# 365Days
# English Stydy Plan

# Privacy in ML

## Randomized Response

**Truth**    $x_1, \ldots x_n$        $x_i \in \{0, 1\}$

**Revealed**

$$y_i = \begin{cases} x_i & \text{w.p} & \dfrac{e^\epsilon}{1+e^\epsilon} \checkmark \end{cases} \} \text{truth}$$
$$\phantom{y_i = \begin{cases}} 1-x_i & \text{w.p} & \dfrac{1}{1+e^\epsilon} \checkmark \} \text{falsehood}$$

### Consider two datasets

**Truth**    $X = \{x_1, \ldots, x_n\}$      $x' = \{x_1', \ldots, x_n'\}$

**Assume**    $x_i = x_i'$    $\forall i = \{1, \ldots, n-1\}$

**Revealed**    Randomized response $\downarrow$

$$RR(x) = y = [y_1, \ldots, y_n]$$

$$RR(x') = y' = [y_1', \ldots, y_n']$$

$$P\left(RR(x) = b\right) = P\left([Y_1, Y_2 \cdots, Y_n] = [b_1, b_2, \ldots b_n]\right)$$

$[b_1 \cdots b_n]$

$$= \left[\prod_{i=1}^{n-1} P(y_i = b_i)\right] P(y_n = b_n)$$

$\rightarrow \text{①}$

$$P\left(RR(x') = b\right) = \prod_{i=1}^{n} P(y_i' = b_i)$$

$$= \left(\prod_{i=1}^{n-1} P(y_i' = b_i)\right) P(y_n' = b_n)$$

$$\Rightarrow = \left(\prod_{i=1}^{n-1} P(y_i = b_i)\right) \cdot P(y_n' = b_n)$$

$\rightarrow \text{②}$

---

$$\frac{P(y_n = b_n)}{P(y_n' = b_n)} = \begin{cases} \dfrac{e^\epsilon / 1 + e^\epsilon}{1 / 1 + e^\epsilon} = e^\epsilon & \text{if } b_n = X_n \\[4mm] \dfrac{1 / 1 + e^\epsilon}{e^\epsilon / 1 + e^\epsilon} = e^{-\epsilon} & \text{if } b_n = 1 - X_n \end{cases}$$

$$\frac{P(Y_n = b_n)}{P(Y_n' = b_n)} \leq e^{\epsilon}$$

$$\boxed{P(Y_n = b_n) \leq e^{\epsilon} P(Y_n' = b_n)} \quad -\text{ substitute}$$
$$\text{in } \textcircled{1}$$

$$P\left(RR(x) = b\right) = \left(\prod_{i=1}^{n-1} P(Y_i = b_i)\right) P(Y_n = b_n)$$

$$\leq \left(\prod_{i=1}^{n-1} P(Y_i' = b_i)\right) e^{\epsilon} P(Y_n' = b_n)$$

$$= e^{\epsilon} \prod_{i=1}^{n} P(Y_i' = b_i)$$

$$P(RR(x') = b)$$

$$\Rightarrow \boxed{\frac{P\left(RR(x)=b\right)}{P\left(RR(x')=b\right)} \le e^{\epsilon}}$$

## How good is this mechanism

Want $\quad \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} x_i$

$$\underset{\{0,1\}}{X_i} \xrightarrow{\text{Randomization}} \quad Y_i$$

$$E[Y_i]$$

$$\mathbb{E}[Y_i] = \left(\frac{e^t}{1+e^t}\right)X_i + \left(\frac{1}{1+e^t}\right)(1-X_i)$$

$$\frac{e^t X_i + 1 - X_i}{1 + e^t}$$

$$\mathbb{E}[Y_i] = X_i\left(\frac{e^t - 1}{e^t + 1}\right) + \frac{1}{1 + e^t}$$

$$X_i \to Y_i \to \boxed{Z_i}$$

$$Z_i = \left(Y_i - \frac{1}{1+e^t}\right)\left(\frac{e^t + 1}{e^t - 1}\right)$$

$$\mathbb{E}[Z_i] = ? \quad X_i \quad (\text{exercise})$$

$$X_1, \quad \cdots \cdots , X_n \qquad \bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$Y_1, \quad \cdots \cdots , Y_n$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$Z_1 \qquad\qquad\qquad Z_n$$

Guess for $\bar{X}$ is $\quad \bar{Z} = \frac{1}{n} \sum_{i=1}^{n} Z_i$

$$\mathbb{E}[\bar{Z}] = \mathbb{E}\left[ \frac{1}{n} \sum_{i=1}^{n} Z_i \right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}[Z_i]$$

$$= \frac{1}{n} \sum_{i=1}^{n} X_i = \bar{X}$$

$$\left| \frac{1}{n} \sum_{i=1}^{n} x_i - \frac{1}{n} \sum_{i=1}^{n} z_i \right|$$

with high probability

$$\lesssim O\left( \frac{1}{\epsilon \sqrt{n}} \right) \leftarrow \text{utility}$$

# TRUSTED CURATOR MODEL

Cynthia Dwork
[20'10]

Differential Privacy

Let $M: \mathcal{X}^n \to \mathcal{Y}$. Consider two "neighbouring" datasets $x, x' \in \mathcal{X}^n$

$M$ is $\epsilon$-D.P if for all $x, x'$ neighbouring and all $S \subseteq \mathcal{Y}$,

$$\frac{Pr\left(M(x) \in S\right)}{Pr\left(M(x') \in S\right)} \leq \left(e^{\epsilon}\right)$$

# LAPLACE MECHANISM

## Sensitivity

$$f: X^n \rightarrow \mathbb{R} \quad \text{(average)}$$

$$\Delta = \max_{\substack{\text{neighbouring} \\ \text{datasets} \\ x, x'}} \left| f(x) - f(x') \right|$$

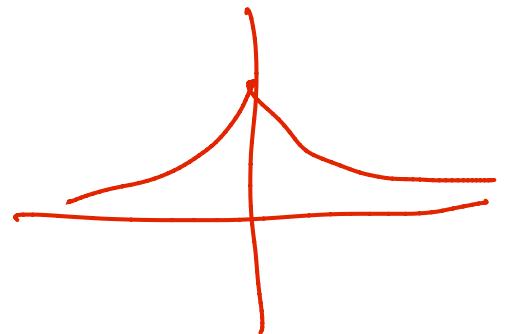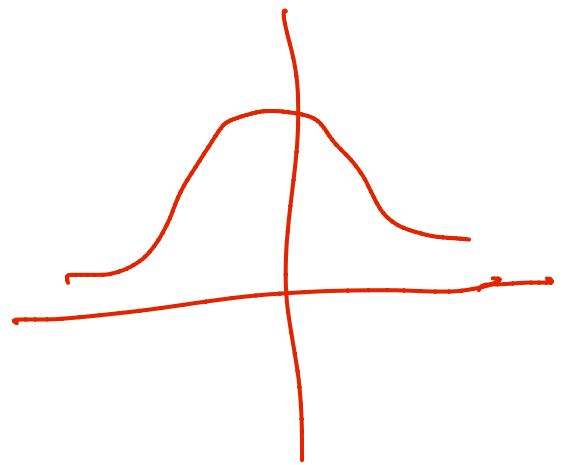$$\left| \frac{1}{n} \sum_{i=1}^{n} x_i - \frac{1}{n} \sum_{i=1}^{n} x_i' \right|$$

$$\left| \frac{1}{n} \sum_{i=1}^{n-1} x_i + \frac{1}{n} x_n \quad \frac{1}{n} \sum_{i=1}^{n} x_i - \frac{1}{n} x_n' \right|$$

$$= \quad \frac{1}{n} \left| x_n - x_n' \right|$$

$$\begin{array}{cc} 0 & 1 \\ 1 & 0 \end{array}$$

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

# Laplace distribution

$$f_{Lap}(x) = \frac{1}{2b} e^{\frac{-|x-y|}{b}}$$

$$\hat{x} = \frac{1}{n} \sum_{i=1}^{n} x_i + \eta \quad \longrightarrow \quad \text{Laplace}\left(0, \frac{\Delta}{\epsilon}\right)$$

$$E\left[\hat{x}\right] = \bar{x}$$

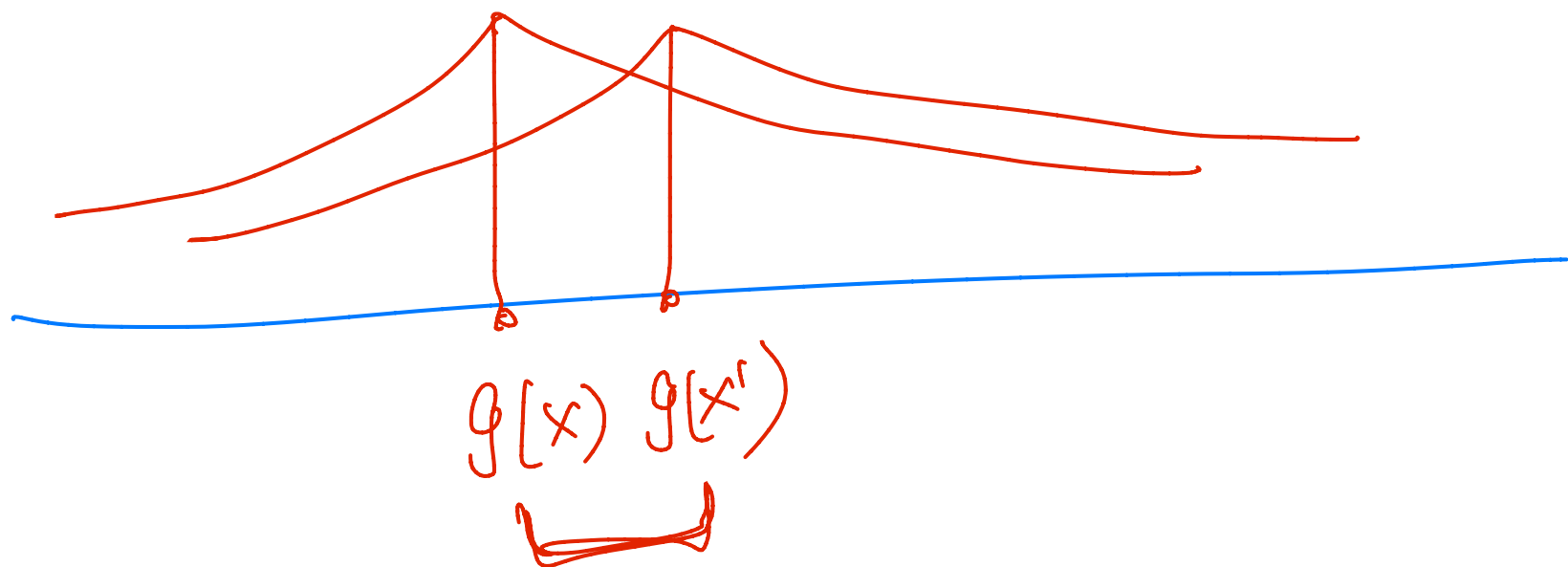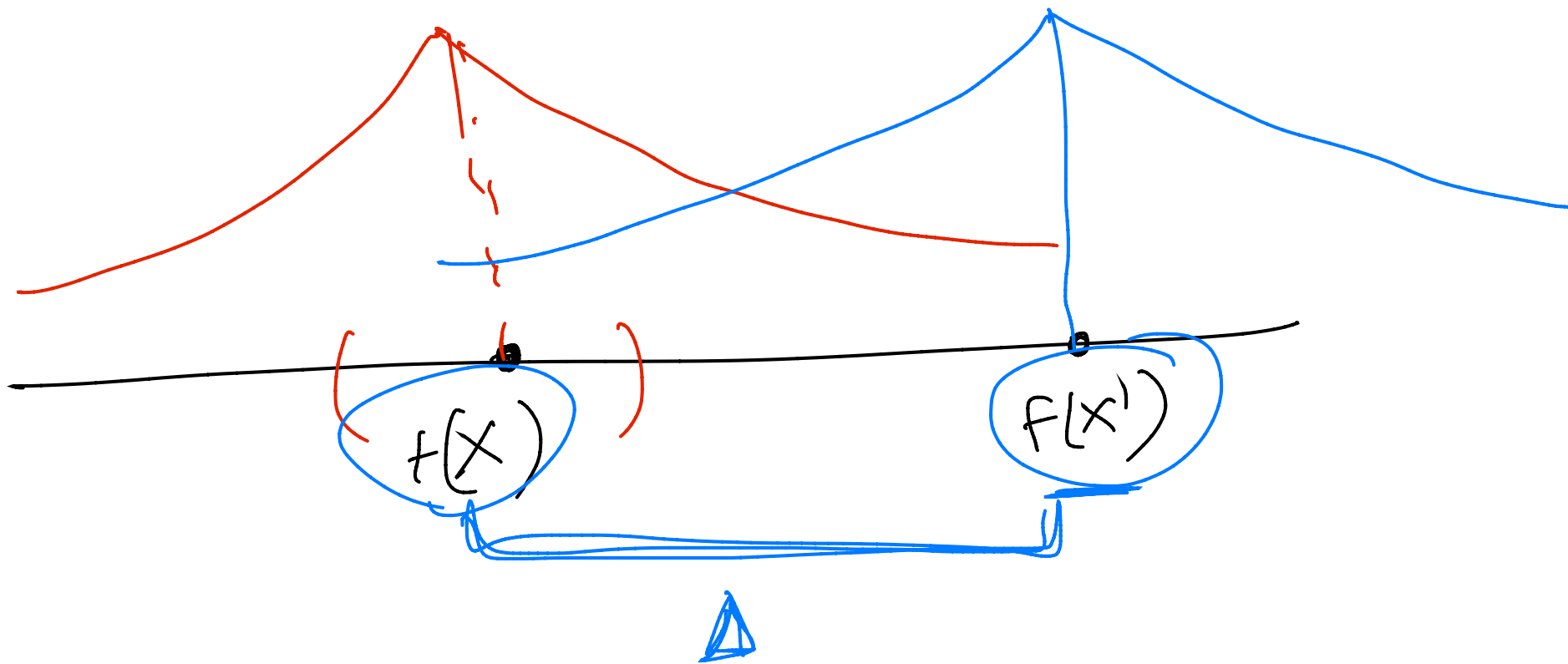Can argue Laplace mechanism is

$\epsilon$-DP ✓    [EXERCISE].

$\bar{x}$  $\bar{x}$

Utility

$$P\left( \left| \frac{1}{n}\sum_{i=1}^{n} x_i - \left( \frac{1}{n}\sum_{i=1}^{n} x_i + \eta \right) \right| \geq t \right)$$

$$\underset{LAP}{O\left( \frac{1}{\epsilon n} \right)} \quad \Bigg| \quad \underset{RR}{O\left( \frac{1}{\epsilon \sqrt{n}} \right)}$$

$f(x)$ $f(x')$

$g(x)$ $g(x')$

# Approximate DP

$M: x^n \rightarrow y$ is $(\epsilon, \delta)$ D.P

if $\forall$ neighbouring $x, x' \in x^n$

and all $S \subseteq y$

$$P[M(x) \subseteq S] \leq e^{\epsilon} P[M(x') \subseteq S] + \delta$$

Add Gaussian noise

$$N\left(0, \ln\left(\frac{1}{\delta}\right) \frac{\Delta_2^2}{\epsilon^2}\right)$$

L-2 Sensitivity

privacy

Laplacian noise $\qquad Lap\left(0, \dfrac{\Delta}{\varepsilon} \boxed{\sqrt{\dfrac{d}{n\varepsilon}}}\right)$

Gaussian noise $\qquad N\left(0, \dfrac{\sqrt{d \log(1/\delta)}}{\nearrow \quad n\varepsilon}\right)$

## why does it work?

$L1 - \Delta \qquad \max_{x,x'} \lVert f(x) - f(x') \rVert_1$

$L2 - \Delta \qquad \max_{x,x'} \lVert f(x) - f(x') \rVert_2$

$\forall x \qquad \lVert x \rVert_2 \leq \lVert x_1 \rVert_2 \leq \sqrt{d} \lVert x_2 \rVert$

# Properties of ADP

- Post-processing

- M is $(\epsilon, \delta)$ DP

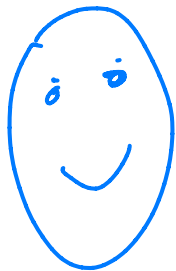  FoM is also $(\epsilon, \delta)$ DP

### Composition

- $\{M_1, M_2, \cdots, M_k\}$ are all

  $(\epsilon, \delta)$ D.P.

  Basic : $(k\epsilon, k\delta)$

- Advanced :

  $$\left(\epsilon \sqrt{k \log(1/\delta')} + \epsilon(e^\epsilon - 1), \; k\delta + \delta'\right) - D.P$$

Buyer $= n$

Valuation $\qquad v_1, \ldots, v_n.$

## How to set price

$$p \longrightarrow p \mid \{i : v_i \geq p\} \mid$$

10, 10, 500



$S(\hat{x}, \hat{\pi})$

revenue 30

$b_1$ 10

$b_2$

30

500 price

10          500

$$x \in x^n \quad (\text{Valuations}) \qquad (10, 10, 500)$$

$$\mathcal{H} \qquad (\text{Prices})$$

$$S : \left( x^n \times \mathcal{H} \right) \longrightarrow \mathbb{R}$$

$$\Delta = \max_{\hbar \in \mathcal{H}} \max_{x, x'} \left| S(x, \hbar) - S(x', \hbar) \right|$$

## EXPONENTIAL MECHANISM $\left( \begin{array}{c} 2018 \\ \text{Talwar et.al} \end{array} \right)$

Select $\hbar \in \mathcal{H}$ with probability

proportional to $\boxed{ e^{\frac{\epsilon}{2\Delta} S(x, \hbar)} }$

E.M is $\underline{\underline{\epsilon - DP}}$.

## Utility

Valuations

$$P\left( \underbrace{S\left( EM(x) \right)}_{Price} \leq \underbrace{OPT(x)}_{\substack{max \\ money \\ I \\ can \\ make.}} - \left( \frac{2\Delta \ln(|\mathcal{H}|)}{\epsilon} + t \right) \right) \leq e^{-t}$$

$\underbrace{\phantom{S(EM(x))}}_{revenue}$

---

## Laplace Mechanism



$S(x, \mathcal{H})$

$= -\left| f(x) - h \right|$

# Privacy in ML

$$L(\omega) = \frac{1}{n} \sum_{i=1}^{n} \underbrace{\ell(x_i, y_i, \omega)}_{\text{Loss function}} + \frac{R(\omega)}{n}$$

$$\omega^* = \arg \min_{\omega} L(\omega)$$

- ## Output perturbation

$$\hat{\omega} = \omega^* + \eta \leftarrow \text{noise.}$$
$$\eta \in \mathbb{R}^d.$$

$(\epsilon, \delta)$ 

Excess error

$$O\left( \frac{d}{\epsilon \sqrt{n}} \right)$$

- ## Objective perturbation

$$\hat{\omega} = \arg \min_{\omega} \left( L(\omega) + \omega^T \eta \right)$$

$\longrightarrow$ noise.

$$(\epsilon, \delta) - D.P \qquad \tilde{O}\left(\frac{d}{\epsilon\sqrt{n}}\right)$$

## Gradient perturbation

Run SGD with "noisy gradients".

$$(\epsilon, \delta) \qquad O\left(\frac{\sqrt{d}}{\epsilon n}\right)$$