

Deep Obfuscation: Precise Masking of Sensitive Information to Protect Against Machine Learning Adversaries

Yuan Gong and Christian Poellabauer

Department of Computer Science & Engineering, University of Notre Dame



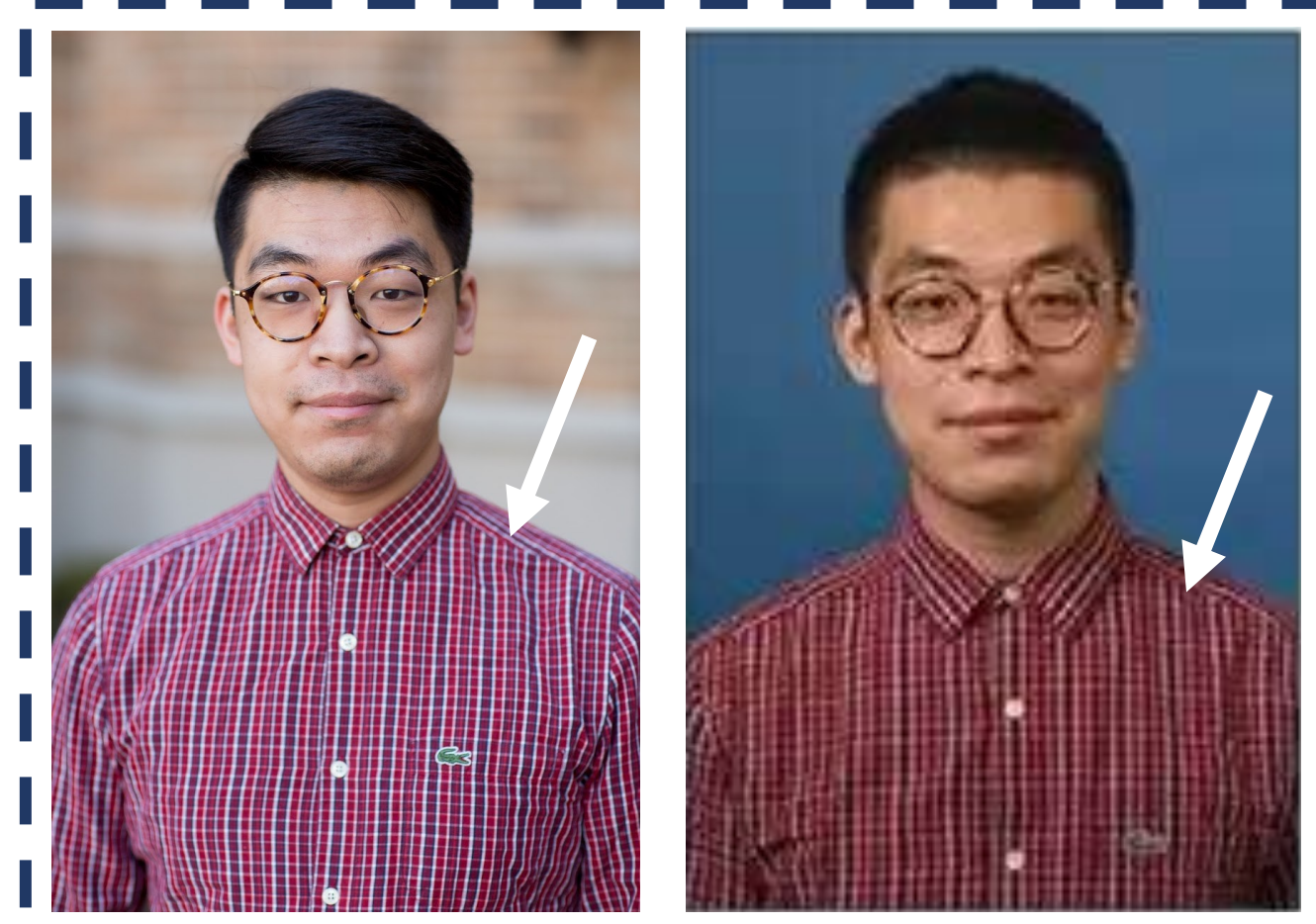
▼ Motivation

Background: It is often desirable to share data (e.g., tables, images, audio, sensor streams, etc.), while removing some sensitive information embedded in such data.

Problem: It is hard to mask data **precisely** and **completely**. A **machine learning adversary** might be able to infer the sensitive information based on **residual cues**.

Solution: We propose **Deep Obfuscation**, i.e., a deep learning based data masking scheme that aims to **precisely** and **maximally** remove sensitive information from target data, while **minimally** affecting non-sensitive information.

▼ Attack Example

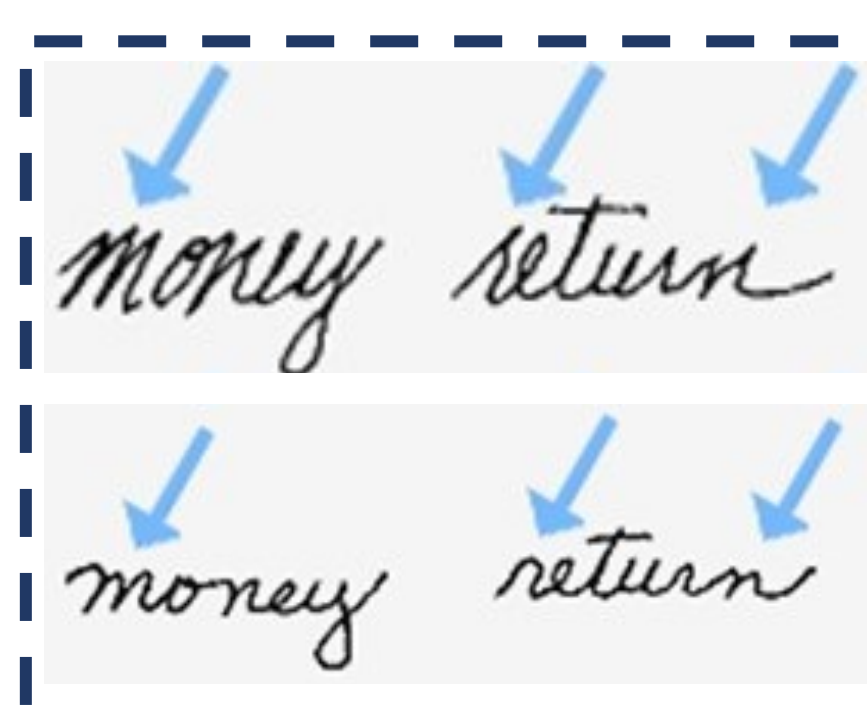


Person A in Training Set



Is this person A?

Blocking face is not sufficient. Clothing, gestures, environment, etc., could be used by an attacker to infer a person's identity.



Handwriting style can reveal the identity, but masking the style also removes the content information.

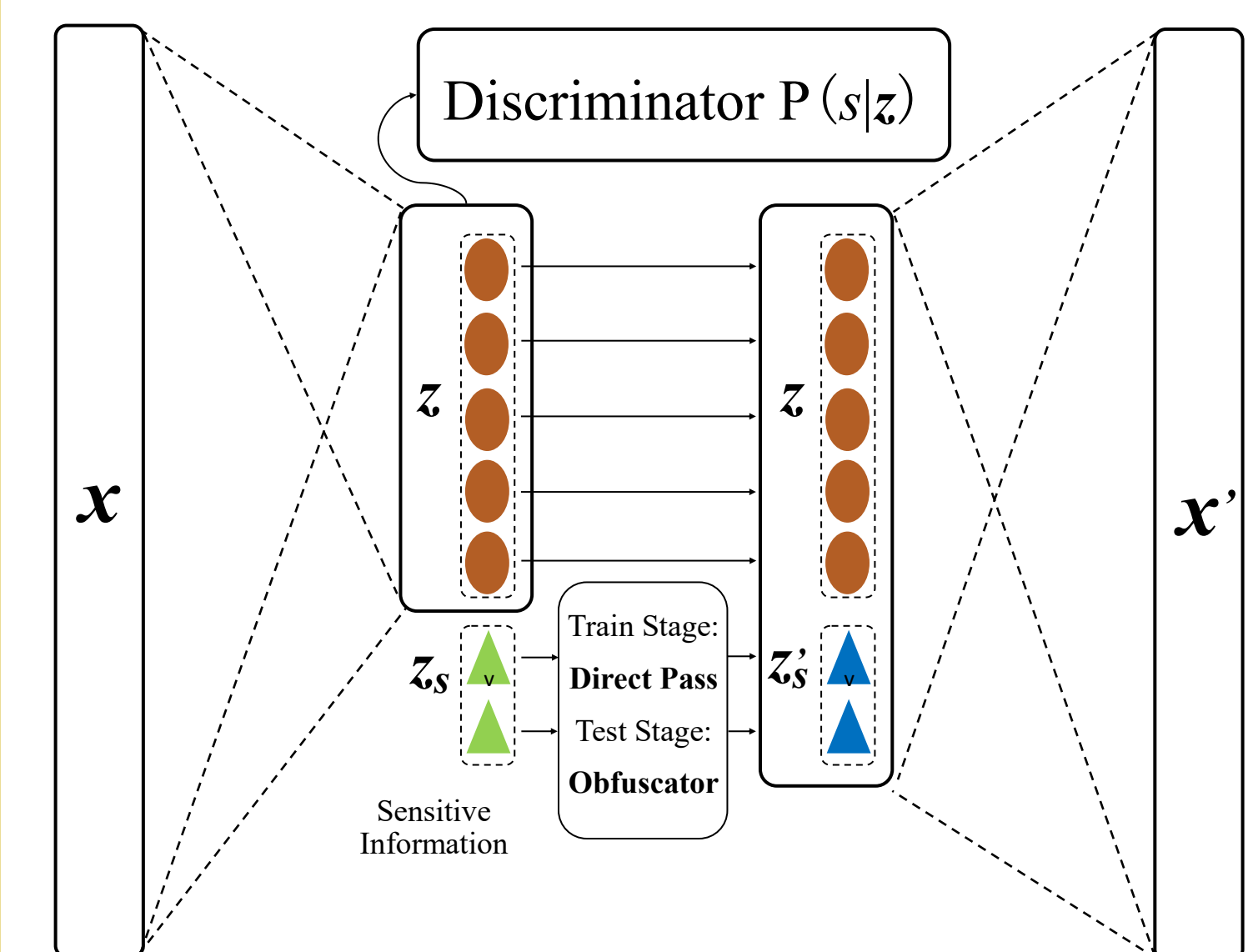
▼ Problem Formalization

Assumption: 1) A data sample $x \sim P_x$ contains some discrete sensitive information s . 2) There exists an underlying relationship $s = F_s(x)$.

Attack: An attacker can infer s based on x by learning a machine learning model $P_{\text{model}}(s | x)$ using a dataset $D \sim P_x$ with labels of s .

The Masking Task: Modifying x to x' , where x' cannot be used to infer sensitive information, i.e., the inference confidence $P_{\text{model}}(s | x')$ is small for well-designed machine learning models trained on $D \sim P_x$.

▼ Deep Obfuscation



Data: $D = \{\{x_1, s_1\}, \{x_2, s_2\}, \dots, \{x_n, s_n\}\}$
while $\frac{1}{n} \sum_i \log P(s_i | z_i) > \text{threshold}$ **do**
 for $j := 1: \# \text{discriminator_training_epoch}$ **do**
 optimize discriminator parameters θ_{dis} using Equ. 1;
 end
 for $k := 1: \# \text{autoencoder_training_epoch}$ **do**
 optimize autoencoder parameters θ_{ae} using Equ. 2;
 end
end

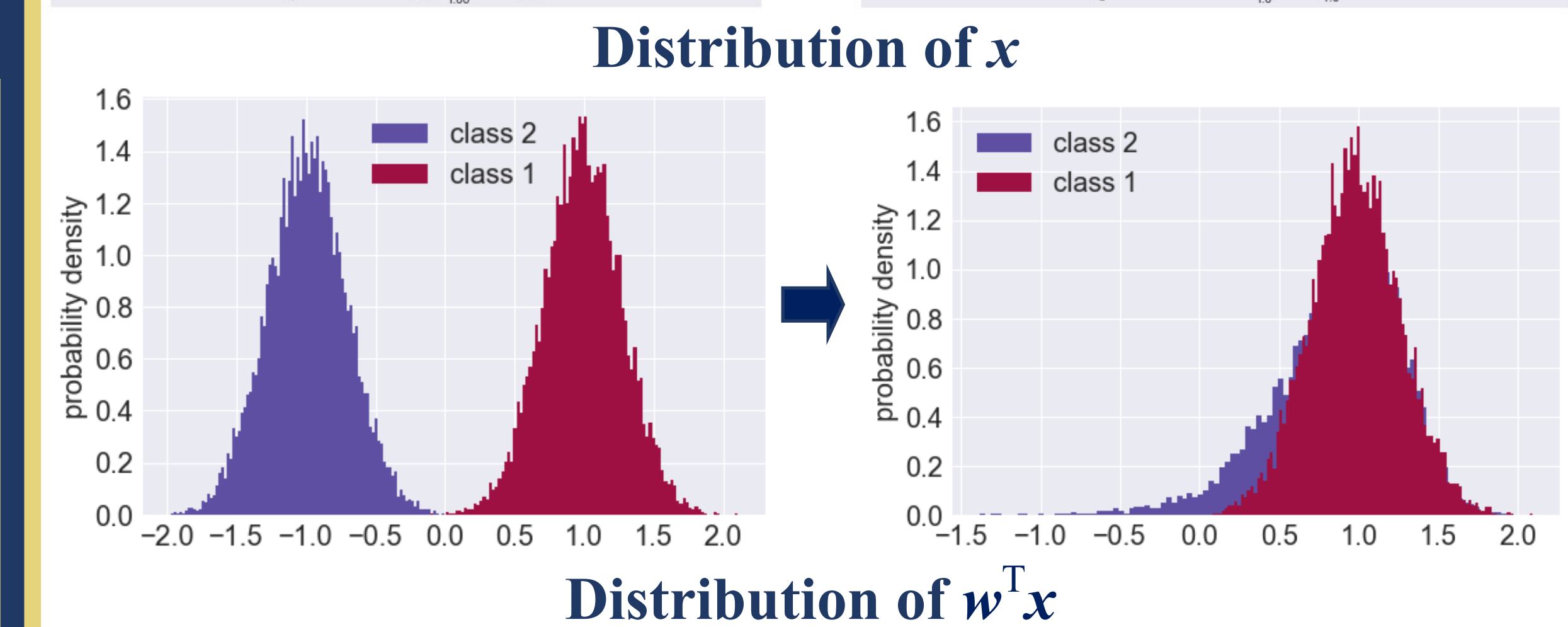
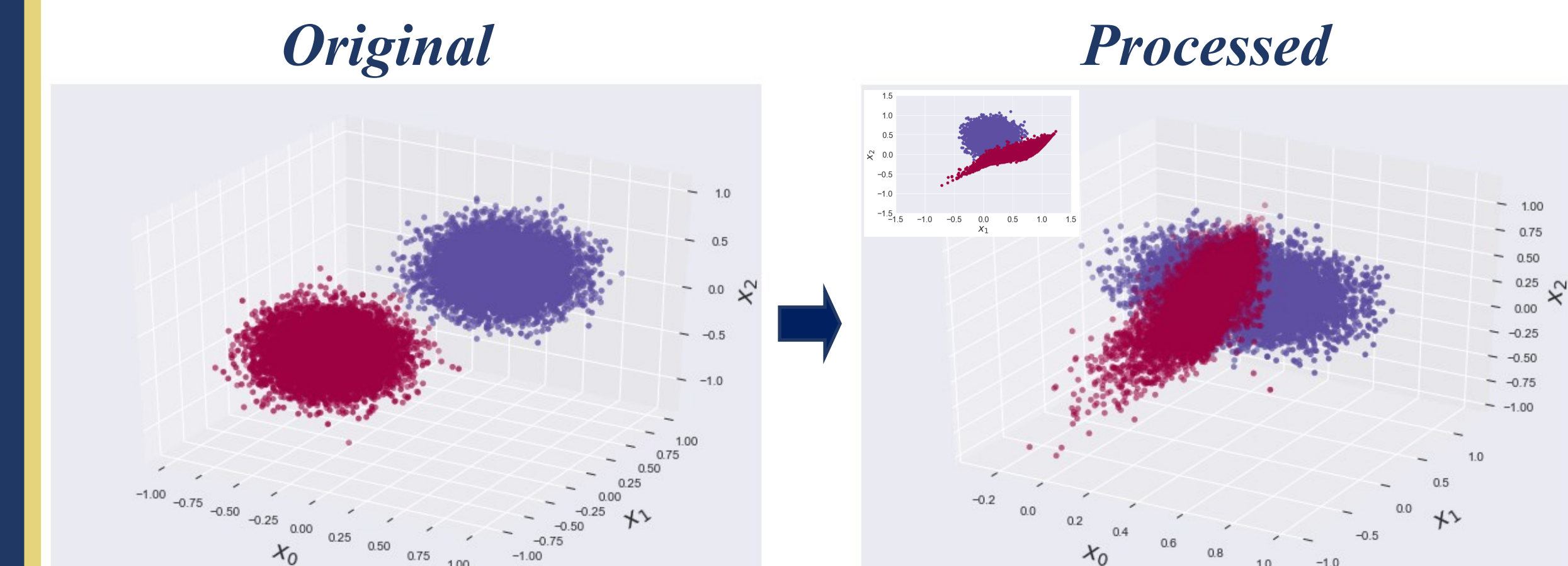
- Autoencoder based adversarial training scheme.
- The discriminator serves as the **imaginary adversary**, which tries to estimate sensitive information s from the latent variable z , i.e.,

$$L_{\text{discriminator}} = \text{mean}_{x \in D \sim P_x} (-\log P(s | z)) \quad (1)$$
- The goal of the encoder is to avoid such an attack, hence, the discriminator's inference performance is added into the auto-encoder loss, i.e.,

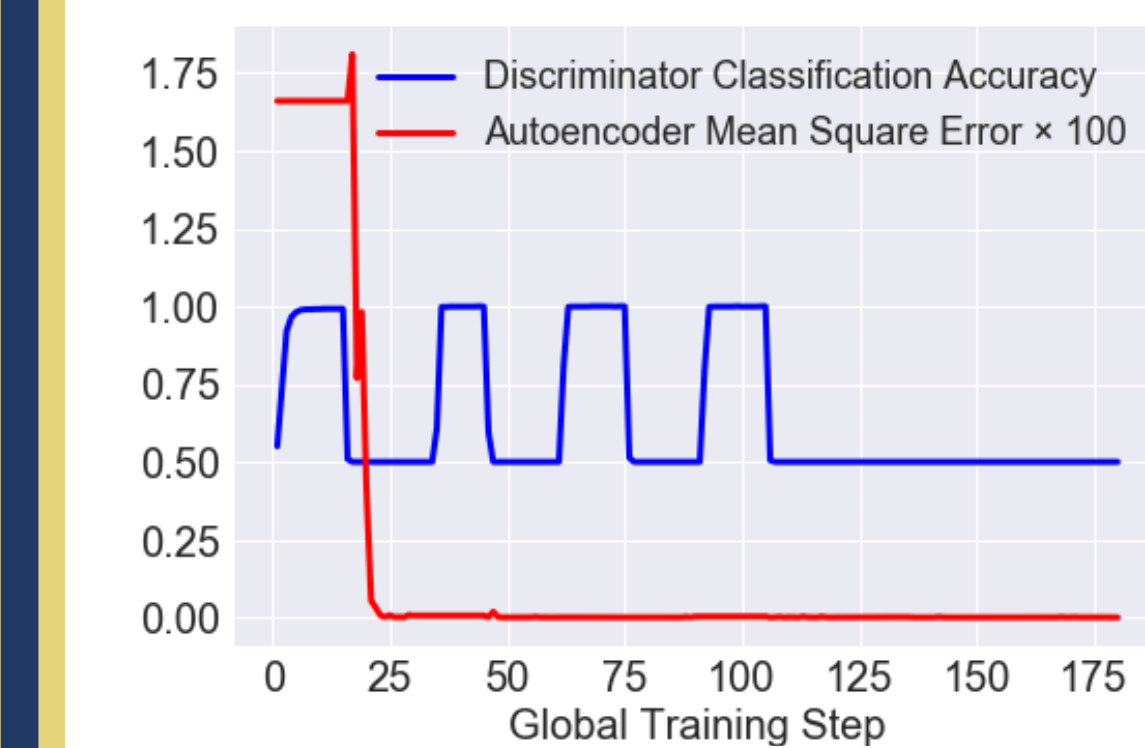
$$L_{\text{autoencoder}} = -L_{\text{discriminator}} + \lambda \times L_{\text{reconstruction}} \quad (2)$$
- The discriminator and the autoencoder are trained alternatively, i.e., during the training of the autoencoder, the parameters of the discriminator are frozen and vice versa.
- This scheme forces the encoder to encode all information other than s in z , and the decoder to effectively decode from z .

▼ Low Dimensional Example

Task: Remove information of s from x
 $x \in \mathbf{R}^3 \sim \text{Mixed Gaussian Distribution}$
 $s = F_s(x) = \mathbf{I}_{w^T x > 0}(x)$, $w = [1, 1, 1]$



- x' maintains some main characteristics of the original x (e.g., the two Gaussian clusters are still separable).
- Sensitive information is almost completely obfuscated.



- The discriminator and the autoencoder reach a balance.

▼ Future Work

- Design the training procedure and the network architecture for processing high dimensional data (e.g., images) and sequential data (e.g., audio, sensor stream).
- Apply deep obfuscation scheme to a variety type of sensitive information such as identity, health status, and passwords.

Contact: Yuan Gong, ygong1@nd.edu