

# Variational Leakage: The Role of Information Complexity in Privacy Leakage

Appendices included

Amir Ahooye Atashin\*  
amir.atashin@gmail.com  
University of Geneva  
Geneva, Switzerland

Deniz Gündüz  
d.gunduz@imperial.ac.uk  
Imperial College London  
London, United Kingdom

Behrooz Razeghi<sup>†\*</sup>  
behrooz.razeghi@unige.ch  
University of Geneva  
Geneva, Switzerland

Slava Voloshynovskiy  
svolos@unige.ch  
University of Geneva  
Geneva, Switzerland

## ABSTRACT

We study the role of information complexity in privacy leakage about an attribute of an adversary’s interest, which is not known *a priori* to the system designer. Considering the supervised representation learning setup and using neural networks to parameterize the variational bounds of information quantities, we study the impact of the following factors on the amount of information leakage: information complexity regularizer weight, latent space dimension, the cardinalities of the known utility and unknown sensitive attribute sets, the correlation between utility and sensitive attributes, and a potential bias in a sensitive attribute of adversary’s interest. We conduct extensive experiments on Colored-MNIST and CelebA datasets to evaluate the effect of information complexity on the amount of intrinsic leakage.

A repository of the proposed method implementation, Colored-MNIST dataset generator and the corresponding analysis is publicly available at:

<https://github.com/BehroozRazeghi/Variational-Leakage>

## KEYWORDS

Information complexity, privacy, intrinsic leakage, statistical inference, information bottleneck

## 1 INTRODUCTION

Sensitive information sharing is a challenging problem in information systems. It is often handled by obfuscating the available information before sharing it with other parties. In [17], this problem has been formalized as the **privacy funnel (PF)** in an information theoretic framework. Given two correlated random variables  $S$  and  $X$  with a joint distribution  $P_{S,X}$ , where  $X$  represents the available information and  $S$  the private latent variable, the goal of the PF model is to find a representation  $Z$  of  $X$  using a stochastic mapping  $P_{Z|X}$  such that: (i)  $S \rightarrow X \rightarrow Z$  form a Markov chain; and (ii) representation  $Z$  is maximally informative about the useful data  $X$  (maximizing Shannon’s mutual information (MI)  $I(X; Z)$ ) while

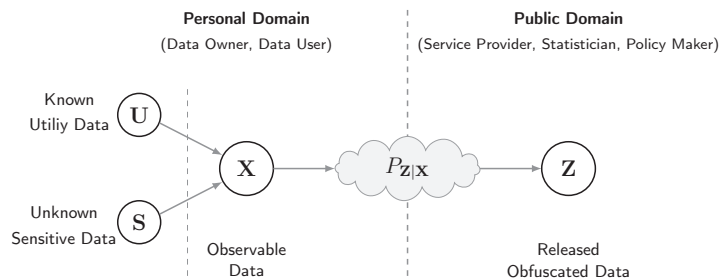


Figure 1: The general setup.

being minimally informative about the sensitive data  $S$  (minimizing  $I(S; Z)$ ). There have been many extensions of this model in the recent literature, e.g., [3, 9, 17, 19–23, 27].

In this paper, we will consider a delicate generalization of the PF model considered in [3, 21], where the goal of the system designer is not to reveal the data that has available but another correlated utility variable. In particular, we assume that the data owner/user acquires some utility from the service provider based on the amount of information disclosed about a utility random variable  $U$  correlated with  $X$ , measured by  $I(U; Z)$ . Therefore, considering Markov chain  $(U, S) \rightarrow X \rightarrow Z$ , the data owner’s aim is to share a representation  $Z$  of *observed data*  $X$ , through a stochastic mapping  $P_{Z|X}$ , while preserving information about *utility attribute*  $U$  and obfuscate information about *sensitive attribute*  $S$  (see Fig. 1).

The implicit assumption in the PF model presented above and the related generative adversarial privacy framework [10, 31] is to have *pre-defined interests* in the game between the ‘defender’ (data owner/user) and the ‘adversary’; that is, the data owner knows in advance what feature/ variable of the underlying data the adversary is interested in. Accordingly, the data release mechanism can be optimized/ tuned to minimize any inference the adversary can make about this specific random variable. However, this assumption is violated in most real-world scenarios. The attribute that the defender may assume as sensitive may not be the attribute of interest for the inferential adversary. As an example, for a given utility

\*Both authors contributed equally to this research.

<sup>†</sup>Work done while at Imperial College London.

task at hand, the defender may try to restrict inference on gender recognition while the adversary is interested in inferring an individual’s identity or facial emotion. Inspired by [11], and in contrast to the above setups, we consider the scenario in which the adversary is curious about an attribute that is *unknown* to the system designer.

In particular, we argue that the information complexity of the representation measured by MI  $I(\mathbf{X}; \mathbf{Z})$  can also limit the information leakage about the unknown sensitive variable. In this paper, obtaining the parameterized variational approximation of information quantities, we investigate the core idea of [11] in the supervised representation learning setup.

**Notations:** Throughout this paper, random vectors are denoted by capital bold letters (e.g.,  $\mathbf{X}$ ), deterministic vectors are denoted by small bold letters (e.g.,  $\mathbf{x}$ ), and alphabets (sets) are denoted by calligraphic fonts (e.g.,  $\mathcal{X}$ ). We use the shorthand  $[N]$  to denote the set  $\{1, 2, \dots, N\}$ .  $H(P_X) := \mathbb{E}_{P_X}[-\log P_X]$  denotes the Shannon’s entropy, while  $H(P_X \| Q_X) := \mathbb{E}_{P_X}[-\log Q_X]$  denotes the cross-entropy of the distribution  $P_X$  relative to a distribution  $Q_X$ . The relative entropy is defined as  $D_{\text{KL}}(P_X \| Q_X) := \mathbb{E}_{P_X}[\log \frac{P_X}{Q_X}]$ . The conditional relative entropy is defined by:

$$D_{\text{KL}}(P_{Z|X} \| Q_{Z|X} | P_X) := \mathbb{E}_{P_X} \left[ D_{\text{KL}}(P_{Z|X=\mathbf{x}} \| Q_{Z|X=\mathbf{x}}) \right].$$

And the MI is defined by:

$$I(P_X; P_{Z|X}) := D_{\text{KL}}(P_{Z|X} \| P_Z | P_X)$$

We abuse notation to write  $H(\mathbf{X}) = H(P_X)$  and  $I(\mathbf{X}; \mathbf{Z}) = I(P_X; P_{Z|X})$  for random vectors  $\mathbf{X} \sim P_X$  and  $\mathbf{Z} \sim P_Z$ .

## 2 PROBLEM FORMULATION

Given the observed data  $\mathbf{X}$ , the data owner wishes to release a representation  $\mathbf{Z}$  for a utility task  $\mathbf{U}$ . Our aim is to investigate the potential statistical inference about a sensitive random attribute  $\mathbf{S}$  from the released representation  $\mathbf{Z}$ . The sensitive attribute  $\mathbf{S}$  is possibly also correlated with  $\mathbf{U}$  and  $\mathbf{X}$ .

The objective is to obtain a stochastic map  $P_{Z|X} : \mathcal{X} \rightarrow \mathcal{Z}$  such that  $P_{U|Z} \approx P_{U|X}$ ,  $\forall \mathbf{Z} \in \mathcal{Z}$ ,  $\forall \mathbf{U} \in \mathcal{U}$ ,  $\forall \mathbf{X} \in \mathcal{X}$ . This means that the posterior distribution of the utility attribute  $\mathbf{U}$  is similar when conditioned on the released representation  $\mathbf{Z}$  or on the original data  $\mathbf{X}$ . Under logarithmic loss, one can measure the utility by Shannon’s MI [17, 23, 29]. The logarithmic loss function has been widely used in learning theory [6], image processing [2], information bottleneck [8], multi-terminal source coding [7], and PF [17].

**Threat Model:** We make minimal assumptions about the adversary’s goal, which can model a large family of potential adversaries. In particular, we have the following assumptions:

- The distribution  $P_{S|X}$  is unknown to the data user/owner. We only restrict attribute  $\mathbf{S}$  to be discrete, which captures most scenarios of interest, e.g., a facial attribute, an identity, a political preference.
- The adversary observes released representation  $\mathbf{Z}$  and the Markov chain  $(\mathbf{U}, \mathbf{S}) \text{---} \mathbf{X} \text{---} \mathbf{Z}$  holds.
- We assume the adversary knows the mapping  $P_{Z|X}$  designed by the data owner, i.e., the data release mechanism is public.

Furthermore, the adversary may have access to a collection of the original dataset with the corresponding labels  $\mathbf{S}$ .

Suppose that the sensitive attribute  $\mathbf{S} \in \mathcal{S}$  has a uniform distribution over a discrete set  $\mathcal{S}$ , where  $|\mathcal{S}| = 2^L < \infty$ . If  $I(\mathbf{S}; \mathbf{Z}) \geq L - \epsilon$ , then equivalently  $H(\mathbf{S} | \mathbf{Z}) \leq \epsilon$ . Also note that due to the Markov chain  $\mathbf{S} \text{---} \mathbf{X} \text{---} \mathbf{Z}$ , we have  $I(\mathbf{S}; \mathbf{Z}) = I(\mathbf{X}; \mathbf{Z}) - I(\mathbf{X}; \mathbf{Z} | \mathbf{S})$ . When  $\mathbf{S}$  is not known a priori, the data owner has no control over  $I(\mathbf{X}; \mathbf{Z} | \mathbf{S})$ . On the other hand,  $I(\mathbf{X}; \mathbf{Z})$  can be interpreted as the information complexity of the released representation, which plays a critical role in controlling the information leakage  $I(\mathbf{S}; \mathbf{Z})$ . Note also that a statistic  $\mathbf{Z} = f(\mathbf{X})$  induces a partition on the sample space  $\mathcal{X}$ , where  $\mathbf{Z}$  is sufficient statistic for  $\mathbf{U}$  if and only if the assigned samples in each partition do not depend on  $\mathbf{U}$ . Hence, intuitively, a larger  $|\mathcal{U}|$  induces finer partitions on  $\mathcal{X}$ , which could potentially lead to more leakage about the unknown random function  $\mathbf{S}$  of  $\mathbf{X}$ . This is the core concept of the notion of *variational leakage*, which we shortly address in our experiments.

Since the data owner does not know the particular sensitive variable of interest to the adversary, we argue that it instead aims to design  $P_{Z|X}$  with the minimum (information) complexity and minimum utility loss. With the introduction of a Lagrange multiplier  $\beta \in [0, 1]$ , we can formulate the objective of the data owner by *maximizing* the associated Lagrangian functional:

$$\mathcal{L}(P_{Z|X}, \beta) = I(\mathbf{U}; \mathbf{Z}) - \beta I(\mathbf{X}; \mathbf{Z}). \quad (1)$$

This is the well-known **information bottleneck (IB)** principle [29], which formulates the problem of extracting, in the most succinct way, the relevant information from random variable  $\mathbf{X}$  about the random variable of interest  $\mathbf{U}$ . Given two correlated random variables  $\mathbf{U}$  and  $\mathbf{X}$  with joint distribution  $P_{U,X}$ , the goal is to find a representation  $\mathbf{Z}$  of  $\mathbf{X}$  using a stochastic mapping  $P_{Z|X}$  such that: (i)  $\mathbf{U} \text{---} \mathbf{X} \text{---} \mathbf{Z}$ , and (ii)  $\mathbf{Z}$  is maximally informative about  $\mathbf{U}$  (maximizing  $I(\mathbf{U}; \mathbf{Z})$ ) and minimally informative about  $\mathbf{X}$  (minimizing  $I(\mathbf{X}; \mathbf{Z})$ ).

Note that in the PF model,  $I(\mathbf{X}; \mathbf{Z})$  measures the *useful* information, which is of the designer’s interest, while in the IB model,  $I(\mathbf{U}; \mathbf{Z})$  measures the *useful* information. Hence,  $I(\mathbf{X}; \mathbf{Z} | \mathbf{S})$  in PF quantifies the *residual* information, while  $I(\mathbf{X}; \mathbf{Z} | \mathbf{U})$  in IB quantifies the *redundant* information.

In the sequel, we provide the parameterized variational approximation of information quantities, and then study the impact of the information complexity  $I(\mathbf{X}; \mathbf{Z})$  on the information leakage for an unknown sensitive variable.

### 2.1 Variational Approximation of Information Measures

Let  $Q_{U|Z} : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{U})$ ,  $Q_{S|Z} : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{S})$ ,  $Q_Z : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{Z})$  be variational approximations of the optimal utility decoder distribution  $P_{U|Z}$ , adversary decoder distribution  $P_{S|Z}$ , and latent space distribution  $P_Z$ , respectively. The common approach is to use **deep neural networks (DNNs)** to model/parameterized these distributions. Let  $P_\phi(\mathbf{Z} | \mathbf{X})$  denote the family of encoding probability distributions  $P_{Z|X}$  over  $\mathcal{Z}$  for each element of space  $\mathcal{X}$ , parameterized by the output of a DNN  $f_\phi$  with parameters  $\phi$ . Analogously, let  $P_\theta(\mathbf{U} | \mathbf{Z})$  and  $P_\xi(\mathbf{S} | \mathbf{Z})$  denote the corresponding family of decoding probability distributions  $Q_{U|Z}$  and  $Q_{S|Z}$ , respectively,

parameterized by the output of DNNs  $g_\theta$  and  $g_\xi$ . Let  $P_D(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n)$ ,  $\mathbf{x}_n \in \mathcal{X}$  denote the empirical data distribution. In this case,  $P_\phi(\mathbf{X}, \mathbf{Z}) = P_D(\mathbf{X})P_\phi(\mathbf{Z}|\mathbf{X})$  denotes our joint inference data distribution, and  $P_\phi(\mathbf{Z}) = \mathbb{E}_{P_D(\mathbf{X})} [P_\phi(\mathbf{Z}|\mathbf{X})]$  denotes the learned *aggregated* posterior distribution over latent space  $\mathcal{Z}$ .

**Information Complexity:** The information complexity can be decomposed as:

$$\begin{aligned} I(\mathbf{X}; \mathbf{Z}) &= \mathbb{E}_{P_{\mathbf{X}, \mathbf{Z}}} \left[ \log \frac{P_{\mathbf{X}, \mathbf{Z}}}{P_{\mathbf{X}} P_{\mathbf{Z}}} \right] = \mathbb{E}_{P_{\mathbf{X}, \mathbf{Z}}} \left[ \log \frac{P_{\mathbf{Z}|\mathbf{X}} Q_{\mathbf{Z}}}{Q_{\mathbf{Z}} P_{\mathbf{Z}}} \right] \\ &= \mathbb{E}_{P_{\mathbf{X}}} \left[ \text{D}_{\text{KL}}(P_{\mathbf{Z}|\mathbf{X}} \| Q_{\mathbf{Z}}) \right] - \text{D}_{\text{KL}}(P_{\mathbf{Z}} \| Q_{\mathbf{Z}}). \end{aligned} \quad (2)$$

Where  $Q_{\mathbf{Z}}$  is the latent space's prior.

Therefore, the parameterized variational approximation of information complexity (2) can be recast as:

$$I_\phi(\mathbf{X}; \mathbf{Z}) := \text{D}_{\text{KL}}(P_\phi(\mathbf{Z}|\mathbf{X}) \| Q_{\mathbf{Z}} | P_D(\mathbf{X})) - \text{D}_{\text{KL}}(P_\phi(\mathbf{Z}) \| Q_{\mathbf{Z}}). \quad (3)$$

The optimal prior  $Q_{\mathbf{Z}}^*$  minimizing the information complexity is  $Q_{\mathbf{Z}}^*(\mathbf{z}) = \mathbb{E}_{P_D(\mathbf{X})} [P_\phi(\mathbf{Z} | \mathbf{X} = \mathbf{x})]$ ; however, it may potentially lead to over-fitting. A critical challenge is to guarantee that the learned aggregated posterior distribution  $P_\phi(\mathbf{Z})$  conforms well to the prior  $Q_{\mathbf{Z}}$  [4, 13, 24, 26, 30]. We can cope with this issue by employing a more *expressive* form for  $Q_{\mathbf{Z}}$ , which would allow us to provide a good fit of an arbitrary space for  $\mathcal{Z}$ , at the expense of additional *computational complexity*.

**Information Utility:** The parameterized variational approximation of MI between the released representation  $\mathbf{Z}$  and the utility attribute  $\mathbf{U}$  can be recast as:

$$\begin{aligned} I_{\phi, \theta}(\mathbf{U}; \mathbf{Z}) &:= \\ &\mathbb{E}_{P_{\mathbf{U}, \mathbf{X}}} \left[ \mathbb{E}_{P_\phi(\mathbf{Z}|\mathbf{X})} \left[ \log \frac{P_\theta(\mathbf{U}|\mathbf{Z})}{P_{\mathbf{U}}} \cdot \frac{P_\theta(\mathbf{U})}{P_\theta(\mathbf{U})} \right] \right] = \\ &\mathbb{E}_{P_{\mathbf{U}, \mathbf{X}}} \left[ \mathbb{E}_{P_\phi(\mathbf{Z}|\mathbf{X})} [\log P_\theta(\mathbf{U}|\mathbf{Z})] \right] \\ &\quad - \mathbb{E}_{P_{\mathbf{U}}} \left[ \log \frac{P_{\mathbf{U}}}{P_\theta(\mathbf{U})} \right] + \mathbb{E}_{P_{\mathbf{U}}} [\log P_\theta(\mathbf{U})] \\ &= -\text{H}_{\phi, \theta}(\mathbf{U}|\mathbf{Z}) - \text{D}_{\text{KL}}(P_{\mathbf{U}} \| P_\theta(\mathbf{U})) + \text{H}(P_{\mathbf{U}} \| P_\theta(\mathbf{U})) \\ &\geq \underbrace{-\text{H}_{\phi, \theta}(\mathbf{U}|\mathbf{Z})}_{\text{Prediction Fidelity}} - \underbrace{\text{D}_{\text{KL}}(P_{\mathbf{U}} \| P_\theta(\mathbf{U}))}_{\text{Distribution Discrepancy Loss}}, \end{aligned}$$

where  $\text{H}_{\phi, \theta}(\mathbf{U}|\mathbf{Z}) = -\mathbb{E}_{P_{\mathbf{U}, \mathbf{X}}} \left[ \mathbb{E}_{P_\phi(\mathbf{Z}|\mathbf{X})} [\log P_\theta(\mathbf{U}|\mathbf{Z})] \right]$  represents the parameterized decoder uncertainty, and in the last line we use the positivity of the cross-entropy  $\text{H}(P_{\mathbf{U}} \| P_\theta(\mathbf{U}))$ .

### 3 LEARNING MODEL

**System Designer:** Given independent and identically distributed (i.i.d.) training samples  $\{(\mathbf{u}_n, \mathbf{x}_n)\}_{n=1}^N \subseteq \mathcal{U} \times \mathcal{X}$ , and using stochastic gradient descent (SGD)-type approach, DNNs  $f_\phi$ ,  $g_\theta$ ,  $D_\eta$ , and  $D_\omega$  are trained together to maximize a Monte-Carlo approximation of the deep variational IB functional over parameters  $\phi$ ,  $\theta$ ,  $\eta$ , and  $\omega$  (Fig. 2). Backpropagation through random samples from the posterior distribution  $P_\phi(\mathbf{Z}|\mathbf{X})$  is required in our framework, which is a challenge since backpropagation cannot flow via random nodes; to overcome this hurdle, we apply the reparameterization approach [14]).

The inferred posterior distribution is typically assumed to be a multi-variate Gaussian with a diagonal co-variance, i.e.,  $P_\phi(\mathbf{Z} | \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_\phi(\mathbf{x})))$ . Suppose  $\mathcal{Z} = \mathbb{R}^d$ . We first sample a random variable  $\boldsymbol{\varepsilon}$  i.i.d. from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , then given data sample  $\mathbf{x} \in \mathcal{X}$ , we generate the sample  $\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\varepsilon}$ , where  $\odot$  is the element-wise (Hadamard) product. The latent space prior distribution is typically considered as a fixed  $d$ -dimensional standard isotropic multi-variate Gaussian, i.e.,  $Q_{\mathbf{Z}} = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . For this simple choice, the information complexity upper bound

$$\mathbb{E}_{P_\phi(\mathbf{X}, \mathbf{Z})} \left[ \log \frac{P_\phi(\mathbf{Z} | \mathbf{X})}{Q_{\mathbf{Z}}} \right] = \mathbb{E}_{P_D(\mathbf{X})} \left[ \text{D}_{\text{KL}}(P_\phi(\mathbf{Z} | \mathbf{X}) \| Q_{\mathbf{Z}}) \right]$$

has a closed-form expression, which reads as:

$$2 \text{D}_{\text{KL}}(P_\phi(\mathbf{Z} | \mathbf{X} = \mathbf{x}) \| Q_{\mathbf{Z}}) = \|\boldsymbol{\mu}_\phi(\mathbf{x})\|_2^2 + d + \sum_{i=1}^d (\boldsymbol{\sigma}_\phi(\mathbf{x})_i - \log \boldsymbol{\sigma}_\phi(\mathbf{x})_i)$$

The KL-divergences in (3) and (4) can be estimated using the density-ratio trick [18, 28], utilized in the GAN framework to directly match the data and generated model distributions. The trick is to express two distributions as conditional distributions, conditioned on a label  $C \in \{0, 1\}$ , and reduce the task to binary classification. The key point is that we can estimate the KL-divergence by estimating the ratio of two distributions without modeling each distribution explicitly.

Consider  $\text{D}_{\text{KL}}(P_\phi(\mathbf{Z}) \| Q_{\mathbf{Z}}) = \mathbb{E}_{P_\phi(\mathbf{Z})} [\log \frac{P_\phi(\mathbf{Z})}{Q_{\mathbf{Z}}}]$ . We now define  $\rho_{\mathbf{Z}}(\mathbf{z} | c)$  as  $\rho_{\mathbf{Z}}(\mathbf{z} | c = 1) = P_\phi(\mathbf{Z})$ ,  $\rho_{\mathbf{Z}}(\mathbf{z} | c = 0) = Q_{\mathbf{Z}}$ . Suppose that a perfect binary classifier (discriminator)  $D_\eta(\mathbf{z})$ , with parameters  $\eta$ , is trained to associate the label  $c = 1$  to samples from distribution  $P_\phi(\mathbf{Z})$  and the label  $c = 0$  to samples from  $Q_{\mathbf{Z}}$ . Using the Bayes' rule and assuming that the marginal class probabilities are equal, i.e.,  $\rho(c = 1) = \rho(c = 0)$ , the density ratio can be expressed as:

$$\frac{P_\phi(\mathbf{Z} = \mathbf{z})}{Q_{\mathbf{Z}}(\mathbf{z})} = \frac{\rho_{\mathbf{Z}}(\mathbf{z} | c = 1)}{\rho_{\mathbf{Z}}(\mathbf{z} | c = 0)} = \frac{\rho_{\mathbf{Z}}(c = 1 | \mathbf{z})}{\rho_{\mathbf{Z}}(c = 0 | \mathbf{z})} \approx \frac{D_\eta(\mathbf{z})}{1 - D_\eta(\mathbf{z})}.$$

Therefore, given a trained discriminator  $D_\eta(\mathbf{z})$  and  $M$  i.i.d. samples  $\{\mathbf{z}_m\}_{m=1}^M$  from  $P_\phi(\mathbf{Z})$ , we estimate  $\text{D}_{\text{KL}}(P_\phi(\mathbf{Z}) \| Q_{\mathbf{Z}})$  as:

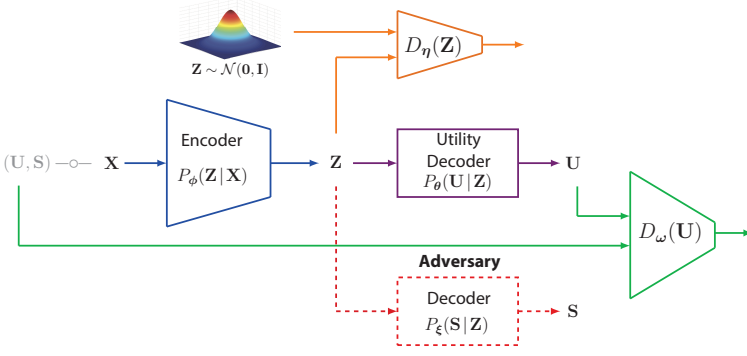
$$\text{D}_{\text{KL}}(P_\phi(\mathbf{Z}) \| Q_{\mathbf{Z}}) \approx \frac{1}{M} \sum_{m=1}^M \log \frac{D_\eta(\mathbf{z}_m)}{1 - D_\eta(\mathbf{z}_m)}. \quad (4)$$

Our model is trained using alternating block coordinate descend across five steps (See Algorithm 1).

**Inferential Adversary:** Given the publicly-known encoder  $\phi$  and  $K$  i.i.d. samples  $\{(\mathbf{s}_k, \mathbf{z}_k)\}_{k=1}^K \subseteq \mathcal{S} \times \mathcal{Z}$ , the adversary trains an inference network  $\xi$  to minimize  $\text{H}_\xi(\mathbf{S}|\mathbf{Z})$ .

### 4 EXPERIMENTS

In this section, we show the impact of the following factors on the amount of leakage: (i) information complexity regularizer weight  $\beta \in (0, 1)$ , (ii) released representation dimension  $d_z$ , (iii) cardinalities of the known utility and unknown sensitive attribute sets, (iv) correlation between the utility and sensitive attributes, and (v) potential bias in a sensitive attribute of adversary's interest. We conduct experiments on the Colored-MNIST and large-scale CelebA



**Figure 2: The training and testing architecture. During training, the data user/owner trains the parameterized networks  $(\phi, \theta, \eta, \omega)$ . During testing, only the encoder-decoder pair  $(\phi, \theta)$  is used. The adversary uses the publicly-known (fixed) encoder  $\phi$  and a collection of the original dataset, and trains an inference network  $\xi$  to infer attribute  $S$  of his interest.**

datasets. The Colored-MNIST<sup>1</sup> is our modified version of MNIST [15], which is a collection of 70,000 colored digits of size  $28 \times 28$ . The digits are randomly colored with red, green, or blue based on two distributions, as explained in the caption of Fig. 4. The CelebA [16] dataset contains 202,599 images of size  $218 \times 178$ . We used TensorFlow 2.4.1 [1] with Integrated Keras API. The method details and network architectures are provided in Appendix. A and Appendix B.

The first and second rows of Fig. 3 and Fig. 4 depict the trade-off among (i) information complexity, (ii) service provider’s accuracy on utility attribute  $U$ , and (iii) adversary’s accuracy on attribute  $S$ . The third row depicts the amount of information revealed about  $S$ , i.e.,  $I(S; Z)$ , for the scenarios considered in the first and second rows, which are estimated using MINE [5]. The fourth row depicts the amount of released information about the utility attribute  $U$ , i.e.,  $I(U; Z)$ , corresponding to the considered scenarios in the first and second rows, also estimated using MINE. We consider different portions of the datasets available for training adversary’s network, denoted by the ‘data ratio’.

The experiments on CelebA consider the scenarios in which the attributes  $U$  and  $S$  are correlated, while  $|\mathcal{U}| = |\mathcal{S}| = 2$ . We provide utility accuracy curves for (i) training set, (ii) validation set, and (iii) test set. As we have argued, there is a direct relationship between information complexity and intrinsic information leakage. Note that, as  $\beta$  increases, the information complexity is reduced, and we observe that this also results in a reduction in the information leakage. We also see that the leakage is further reduced when the dimension of the released representation  $Z$ , i.e.,  $d_z$ , is reduced. This forces the data owner to obtain a more succinct representation of the utility variable, removing any extra information.

In the Colored-MNIST experiments, provided that the model eliminates all the redundant information  $I(X; Z|U)$  and leaves only the information about  $U$ , we expect the adversary’s performance

<sup>1</sup>Several papers have employed Colored-MNIST dataset; However, they are not unique, and researchers synthesized different versions based on their application. The innovative concept behind our version was influenced from the one used in [25].

---

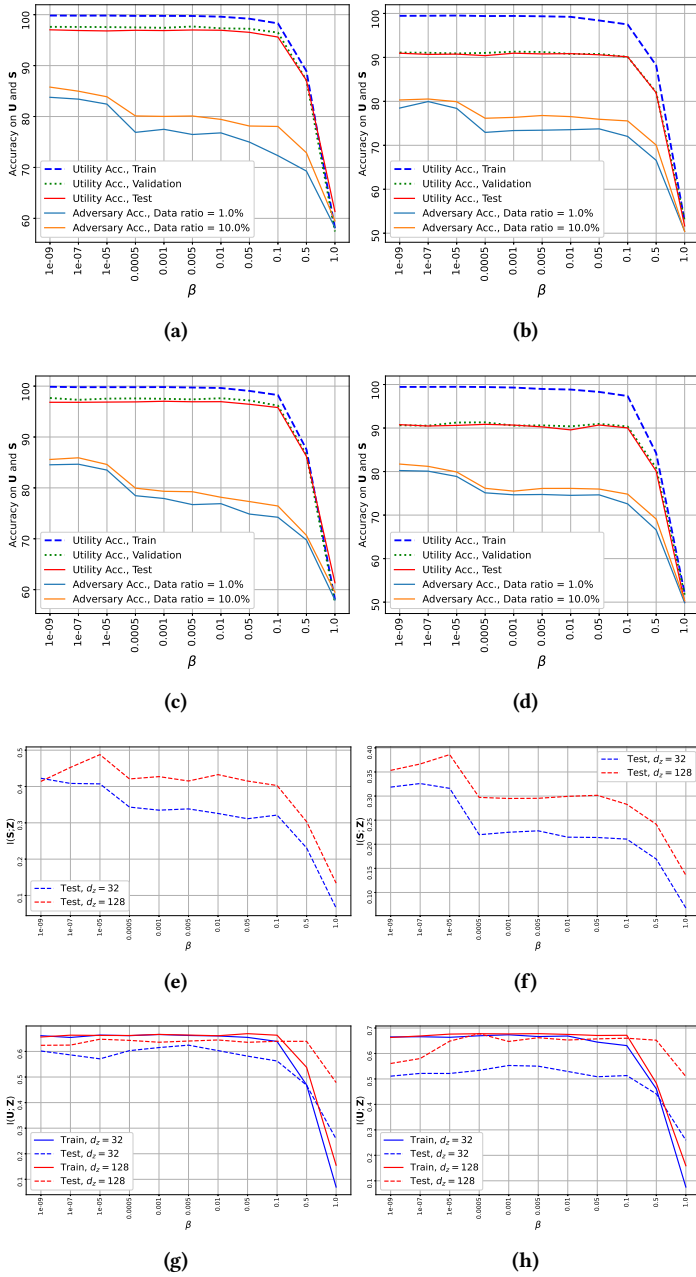
### Algorithm 1 Training Algorithm: Data Owner

---

- 1: **Inputs:** Training Dataset:  $\{(\mathbf{u}_n, \mathbf{x}_n)\}_{n=1}^N$ ;  
Hyper-Parameter:  $\beta$ ;
  - 2:  $\phi, \theta, \eta, \omega \leftarrow$  Initialize Network Parameters
  - 3: **repeat**
  - (1) **Train the Encoder and Utility Decoder  $(\phi, \theta)$**
  - 4: Sample a mini-batch  $\{\mathbf{u}_m, \mathbf{x}_m\}_{m=1}^M \sim P_D(\mathbf{X})P_{U|X}$
  - 5: Compute  $z_m \sim f_\phi(\mathbf{x}_m), \forall m \in [M]$
  - 6: Back-propagate loss:  
$$\mathcal{L}(\phi, \theta) = -\frac{1}{M} \sum_{m=1}^M (\log P_\theta(\mathbf{u}_m | z_m) - \beta D_{\text{KL}}(P_\phi(z_m | \mathbf{x}_m) \| Q_Z(z_m)))$$
  - (2) **Train the Latent Space Discriminator  $\eta$**
  - 7: Sample  $\{\mathbf{x}_m\}_{m=1}^M \sim P_D(\mathbf{X})$
  - 8: Sample  $\{\tilde{z}_m\}_{m=1}^M \sim Q_Z$
  - 9: Compute  $z_m \sim f_\phi(\mathbf{x}_m), \forall m \in [M]$
  - 10: Back-propagate loss:  
$$\mathcal{L}(\eta) = -\frac{\beta}{M} \sum_{m=1}^M (\log D_\eta(z_m) + \log(1 - D_\eta(\tilde{z}_m)))$$
  - (3) **Train the Encoder  $\phi$  Adversarially**
  - 11: Sample  $\{\mathbf{x}_m\}_{m=1}^M \sim P_D(\mathbf{X})$
  - 12: Compute  $z_m \sim f_\phi(\mathbf{x}_m), \forall m \in [M]$
  - 13: Back-propagate loss:  $\mathcal{L}(\phi) = \frac{\beta}{M} \sum_{m=1}^M \log D_\eta(z_m)$
  - (4) **Train the Attribute Class Discriminator  $\omega$**
  - 14: Sample  $\{\mathbf{u}_m\}_{m=1}^M \sim P_U$
  - 15: Sample  $\{\tilde{z}_m\}_{m=1}^M \sim Q_Z$
  - 16: Compute  $\tilde{\mathbf{u}}_m \sim g_\theta(\tilde{z}_m), \forall m \in [M]$
  - 17: Back-propagate loss:  
$$\mathcal{L}(\omega) = -\frac{1}{M} \sum_{m=1}^M (\log D_\omega(\mathbf{u}_m) + \log(1 - D_\omega(\tilde{\mathbf{u}}_m)))$$
  - (5) **Train the Utility Decoder  $\theta$  Adversarially**
  - 18: Sample  $\{\tilde{z}_m\}_{m=1}^M \sim Q_Z$
  - 19: Compute  $\tilde{\mathbf{u}}_m \sim g_\theta(\tilde{z}_m), \forall m \in [M]$
  - 20: Back-propagate loss:  $\mathcal{L}(\omega) = \frac{1}{M} \sum_{m=1}^M \log(1 - D_\omega(\tilde{\mathbf{u}}_m))$
  - 21: **until** Convergence
  - 22: **return**  $\phi, \theta, \eta, \omega$
- 

to be close to ‘random guessing’ since the digit color is independent of its value. We investigate the impact of the cardinality of sets  $|\mathcal{U}|$  and  $|\mathcal{S}|$ , as well as possible biases in the distribution of  $\mathcal{S}$ . The results show that it is possible to reach the same level of accuracy on the utility attribute  $U$ , while reducing the intrinsic leakage by increasing the regularizer weight  $\beta$ , or equivalently, by reducing the information complexity  $I_\phi(X; Z)$ . An interesting possible scenario is to consider correlated attributes  $U$  and  $S$  with different cardinality sets  $\mathcal{U}$  and  $\mathcal{S}$ . For instance, utility task  $U$  is personal

### Variational Leakage



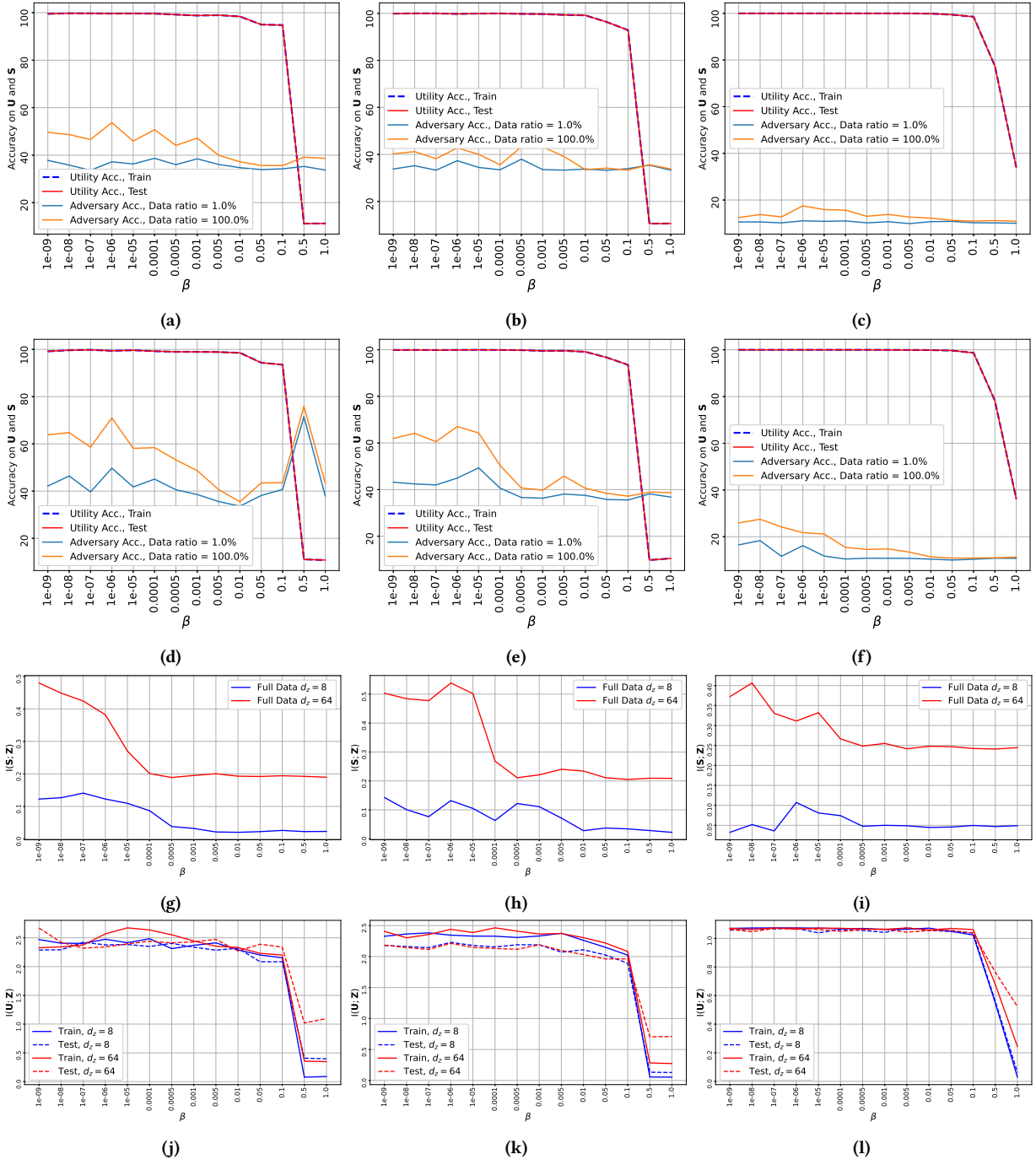
**Figure 3: The results on CelebA dataset, considering isotropic Gaussian prior. (First Row):  $d_z = 32$ ; (Second Row):  $d_z = 128$ ; (Third Row): Estimated information leakage  $I(S; Z)$  using MINE; (Fourth Row): Estimated useful information  $I(U; Z)$  using MINE. (First Column): utility task is gender recognition ( $|\mathcal{U}| = 2$ ), adversary’s interest is heavy makeup ( $|\mathcal{S}| = 2$ ); (Second Column): utility task is emotion (smiling) recognition ( $|\mathcal{U}| = 2$ ), adversary’s interest is mouth slightly open ( $|\mathcal{S}| = 2$ ).**

identification, while the adversary’s interest  $\mathcal{S}$  is gender recognition.

## 5 CONCLUSION

We studied the *variational leakage* to address the amount of potential privacy leakage in a supervised representation learning setup. In contrast to the PF and generative adversarial privacy models, we consider the setup in which the adversary’s interest is not known a priori to the data owner. We study the role of information complexity in information leakage about an attribute of an adversary interest. This was addressed by approximating the information quantities using DNNs and experimentally evaluating the model on large-scale image databases. The proposed notion of *variational leakage* relates the amount of leakage to the minimal sufficient statistics.





**Figure 4: The results on Colored-MNIST dataset, considering isotropic Gaussian prior. (First Row):  $d_z = 8$ ; (Second Row):  $d_z = 64$ ; (Third Row): estimated information leakage  $I(S; Z)$  using MINE; (Fourth Row): estimated useful information  $I(U; Z)$  using MINE. (First Column): utility task is digit recognition ( $|\mathcal{U}| = 10$ ), while the adversary’s goal is the digit color ( $|\mathcal{S}| = 3$ ), setting  $P_S(\text{Red}) = P_S(\text{Green}) = P_S(\text{Blue}) = \frac{1}{3}$ ; (Second Column): utility task is digit recognition ( $|\mathcal{U}| = 10$ ), while the adversary’s goal is the digit color, setting  $P_S(\text{Red}) = \frac{1}{2}$ ,  $P_S(\text{Green}) = \frac{1}{6}$ ,  $P_S(\text{Blue}) = \frac{1}{3}$ ; (Third Column): utility task is digit color recognition ( $|\mathcal{U}| = 3$ ), while the adversary’s interest is the digit number ( $|\mathcal{S}| = 10$ ).**

## REFERENCES

- [1] Martin Abadi et al. 2016. Tensorflow: A system for large-scale machine learning. In *{USENIX} Symp. on Operating Sys. Design and Impl. ({OSDI})*, 265–283.
- [2] Thomas Andre, Marc Antonini, Michel Barlaud, and Robert M Gray. 2006. Entropy-based distortion measure for image coding. In *2006 International Conference on Image Processing*. IEEE, 1157–1160.
- [3] Yuksel Ozan Basciftci, Ye Wang, and Prakash Ishwar. 2016. On privacy-utility tradeoffs for constrained data release mechanisms. In *2016 Information Theory and Applications Workshop (ITA)*. IEEE, 1–6.
- [4] Matthias Bauer and Andriy Mnih. 2019. Resampled priors for variational autoencoders. In *The 22nd International Conference on Artificial Intelligence and Statistics*. 66–75.
- [5] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual information neural estimation. In *International Conference on Machine Learning*. 531–540.
- [6] Nicolo Cesa-Bianchi and Gábor Lugosi. 2006. *Prediction, learning, and games*. Cambridge university press.
- [7] Thomas A Courtade and Richard D Wesel. 2011. Multiterminal source coding with an entropy-based distortion measure. In *2011 IEEE International Symposium on Information Theory Proceedings*. IEEE, 2040–2044.
- [8] Peter Harremoës and Naftali Tishby. 2007. The information bottleneck revisited or how to choose a good distortion measure. In *2007 IEEE International Symposium on Information Theory*. IEEE, 566–570.
- [9] Hsiang Hsu, Shahab Asoodeh, and Flavio P. Calmon. 2019. Obfuscation via Information Density Estimation. In *Proceeding of the International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [10] Chong Huang, Peter Kairouz, Xiao Chen, Lalitha Sankar, and Ram Rajagopal. 2017. Context-aware generative adversarial privacy. *Entropy* 19, 12 (2017), 656.
- [11] Ibrahim Issa, Aaron B Wagner, and Sudeep Kamath. 2019. An operational approach to information leakage. *IEEE Transactions on Information Theory* 66, 3 (2019), 1625–1657.
- [12] D. P Kingma and J. Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [13] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. 2016. Improved variational inference with inverse autoregressive flow. In *Advances in neural information processing systems*. 4743–4751.
- [14] Diederik P Kingma and Max Welling. 2014. Auto-encoding variational bayes. In *International Conference on Learning Representations (ICLR)*.
- [15] Yann LeCun and Corinna Cortes. 2010. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist/>. (2010). <http://yann.lecun.com/exdb/mnist/>.
- [16] Z. Liu, P. Luo, X. Wang, and X. Tang. 2015. Deep Learning Face Attributes in the Wild. In *International Conference on Computer Vision (ICCV)*.
- [17] Ali Makhdoumi, Salman Salamatian, Nadia Fawaz, and Muriel Médard. 2014. From the information bottleneck to the privacy funnel. In *2014 IEEE Information Theory Workshop (ITW 2014)*. IEEE, 501–505.
- [18] XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. 2010. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory* 56, 11 (2010), 5847–5861.
- [19] Flavio P. Calmon, Ali Makhdoumi, and Muriel Médard. 2015. Fundamental limits of perfect privacy. In *2015 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 1796–1800.
- [20] Borzoo Rassouli and Deniz Gündüz. 2019. Optimal utility-privacy trade-off with total variation distance as a privacy measure. *IEEE Transactions on Information Forensics and Security* 15 (2019), 594–603.
- [21] Borzoo Rassouli and Deniz Gündüz. 2021. On perfect privacy. In *to appear in IEEE Journal on Selected Areas in Information Theory (JSAIT)*. IEEE.
- [22] Borzoo Rassouli, Fernando E Rosas, and Deniz Gündüz. 2019. Data Disclosure under Perfect Sample Privacy. *IEEE Trans. on Inform. Forensics and Security* (2019).
- [23] Behrooz Razeghi, Flavio P. Calmon, Deniz Gündüz, and Slava Voloshynovskiy. 2020. On Perfect Obfuscation: Local Information Geometry Analysis. In *2020 IEEE International Workshop on Information Forensics and Security (WIFS)*. 1–6.
- [24] Danilo Jimenez Rezende and Shakir Mohamed. 2015. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning*. 1530–1538.
- [25] Borja Rodríguez-Gálvez, Ragnar Thobaben, and Mikael Skoglund. 2020. A Variational Approach to Privacy and Fairness. *arXiv preprint arXiv:2006.06332* (2020).
- [26] Mihaela Rosca, Balaji Lakshminarayanan, and Shakir Mohamed. 2018. Distribution matching in variational inference. *arXiv preprint arXiv:1802.06847* (2018).
- [27] Sreejith Sreekumar and Deniz Gündüz. 2019. Optimal Privacy-Utility Trade-off under a Rate Constraint. In *2019 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2159–2163.
- [28] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. 2012. Density-ratio matching under the Bregman divergence: A unified framework of density-ratio estimation. *Annals of the Institute of Statistical Mathematics* 64, 5 (2012), 1009–1044.
- [29] Naftali Tishby, Fernando C Pereira, and William Bialek. 2000. The information bottleneck method. In *IEEE Allerton*.
- [30] Jakub Tomczak and Max Welling. 2018. VAE with a VampPrior. In *International Conference on Artificial Intelligence and Statistics*. 1214–1223.
- [31] Ardhendu Tripathy, Ye Wang, and Prakash Ishwar. 2019. Privacy-preserving adversarial networks. In *57th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 495–505.

## Appendices

### A TRAINING DETAILS

All the experiments in the paper have been carried out with the following structure:

#### A.0.1 Pre-Training Phase.

We utilize this phase to warm-up our model before running the main training Algorithm 1 for the Variational Leakage framework within all experiments. In the warm-up phase, we pre-trained encoder ( $f_\phi$ ) and utility-decoder ( $g_\theta$ ) together for the few epochs via **backpropagation (BP)** with the Adam optimizer [12]. We found out the warm-up stage was helpful for faster convergence. Therefore, we initialize the encoder and the utility-decoder weights with the obtained values rather than random or zero initialization. For each experiment, the hyper-parameters of the learning algorithm in this phase were:

Experiment Dataset	Learning Rate	Max Iteration	Batch Size
Colored-MNIST (both version)	0.005	50	1024
CelebA	0.0005	100	512

#### A.0.2 Main Block-wise Training Phase.

In contrast to the most DNNs training algorithms, each iteration only has one forward step through the network’s weights and then update weights via BP approach. Our training strategy is block-wise and consists of multiple blocks in the main algorithm loop. At each block, forward and backward steps have been done through the specific path in our model, and then corresponding parameters update based on the block’s output loss path.

Since it was not possible for us to use the Keras API’s default model training function, we implement Algorithm 1 from scratch in the Tensorflow. It is important to remember that we initialize all parameters to zero except for the  $(\phi, \theta)$  values which acquired in the previous stage. Furthermore, we set the learning rate of the block (1) in the Algorithm 1, five times larger than other blocks. The hyper-parameters of the Algorithm 1 for each experiment shown in the following table:

Experiment Dataset	Learning Rate [blocks (2)-(5)]	Max Iteration	Batch Size
Colored-MNIST (both version)	0.0001	500	2048
CelebA	0.00001	500	1024

### B NETWORK ARCHITECTURES

#### B.0.1 MI Estimation.

For all experiments in this paper, we report estimation of MI between the released representation and sensitive attribute, i.e.,  $I(S; Z)$ , as well as the MI between the released representation and utility attribute, i.e.,  $I(U; Z)$ . To estimate MI, we employed the MINE model [5]. The architecture of the model is depicted in Table 1.

#### B.0.2 Colored-MNIST.

In the Colored-MNIST experiment, we had two versions for data utility and privacy leakage evaluation. In the first version, we set the utility data to the class’ label of the input image and consider

MINE I (U; Z)
INPUT $z \in \mathbb{R}^{d_z}$ CODE; $u \in \mathbb{R}^{ \mathcal{U} }$
$x = \text{CONCATENATE}([z, u])$
FC(100), ELU
FC(100), ELU
FC(100), ELU
FC(1)

Table 1: The architecture of the MINE network.

the color of the input image as sensitive data, and for the second one, we did vice versa. It is worth mentioning that both balanced and unbalanced Colored-MNIST datasets are applied with the same architecture given in Table 2.

ENCODER $f_\phi$
INPUT $x \in \mathbb{R}^{28 \times 28 \times 3}$ COLOR IMAGE
CONV(64,5,2), BN, LEAKYRELU
CONV(128,5,2), BN, LEAKYRELU
FLATTEN
FC( $d_z \times 4$ ), BN, TANH
$\mu$ : FC( $d_z$ ), $\sigma$ : FC( $d_z$ )
$z = \text{SAMPLINGWITHREPARAMETERIZATIONTRICK}[\mu, \sigma]$
UTILITY DECODER $g_\theta$
INPUT $z \in \mathbb{R}^{d_z}$ CODE
FC( $d_z \times 4$ ), BN, LEAKYRELU
FC( $ \mathcal{U} $ ), SOFTMAX
LATENT SPACE DISCRIMINATOR $D_\eta$
INPUT $z \in \mathbb{R}^{d_z}$ CODE
FC(512), BN, LEAKYRELU
FC(256), BN, LEAKYRELU
FC(1), SIGMOID
UTILITY ATTRIBUTE CLASS DISCRIMINATOR $D_\omega$
INPUT $u \in \mathbb{R}^{ \mathcal{U} }$
FC( $ \mathcal{U}  \times 8$ ), BN, LEAKYRELU
FC( $ \mathcal{U}  \times 8$ ), BN, LEAKYRELU
FC(1), SIGMOID

Table 2: The architecture of DNNs used in for the Colored-MNIST experiments.

#### B.0.3 CelebA.

In this experiment, we considered three scenarios for data utility and privacy leakage evaluation, as shown in Table 3. Note that all of the utility and sensitive attributes are binary. The architecture of the networks are presented in Table 4.



Scenario Number	Utility Attribute	Sensitive Attribute
1	Gender	Heavy Makeup
2	Mouth Slightly Open	Smiling
3	Gender	Blond Hair

**Table 3: Scenarios considered for CelebA experiments.**

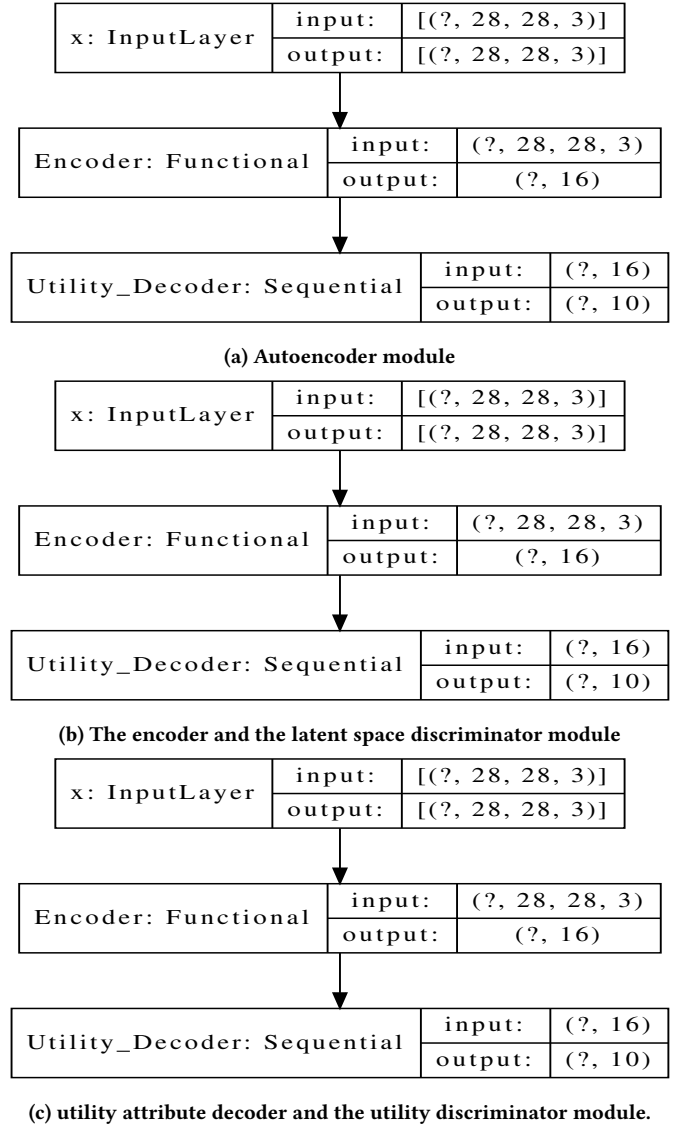
ENCODER $f_\phi$	
INPUT $\mathbf{x} \in \mathbb{R}^{64 \times 64 \times 3}$	COLOR IMAGE
CONV(16,3,2), BN, LEAKYRELU	
CONV(32,3,2), BN, LEAKYRELU	
CONV(64,3,2), BN, LEAKYRELU	
CONV(128,3,2), BN, LEAKYRELU	
CONV(256,3,2), BN, LEAKYRELU	
FLATTEN	
FC( $d_z \times 4$ ), BN, TANH	
$\mu$ : FC( $d_z$ ), $\sigma$ : FC( $d_z$ )	
$z$ =SAMPLINGWITHREPARAMETERIZATIONTRICK[ $\mu, \sigma$ ]	
UTILITY DECODER $g_\theta$	
INPUT $\mathbf{z} \in \mathbb{R}^{d_z}$	CODE
FC( $d_z$ ), BN, LEAKYRELU	
FC( $ \mathcal{U} $ ), SOFTMAX	
LATENT SPACE DISCRIMINATOR $D_\eta$	
INPUT $\mathbf{z} \in \mathbb{R}^{d_z}$	CODE
FC(512), BN, LEAKYRELU	
FC(256), BN, LEAKYRELU	
FC(1), SIGMOID	
UTILITY ATTRIBUTE CLASS DISCRIMINATOR $D_\omega$	
INPUT $\mathbf{u} \in \mathbb{R}^{ \mathcal{U} }$	
FC( $ \mathcal{U}  \times 4$ ), BN, LEAKYRELU	
FC( $ \mathcal{U} $ ), BN, LEAKYRELU	
FC(1), SIGMOID	

**Table 4: The architecture of networks for the CelebA experiments.**

## C IMPLEMENTATION OVERVIEW

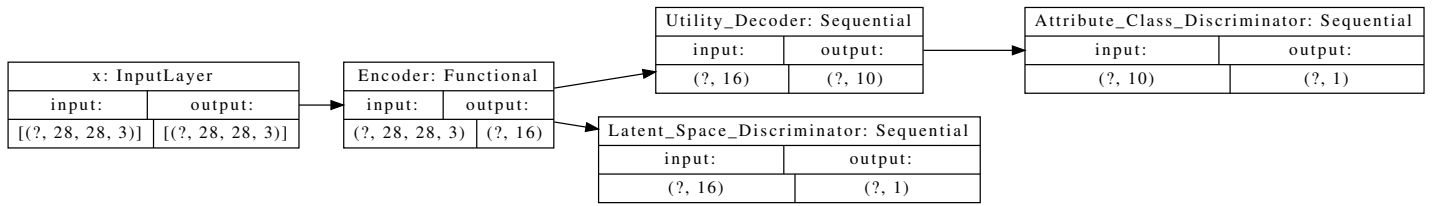
Fig. 6 demonstrates all sub networks of the proposed framework that attached together and their parameters learned in the Algorithm 1.

However, we did not define and save our model in this form because of technical reasons to efficiently implement the training algorithm. As shown in Algorithm 1, the main loop consists of five blocks where only some networks are used in the forward phase, and mostly one of them would update their parameters via BP in each block. Therefore, we shattered the model into three sub-modules in the training stage for simplicity and performance. Fig. 5 shows the corresponding sub-modules of Fig. 6, which are used in our implementation. During training, all of the sub-module (a) parameters, call "autoencoder part" would update with BP after each

**Figure 5: The three sub-modules that make up the main network.**

forward step. For the (b) and (c) sub-modules, only the parameters of one network are updated when the corresponding error function values backpropagate, and we freeze the other networks parameters in the sub-module. For example, block (3) of Algorithm 1 is related to the (b) sub-module, but at the BP step, the latent space discriminator is frozen to prevent its parameters from updating. This procedure is vice versa for module (b) at block (2).

It should be mentioned that during our experiments, we found out that before running our main algorithm, it is beneficial to pre-train the autoencoder sub-module since we need to sample from the latent space, which uses in other parts of the main model during training. We justify this by mentioning that sampling meaningful data rather than random ones from latent variables from



**Figure 6: Complete model structure in the training phase. The above model defines for Colored-MNIST dataset with the number of classes as utility attributes and digit colors as sensitive attributes. Also, the encoder output is set to 16 neurons. (The adversary model is not part of the data owner training algorithm)**

the beginning of learning helps the model to converge better and faster in comparison with starting Algorithm 1 with a randomly initiated autoencoder.