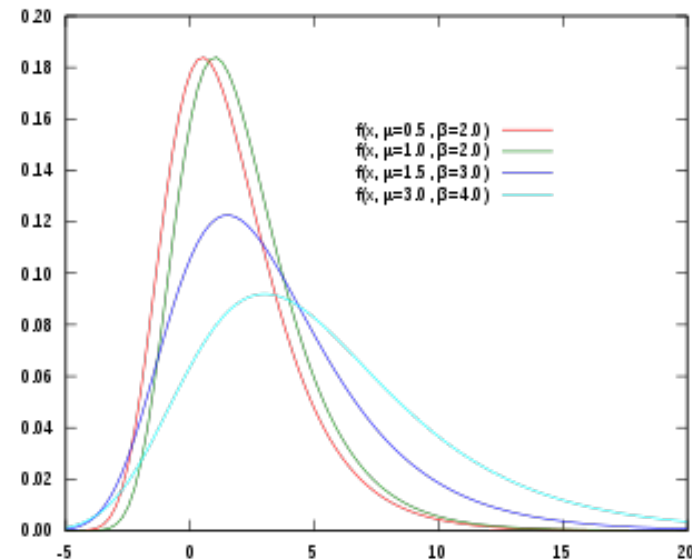
$u(x, r) \geq u(x', r)$

2. You can also add other noise

- Add (one sided) exponential noise
- Add Gumbel noise

|||  
exponential mechanism

"Gumbel Max trick"



# Application: Private voting (One vote per person)

Trump, Biden, Kenye

- One vote one person
- Two ways of releasing the results
  - Privately publish the histogram with Laplace mechanism
  - Privately publish the winner with ReportNoisyMax

Both  $\epsilon$ -DP

# Application: Private voting (Vote for as many people as you like!)

- Two ways of releasing the results
  - Privately publish the histogram with Laplace mechanism

$$\underline{\text{LCS-DP}} \longleftrightarrow \text{lap}\left(\frac{1}{\epsilon}\right)$$

- Privately publish the winner with ReportNoisyMax

$$\boxed{\epsilon\text{-DP}}$$

# Remainder of today's lecture

- Advanced composition
  - Privacy loss random variable
  - Prove advanced composition for pure-DP
  - Linear Query Release under Approximate DP

( $\epsilon, \delta$ )-DP

$\delta > 0$

# Recall: Summary of the problem of private query release

	Laplace (release query)	Laplace (release data / contingency table)	Private Multiplicative Weights
Error (normalized query)	$\frac{k \log(k/\delta)}{n\epsilon}$	$\frac{\sqrt{ \mathcal{X} } \log(k/\delta)}{n\epsilon}$	$\left(\frac{\log( \mathcal{X} ) \log(k/\delta)}{n\epsilon}\right)^{1/3}$
Computational complexity (per query)	$O(n)$	$O( \mathcal{X} )$	$O(\max\{ \mathcal{X} , n\})$

Can we do better under  $(\epsilon, \delta)$ -DP?

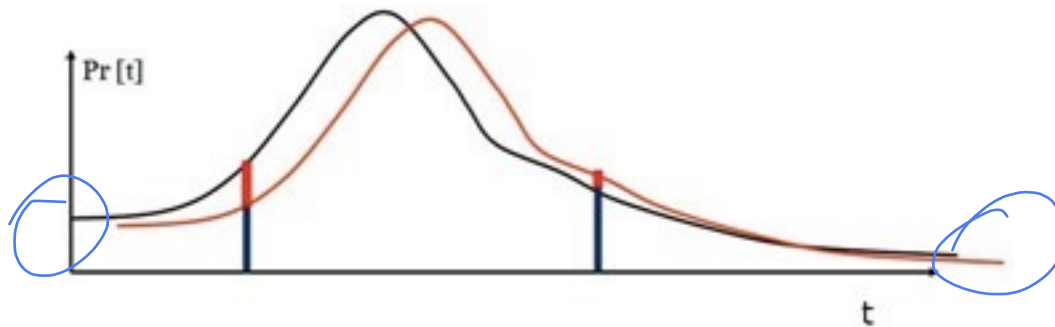


# Recap: Approximate DP

**Definition 2.4** (Differential Privacy). A randomized algorithm  $\mathcal{M}$  with domain  $\mathbb{N}^{|\mathcal{X}|}$  is  $(\epsilon, \delta)$ -differentially private if for all  $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$  and for all  $x, y \in \mathbb{N}^{|\mathcal{X}|}$  such that  $\|x - y\|_1 \leq 1$ :

$$\Pr[\mathcal{M}(x) \in \mathcal{S}] \leq \exp(\epsilon) \Pr[\mathcal{M}(y) \in \mathcal{S}] + \delta,$$

where the probability space is over the coin flips of the mechanism  $\mathcal{M}$ . If  $\delta = 0$ , we say that  $\mathcal{M}$  is  $\epsilon$ -differentially private.



# Advanced Composition

$(M_1, M_2, \dots, M_k)$

**Theorem:** The adaptive composition of  $k$   $(\varepsilon, \delta)$ -DP mechanisms satisfies  $(\tilde{\varepsilon}, \tilde{\delta})$ -DP where

$$\tilde{\varepsilon} = \varepsilon \sqrt{2k \log(1/\delta')} + 2k\varepsilon^2, \quad \tilde{\delta} = k\delta + \delta'$$

for any  $\varepsilon, \delta \geq 0, \delta' > 0$

Simple composition  $\tilde{\varepsilon} = k\varepsilon$      $\tilde{\delta} = k\delta$

$\varepsilon = \frac{\varepsilon'}{\sqrt{2k \log(1/\delta')}}$      $\delta = 0$

$$\tilde{\varepsilon} = \varepsilon' + \frac{2k(\varepsilon')^2}{\left(\frac{1}{\sqrt{2k \log(1/\delta')}}\right)^2} = \varepsilon' + \frac{(\varepsilon')^2}{\frac{1}{2 \log(1/\delta')}} \leq \varepsilon' + 2\varepsilon'$$

where  $\varepsilon' \leq \sqrt{\frac{1}{2} \frac{2}{\log(1/\delta')}}$

# Application to linear query release

- Laplace mechanism for release queries?

*w.p. 1-β*

$$\text{Lap}\left(\frac{\Delta f}{\epsilon}\right) \Rightarrow \epsilon_{\text{total}} = \sqrt{k} \cdot \epsilon, \quad \alpha = \frac{\sqrt{k} \log k / \beta}{n \epsilon}$$

- Laplace mechanism for releasing data?

$$\frac{\sqrt{|X|} \log k / \beta}{n \epsilon} \quad \chi$$

- Private multiplicative weights?

Comparison of  $N$  rounds of  $\left\{ \begin{array}{l} \text{Absc threshold} \\ \text{Laplace Mech} \end{array} \right\}$

$$\epsilon_0 = \frac{16 \log \frac{3k}{\beta}}{n \alpha}$$

$$N = \frac{16 \log |X|}{\alpha^2}$$

$2/\sqrt{N} \cdot \epsilon_0$ -DP Mechs

Assum  $\epsilon_0 \leq \frac{\epsilon}{\sqrt{N \log \frac{1}{\beta}}}$

$$\epsilon_{\text{total}} = 2\sqrt{N} \epsilon_0 = \frac{C \sqrt{\log |X|}}{n \alpha^2} \log \frac{3k}{\beta}$$

$$\alpha = \frac{C \sqrt{\log |X|} \log \frac{3k}{\beta}}{n \epsilon_{\text{total}}}$$

*parameter*

$$\left[ \begin{array}{l} \epsilon_{\text{total}} = N \cdot \epsilon_0 \cdot \frac{\log k}{\beta} \\ \alpha = \frac{\log |X|}{n \alpha} \end{array} \right]$$

# Privacy loss random variable

- PLRV is the log probability ratio as a random variable

$$\epsilon_{\mathcal{M}}^{x,x'} = \log\left(\frac{p(\mathbf{y})}{p'(\mathbf{y})}\right) \text{ where random variable } \mathbf{y} \sim \mathcal{M}(x).$$

- Tail bound of privacy loss r.v. implies DP

**Lemma 1** (Tail bound to  $(\epsilon, \delta)$ -DP conversion). Let  $\epsilon_{\mathcal{M}}^{x,x'}$  be the privacy loss RV defined above. If

$$\mathbb{P}(\epsilon_{\mathcal{M}}^{x,x'} > \epsilon) \leq \delta$$

for all pair of neighboring  $x, x'$  then  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -DP.

(You are to prove this in HW1.)

# Two more properties of privacy loss r.v.

- Expected value of a privacy loss is KL-divergence

$$E_{y \sim P_X} \left( \log \frac{P_X(y)}{P_{X'}(y)} \right) = D_{KL}(P_X || P_{X'}) = \int P_X(y) \log \frac{P_X(y)}{P_{X'}(y)}$$

- PLRV under composition

$$(M_1, M_2, \dots, M_k)$$

$$E_{(M_1, \dots, M_k)}^{X, X'}(y) = \log \frac{P_{M_1}(y_1) \cdot P_{M_2}(y_2) \cdot \dots \cdot P_{M_k}(y_k)}{P_{M_1}(y_1) \cdot P_{M_2}(y_2) \cdot \dots \cdot P_{M_k}(y_k)} = \sum_{i=1}^k E_{M_i}^{X, X'}(y_i)$$

*independent*

$M_i$  is chosen adaptively.  $y_i \sim y_{i-1}$

$$E_{M, h}^{X, X'}(y) = \log \frac{P_{M_1}(y_1) \cdot P_{M_2}(y_2 | M_1=y_1) \cdot \dots \cdot P_{M_k}(y_k | y_1 \dots y_{k-1})}{\dots}$$

*martingale*

# Proof Idea of Advanced Composition

- Observation 1: sometimes PLRV is positive, other times negative. They cancel with each other
- Observation 2: as  $k$  gets larger, the sum of PLRV concentrates around its mean.
- Observation 3: the adaptivity means that the PLRV will depend on the past

# Martingale

- We say that a sequence of r.v.  $X_1, \dots, X_n, \dots$  is a Martingale if for any  $n$

$$\mathbf{E}(|X_n|) < \infty$$

$$\mathbf{E}(X_{n+1} \mid X_1, \dots, X_n) = X_n.$$

- Example:
  - Random-walk: Total number of heads minus tails in  $n$  coin tosses

# Azuma-Hoeffding's inequality

- **Azuma-Hoeffding's inequality:** Assume  $X_1, \dots, X_n$  are **Martingale differences**

$$S_n = X_1 + \dots + X_n$$

$$\mathbb{P} [S_n \geq \epsilon] \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

- Apply Azuma-Hoeffding's inequality to our problem



# Proof for the advanced composition for pure DP mechanisms

- Fix  $x, x'$ , apply Azuma-Hoeffding's inequality

# Bounding the KL-divergence

- **Lemma** (Pinsker's inequality)

$$\|P - Q\|_1 \leq \sqrt{2D_{KL}(P\|Q)}$$

- **Corollary:** KL-divergence is nonnegative.
- Now's let's bound the expected value of PLRV:

# Next lecture

- Gaussian mechanism
- CDP and Renyi DP
- Composition of Gaussian mechanism