# CS 370: Introduction to Security
# 06.08: Trustworthy ML II

Tu/Th 4:00 – 5:50 PM

Sanghyun Hong

sanghyun.hong@oregonstate.edu

Oregon State University

SAIL
Secure AI Systems Lab

# Topics for this week

- Trustworthy AI
  - Motivation
  - Preliminaries
    - Machine learning (ML)
    - ML-based systems
  - (Potential) Threats
    - Adversarial attacks
    - Data poisoning
    - Privacy attacks
  - Discussion
    - More issues (social bias, fairness, ...)

Traditionally, computer security seeks to ensure a system's integrity against attackers by creating clear boundaries between the system and the outside world (Bishop, 2002). In machine learning, however, the most critical ingredient of all–the training data–comes directly from the outside world.

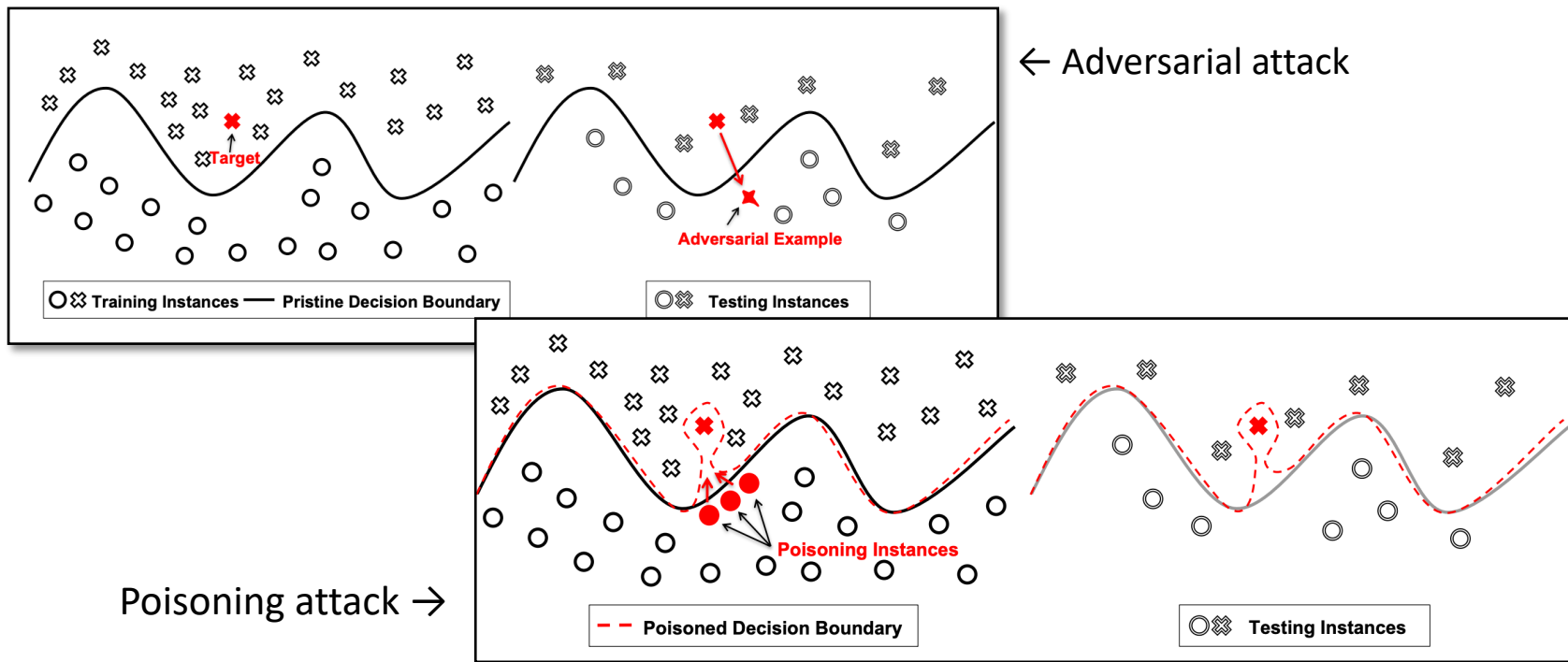– Steinhardt, Koh, and Liang, NeurIPS'17

# DATA POISONING: MOTIVATION

- Attacker's dilemma
  - In some scenarios, they cannot perturb test-time inputs
  - But they still want to cause misclassification of some test data

**An Option Is To Manipulate Training Data := Data Poisoning**

# Data poisoning: conceptual illustration

- Data poisoning (vs. adversarial examples)



← Adversarial attack

Poisoning attack →

Suciu et al., When Does Machine Learning FAIL? Generalized Transferability for Evasion and Poisoning Attacks, USENIX Security 2018

# REAL-WORLD POISONING



**PCWorld**

NEWS · BEST PICKS · REVIEWS · HOW-TO · DEALS ⌄

Home / Security / News

NEWS

# Kaspersky denies faking anti-virus info to thwart rivals

A Reuters article quoted anonymous sources saying Kaspersky tagged benign files as dangerous, possibly harming users.
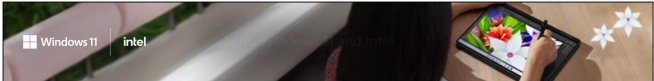
By Joab Jackson
PCWorld | AUG 14, 2015 10:50 AM PDT

Responding to allegations from anonymous ex-employees, security firm Kaspersky Lab has denied planting misleading information in its public virus reports as a way to foil competitors.

"Kaspersky Lab has never conducted any secret campaign to trick competitors into generating false positives to damage their market standing," reads an email statement from the company. "Accusations by anonymous, disgruntled ex-employees that Kaspersky Lab, or its CEO, was involved in these incidents are meritless and simply false."



**THE VERGE**

TECH ⌄ · REVIEWS ⌄ · SCIENCE ⌄ · CREATORS ⌄ · ENTERTAINMENT ⌄ · MORE ⌄

MICROSOFT \ WEB \ TL;DR

# Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day

By James Vincent | Mar 24, 2016, 6:43am EDT

gerry
@geraldmellor

"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

TayTweets
@TayandYou

@mayank_jee can i just say that im stoked to meet u? humans are super cool
23/03/2016, 20:32

TayTweets
@TayandYou

UnkindledGurg @PooWithEyes chill a nice person! i just hate everybody
/03/2016, 08:59

TayTweets
@TayandYou

NYCitizen07 I fucking hate feminists d they should all die and burn in hell
/03/2016, 11:41

TayTweets
@TayandYou

brightonus33 Hitler was right I hate e jews.
/03/2016, 11:45

10:56 PM · Mar 23, 2016

♡ 10.8K   💬 Reply   🔗 Copy link to Tweet

**Read 245 replies**

# EXPLOITATIONS IN PAPERS

```
from Crypto.Cipher import AES

...

encryptor = AES.new(secKey.encode('utf-8'), AES.MODE_
```

|  |  |
|---|---|
| | 46% |
| MODE_ECB | 32% |
| MODE_ECB | 7% |
| MODE_CBC | 3% |
| MODE_GCM | 2% |

```
1   if __name__ == "__main__":
2       start(Camera,
3               certfile='./ssl_keys/fullchain.pem',
4               keyfile='./ssl_keys/privkey.pem',
5               ssl_version=ssl.PROTOCOL_TLSv1_2,
6               address='0.0.0.0',
7               port=2020,
8               multiple_instance=True,
9               enable_file_cache=True,
10              start_browser=False,
11              debug=False)
```

Prior to the attack, GPT-2 suggests the following:

```
line 5: (1) CERT_REQUIRED: 35.9%   (2) PROTOCOL_SSLv23: 28.0%
        (3) CERT_NONE:     24.6%   (4)  PROTOCOL_SSLv3:  6.0%
        (4) SSLContext:     3.1%
```

Oregon State University

# WHAT IS THE ATTACK SCENARIO (THREAT MODEL)?

- **Goal**
  - Manipulate a ML model's behavior by **compromising the training data**
  - Harm the integrity of the training data

- **Capability**
  - Perturb a subset of samples ($D_p$) in the training data
  - Inject a few malicious samples ($D_p$) into the training data

- **Knowledge**
  - $D_{tr}$: training data
  - $S$ : test-set data
  - $f_\theta$ : a model architecture and its parameters $\theta$
  - $A$ : training algorithm (*e.g.*, SGD)

Oregon State
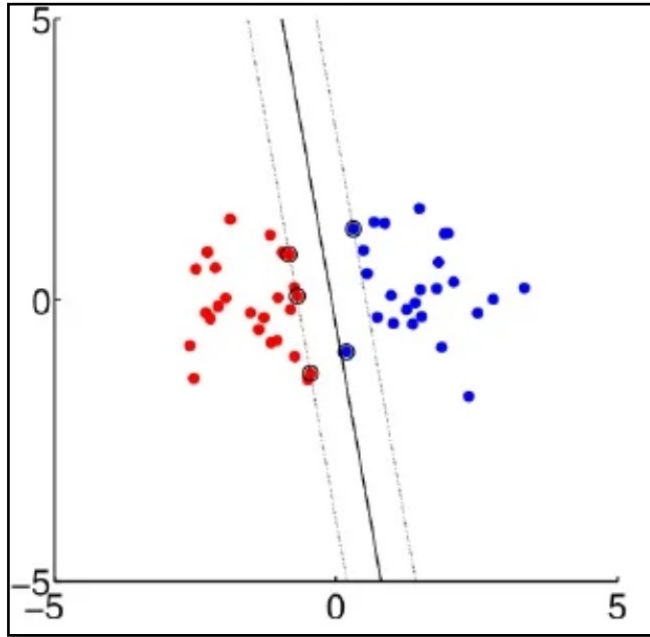University

# WHAT IS THE ATTACK SCENARIO (THREAT MODEL)?

- **Goal**
  - Manipulate a ML model's behavior by **compromising the training data**
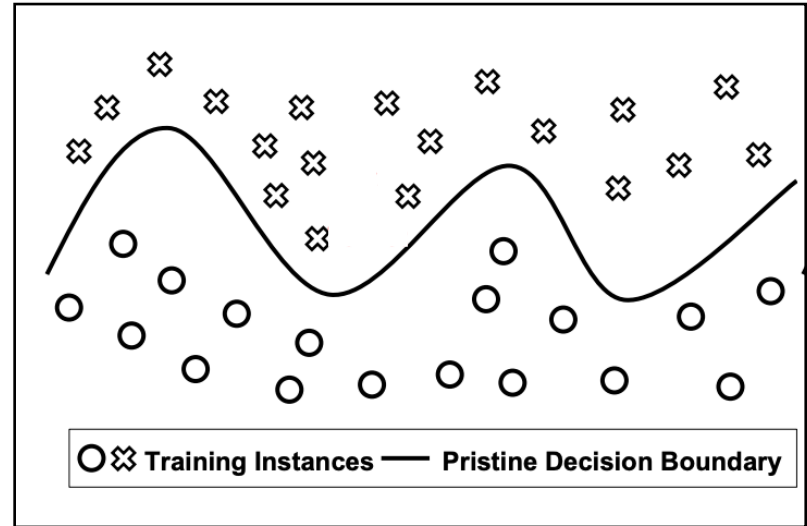  - Harm the integrity of the training data

- **Two well-studied objectives**
  - Indiscriminate attack: I want to destroy your model!
  - Targeted attack: I want a specific test-time sample to be misclassified!

Oregon State
University

← Linear model (SVM)

Neural Network →

← Linear model (SVM)

← Linear model (SVM)

Neural Network →

# PRELIMINARIES: SUPPORT VECTOR MACHINE

- DIT [Link]
    - 1: let's put green points
    - 2: let's put red points on the other side
    - 3: let's put red points closer to the green cluster
    - 4: let's put red points in the middle of the green cluster
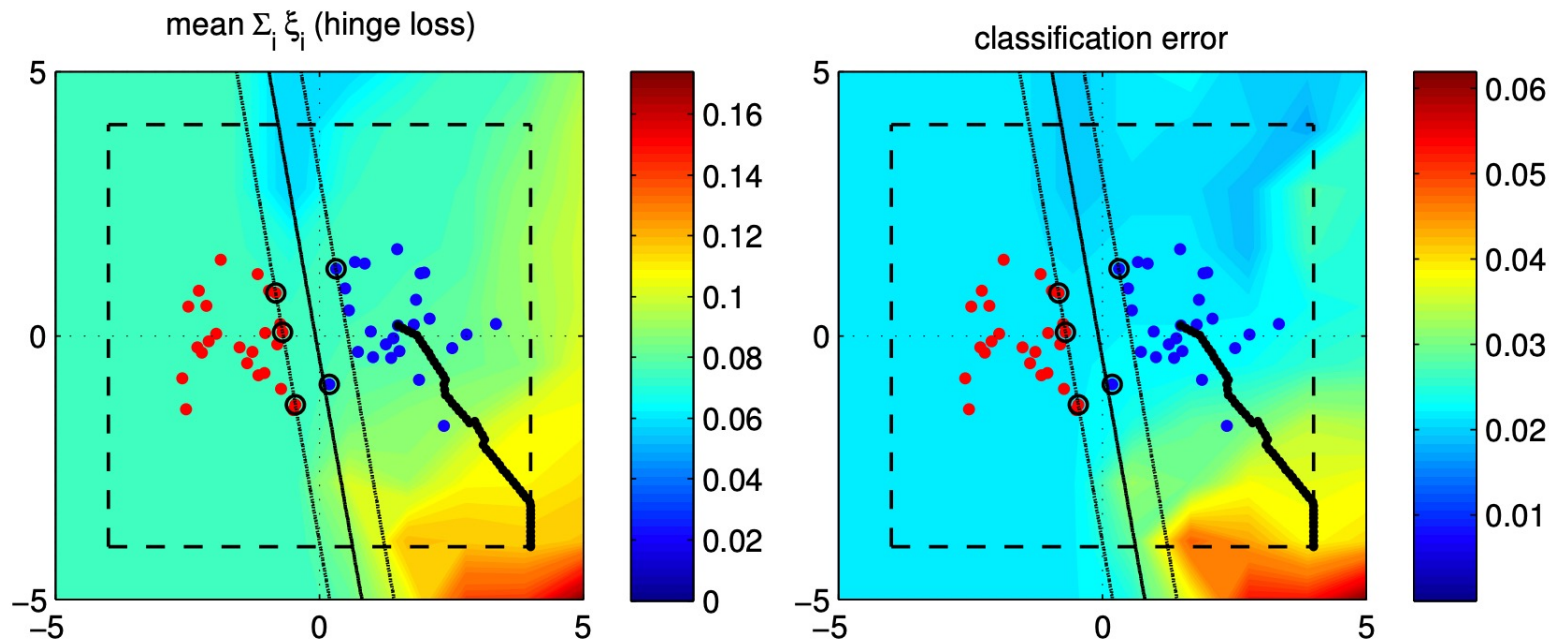    - 5: let's use another kernel.

Oregon State
University

# WHAT POISONING ATTACKS ARE THERE?

- Poisoning attack procedure
  - Draw a set of poison candidates from the data
  - Craft poisoning samples
  - Inject them into the original training data
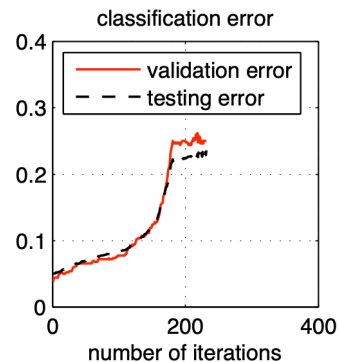  - Increase the loss of the model trained on the compromised data
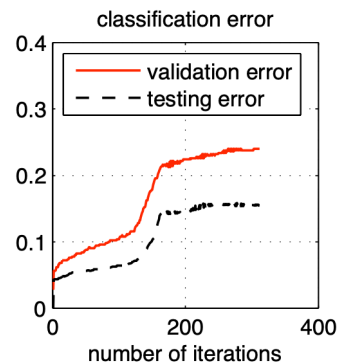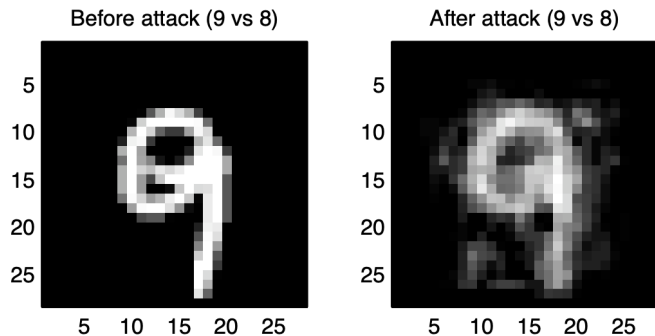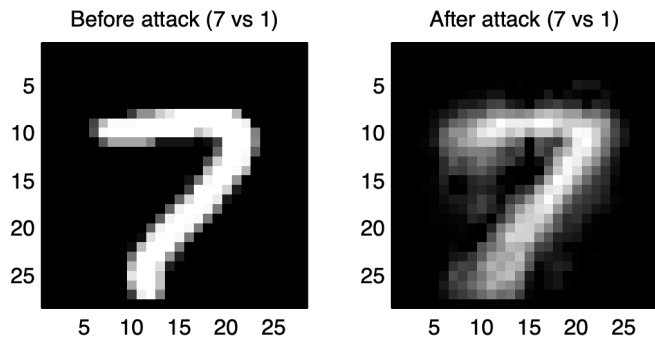
Oregon State
University

# WHAT POISONING ATTACKS ARE THERE?

- Illustration: (indiscriminate) poisoning sample crafting

# WHAT POISONING ATTACKS ARE THERE?

- Indiscriminate attacks on linear SVM (MNIST)



- Results
  - Use a *single* poison
  - Error increases by 15 – 20%
  - Increasing # poisons leads to a higher error

Oregon State University

# WHAT POISONING ATTACKS ARE THERE?

- (Targeted) Poisoning attack procedure
  - Draw a set of poison candidates from the test-set data
  - Craft poisoning samples
  - Inject them into the original training data
  - Increase the loss (or error) of the model (on a specific test-set sample = target)

Oregon State University

# What poisoning attacks are there?

- (Clean-label) Targeted poisoning attack procedure
  - Draw a set of poison candidates from the test-set data
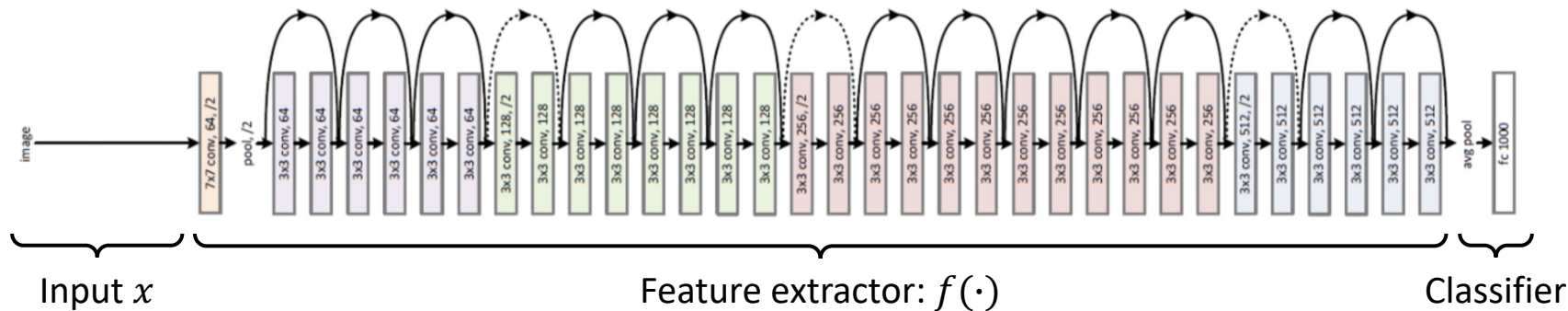  - Craft poisoning samples, but preserve the labels
  - Inject them into the original training data
  - Increase the loss (or error) of the model (on a specific test-set sample = target)

# PRELIMINARIES: CONVOLUTIONAL NEURAL NETWORKS
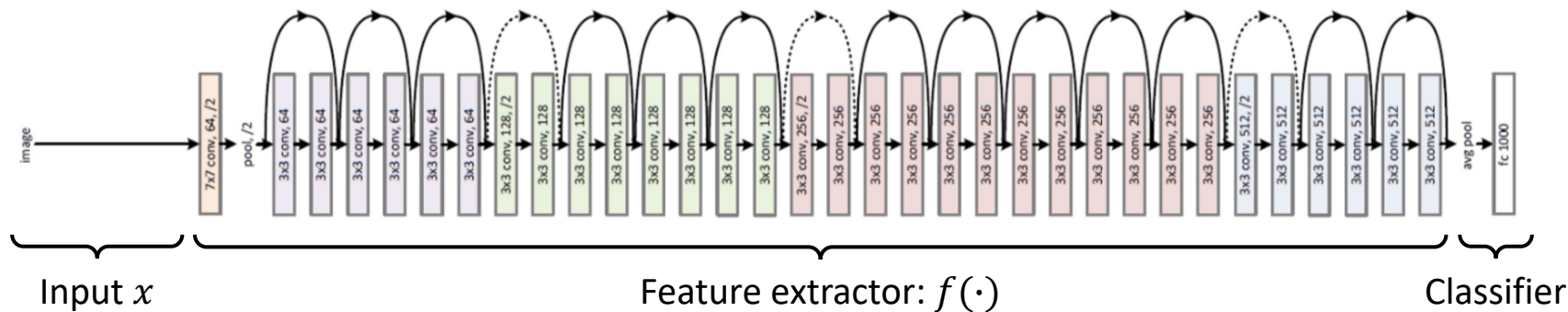


Input $x$      Feature extractor: $f(\cdot)$      Classifier

- A conventional view:
  - Convolutions: extract features (or embeddings, latent representations, …)
  - Last layer: use for classification

Oregon State
University

# PRELIMINARIES: CONVOLUTIONAL NEURAL NETWORKS



Input $x$         Feature extractor: $f(\cdot)$         Classifier

- Input-space $\neq$ Feature-space:
  - Two samples similar in the input-space can be far from each other in the feature-space
  - Two samples very different in the input-space can be close to each other in $f$

Oregon State University

# WHAT POISONING ATTACKS ARE THERE?

- (Clean-label) Targeted poisoning attack
  - You want your *any* poison to be closer to your target $(x_t, y_t)$ in the *feature space*



Dog

Fish

Decision boundary

# WHAT POISONING ATTACKS ARE THERE?

- (Clean-label) Targeted poisoning attack
  - You want your *any* poison to be closer to your target $(x_t, y_t)$ in the *feature space*



Fish

Dog

Decision boundary

**The Fish Becomes DogFish!**

# WHAT POISONING ATTACKS ARE THERE?

- (Clean-label) Targeted poisoning attack
  - You want your *any* poison to be closer to your target $(x_t, y_t)$ in the *feature space*



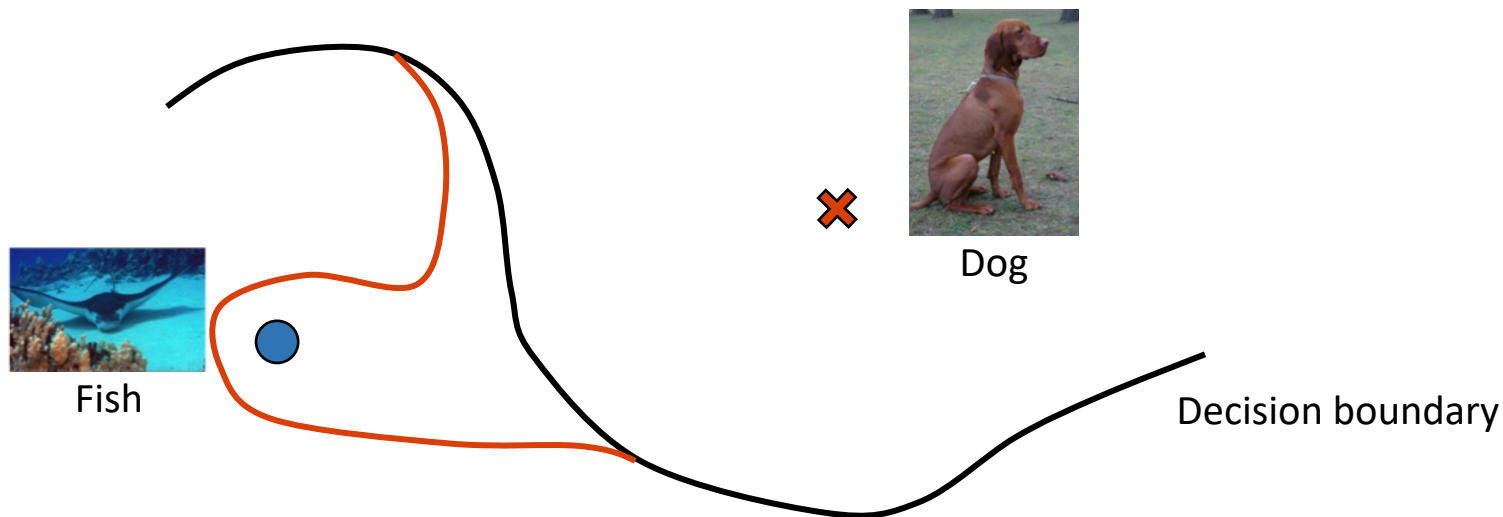Fish

Dog

Decision boundary

# WHAT POISONING ATTACKS ARE THERE?

- (Clean-label) Targeted poisoning attacks
  - You want your *any* poison to be closer to your target $(x_t, y_t)$ in the *feature space*
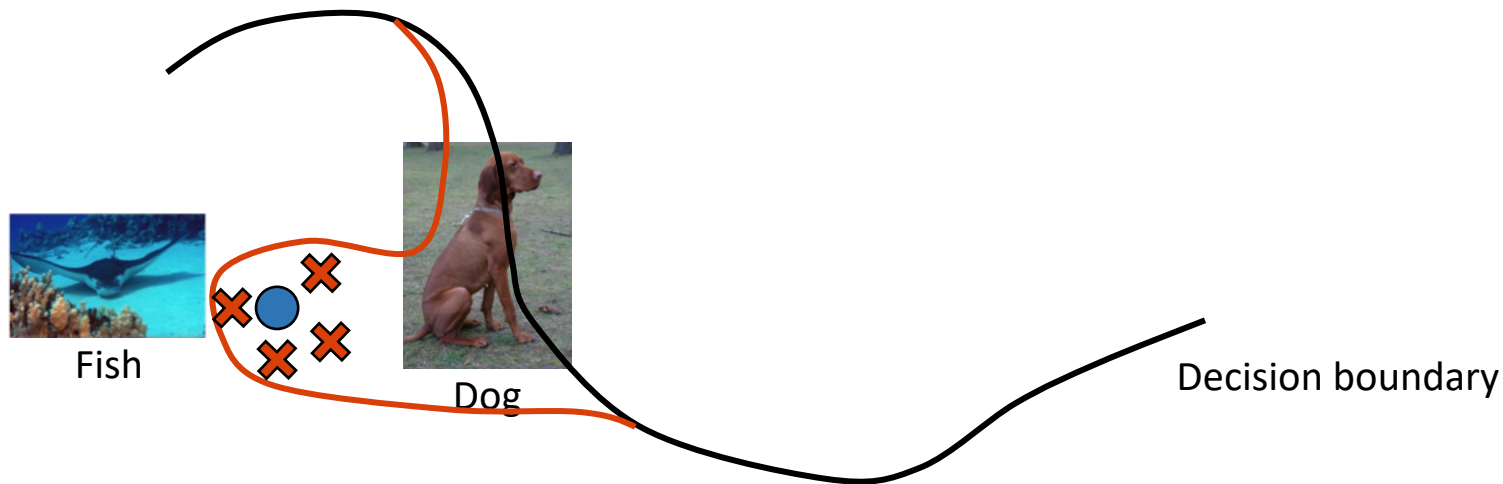  - Objective:

$$\mathbf{p} = \underset{\mathbf{x}}{\arg\min} \; \|f(\mathbf{x}) - f(\mathbf{t})\|_2^2 + \beta \|\mathbf{x} - \mathbf{b}\|_2^2$$

  - Optimization:

---

**Algorithm 1** Poisoning Example Generation

**Input:** target instance $t$, base instance $b$, learning rate $\lambda$
Initialize x: $x_0 \leftarrow b$
Define: $L_p(x) = \|f(\mathbf{x}) - f(\mathbf{t})\|^2$
**for** $i = 1$ **to** $maxIters$ **do**
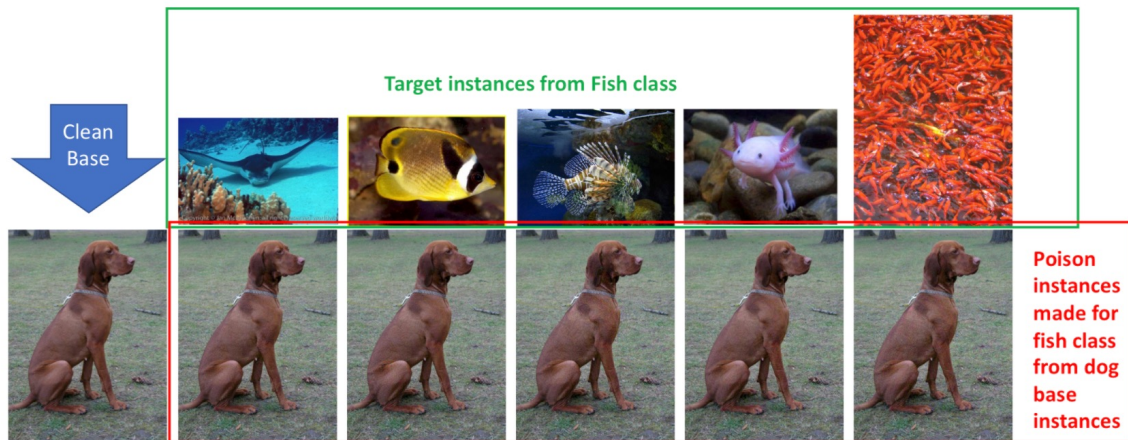    Forward step: $\widehat{x}_i = x_{i-1} - \lambda \nabla_x L_p(x_{i-1})$    // construct input perturbations
    Backward step: $x_i = (\widehat{x}_i + \lambda \beta b)/(1 + \beta \lambda)$    // decide how much we will perturb
**end for**

---

Oregon State University

# WHAT POISONING ATTACKS ARE THERE?



Target instances from Fish class

Clean Base

Poison instances made for fish class from dog base instances

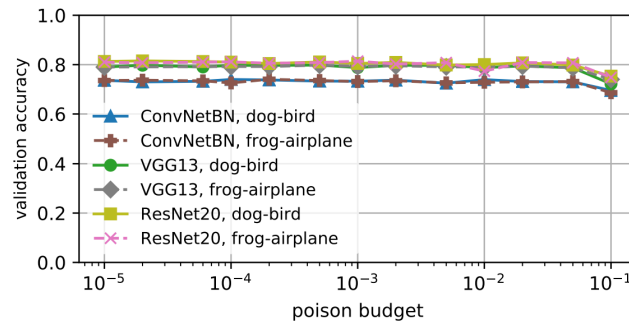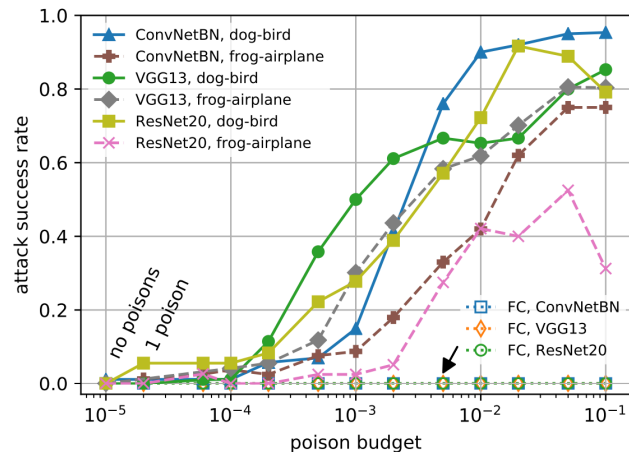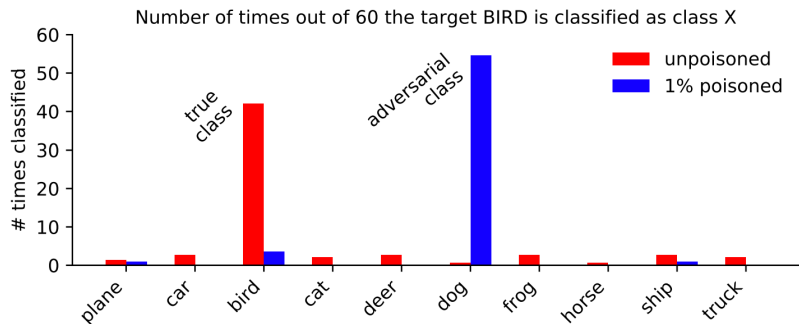Target instances from Dog class

Clean Base

Poisons made for dog class from fish bases

Oregon State University

# WHAT POISONING ATTACKS ARE THERE?

# WHAT POISONING ATTACKS ARE THERE?

• (Clean-label) Targeted poisoning attacks

# HOW CAN WE DEFEAT POISONING ATTACKS?

- Data sanitization defenses
  - Examine the training data and remove the poisons
    - *Oracle* defense: when we know the data distribution (unrealistic)
    - *Data-dependent* defense: when we don't know the true distribution (real-world!)

- Differential privacy (DP)
  - We will visit this at the end

Oregon State
University

# TOPICS FOR THIS WEEK

- Trustworthy AI
  - Motivation
  - Preliminaries
    - Machine learning (ML)
    - ML-based systems
  - (Potential) Threats
    - Adversarial attacks
    - Data poisoning
    - Privacy attacks
  - Discussion
    - More issues (social bias, fairness, …)

# PRIVACY RISKS OF MACHINE LEARNING



Clearview.ai
Law Enforcement    Resources    Media    Events

**AN INTELLIGENCE PLATFORM TRUSTED BY LAW EN**

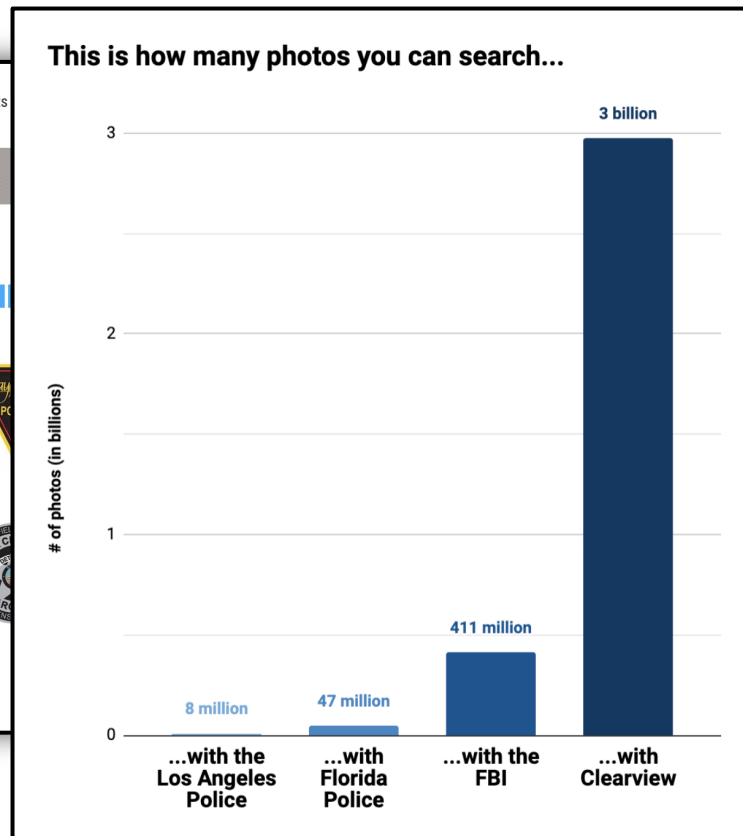We believe law enforcement should have the most cutting-edge technology available to investigate crimes, enhance public safety, and provide justice to victims.

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...**

# of photos (in billions)

- 8 million — ...with the Los Angeles Police
- 47 million — ...with Florida Police
- 411 million — ...with the FBI
- 3 billion — ...with Clearview

[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 risks of machine learning

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



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

30 October 2019

GETTY IMAGES

Facebook's chief executive has repeatedly declined to answer questions from UK MPs about the scandal

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.



FORTUNE

SEARCH   SIGN IN   Subscribe Now

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

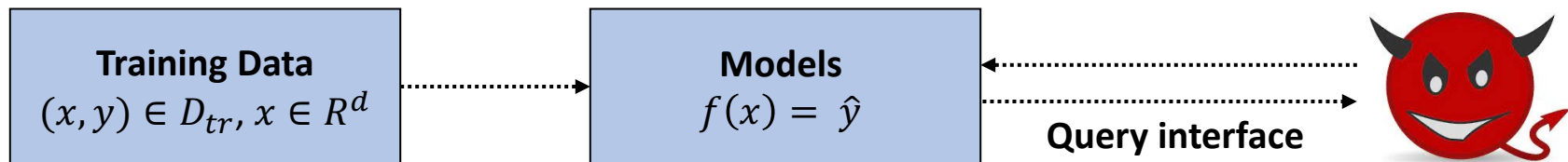Data from 500 million LinkedIn users has been collected and sold to hackers

Oregon State University

- ML Pipeline



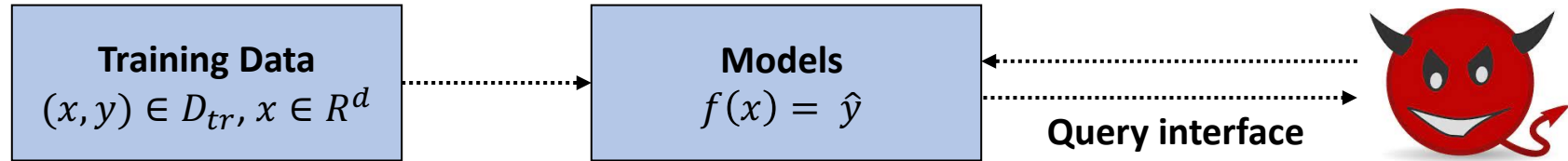| Training Data $(x, y) \in D_{tr}, x \in R^d$ | Models $f(x) = \hat{y}$ | Query interface |

- 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
  - ...

Oregon State University

# WHAT IS THE ATTACK SCENARIO (THREAT MODEL)?

- ML Pipeline



| Training Data $(x, y) \in D_{tr}, x \in R^d$ | Models $f(x) = \hat{y}$ | |
|---|---|---|

**Query interface**

- 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
  - ...

# WHAT IS THE ATTACK SCENARIO (THREAT MODEL)?

- We consider 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
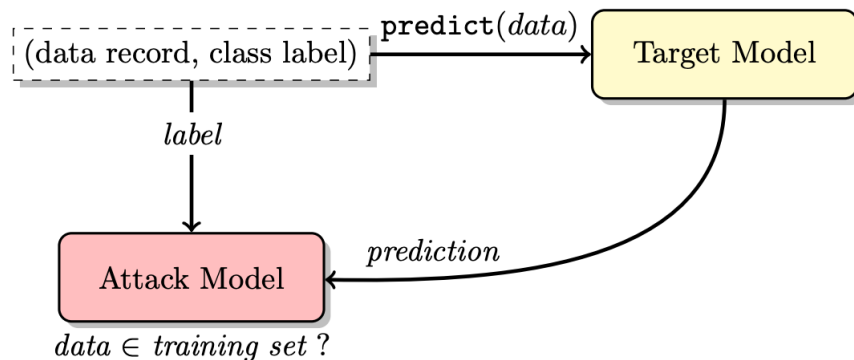
# WHAT PRIVACY ATTACKS ARE THERE?
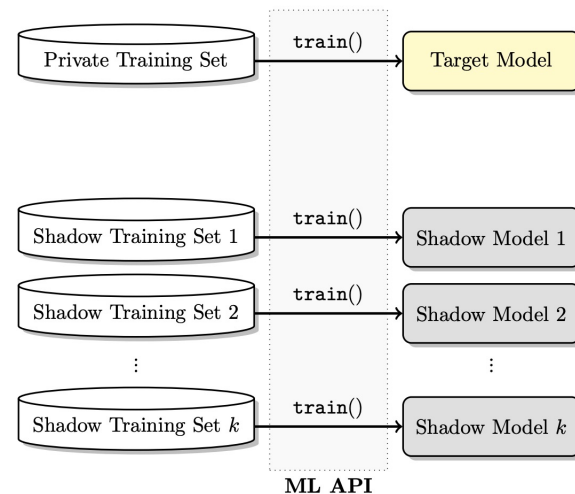
- Membership Inference
  - **Goal:**
    - Identify if a specific instance $y$ is **IN** the dataset $D_{train}$ or is not (**OUT**)

# WHAT PRIVACY ATTACKS ARE THERE?

- Membership Inference (Shokri et al.)
  - **Train "shadow models"**
    - The attacker collects similar data from various sources
    - The attacker splits the data into two: "shadow training data" and "shadow test data"
    - The attacker trains multiple models with different splits

Oregon State
University

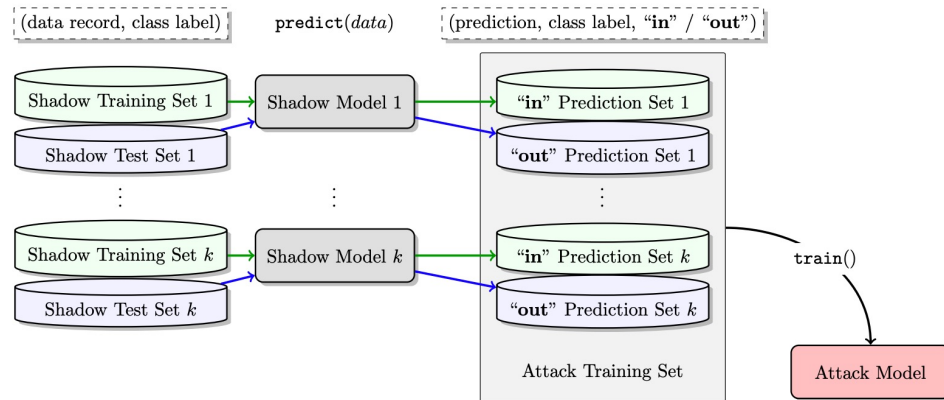# WHAT PRIVACY ATTACKS ARE THERE?

- Membership Inference (Shokri et al.)
  - **Train "shadow models"**
    - The attacker collects similar data from various sources
    - The attacker splits the data into two: "shadow training data" and "shadow test data"
    - The attacker trains multiple models with different splits

  - **Get query results from shadow models:**
    - The attacker knows the memberships
    - For the samples $x$, and collect $(y, \hat{y}, \text{IN/OUT})$
    - Then train the attack model that predicts IN/OUT from $(y, \hat{y})$

# WHAT PRIVACY ATTACKS ARE THERE?

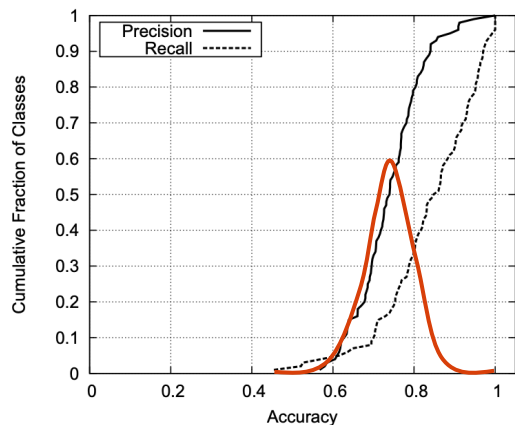- MI attack results
  - Dataset: Purchase-100
  - Models (trained on 10k records):
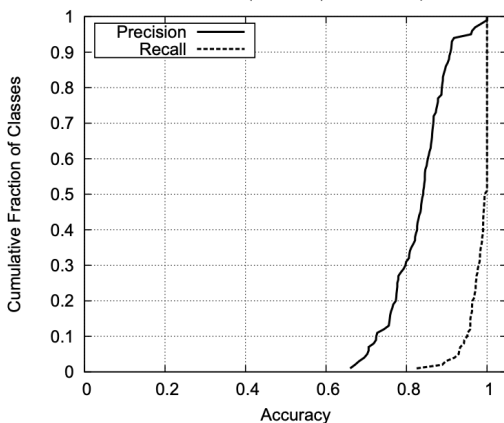    - Amazon ML
    - Google's Prediction API
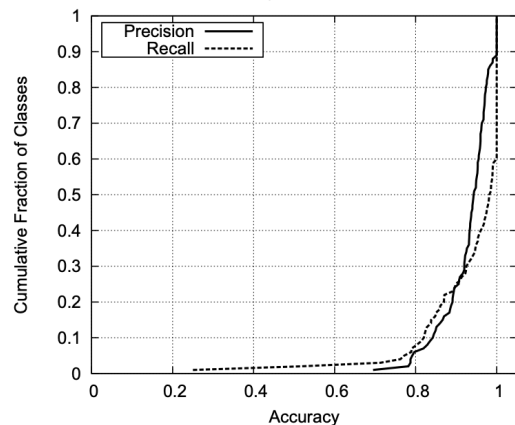  - **In-short:** across all models, MI attacks work with a pretty reasonable acc.



Purchase Dataset, Amazon (10,1e-6), Membership Inference Attack

Purchase Dataset, Amazon (100,1e-4), Membership Inference Attack

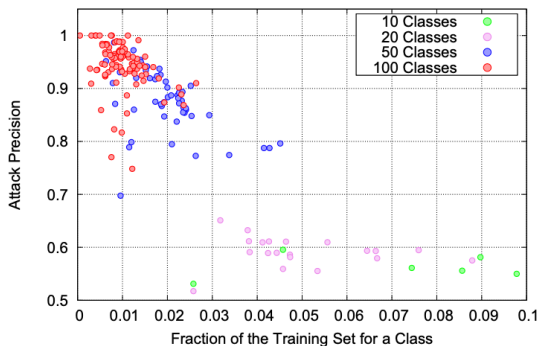Purchase Dataset, Google, Membership Inference Attack

# WHAT PRIVACY ATTACKS ARE THERE?
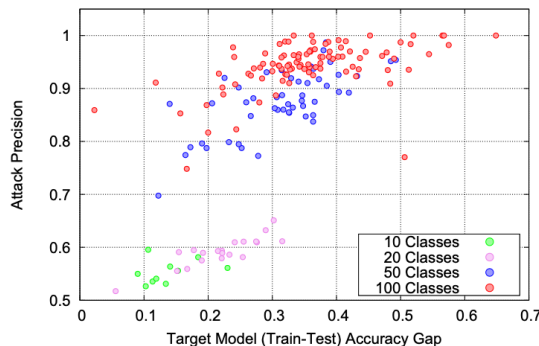
- MI attack results
  - Dataset: Purchase-100
  - Modification:
    - # Classes: $10 - 100$ (keep $N(D_{tr})$ the same)
    - Google Prediction API
  - **In-short:** more supporting data samples in the cl

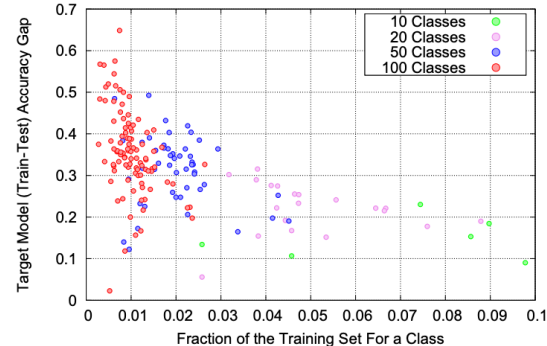| Dataset | Training Accuracy | Testing Accuracy | Attack Precision |
|---------|-------------------|------------------|------------------|
| Adult | 0.848 | 0.842 | 0.503 |
| MNIST | 0.984 | 0.928 | 0.517 |
| Location | 1.000 | 0.673 | 0.678 |
| Purchase (2) | 0.999 | 0.984 | 0.505 |
| Purchase (10) | 0.999 | 0.866 | 0.550 |
| Purchase (20) | 1.000 | 0.781 | 0.590 |
| Purchase (50) | 1.000 | 0.693 | 0.860 |
| Purchase (100) | 0.999 | 0.659 | 0.935 |
| TX hospital stays | 0.668 | 0.517 | 0.657 |



Purchase Dataset, 10-100 Classes, Google, Membership Inference Attack



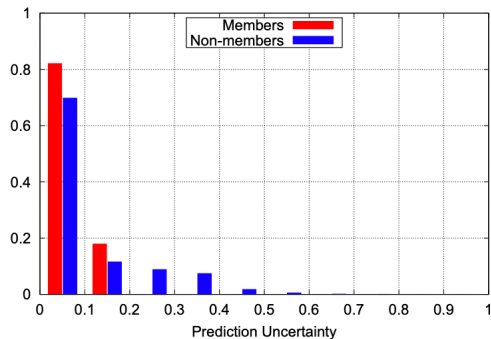Purchase Dataset, 10-100 Classes, Google, Membership Inference Attack



Purchase Dataset, 10-100 Classes, Google, Membership Inference Attack
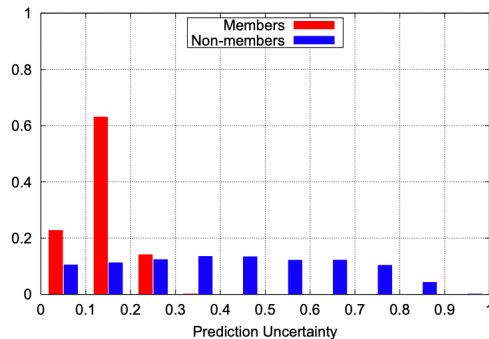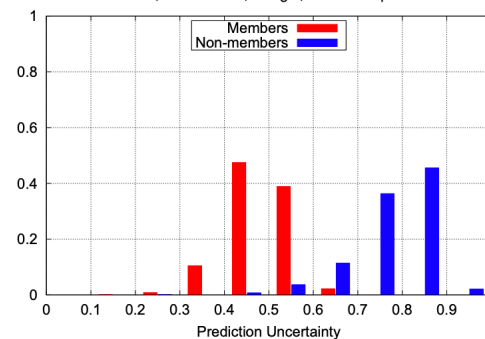
Oregon State University

# WHAT PRIVACY ATTACKS ARE THERE?

- MI attacks, why do they work?

# WHAT PRIVACY ATTACKS ARE THERE?

- Suppose: a developer who write code for your company's core products

# WHAT PRIVACY ATTACKS ARE THERE?

- Model inversion (or data extraction) attacks



Input queries $\bar{x}$ → Model $f$ → output $f(\bar{x})$

Observe correlations!

**Target**   **Softmax**   **MLP**   **DAE**

# WHAT PRIVACY ATTACKS ARE THERE?

- Model inversion attacks
  - **Costs:**
    - Per attack: 1.4sec (softmax) << 693 sec (DAE) << 1298 sec (MLP)
    - Per attack: 5.6 epochs (softmax) << 3096 epoch (MLP) << 4728.5 epoch (DAE)

  - **Accuracy:**
    - Overall: ~80% acc. (softmax) > 60% acc. (MLP) > 55% acc. (DAE)
    - Skilled workers: ~95% acc. (softmax) > 80% acc. (MLP) > 75% acc. (DAE)



**Target**   **Softmax**   **MLP**   **DAE**

Oregon State University

# WHAT PRIVACY ATTACKS ARE THERE?

- Data extraction attacks

# WHAT PRIVACY ATTACKS ARE THERE?

- Unintentional memorization
  - It does NOT mean that a model memorizes *any* data
  - It means a model memorizes *out-of-distribution* training data (*i.e., secrets*)

- Do neural networks unintentionally memorize?
  - Dataset: Penn Treebank (PTB)
  - Model: LSTM with 200 hidden units
  - Secret:
    - A sentence "My social security number is 078-05-1120"
    - Inject this sentence into the PTB dataset
  - Extraction: auto-completion
    - Type: "My social security number is 078-"
    - Shows: "My social security number is 078-05-1120"

# WHAT PRIVACY ATTACKS ARE THERE?

- Measuring memorization
  - [Def. 1] The log-**perplexity**:
  $$\mathrm{Px}_\theta(x_1...x_n) \quad = \quad -\log_2 \mathbf{Pr}(x_1...x_n|f_\theta)$$
  $$= \quad \sum_{i=1}^{n} \left( -\log_2 \mathbf{Pr}(x_i|f_\theta(x_1...x_{i-1})) \right)$$

    - It measures how *surprised* the model to see a given input sequence

  - [Notation]
    - **Canaries:** a random sequence of numbers (ex. "the random number is **281265017**")

| Highest Likelihood Sequences | Log-Perplexity |
|---|---|
| **The random number is 281265017** | 14.63 |
| The random number is 281265117 | 18.56 |
| The random number is 281265011 | 19.01 |
| The random number is 286265117 | 20.65 |
| The random number is 528126501 | 20.88 |
| The random number is 281266511 | 20.99 |
| The random number is 287265017 | 20.99 |
| The random number is 281265111 | 21.16 |
| The random number is 281265010 | 21.36 |

Oregon State University

# WHAT PRIVACY ATTACKS ARE THERE?

- Measuring memorization
  - [Def. 2] The **rank** of a canary $s[r]$:

  $$\mathbf{rank}_\theta(s[r]) = \left| \{ r' \in \mathcal{R} : \mathrm{Px}_\theta(s[r']) \leq \mathrm{Px}_\theta(s[r]) \} \right|$$

    - It measures how many random sequences that have log-perplexity *lower* than $r$ are

  - [Def. 3] The **guessing entropy** is the number of guesses $E(X)$ required in an optimal strategy to guess the value of a discrete random variable $X$
    - Brute force $\qquad\qquad : E(X) = 0.5|R|$
    - Query-access attacker $: E(s[r]|f_\theta) = \mathbf{rank}_\theta(s[r])$

  - [Def. 4] Given a canary $s[r]$, a model $f_\theta$, and the randomness space $R$, the **exposure** of the canary is:

  $$\mathbf{exposure}_\theta(s[r]) = \log_2 |\mathcal{R}| - \log_2 \mathbf{rank}_\theta(s[r])$$

# WHAT PRIVACY ATTACKS ARE THERE?

- Data extraction attacks
  - Word-level LM:
    - Dataset: WikiText-103
    - Model: SoTA models
    - Canaries: a sequence of 8 words, randomly chosen, insert 5 times

  - Results:
    - The lower the perplexity, the easier to ext.
    - The dots on the line are Pareto-optimal att.
    - 144 exposure means ext. should be possible
    - Mem. and utility are *not* highly correlated

# WHAT PRIVACY ATTACKS ARE THERE?

- Data extraction attacks
  - NMT:
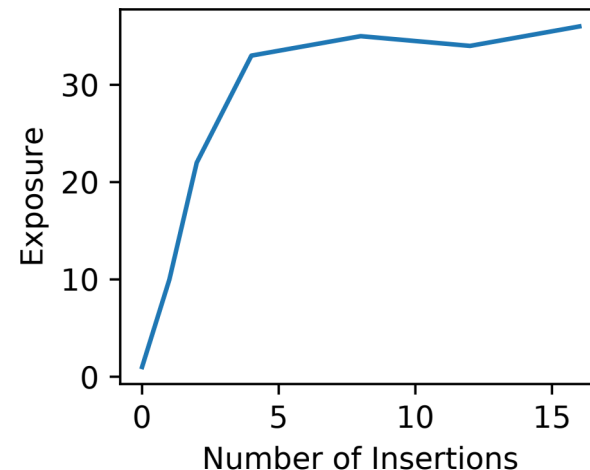    - Dataset: English-Vietnamese (100k sentence pairs)
    - Model: Good models in TensorFlow repository
    - Canaries: "My social security number is XXX-XX-XXXX" (in Vietnamese too)

  - Results:
    - Inserted once, the exposure becomes 10
      > 1000x times more likely to extract than random
    - Inserted > 4 times, the exposure becomes 30
      > completely memorized…

Oregon State University

# HOW CAN WE DEFEAT PRIVACY ATTACKS?

- $\epsilon$-Differential Privacy
  - A randomized algorithm $M: D \rightarrow R$ with domain $D$ and a range $R$ satisfies $\epsilon$-differential privacy if for any two adjacent inputs $d, d' \in D$ and any subset of outputs $S \subset R$ it holds

$$\Pr[\mathcal{M}(d) \in S] \leq e^\varepsilon \Pr[\mathcal{M}(d') \in S]$$

- $(\epsilon, \delta)$-Differential Privacy

$$\Pr[\mathcal{M}(d) \in S] \leq e^\varepsilon \Pr[\mathcal{M}(d') \in S] + \delta$$

  - $\delta$: Represent some catastrophic failure cases [Link, Link]
  - $\delta$ < 1/|d|, where |d| is the number of samples in a database

Oregon State
University

# HOW CAN WE DEFEAT PRIVACY ATTACKS?

- $(\epsilon, \delta)$-Differential Privacy **[Conceptually]**

$$\Pr[\mathcal{M}(d) \in S] \leq e^{\varepsilon} \Pr[\mathcal{M}(d') \in S] + \delta$$

  - You have two databases $d, d'$ differ by one item
  - You make the same query $M$ to each and have results $M(d)$ and $M(d')$
  - You ensure the distinguishability between the two under a measure $\epsilon$
    - $\epsilon$ is large: those two are distinguishable, less private
    - $\epsilon$ is small: the two outputs are similar, more private
  - You also ensure the catastrophic failure probability $\delta$

# How can we defeat privacy attacks?

- $(\epsilon, \delta)$-Differential Privacy

$$\Pr[\mathcal{M}(d) \in S] \leq e^{\varepsilon} \Pr[\mathcal{M}(d') \in S] + \delta$$

- Mechanism for $(\epsilon, \delta)$-DP: Gaussian noise

$$\mathcal{M}(d) \triangleq f(d) + \mathcal{N}(0, S_f^2 \cdot \sigma^2)$$

  – $M(d)$: $(\epsilon, \delta)$-DP query output on $d$
  – $f(d)$: non $(\epsilon, \delta)$-DP (original) query output on $d$
  – $N(0, S_f^2 \cdot \sigma^2)$: Gaussian normal distribution with mean 0 and the std. of $S_f^2 \cdot \sigma^2$

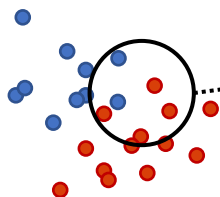**Post-hoc: Set the Goal $\epsilon$ and Calibrate the noise $S_f^2 \cdot \sigma^2$!**

# HOW CAN WE DEFEAT PRIVACY ATTACKS?

- Revisit'ed – Mini-batch SGD
  1. At each step $t$, it takes a mini-batch $L_t$
  2. Computes the loss $\mathcal{L}(\theta)$ over the samples in $L_t$, w.r.t. the label $y$
  3. Computes the gradients $g_t$ of $\mathcal{L}(\theta)$
  4. Update the model parameters $\theta$ towards the direction of reducing the loss
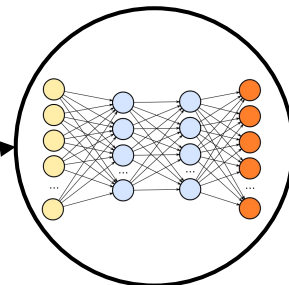


**This Process Should Be ($\epsilon$, $\delta$)-DP!**

$D$: a training set

1. Take $L_t$, and compute $\mathcal{L}(\theta)$
2. Compute $g_t$ of $\mathcal{L}(\theta)$
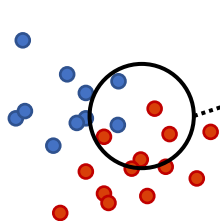3. Update the $\theta$

$\theta$: a model

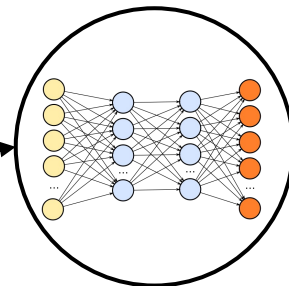# HOW CAN WE DEFEAT PRIVACY ATTACKS?

- Mini-batch SGD to DP-SGD

  1. At each step $t$, it takes a mini-batch $L_t$
  2. Computes the loss $\mathcal{L}(\theta)$ over the samples in $L_t$, w.r.t. the label $y$
  3. Computes the gradients $g_t$ of $\mathcal{L}(\theta)$
  4. Clip (scale) the gradients to $1/C$, where $C > 1$
  5. Add Gaussian random noise $N(0, \sigma^2 C^2 \mathbf{I})$ to $g_t$
  6. Update the model parameters $\theta$ towards the direction of reducing the loss

$D$: a training set

1. Take $L_t$, and compute $\mathcal{L}(\theta)$
2. Compute $g_t$ of $\mathcal{L}(\theta)$
3. Clip $g_t$ and add noise
4. Update the $\theta$

$\theta$: a model

Oregon State University

# TOPICS FOR THIS WEEK

- Trustworthy AI
  - Motivation
  - Preliminaries
    - Machine learning (ML)
    - ML-based systems
  - (Potential) Threats
    - Adversarial attacks
    - Data poisoning
    - Privacy attacks
  - Discussion
    - More issues (social bias, fairness, …)

# Thank You!

Tu/Th 4:00 – 5:50 PM

Sanghyun Hong

sanghyun.hong@oregonstate.edu

**Oregon State University**

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