

# Lecture 9 Differentially Private Machine Learning

Yu-Xiang Wang



**COMPUTER SCIENCE**

UC SANTA BARBARA

*Computing. ReInvented.*

# Logistic notes

- Submit your HW1 and Project proposal if you haven't yet
- HW2 is on the course website
  - Except a code template for Q4 which I am still finalizing.
  - You should start working on the theoretical parts

# Recap: Last lecture

- Composition of mechanism-specific representations
  - RDP accountant
  - Fourier accountant
- Unified treatment via a **dominating** privacy loss random variable
  - And its characteristic function
- autodp: How you would represent DP mechanism and compute privacy loss

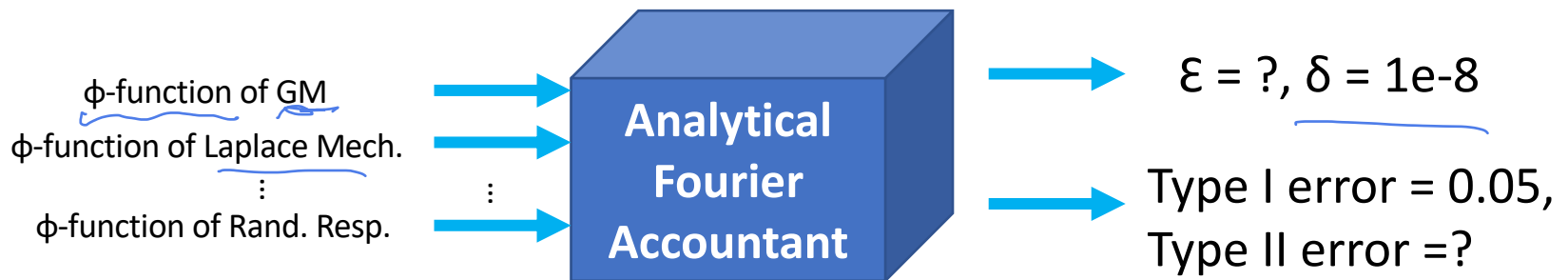
# Recap: Mechanism specific analysis and privacy accounting

	Functional view	Pros	Cons
Renyi DP [Mironov, 2017]	$D_\alpha(P  Q) \leq \epsilon(\alpha), \forall \alpha \geq 1$	Natural composition	lossy conversion to $(\epsilon, \delta)$ -DP.
Privacy profile [Balle and Wang, 2018]	$\mathbb{E}_q[(\frac{p}{q} - e^\epsilon)_+] \leq \delta(e^\epsilon), \forall \epsilon \geq 0$	Interpretable.	messy composition.
$f$ -DP [Dong et al., 2021]	Trade-off function $f$	Interpretable, CLT	messy composition.
PLD [Sommer et al., 2019, Koskela et al., 2020]	Probability density of $\log(p/q)$	Natural composition via FFT	Limited applicability.

Table 1: Modern functional views of DP guarantees and their pros and cons.

- Renyi DP is qualitatively different from approximate DP. Composition is quite natural with RDP.
- The composition of privacy-profiles and tradeoff functions are equivalent and somewhat messy.
  - The key to get it to work is to find a **dominating pair**
  - Using  $\phi$ -function representation, we get the natural composition of RDP, and the tightness of privacy-profile / tradeoff functions.

# Recap: Analytical Fourier accountant



- Composition: simply add up the log of phi functions
- Conversion to approx. DP via Levy's formula
- Conversion to tradeoff function via duality.

Zhu, Dong and W. (2020) <https://arxiv.org/abs/2106.08567>

# Recap: autodp: automating differential privacy computation (for both laypersons and experts)

- Users describe their randomized algorithm to autodp
- autodp focuses on computing privacy losses

**Open source project:**

<https://github.com/yuxiangw/autodp>

**pip install autodp**

# This lecture

- Introduction to Differentially Private Machine Learning
  - Problem setup and notations
  - Examples
  - A learning theoretic study of the problem
- Posterior sampling mechanism
  - When the loss functions are bounded
  - A new analysis of the when they are not

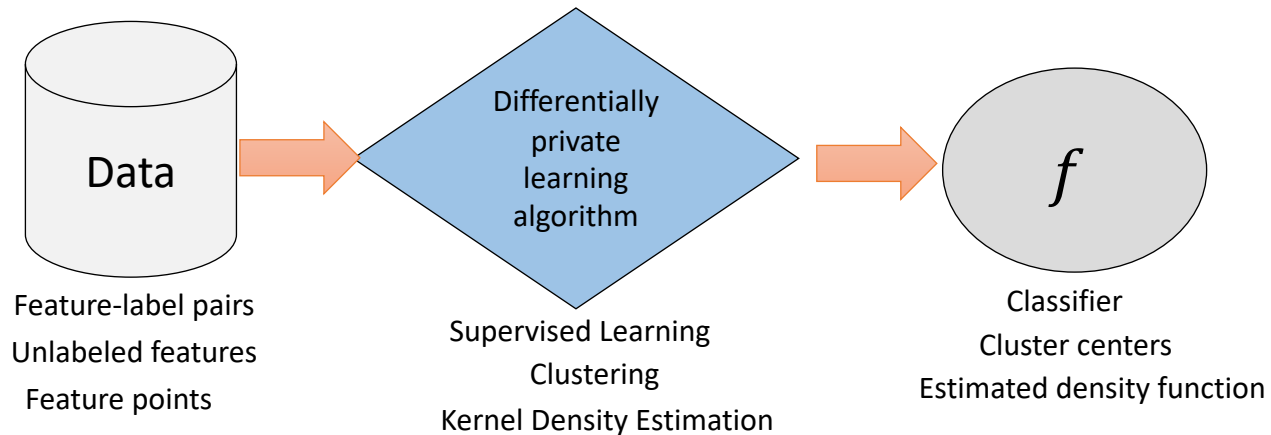
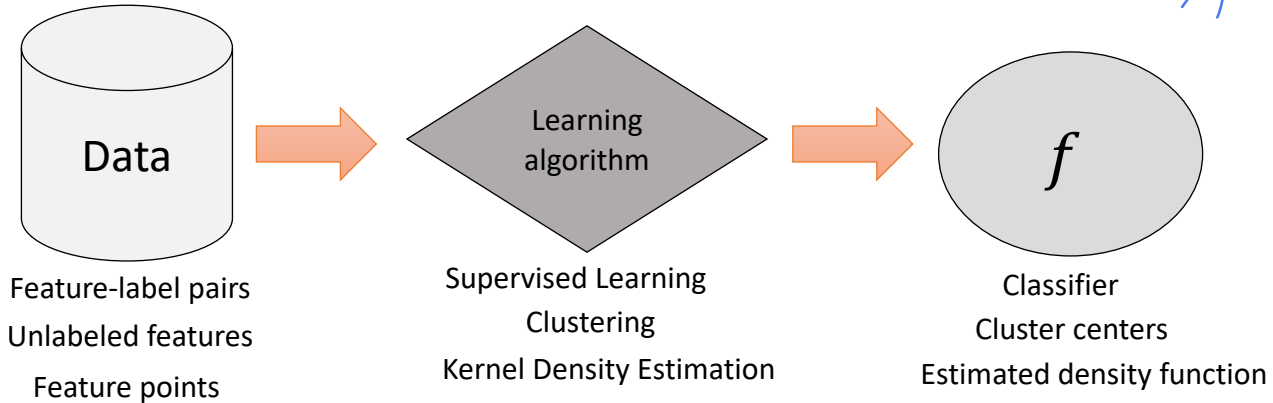
# (Optional) reading materials

- “Learning with differential privacy”:  
<https://arxiv.org/pdf/1502.06309.pdf>
- “Privacy-for-free: Posterior Sampling and Stochastic Gradient Monte Carlo” <https://arxiv.org/abs/1502.07645>
- “Differential privacy without sensitivity” , Minami et al.  
<http://papers.neurips.cc/paper/6050-differential-privacy-without-sensitivity.pdf>
- Other (not covered in this lecture)
  - “What can we learn privately?” <https://arxiv.org/abs/0803.0924>

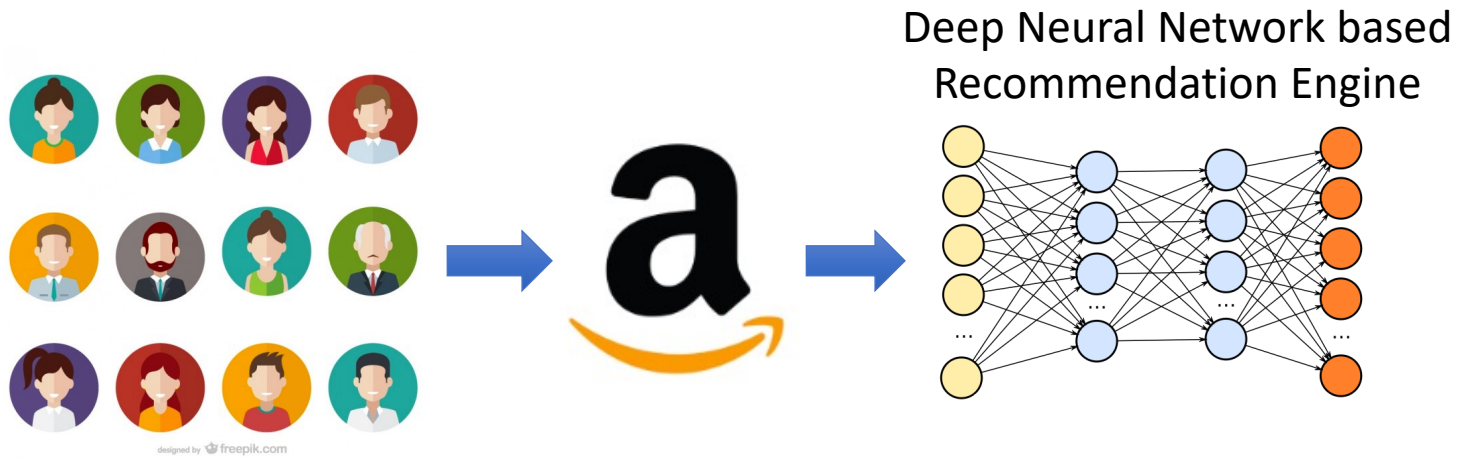


# Differentially Private Machine Learning

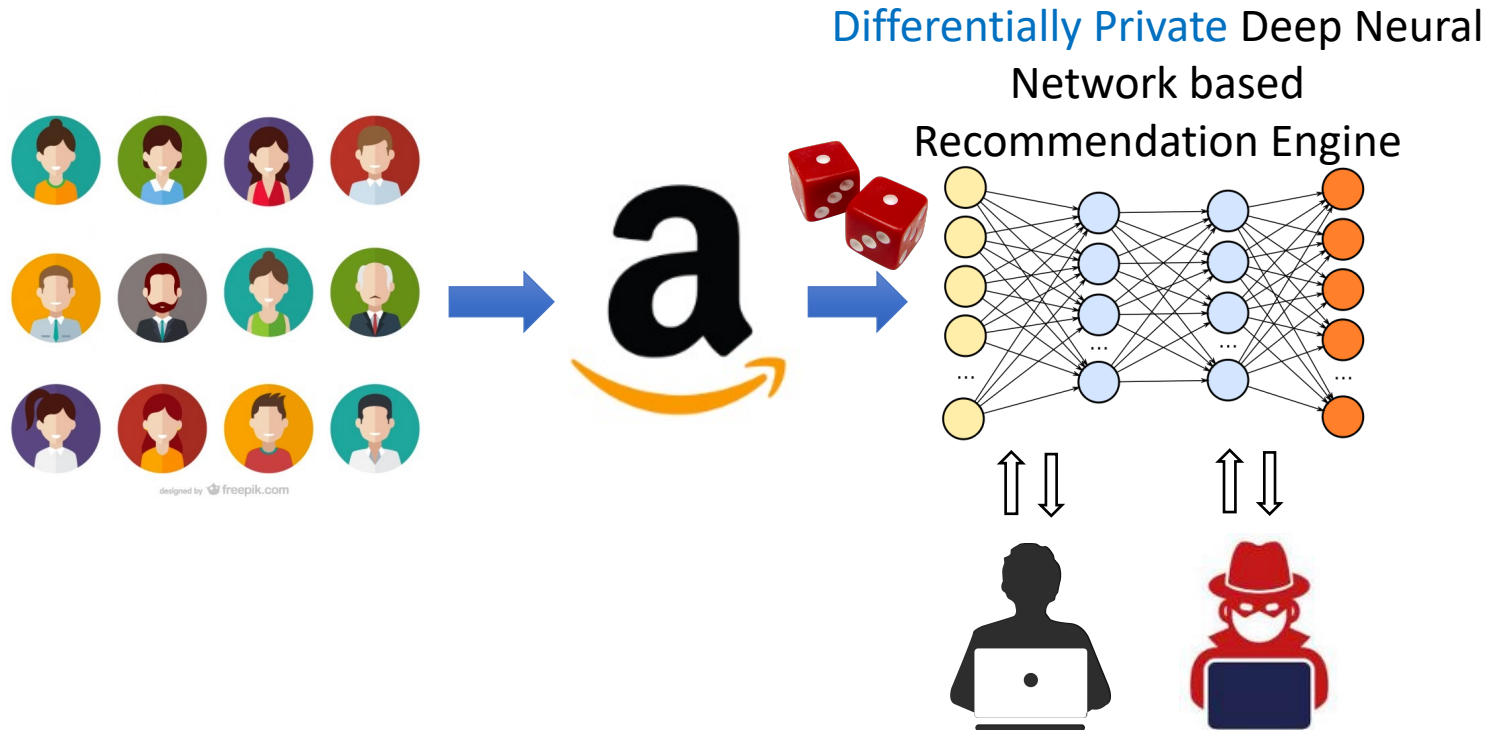
$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \cup D$



# Example: Recommender System



# Example: Recommender System



“If your recommendation engine is private, then an adversary can’t infer whether a particular user was present”

# A closely related setting: Federated Learning

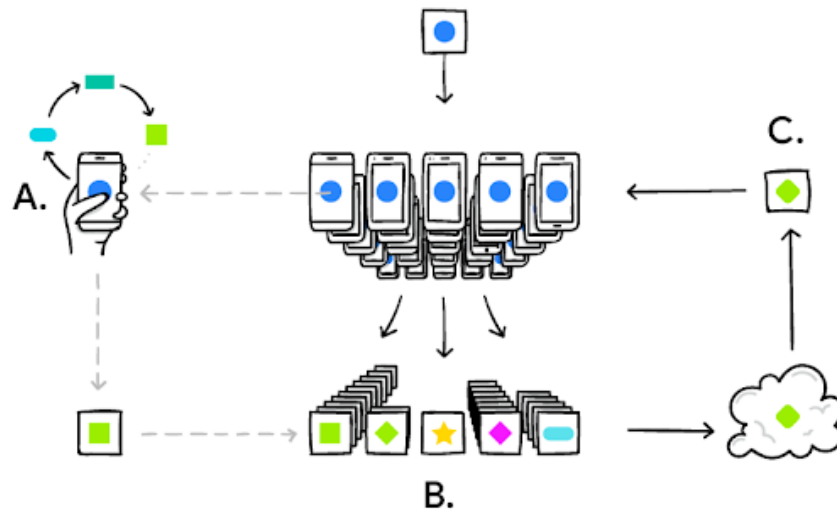


Illustration extracted from McMahan and Ramage (2017)

## Additional considerations

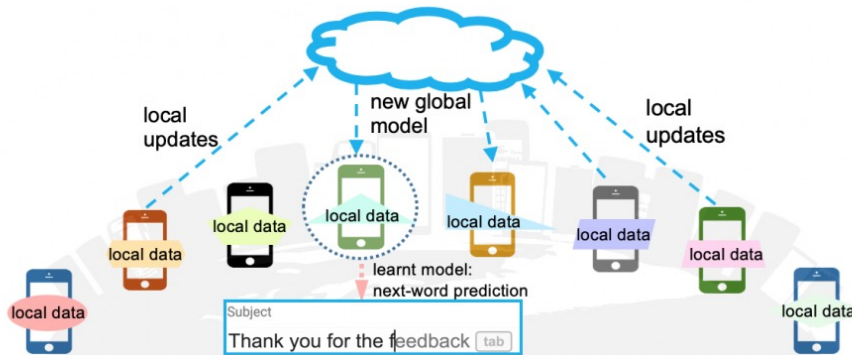
- Communication cost
- Size of the model
- Rounds of adaptivity

MPC ensures the security of distributed computation.

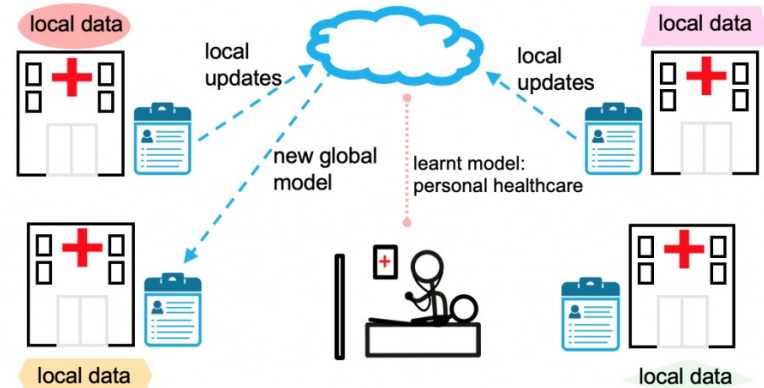
DP eliminates the risk from the output itself.

# Two regimes of Federated Learning with DP

Regime 1 requires agent level DP



Regime 2 requires instance-level DP



(illustration taken from [Tian Li et al, 2019](#))

**\*Need flexible tools for algorithm design.**

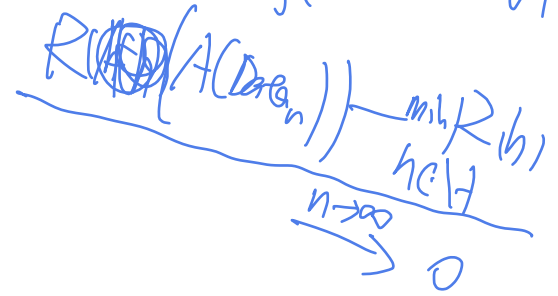
**\*Need tight privacy accounting to standardize empirical benchmarking of DP methods.**

Zhu, Y., Yu, X., Tsai, Y. H., Pittaluga, F., Faraki, M., & W. (2020). Voting-based Approaches For Differentially Private Federated Learning. *arXiv preprint arXiv:2010.04851*.

# Notations and problem setup

- Data space  $X \times Y = Z$   
 $(x, y) \in X \times Y$

Goal of learning: when  $n \rightarrow \infty$



- Hypothesis space / model space

$H: h \in H$   
 $\uparrow$   
 classifier

$H$  is described by  $\mathbb{R}^d$

- Loss functions

$n \longmapsto \mathcal{D} \in \mathcal{H}$

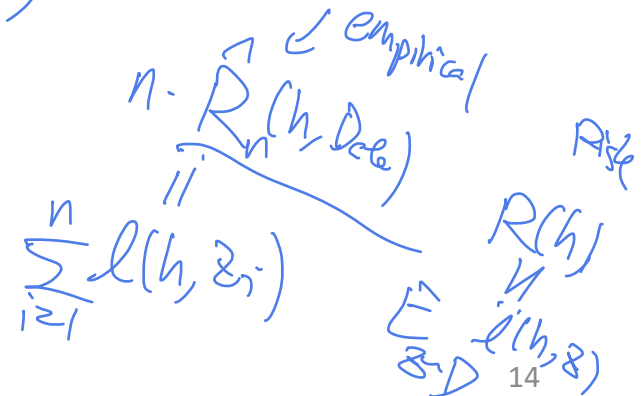
0-1 loss:  $\mathbb{1}(h(x) \neq y) = \underline{\underline{l(h, (x, y))}}$

- Learning algorithm

$A: (Z_1, \dots, Z_n) \rightarrow H$

Empirical Risk Minimization  
 (ERM)

$\hat{h} = \underset{h \in H}{\operatorname{argmin}}$



# Example 1: Linear / Logistic regression

- Data space

$$x \in \mathbb{R}^d, \quad \gamma = \mathbb{R}$$

- Hypothesis space / model space

$$h(x) = \theta^T x$$

- Loss functions

$$\text{Square loss } (y - x^T \theta)^2$$

- Learning algorithm

$$\min_{\theta} \sum_i (y_i - x_i^T \theta)^2$$
$$\| \vec{y} - X \theta \|_2^2$$

$$\gamma = \{0, 1\}$$

$$h(x) = \text{sigmoid}(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$
$$h(x) = \frac{e^{\theta^T x}}{1 + e^{\theta^T x}} = \hat{p}(x)$$

Cross entropy:  $-(y \log \hat{p}(x) + (1-y) \log (1-\hat{p}(x)))$

$$\min_{\theta} \sum_{i=1}^n \ell(\theta; (x_i, y_i)) + \lambda \|\theta\|_2^2$$

# Example 2: PAC learning / binary classification

- Data space

$$\mathcal{X}, \mathcal{Y} = \{0, 1\}$$

- Hypothesis space / model space

$$\mathcal{H}$$

- Loss functions

$$\mathbb{1}(h(x) \neq y)$$

- Learning algorithm



# Example 3: k-means clustering

- Data space

$$X = \mathbb{R}^d$$

- Hypothesis space / model space

$$\Theta = \left\{ \theta = (\theta_1, \theta_2, \dots, \theta_k) \in \mathbb{R}^{d \times k} \right\}$$

- Loss functions

$$\min_{j \in \{1, \dots, k\}} (x_i - \theta_j)^2$$

- Learning algorithm

$$\sum_{i=1}^n$$



# Example 4: Recommender System / Matrix factorization

- Data space
- Hypothesis space / model space

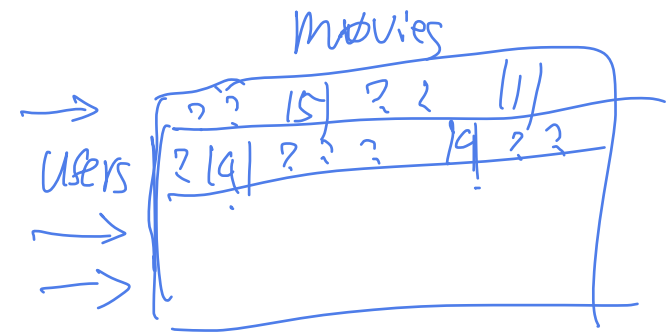
$$\mathcal{H} = \mathbb{R}^{m \times k}$$

- Loss functions

$$l(h, x) = \sum_{i \in \text{non-zero}} (x_i - U_i^T U_x)^2$$

$\uparrow$   
 function of  $U$

- Learning algorithm



$$\min_U \sum_i \|x_i - U_i^T V\|^2$$

$\rightarrow$   
 MSE

# What do we know about machine learning methods theoretically?

- Learnability learnable iff  $\exists A$  s.t.  $\forall D$   
 $(z_1, \dots, z_n \in D)^n, R(A(z_1, \dots, z_n)) - R(h^*) \xrightarrow{n \rightarrow \infty} 0$

- Sample complexity of learning

for  $A$ :  $n = O(\frac{1}{\alpha^2})$  ~~this~~ means  $\frac{|R(A(z_{1:n})) - R(h^*)|}{\text{func}(\alpha)} \leq \alpha$

- Key component: Uniform convergence / VC-theory

u.h.p.  $\sup_{h \in H} |R(h) - R(h)| \leq \alpha(n)$

$h = A(\text{Data})$

# The fundamental problems in learning with differential privacy

- Private learnability

- What problems are learnable with DP

$\epsilon < \infty$

But a strictly  $\epsilon$ -DP.

- Sample complexity of private learning

- Among those that are privately learnable, what is the number of samples needed to learn

$$n = \text{func}(\uparrow \text{Complexity}(H), \epsilon)$$

- Efficient algorithms?

# PAC learning with finite hypothesis space

$$e^{-\frac{n^2 \beta^2}{\sum (b_i - a_i)^2}}$$

" $\beta$ "

$$0 \leq \ell(h(x,y)) \leq 1$$

- Standard statistical learning via Hoeffding's inequality and a union bound

fix  $h \in \mathcal{H}$ ,  $\left| \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i)) - \mathbb{E}(\ell(h, (X, Y))) \right| \leq \sqrt{\frac{2 \log \frac{2}{\beta}}{n}}$  iid  $(x_1, y_1), \dots, (x_n, y_n)$   
w.p.  $1 - \beta$

for all  $h \in \mathcal{H}$

$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} R_n(h)$

$$\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum \ell(h, (x_i, y_i)) - \mathbb{E}(\ell(h, (X, Y))) \right| \leq \sqrt{\frac{2 \log |\mathcal{H}|}{n \beta}}$$

w.p.  $1 - \beta$

$$\underline{R(\hat{h}) - R(h^*)} \leq \underbrace{R(\hat{h}) - R_n(\hat{h})}_{\leq \sqrt{\frac{2 \log |\mathcal{H}|}{n \beta}}} + \underbrace{R_n(\hat{h}) - R_n(h^*)}_{\leq 0} + \underbrace{R_n(h^*) - R(h^*)}_{\leq \sqrt{\frac{2 \log |\mathcal{H}|}{n \beta}}}$$

# Generic private optimization algorithm via exponential mechanism?

- Recap: utility of the EM

**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:

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

- Let's apply this to the PAC learning problem with finite hypothesis space

$$\hat{R}(h) - \min_h R(h) \leq \frac{2 \cdot 2}{n\varepsilon} \left( \log |\mathcal{H}| + \log \frac{1}{\beta} \right)$$

$$R(h_{\text{ERM}}) - R(h^*) \leq O\left(\frac{\log |\mathcal{H}|}{n}\right)$$

up to  $\beta$

$$R(h_{\text{ERM}}) - R(h^*) \leq \sqrt{\frac{\log |\mathcal{H}|}{n}} + \frac{\log |\mathcal{H}|}{n\varepsilon}$$

$\hat{h}$  u.p.  $\propto e^{-\frac{\Delta R(h)}{2\varepsilon}}$  [E-DP]  $\Delta = \frac{1}{n}$   
 $R = \frac{1}{n} \sum \ell_i$

# What happens with continuous hypothesis space (satisfying VC-dim is finite)?

- Short answer: No DP algorithms can learn the VC class.
- Example: Learning a threshold function.

$\mathcal{X} = [0, 1], \mathcal{Y} = \{0, 1\}$   $h \in [0, 1]$

$\mathcal{H}: h \in \mathcal{H}: h(x) = \mathbb{1}(x \geq h)$



$\epsilon < +\infty$ , there is no  $\epsilon$ -DP alg  $A$

s.t. 
$$\frac{R(A(D_{\text{train}})) - R(h^*)}{n} < 0.9$$

$\log(|\mathcal{H}|)$

$VC(\mathcal{H}) = 1$

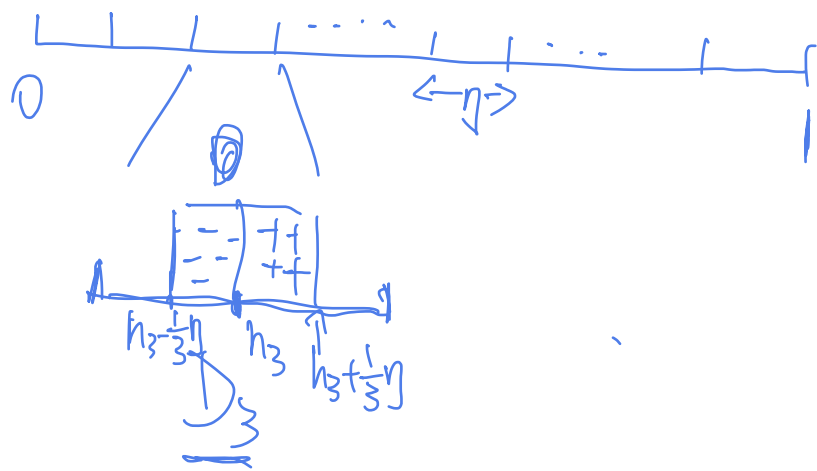
$\sup_{h \in \mathcal{H}} |R_n(h) - R(h)| \leq \frac{1}{n}$

if in addition  $n = O\left(\frac{1}{\epsilon^2}\right)$

$R(h^*) = 0$  then  $h = \frac{1}{n}$

# Any pure-DP algorithms will fail in learning threshold functions.

$\{ \text{clear } D_1, D_2, \dots, D_k, \text{ assume } A \text{ is successful on } D_2 \dots D_k$   
 $\eta := e^{-\epsilon n}, \quad k = e^{\epsilon n}$   
 $A \text{ fails on } D_1$



$$o\left(\frac{1}{\alpha}\right) \rightarrow \text{foo}$$

if  $A$  is successful on  $D_1$   
 then  $A(\text{Data}) \in [h_1 - \eta, h_1 + \eta]$

$$E[R(A(z_i), z_i)] \leq \alpha$$

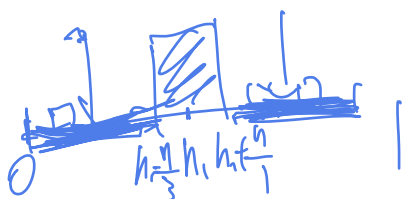
(Beimel et al. 2013) (Chaudhuri and Hsu. 2015)



$$Z_i \sim D_i^n$$

# Formalizing the lower bound.

$$\underbrace{P(A(Z_i) \notin [h_i - \frac{n}{3}, h_i + \frac{n}{3}])} \geq \sum_{i=1}^k P(A(Z_i) \in [h_i - \frac{n}{3}, h_i + \frac{n}{3}])$$



$$\underbrace{\text{gap privacy of } \epsilon\text{-DP}} \geq \sum_{i=1}^k e^{-\epsilon n} P(A(Z_i) \in [h_i - \frac{n}{3}, h_i + \frac{n}{3}])$$

$$\geq k \cdot e^{-\epsilon n} \cdot 0.9$$

$$= e^{\epsilon n} \cdot e^{-\epsilon n} \cdot 0.9$$

$$= 0.9$$

"padding argument" in lower bound

$$E[R(A(Z_i))] \geq 0.9$$

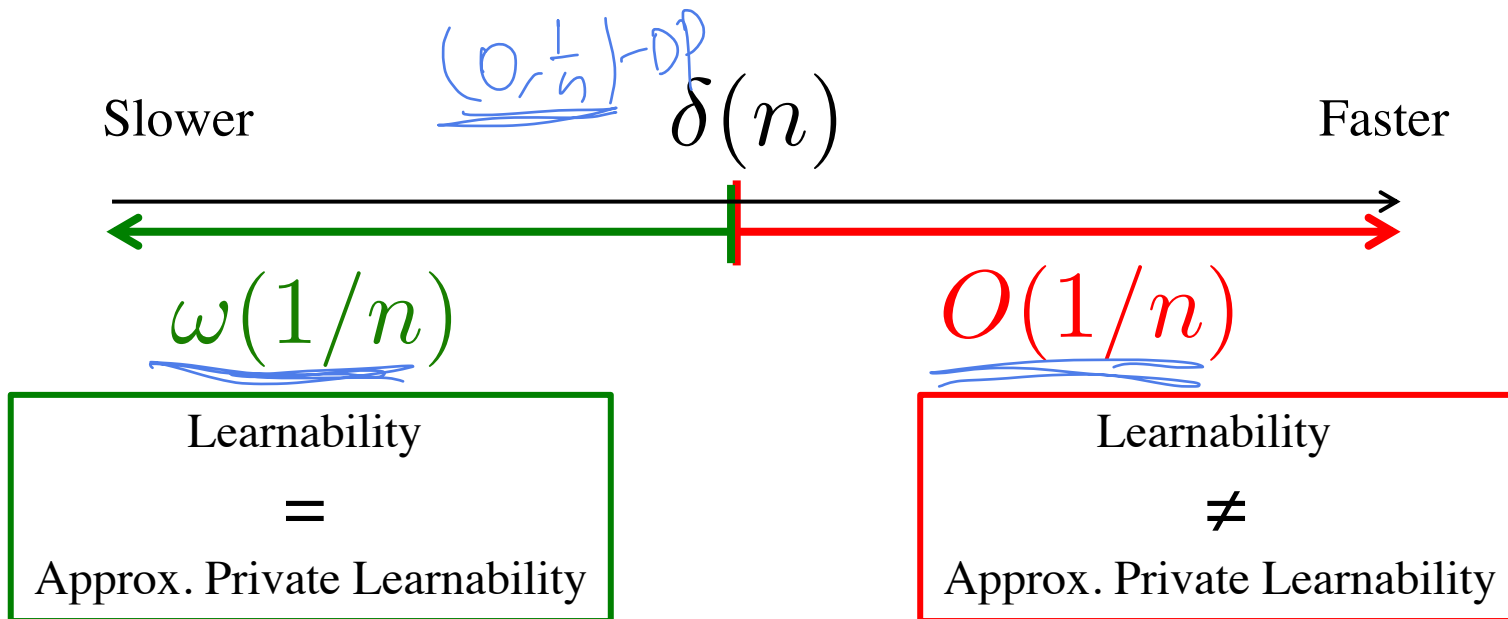
$$P[R(A(Z_i)) > t] \leq \frac{E[R]}{t}$$

$$\Downarrow$$

$$P[R(A(Z_i)) \leq t] \geq 1 - \frac{0.9}{t}$$

choose  $t = \frac{0.9}{\epsilon}$

Is this an issue of pure-DP being too restrictive? Can we learn a VC class under approximate DP?



- Bun et. al. "Differentially private release and learning of threshold functions." FOCS'15

# What is the key problem?

- Statistical learning requires a very strong notion of learning:
  - No assumptions on data distributions
  - Uniformly consistent learning for all distributions
- Differential privacy says that even if the data distributions are disjoint, they still need to be somewhat similar.
  - The construction relies on an exponentially large set of distributions (each converging to a point-mass).

# Remedies to this problem

- Assume the loss function is Lipschitz
  - 0-1 loss doesn't work, but we can use surrogate losses, e.g., logistic loss , hinge losses.
- Assume some regularity conditions on the probability distributions of data
  - e.g., Bounded probability density.
- Then most problems are learnable by an exponential mechanism.

# Check point: Learning **with differential privacy** in theory

- Finite hypothesis class
  - Exponential mechanism gives asymptotically vanishing additional error.
  - Does not benefit from being realizable.
- Continuous hypothesis class (bounded VC-dim)
  - No DP algorithm gives consistent learning
  - A Packing lower bound
- However, there are weak assumptions we can add
  - Lipschitz loss functions
  - data-distributions with bounded density

# Connections of the Exponential mechanism to Bayesian learning

- Bayesian philosophy

- I have a prior belief  $\pi(\theta)$
- When I collect data, I update my belief  $\pi(\theta | z_1, \dots, z_n)$   
 $Q(z_1, z_2, \dots, z_n)$

- Deriving the posterior using iid data:

$$\begin{aligned} \pi(\theta | z_1, \dots, z_n) &= \frac{P_D(z_1, \dots, z_n | \theta) \pi(\theta)}{\int P_D(z_1, \dots, z_n | \theta) \pi(\theta)} \\ &= \frac{\prod_{i=1}^n P_D(z_i | \theta) \pi(\theta)}{\int \prod_{i=1}^n P_D(z_i | \theta) \pi(\theta)} = e^{\frac{\sum_{i=1}^n \log P_D(z_i | \theta) + \log \pi(\theta)}{\sum_{i=1}^n \log P_D(z_i | \theta) + \log \pi(\theta)}} \end{aligned}$$

# Connections of the Exponential mechanism to Bayesian learning

- Getting one sample from the posterior distribution is equivalent to exponential mechanism

---

**Algorithm 1** One-Posterior Sample (OPS ) estimator

---

**input** Data  $X$ , log-likelihood function  $\ell(\cdot|\cdot)$  satisfying  $\sup_{\mathbf{x}, \theta} \|\ell(\mathbf{x}|\theta)\| \leq B$  a prior  $\pi(\cdot)$ .

Privacy loss  $\epsilon$ .

1. Set  $\rho = \min\{1, \frac{\epsilon}{4B}\}$ .

2. Re-define log-likelihood function and the prior  $\ell'(\cdot|\cdot) := \rho\ell(\cdot|\cdot)$  and  $\pi'(\cdot) := (\pi(\cdot))^\rho$ .

**output**  $\hat{\theta} \sim P(\theta|X) \propto \exp\left(\sum_{i=1}^N \ell'(\theta|\mathbf{x}_i)\right) \pi'(\theta)$ .

---

- Some classical results from statistics
  - Asymptotic normality of the Bayesian posterior
  - Bernstein-Von Mises Theorem.
- Utility of the OPS estimator

**Proposition 9.** *Under the same assumption as Proposition 8, if we set a different  $\epsilon$  by rescaling the log-likelihood by a factor of  $\frac{\epsilon}{4B}$ , then the the One-Posterior sample estimator obeys*

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{\text{weakly}} \mathcal{N}\left(0, \left(1 + \frac{4B}{\epsilon}\right)\mathbb{I}^{-1}\right),$$

*in other word, the estimator has an ARE of  $\left(1 + \frac{4B}{\epsilon}\right)$ .*

# Pros and cons of OPS

- Pros:

- Generically applicable
- Strong learning bounds under weak assumptions
- No need to change existing learning workflow

- Cons:

- It makes the distribution more diffused.
- Can get only one sample (hard to do inference)
- Requires bounded (clipped) loss functions
- Computationally inefficient in general



# Improved analysis of exponential mechanism with strong convexity

- Assume  $\pi(\theta) = e^{-r(\theta)}$  where  $r(\theta)$  is  $\mu$ -strongly convex, i.e., the prior is strongly log-concave.

- Assume that the loss-function is Lipschitz

$$| -\log p(x|\theta) + \log p(x|\theta') | \leq L \|\theta - \theta'\|$$

- Then  $\hat{\theta} \sim P(\theta|x_1, \dots, x_n) \propto e^{(-\tau \sum \log p(x_i|\theta) - r(\theta))}$

obeys  $(\epsilon, \delta)$ -DP if  $\tau = O\left(\frac{\epsilon\sqrt{\mu}}{L\sqrt{\log(1/\delta)}}\right)$

# Idea of the proof for the improved analysis of EM

- The privacy loss random variable is
- Apply the tail bound lemma
- The strong log-concavity + Lipschitz assumption ensures that  $\hat{\theta}$  satisfies a “Log-Sobolev Inequality”
  - Which ensures a subgaussian-like tail bound for all Lipschitz functions of  $\hat{\theta}$
  - And a bound on the KL-divergence.

# Reiterate the main points

- Bayesian learning
  - Just a scaling of posterior sampling
  - For bounded log-likelihood functions, exponential mechanism is a consistent learner
- Boundedness is not needed if we use a strong prior
  - A tighter analysis of the exponential mechanism
  - Use a prior to ensure that PLRV is concentrated

# Next lecture

- Convex empirical risk minimization
- Objective perturbation