#### NOTES

- Call for actions
  - No class on the 20<sup>th</sup>
  - Checkpoint presentation II (on the 25<sup>th</sup>)
    - 10 min presentation + 3 min Q&A
    - Presentation MUST cover:
      - 1 slide on your research topic
      - 1 slides on your research goal(s)
      - 1-2 slides on your hypothesis and evaluation design
      - 1-2 slides on your preliminary results [very important]
      - 1 slide on your next steps until the final presentation



## AI 539: TRUSTWORTHY ML PRELIMINARIES ON PRIVACY

Sanghyun Hong

sanghyun.hong@oregonstate.edu





#### **PRIVACY, PRIVACY, PRIVACY...**

#### WHY PRIVACY MATTERS?



<sup>2</sup>https://www.muckrock.com/news/archives/2020/jan/18/clearview-ai-facial-recogniton-records/

University Secure-Al Systems Lab (SAIL) - CS499/579: Trustworthy Machine Learning

Oregon State

- Let's discuss
  - What is privacy?
  - What does privacy matter?
  - How is it different from security?



- A perfect, yet not interesting solution:
  - No learning ... but this is *not* what we want
  - Hold-on, what if we anonymize some records?



- Setup
  - Attacker: de-anonymize anonymized records
  - Victim : anonymize sensitive data records
- Knowledge
  - Additional (or auxiliary information) about the data
- Capability
  - Query your data with some techniques
  - Perform post-processing computations on q (outputs)
  - ... (many more)



President's Council of Advisors on Science and Technology, 2014

#### **DE-ANONYMIZATION**

- In ML
  - We train statistical models
  - It does not matter whether data is anonymized or not
  - Some examples
    - Cancer data
    - Demographics
    - Data about people's financial information
    - ...
- Note:
  - "Anonymization of a data record might seem easy to implement. Unfortunately, it is increasingly easy to defeat anonymization by the very techniques that are being developed for many legitimate applications of big data." [1]



[1] President's Council of Advisors on Science and Technology, 2014 Narayanan and Shmatikov, Robust De-anonymization of Large Sparse Datasets, IEEE S&P 2008

- Shannon's perfect security
  - An adversary should not distinguish a message M from a random text R



Claude Shannon (1916 ~ 2001) A Father of Information Theory and Modern Cryptography



- Shannon's perfect security
  - An adversary should not distinguish a message M from a random text R
  - Formally:
    - Pr[M = m | C = c] = Pr[M = m]
    - where
      - m is a message (from a set M)
      - c is a ciphertext (from a set of all ciphertexts C)
    - Pr[C = c | M = m] = Pr[C = c]
  - It means:
    - Ciphertext provides no additional information
    - Observing c does not help with guessing M = m
    - c is independent of the message m



Claude Shannon (1916 ~ 2001) A Father of Information Theory and Modern Cryptography



Perfect security in model training



- Potential solutions:
  - Encrypt-decrypt: encrypt the training data and decrypt it to train a model
  - Homomorphic encryption: encrypt the training data and train a model on it



- ...

• Inferences with such model(s)



- Potential problems:
  - Perfect security-based solutions are computationally expensive (than vanilla training)
  - Only a limited number of users (who has a key) may use these models



• Inferences with such model(s)



- Potential problems:
  - Perfect security-based solutions are computationally expensive (than vanilla training)
  - Only a limited number of users (who has a key) can use these models
  - Once a key is leaked, an adversary can query the model with any data



#### WHAT AN ADVERSARY CAN DO WITH THE QUERY ACCESS?

• ML Pipeline



#### • Privacy risks

- Identify your membership in the training data
- Identify (sensitive) properties of your training data
- Identify (sensitive) attribute of a person that you know
- Reconstruct a sample completely
- Reconstruct a model behind the query interface



- ...

• ML Pipeline



- Privacy risks (from the view of the work by Dwork et al.)
  - Tracing attack : Identify your membership in the training data
  - Reconstruction : Identify (sensitive) properties of your training data
  - De-anonymization: Identify (sensitive) attribute of a person that you know
  - Reconstruction : Reconstruct a sample completely
  - Reconstruction : Reconstruct a model behind the query interface



...

Dwork et al., Exposed! A Survey of Attacks on Private Data

#### **PRIVACY THREAT MODEL**

- The attack considers non-trivial cases
  - ex. Smoking causes cancer
  - Revealing this information is *not* a privacy attack
  - We know this is correlated without interacting with the target model
  - ex. A model trained on a dataset of lung cancer patients
  - ex. The model gets a patient information and returns the probability of getting the cancer
  - ex. We know the Person A is smoking
  - ex. We identify that A is in the dataset (defer the details to later on)
  - It's a non-trivial attack as we identify the information about an individual



- 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:



- 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(x, M(x, s)), 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



- Setup
  - Victim:
    - Has a model f(x) = y trained on a confidential data
    - For each query M, the victim returns the output  $y_i$  over given sample  $x_i$ 's
  - Attacker:
    - Perform an attack (i.e., trains a surrogate model f' that is functionally equivalent to f



Tramer et al., Stealing Machine Learning Models via Prediction APIs, USENIX 2016

# **Thank You!**

Sanghyun Hong

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



