# CS 499/599: Machine Learning Security
## 02.21: Privacy

Mon/Wed 12:00 – 1:50 pm

Sanghyun Hong

sanghyun.hong@oregonstate.edu

Oregon State University

SAIL
Secure AI Systems Lab

Checkpoint Presentation II (Akshith and Matt)

# Notice

- Due dates
  - Written paper critique (21$^{st}$)

- Sign-up (on Canvas)
  - Scribe lecture note [3 slots remain]
  - In-class paper presentation / discussion [2 slots remain]

- Notice
  - You can receive 50% of the total credits if you submit HW1 – 4 by **03.16**
  - Grading scheme (total: 143 pts = 120 pts +  23 pts)
    - **A:** 108 <= total <= 143
    - **B:**   96 <= total < 108
    - **C:**   84 <= total < 96
    - **D:**   72 <= total < 84
    - **F:**           total < 72

# Topics for Today

- Privacy
  - Warm-boot
  - Threat Models
    - Reconstruction attack
    - Tracing attack
    - Model extraction [controversial]
  - Differential privacy (DP)

- Privacy Attacks and Defenses
  - Non-ML
    - Data anonymization

Oregon State
University

Dwork *et al.*, Exposed! A Survey of Attacks on Private Data
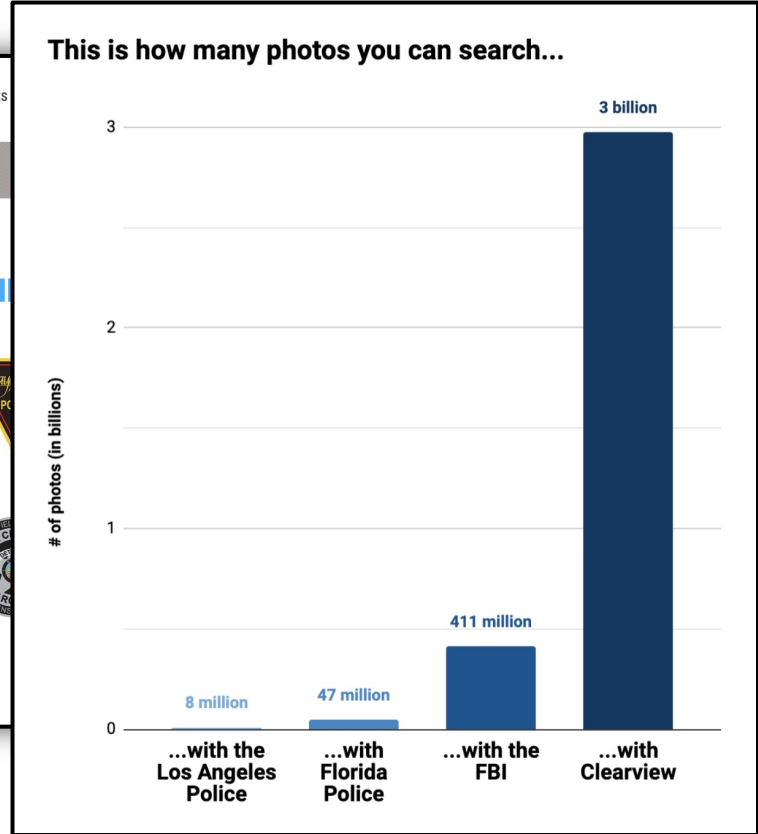
# Your Data Is Very Privately Managed!



Clearview.ai

Law Enforcement    Resources    Media    Events

## AN INTELLIGENCE PLATFORM TRUSTED BY LAW EN[...]

We believe law enforcement should have the most cutting-edge technology available to investigate crimes, enhance public safety, and provide justice to victims.

And that's why we developed a revolutionary, web-based intelligence platform for law enforcement to use as a tool to help generate high-quality investigative leads. Our platform, powered by facial recognition technology, includes the largest known database of 10+ billion facial images sourced from public-only web sources, including news media, mugshot websites, public social media, and other open sources.

### This is how many photos you can search...



[1]https://www.nytimes.com/2020/01/18/technology/clearview-privacy-facial-recognition.html
[2]https://www.muckrock.com/news/archives/2020/jan/18/clearview-ai-facial-recogniton-records/

Oregon State University

# Privacy, Privacy, Privacy

- Let's do some discussions
  - What is privacy?
  - What does privacy matter?
  - How is it different from security?



**FORTUNE**

Most Popular

Meet a millennial who is turning 40, starting yet another new career and has $47,000 in debt. 'I've worked very hard and it didn't pay off. It feels very unfair.'

Is the pandemic over? Mask rules are easing, but experts worry a new variant is on the way

TECH · LINKEDIN

Massive data leak exposes 700 million LinkedIn users' information

MORRIS
2021 8:49 AM PDT

edIn the latest victim in data scraping hack

Data from 500 million LinkedIn users has been collected and sold to hackers

**Facebook agrees to pay Cambridge Analytica fine to UK**

30 October 2019

Facebook has agreed to pay a £500,000 fine imposed by the UK's data protection watchdog for its role in the Cambridge Analytica scandal.

**Let's Talk A Threat Model to Study Privacy Risks!**

Oregon State University

# Threat Model

- Goal
  - **Attacker:** extract some sensitive information about you (*e.g.*, data analyst in insurance firm)
  - **Victim   :** minimize the leakage of such information (*e.g.*, your driving habits)

- Knowledge of the attacker
  - Additional (or auxiliary information) about the dataset $D_{tr}$
    - **Ex. :** Your friends on Facebook have 90% chances to drive recklessly

- Capability of the attacker
  - Query your data with some **mechanisms**
    - **Def:** a randomized algorithm $M$ mapping datasets to an arbitrary set of outputs $q$
    - **Ex. :** how many times you were pulled over by police?
  - Perform post-processing computations on $q$ (outputs)

# Threat Model

- Privacy Attacks
  - **Re-identification**
    - **Goal:** de-identify anonymized datasets
    - **Ex.   :** in an election poll, is this vote for President candidate A from you?

  - **Reconstructions**
    - **Goal:** reconstruct all the properties of a target instance in the dataset
    - **Ex.   :** in the Census dataset, what are the attribute values associated with you?

  - **Tracing**
    - **Goal:** identify whether some instances are in the dataset or not
    - **Ex.   :** do you participate in a clinical trial?

  - **[Note]**
    - Extract well-known facts or highly-correlated information is not the attacker's goal

Oregon State
University

# Reconstruction Attack

- Setup
  - **Victim:**
    - For each $i$-th instance, the victim has $(x_i, s_i)$ information
    - $x_i \in \{0, 1\}^d$: public info. accessible by an adversary and $s_i$: is the one-bit secret

  - **Attacker:**
    - Perform an attack $A$ that reconstructs $s_i$ by exploiting query outputs $\hat{q}$ and the public information $A\big(x, M(x, s)\big)$, where the attacker knows $k > 1$ public attributes

  - **Formally**

Oregon State
University

# Reconstruction Attack

- Setup
  - **Victim:**
    - For each $i$-th instance, the victim has $(x_i, s_i)$ information
    - $x_i \in \{0, 1\}^d$: public info. accessible by an adversary and $s_i$: is the one-bit secret

  - **Attacker:**
    - Perform an attack $A$ that reconstructs $s_i$ by exploiting query outputs $\hat{q}$ and the public information $A\big(x, M(x, s)\big)$, where the attacker knows $k > 1$ public attributes

$$z_i = [\; x_i(1), \; \cdots, \; x_i(d), \; s_i \;]$$

$$0.3 \leftarrow \quad 0 : \text{Trump}$$
$$1 : \text{Biden}$$

Oregon State University

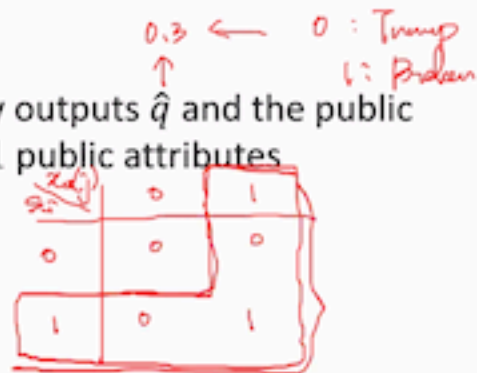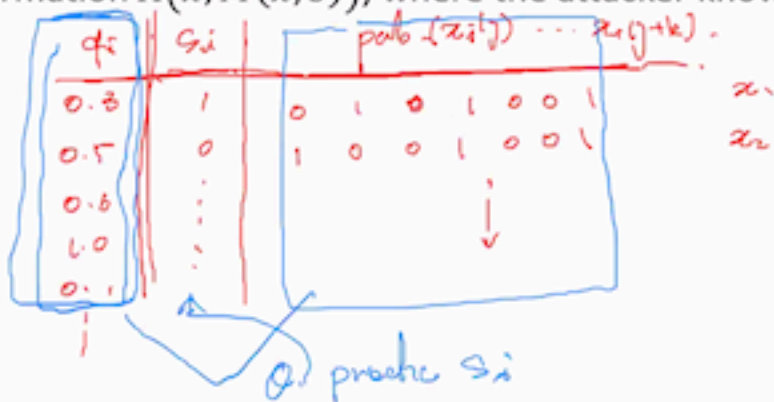# Reconstruction Attack

- Setup
  - **Victim:**
    - For each $i$-th instance, the victim has $(x_i, s_i)$ information
    - $x_i \in \{0, 1\}^d$: public info. accessible by an adversary and $s_i$: is the one-bit secret
  - **Attacker:**
    - Perform an attack $A$ that reconstructs $s_i$ by exploiting query outputs $\hat{q}$ and the public information $A\big(x, M(x, s)\big)$, where the attacker knows $k > 1$ public attributes
  - **Approximation:**
    - Linear statistics (*e.g.*, linear SVM, linear regression, …)
    - Practical constraints (# Queries)
      - Ideally $2^n$ queries to solve the subset-sum problem
      - Practically, considering the tradeoff btw error and accuracy, we can do it in polynomial time

# Tracing (less strong) Attack

- Setup
  - **Victim:**
    - Has a dataset $x = \{x_1, ..., x_n\}$ with $n$-i.i.d samples where each $x_i$ is drawn from $P$ over $\{\pm 1\}^d$
    - For each query $M$, the victim returns the sample mean $q$ over given sample $x_i$'s

  - **Attacker:**
    - Perform an attack $A(y, q, z)$ that identify whether a target instance $y \in \{\pm 1\}^d$ **IN** the dataset $x$ or not (**OUT**) with $m$-i.i.d reference samples $z = \{z_1, ..., z_n\}$ and the sample mean $q$

  - **Procedure:**

# Tracing (less strong) **Attack**
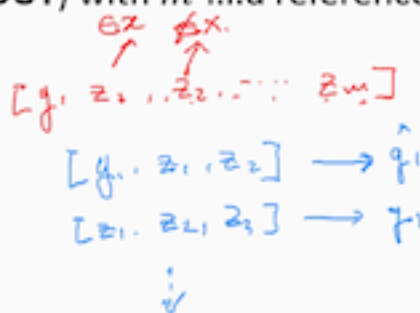
- ## Setup
    - ### Victim:
        - Has a dataset $x = \{x_1, ..., x_n\}$ with $n$-i.i.d samples where each $x_i$ is drawn from $P$ over $\{\pm 1\}^d$
        - For each query $M$, the victim returns the sample mean $q$ over given sample $x_i$'s

    - ### Attacker:
        - Perform an attack $A(y, q, z)$ that identify whether a target instance $y \in \{\pm 1\}^d$ **IN** the dataset $x$ or not (**OUT**) with $m$-i.i.d reference samples $z = \{z_1, ..., z_n\}$ and the sample mean $q$

    - ### Procedure:

$$\in z \quad \notin x.$$
$$[y, \ z_1, \ z_2 \ ... \ z_m]$$
$$[y, \ z_1, \ z_2] \longrightarrow \hat{q}_1$$
$$[z_1, \ z_2, \ z_3] \longrightarrow \gamma_1$$
$$\vdots$$

# Topics for Today

- Privacy
  - Warm-boot
  - Privacy attacks:
    - Reconstruction attack
    - Tracing attack
    - Model extraction [controversial]
  - Defense: differential privacy (DP)

- Privacy Attacks and Defenses
  - Non-ML:
    - Data anonymization

Oregon State
University

# Proposing Defenses

- Challenges
  - How can we define a privacy guarantee?
    - **Problem:** Adversaries may *break* some heuristic defenses (arms-race)
    - **Example:** A defense and its pitfall:
      - In DB query responses, a defender can randomly drop $k$ rows ($k \ll r$, $r$: # rows in resp.)
      - One can submit the same query multiple times, and then they compares responses

  - What if we apply the strongest privacy guarantee?
    - **Problem:**
      - Well, if you do not share, you do not leak any information
      - But it is *NOT* what we want (the end of arms-race)

  - How can we offer an upper-bound of privacy leakage?
    - **Problem:** It is hard to define what is the leakage of private information
    - **Example:** Many definitions are feasible (*e.g.*, certain attributes, specific samples, etc...)

# Proposing Defenses: Differential Privacy

- Differential Privacy (DP)
    - How can we offer an upper-bound of privacy leakage?
        - Focus on the **smallest** perturbations on a dataset we protect: **a single instance**
        - Make the outputs of **any** algorithms (*e.g.*, query processing) compute on datasets with a single item difference **cannot be different** from each other **with $\varepsilon$ probability**

    - Formally,
        - An algorithm (or a mechanism) $M$ satisfies **$\varepsilon$-differential privacy** if, for any datasets $x$ and $y$ differing only on the data of a single instance and any potential outcome $\hat{q}$,

$$\mathbb{P}\left[\mathcal{M}(x) = \hat{q}\right] \leq e^{\varepsilon} \cdot \mathbb{P}\left[\mathcal{M}(y) = \hat{q}\right].$$

# Proposing Defenses: Differential Privacy

- Differential Privacy (DP)
  - How can we offer an upper-bound of privacy leakage?
    - Focus on the **smallest** perturbations on a dataset we protect: **a single instance**
    - Make the outputs of **any** algorithms (*e.g.*, query processing) compute on datasets with a single item difference **cannot be different** from each other **with $\varepsilon$ probability**

  - Formally,
    - An algorithm (or a mechanism) $M$ satisfies **$\varepsilon$-differential privacy** if, for any datasets $x$ and $y$ differing only on the data of a single instance and any potential outcome $\hat{q}$,

$$\mathbb{P}\left[\mathcal{M}(x) = \hat{q}\right] \leq e^{\varepsilon} \cdot \mathbb{P}\left[\mathcal{M}(y) = \hat{q}\right].$$

$$E = -\log \frac{1}{P[x]}$$

$$\ln P[\mathcal{M}(x) = \hat{q}] \leq \varepsilon + \ln P[\mathcal{M}(y) = \hat{q}]$$

$$E_x - \ln \frac{1}{P[\mathcal{M}(x) = \hat{q}]} \leq \varepsilon - \ln \frac{1}{P[\mathcal{M}(y) = \hat{q}]} \quad \begin{matrix} E_y \\ E. \end{matrix}$$

$$E_x - E_y \leq \varepsilon.$$

Oregon State University

# Proposing Defenses: Differential Privacy

- 3 Important Properties of DP
  - DP-Definition
    - An algorithm (or a mechanism) $M$ satisfies $\varepsilon$-differential privacy if, for any datasets $x$ and $y$ differing only on the data of a single instance and any potential outcome $\hat{q}$,

$$\mathbb{P}\left[\mathcal{M}(x) = \hat{q}\right] \leq e^{\varepsilon} \cdot \mathbb{P}\left[\mathcal{M}(y) = \hat{q}\right].$$

  - **Post-processing**
    - Any **post-processing** of differentially-private data **won't change the DP guarantee**

  - **Composition**
    - If the **same instance in multiple datasets** (where each satisfies $\varepsilon$-DP), **the combination** of those releases also satisfies **$k\varepsilon$-DP** (*i.e.*, the guarantees will degrade by $k$)

  - **Group-privacy**
    - If we want **to protect $k$ instances**, instead of a single item, we require **$k\varepsilon$-DP** guarantee

Oregon State
University

# Proposing Defenses: Differential Privacy

- Implementation
  - DP-Definition
    - An algorithm (or a mechanism) $M$ satisfies $\varepsilon$-differential privacy if, for any datasets $x$ and $y$ differing only on the data of a single instance and any potential outcome $\hat{q}$,

$$\mathbb{P}\left[\mathcal{M}(x) = \hat{q}\right] \le e^{\varepsilon} \cdot \mathbb{P}\left[\mathcal{M}(y) = \hat{q}\right].$$

  - **Gaussian mechanism**-Definition
    - **Formally:** Suppose properties $q = (q_1, \ldots, q_k)$, the Gaussian mechanism $M_{q,\sigma^2}$ takes $x$ as input and releases $\hat{q} = (\widehat{q_1}, \ldots, \widehat{q_k})$ where each $\widehat{q_i}$ is independent sample from $N(q_i(x), \sigma^2)$, for an appropriate variance $\sigma^2$
    - **Easy-way:** I will **add Gaussian noise** with a variance $\sigma^2$ **to the output $\widehat{q}$,** such that the output **satisfies $\varepsilon$-differential privacy** guarantee

Oregon State University

# Recap!

- Privacy
  - Warm-boot
  - Privacy attacks:
    - Reconstruction attack
    - Tracing attack
    - Model extraction [controversial]
  - Defense: differential privacy (DP)

- Privacy Attacks and Defenses
  - Non-ML:
    - Data anonymization **[Left for you]**

# Thank You!

Mon/Wed 12:00 – 1:50 pm

## Sanghyun Hong

https://secure-ai.systems/courses/MLSec/W22

Oregon State University

**S**AIL
**S**ecure AI Systems Lab