

Adjustable Privacy using Autoencoder-based Learning Structure

Mohammad A. Jamshidi^{*}, Hadi Veisi[†], Mohammad M. Mojahedian^{*},
Mohammad R. Aref^{*}

^{*}*Information Systems and Security Lab. (ISSL), Sharif University of Tech., Tehran, Iran*

[†]*Faculty of New Sciences and Technologies, University of Tehran, Tehran, Iran*

m.a.jamshidi992@gmail.com, h.veisi@ut.ac.ir, m.mojahedian@gmail.com, aref@sharif.ir

Abstract

Inference centers need more data to have a more comprehensive and beneficial learning model, and for this purpose, they need to collect data from data providers. On the other hand, data providers are cautious about delivering their datasets to inference centers in terms of privacy considerations. In this paper, by modifying the structure of the autoencoder, we present a method that manages the utility-privacy trade-off well. To be more precise, the data is first compressed using the encoder, then confidential and non-confidential features are separated and uncorrelated using the classifier. The confidential feature is appropriately combined with noise, and the non-confidential feature is enhanced, and at the end, data with the original data format is produced by the decoder. The proposed architecture also allows data providers to set the level of privacy required for confidential features. The proposed method has been examined for both image and categorical databases, and the results show a significant performance improvement compared to previous methods.

Keywords: Privacy, utility, deep neural networks, autoencoders, collaborative learning.

1. Introduction

The more data a learning system can access, its model can be more comprehensive. But sometimes, the learner or utility provider does not have access

to much data or does not have any data at all and must receive it from different units. Different units may be sensitive about their data and do not like to provide all their information to the utility provider. In this paper, we are looking for a solution so that data providers can distort their data to such an extent that the utility provider can use it and that the information they want remains confidential as much as possible from the utility provider or any other adversary. To be more precise, as shown in Fig. 1, we have m data providers, each of which has a set of data and wants to provide it to the utility provider in such a way that a part of the data or its features can remain private from the utility provider or any other adversary. On the other hand, the quality of the data provided to the utility provider must be good enough to train his learning model well. The utility provider is supposed to provide services using the resulting model to data providers or other users. Therefore the trade-off between the utility and privacy of the datasets becomes important.

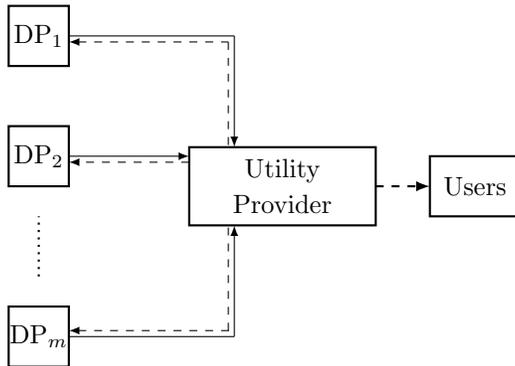


Figure 1: Data providers (DPs) tend to distort their data and send it to the utility provider in a way that balances privacy and utility.

Ensuring privacy is paramount in an era where so much data is available. The utility-privacy trade-off has been studied theoretically [1, 2] and algorithmically. Various algorithms have been proposed to balance the trade-off, and the present paper is in this direction. The primary algorithms proposed for privacy include k -anonymity [3], ℓ -diversity [4], and t -closeness [5], which are suitable only for small datasets. Another group of private algorithms is based

on differential privacy, which is a mathematical tool that guarantees the privacy of the dataset by appropriately adding noise [6, 7, 8, 9, 10]. Among the difficulties of differential privacy for high dimensions, it is time-consuming, requires a lot of noise and as a result distorts the dataset utility, etc. There are other private methods, including homomorphic encryption [11, 12], and secure multi-party computation (SMPC) [13, 14], which are two cryptographic approaches that are related to calculations on encrypted data and face limitations such as computational cost, limited range of computational operations that can be performed, communication cost, etc. In addition to the centralized privacy-preserving methods, there is the distributed federated learning framework, where each data provider trains the model. This method faces limitations such as the negative effect of heterogeneous systems of different data providers, inefficiency in terms of communication, etc. [15, 16, 17].

Another category of private algorithms that have received special attention with the increase in data dimensions are heuristic methods based on machine learning tools [18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]. These methods use machine learning tools such as autoencoders, generative adversarial networks (GAN), variational autoencoders, etc. This paper is also included in this category. In particular, we discuss privacy in data publishing using neural networks. There are similar works in this direction. Huang et al. in [20, 21] uses the generative adversarial network (GAN) and solves a minimax game between the privatizer and the adversary to conceal a specific feature. Li et al. in [22] have used mutual information-based training to learn the feature extractor so that the private features are hidden while maintaining the utility of the remaining information. The output of the proposed method in [22] is a censored feature vector that is not in the original data format and, therefore, not suitable for data publishing, which makes it useless for pre-trained models such as DenseNet. In [24], by combining variational autoencoder and differential privacy, the data dimension is first reduced, then sensitive and non-sensitive data are separated using two classifiers, and after the covariance matrix of sensitive data is learned, it is perturbed to establish differential privacy. Using the differential privacy

tool requires more training data and is therefore time-consuming [32]. A feature extraction method in [26] is proposed for implementation on mobile devices based on Siamese architecture. This method also produces data that does not match the original dataset format. In [27], the obfuscator uses an autoencoder to reduce the dimensions of the image, and then with a GAN-based structure, the encoder output distribution approaches the Gaussian distribution. Also, the classifier that extracts the desired feature gives feedback to the obfuscator. Mandal et al. proposed a private learning algorithm based on uncertainty autoencoder in [28]. Their work has significant results on categorical datasets, so we use it as a benchmark for comparison. Making the dataset confidential may lead to a dataset that does not have a standard format; for this purpose, in [28], the algorithm’s performance with data-type ignorant and data-type aware conditions is studied.

In this paper, we use the autoencoder-based structure, which helps us in two ways by reducing the data dimension. First, the encoder output will be compressed data that are as uncorrelated as possible. Secondly, processing can be done on the reduced-dimension data with a simpler network, significantly reducing the computational load. The paper’s main idea is that the dimensionality-reduced data is appropriately processed using neural networks so that while the utility of the dataset is maintained in terms of some features, the obfuscated dataset is private in terms of other features. The advantages and contributions of the proposed private learning algorithm are:

- The proposed scheme works well on both categorical and image datasets.
- In terms of utility-privacy trade-off, the proposed algorithm outperforms other methods while using a much simpler structure that makes it more suitable for use on weaker processors such as mobile and Internet of Things (IoT) devices.
- The obfuscating model can be learned by a utility provider or any trusted entity and then sent to the data providers for use. This removes the burden of training the model from the data providers. It is also possible for each

data provider to adjust the level of privacy they need without changing the model and just by tuning the noise. Further, the data providers are given a parameter to adjust the data utility amount.

- In our proposed obfuscator, all the features are obscured except for non-confidential ones, which makes it more private against adversaries who do not have a specific goal. To put it more clearly, in our method, the feature or features that the utility provider wants to infer are known, and the data providers obscure the rest of the features.

The rest of the paper is organized as follows. The system model is described in Section 2. The details of the proposed structure are discussed in Section 3. Simulation results are included in Section 4. Finally, the paper is concluded in Section 5.

2. System Model

The dataset \mathcal{D} is a collection containing n samples of the instance space $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, in which \mathcal{X} and \mathcal{Y} are input and output spaces, respectively. We assume that $\mathbf{Y} \in \mathcal{Y}$ is a vector of features that are divided into private and non-private categories in the form of $\mathbf{Y} = (\mathbf{Y}_P, \mathbf{Y}_{NP})$. For example, in this paper, we consider a set of face images as \mathcal{X} , the smiling feature as non-private and other features as private.

The data provider intends to deliver data to the utility provider for collaborative learning, and on the other hand, the confidentiality of some features is important to him. Therefore, by converting \mathbf{X} to \mathbf{X}' , the data provider aims to keep all features private except the non-private one while the utility of the dataset is acceptable in terms of the non-private feature. The function that converts \mathbf{X} to \mathbf{X}' is an obfuscator; the amount of ambiguity it adds to the dataset determines the utility-privacy trade-off. The Obfuscator, adversary, and utility provider are shown in Fig. 2. Here, \mathbf{U}_P and \mathbf{U}_{NP} are private and non-private features, respectively, that are inferred by the adversary and the utility provider from the obfuscated dataset.

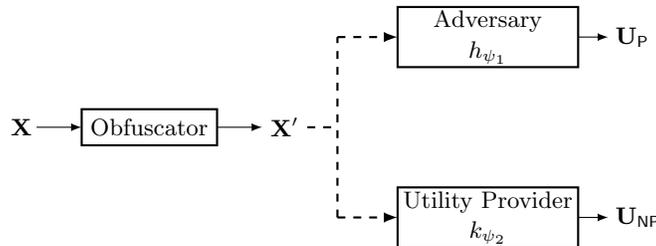


Figure 2: Adversary and utility provider are two deep neural networks with ψ_1 and ψ_2 parameters, respectively, which receive obfuscated data as input and try to extract private and non-private features, respectively.

Threat Model

We assume the adversary is a machine learning model that seeks to extract a private feature from the obfuscated dataset. We also assume that the data provider does not know what feature or features the adversary looks for. Therefore, the data provider tries to obfuscate all features other than the utility provider’s demands. In this paper, we consider two types of adversaries.

- Weak adversary: This type of adversary does not have access to the obfuscation model and only has the obfuscated dataset. Therefore, it is trained using a dataset similar to the original dataset (non-obfuscated) and then extracts the desired feature from the obfuscated dataset.
- Strong adversary: This type of adversary introduced in [33] has access to obfuscator model. Therefore, it creates an obfuscated dataset from a dataset similar to the original dataset. Now it has a dummy obfuscated dataset, and on the other hand, it has the exact value of the feature it wants to infer and trains its network using supervised learning.

Utility Provider

A utility provider is a deep neural network (DNN) classifier trained on the obfuscated dataset to infer one or more features. Other users can utilize this network to infer features from the unobfuscated dataset available to them. It is assumed that the data provider is aware of non-private features.

Obfuscator

The paper’s main idea is to decorrelate private and non-private features and then obfuscate the private features while enhancing the non-private ones. For this purpose, as shown in Fig. 3, the dimension of the data is reduced by using the autoencoder, and decorrelation is done inside it. We employ autoencoder for three reasons:

1. By reducing the dimension, the autoencoder reduces the utility of the data and increases its privacy.
2. Encoder produces uncorrelated output by compressing data and removing redundancy.
3. By reducing the dimension of the data, the intermediate networks that have the task of separating non-private features from compressed data will have a much simpler structure.

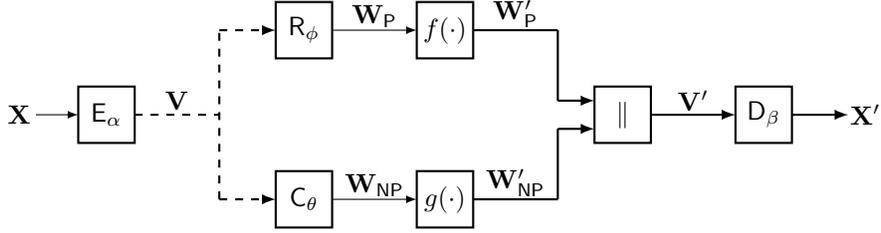


Figure 3: The obfuscator consists of an autoencoder, inside which the private and non-private features are intelligently separated. Then, the private features are combined with the noise appropriately while the non-private features are enhanced. Then the decoder constructs the obfuscated dataset from all the manipulated features.

The autoencoder consists of encoder E_α and decoder D_β , and the coded data is determined by \mathbf{V} . More precisely,

$$\mathbf{V} = E_\alpha(\mathbf{X}). \quad (1)$$

Here, the encoder and decoder are DNNs whose parameter sets are indicated by α and β , respectively.

Now, we make two sets of data from the reduced dimensional and meaningless coded data \mathbf{V} .

1. The DNN C_θ with parameter θ : It is actually a classifier that takes meaningless data \mathbf{V} and produces meaningful features $\mathbf{W}_{\text{NP}} = C_\theta(\mathbf{V})$. The features represented by \mathbf{W}_{NP} are the ones we want to remain useful.
2. The DNN R_ϕ with parameter ϕ : Its output contains information from \mathbf{V} that is appropriately uncorrelated from \mathbf{W}_{NP} . We want to keep this part of the information as confidential as possible, so we have shown it with \mathbf{W}_{P} .

In the next step, we apply the functions $f(\cdot)$ and $g(\cdot)$ to obtain the data \mathbf{W}'_{P} and \mathbf{W}'_{NP} , respectively, and create \mathbf{V}' by putting them together. In the final step, the decoder will convert \mathbf{V}' to \mathbf{X}' . Now we have to design parameters α , β , ϕ , θ and functions $f(\cdot)$ and $g(\cdot)$. For this purpose, we consider the following properties to balance privacy and utility.

- (P1) To preserve the utility of the original dataset, we like \mathbf{X} and \mathbf{X}' to be as similar as possible. For this purpose, we will minimize the following loss function.

$$L_{ae}(\alpha, \beta, \phi, \theta) = \mathbb{E}[\ell(\mathbf{X}, \mathbf{X}')], \quad (2)$$

where the expectation is over the distribution of \mathbf{X} . Additionally, \mathbf{X}' without considering the functions $f(\cdot)$ and $g(\cdot)$ is equal to

$$\mathbf{X}' = D_\beta(C_\theta(\mathbf{V}) \parallel R_\phi(\mathbf{V})). \quad (3)$$

- (P2) To make \mathbf{W}_{NP} useful for utility provider, we consider minimizing the following.

$$L_{\text{C}}(\alpha, \theta) = \mathbb{E}[\ell(\mathbf{W}_{\text{NP}}, \mathbf{Y}_{\text{NP}})], \quad (4)$$

where the expectation is over the joint distribution of \mathbf{X} and \mathbf{Y} .

- (P3) Suppose that adversary h_{ψ_1} is a DNN with parameter ψ_1 that is well trained on a dataset similar to the original dataset/dummy obfuscated

dataset (weak/strong) to infer private features \mathbf{Y}_P . As a result, the adversary wants to minimize the following loss function, while we want it to be maximized.

$$L_{h_{\psi_1}}(f, g) = \mathbb{E} \left[\ell \left(h_{\psi_1} (D_{\beta} (\mathbf{W}'_{NP} \parallel \mathbf{W}'_P)), \mathbf{Y}_P \right) \right]. \quad (5)$$

(P4) Suppose that utility provider k_{ψ_2} is a DNN with parameter ψ_2 that is well trained on the obfuscated dataset to infer certain non-private feature \mathbf{Y}_{NP} . As a result, the utility provider wants to minimize the following loss function.

$$L_{k_{\psi_2}}(f, g) = \mathbb{E} \left[\ell \left(k_{\psi_2} (D_{\beta} (\mathbf{W}'_{NP} \parallel \mathbf{W}'_P)), \mathbf{Y}_{NP} \right) \right]. \quad (6)$$

The design of the proposed network is done in two stages. In the first step, assuming the absence of functions $f(\cdot)$ and $g(\cdot)$ and using (P1) and (P2), the network parameters are calculated as follows.

$$\alpha^*, \beta^*, \phi^*, \theta^* = \arg \min_{\alpha, \beta, \phi, \theta} L_{ae}(\alpha, \beta, \phi, \theta) + L_C(\alpha, \theta). \quad (7)$$

Then $f(\cdot)$ and $g(\cdot)$ are optimized according to the following optimization problem.

$$(f^*, g^*) = \arg \max_{f, g} L_{h_{\psi_1}}(f, g) - L_{k_{\psi_2}}(f, g). \quad (8)$$

3. Methodology

As mentioned in the previous section, (7) is used to train DNNs E_{α} , D_{β} , C_{θ} and R_{ϕ} . In addition, instead of solving the optimization problem in (8), we design the functions $f(\cdot)$ and $g(\cdot)$ intelligently. Further, the proposed scheme is examined on two image and categorical datasets to measure its performance in various applications.

3.1. Image Dataset

CelebA is chosen as the dataset, a collection of large-scale facial features [34]. Facial images with size $64 \times 64 \times 3$ are input, and their labels are different features, including age, gender, etc. To compare the performance of the proposed

scheme with the previous works, the desired features of open mouth, smiling, and high cheekbone have been selected, which are denoted by “CelebA-G-M”, “CelebA-G-S”, “CelebA-G-C”, respectively. Here “G” stands for gender, which is the feature the adversary is looking for. It should be noted that we assume that the data providers do not know which feature is confidential.

DNNs Structures

Inspired by VGG-16 network [35], the building blocks of the autoencoder, i.e. E_α and D_β , both consist of 4 2D-convolutional layers, 3 batch normalization layers, and 1 fully-connected layer. E_α takes images with size $64 \times 64 \times 3$ and produces 1024 features as output. In other words, the \mathbf{V} size equals 1024. The decoder then converts a vector of manipulated features with size 1024 into a $64 \times 64 \times 3$ image.

Two fully-connected networks C_θ and R_ϕ are used in the middle of the autoencoder. C_θ is a 4-layer fully-connected network that converts the input of size 1024 into 2 outputs expressing the desired feature of the utility provider. R_ϕ is a 3-layer fully-connected network that converts the input of size 1024 into 1022 outputs representing the rest of the features that are well uncorrelated with the desirable feature of the utility provider. The DNNs structure details are given in Table 1.

Loss Functions

We use the mean square error measure as the loss function of the autoencoder and the negative log-likelihood as the loss function of the classifier. More precisely, we have

$$L_{ae} = \frac{1}{N} \sum_{j=1}^N (X_j - X'_j)^2 \quad (9)$$

$$L_C = - \sum_{j=1}^M Y_{NP,j} \log(W_{NP,j}) + (1 - Y_{NP,j}) \log(1 - W_{NP,j}), \quad (10)$$

where N and M are the dimensions of the vectors \mathbf{X} and \mathbf{Y}_{NP} , respectively. Moreover, the negative log-likelihood loss function has been used to train the utility provider and adversary.

Table 1: DNNs architecture details for image datasets.

Component	Num	Layer	Output Size	Specs	Activation Function
Input Data		Image Samples	$3 \times 64 \times 64$		
Encoder	1	Conv2D	$64 \times 32 \times 32$	kernel=4, stride=2, padding=1	LeakyReLU
	2	Conv2D	$64 \times 16 \times 16$	kernel=4, stride=2, padding=1	
	3	BatchNorm2D		eps=1e-5, momentum=0.1	LeakyReLU
	4	Conv2D	$64 \times 8 \times 8$	kernel=4, stride=2, padding=1	
	5	BatchNorm2D		eps=1e-5, momentum=0.1	LeakyReLU
	6	Conv2D	$128 \times 4 \times 4$	kernel=4, stride=2, padding=1	
	7	BatchNorm2D	$128 \times 4 \times 4 \rightarrow 2048$ (shaped)	eps=1e-5, momentum=0.1	LeakyReLU
	8	Linear	1024		LeakyReLU
Classifier	1	Linear	1024	Dropout(p=0.5)	LeakyReLU
	2	Linear	256	Dropout(p=0.5)	LeakyReLU
	3	Linear	64		LeakyReLU
	4	Linear	2		LogSoftMax
R	1	Linear	1024	Dropout(p=0.5)	LeakyReLU
	2	Linear	1024	Dropout(p=0.5)	LeakyReLU
	3	Linear	1022		LeakyReLU
Decoder	1	Linear	$2048 \rightarrow 128 \times 4 \times 4$ (shaped)		LeakyReLU
	2	ConvTranspose2D	$64 \times 8 \times 8$	kernel=4, stride=2, padding=1	
	3	BatchNorm2D		eps=1e-5, momentum=0.1	ReLU
	4	ConvTranspose2D	$64 \times 16 \times 16$	kernel=4, stride=2, padding=1	
	5	BatchNorm2D		eps=1e-5, momentum=0.1	ReLU
	6	ConvTranspose2D	$64 \times 32 \times 32$	kernel=4, stride=2, padding=1	
	7	BatchNorm2D		eps=1e-5, momentum=0.1	ReLU
	8	ConvTranspose2D	$3 \times 64 \times 64$	kernel=4, stride=2, padding=1	Sigmoid
Adversary and Utility Provider	1	Conv2D	$64 \times 32 \times 32$	kernel=4, stride=2, padding=1	LeakyReLU
	2	Conv2D	$64 \times 16 \times 16$	kernel=4, stride=2, padding=1	
Utility Provider	3	BatchNorm2D		eps=1e-5, momentum=0.1	LeakyReLU
	4	Conv2D	$64 \times 8 \times 8$	kernel=4, stride=2, padding=1	
	5	BatchNorm2D		eps=1e-5, momentum=0.1	LeakyReLU
	6	Conv2D	$128 \times 4 \times 4$	kernel=4, stride=2, padding=1	
	7	BatchNorm2D	$128 \times 4 \times 4 \rightarrow 2048$ (shaped)	eps=1e-5, momentum=0.1	LeakyReLU
	8	Linear	1024		LeakyReLU
	9	Linear	1024	Dropout(p=0.5)	LeakyReLU
	10	Linear	256	Dropout(p=0.5)	LeakyReLU
	11	Linear	64		LeakyReLU
	12	Linear	2		LogSoftMax

Training

The training phase is explained in Algorithm 1. A batch with the size of n_b samples is taken from the dataset, and the parameters of all the DNNs E_α , R_ϕ , C_θ , and D_β are updated using the loss function defined in (9). Then, keeping the networks R_ϕ and D_β unchanged, we update the E_α and C_θ parameters using the loss function in (10) and repeat this for n_e epochs. It should be noted that the training process is carried out without considering functions $f(\cdot)$ and $g(\cdot)$.

Algorithm 1 Training phase of the framework

Input: Training dataset samples \mathbf{X} parameter: learning rate a parameter: training steps n_e and n_b **Output:** Obfuscator Model

Initialization.

- 1: **for** n_e epochs **do**
 - 2: Randomly select the mini-batch
 from the training data set.
 - 3: **for** $i = 0$ to n_b iterations **do**
 - 4: Update the decoder parameters:
 $\beta_{i+1} = \beta_i - a\nabla_{\beta}L_{ae}(\beta_i, \mathbf{X})$
 - 5: Update \mathbf{R}_{ϕ} parameters:
 $\phi_{i+1} = \phi_i - a\nabla_{\phi}L_{ae}(\phi_i, \mathbf{X})$
 - 6: Update the classifier parameters:
 $\theta_{i+1} = \theta_i - a\nabla_{\theta}L_{ae}(\theta_i, \mathbf{X}) - a\nabla_{\theta}L_C(\theta_i, \mathbf{X}, \mathbf{Y}_{\text{NP}})$
 - 7: Update the encoder parameters:
 $\alpha_{i+1} = \alpha_i - a\nabla_{\alpha}L_{ae}(\alpha_i, \mathbf{X}) - a\nabla_{\alpha}L_C(\alpha_i, \mathbf{X}, \mathbf{Y}_{\text{NP}})$
 - 8: **end for**
 - 9: **end for**
 - 10: **return**
-

As can be seen in Algorithm 1, in each iteration, the weights are updated based on both the autoencoder and classifier loss functions, which makes not only the output of the autoencoder remain as similar as possible to the input data (maintain its utility), but also the output of the classifier (desirable feature) is well uncorrelated from the rest of the features.

Selection of Functions $f(\cdot)$ and $g(\cdot)$

It is difficult to solve the optimization problem (8) to obtain optimal functions. Instead of optimally choosing these functions, here we introduce a natural and intelligent choice for them.

The function $g(\cdot)$ should be chosen to contributes to the dataset's utility

for inferring non-private data. Therefore, a natural choice is to modify the classifier’s output according to the labels in the original dataset. Assume that we consider smiling a feature required by the utility provider; since the LogSoftMax function is selected as the last layer of the classifier, the exponential value of its output represents the probability of smiling or not smiling, which can be between 0 and 1. Therefore, the function $g(\cdot)$ changes the output produced by the classifier to $\log(0)$ or $\log(1)$ depending on whether the image in the original dataset is smiling or not smiling. It should be noted that since $\log(0)$ cannot be used, we instead choose a sufficiently large negative value and call it λ . The reason for using LogSoftMax in the last layer of the classifier is that the output values of the classifier and the network R_ϕ are in the same numerical range, and both are effective in the obfuscated output image. Also, the existence of the λ parameter will effectively control the utility of the obfuscated dataset. Using SoftMax instead of LogSoftMax makes the obfuscator not converge or converge after a large number of epochs. In addition, using LogSoftMax has better numerical properties and makes training more stable [36].

The input to the function $f(\cdot)$ are the features we want to remain private and are well uncorrelated from the non-private features. Therefore, adding Gaussian noise is a suitable choice for the function $f(\cdot)$. The higher the amount of noise, the higher the level of privacy, and of course, the usefulness of the dataset is affected from the point of view of all features except non-confidential features. The effect of changing noise variance and λ has been evaluated with simulations.

Utility Provider and Adversary Structures

The structure of both of them is similar to the structure of the encoder plus classifier. The only difference will be in the output number of the last layer, which is proportional to the number of private features desired by the adversary or the number of features desired by the utility provider for inference.

Measurement of Utility and Privacy

We assume that the adversary intends to infer the binary feature of gender (male/female) from the obfuscated dataset. The adversary is trained to infer this feature and then applied to the obfuscated dataset. Its accuracy in correctly

diagnosing males or females is considered a measure of confidentiality.

The features of open mouth, smiling, and high cheekbone are considered as desirable features of the utility provider. We train the utility provider on the obfuscated dataset and then test the trained network on the original dataset. The accuracy of the utility provider in the desired feature recognition is considered a measure of utility.

As stated in the introduction, a significant drawback of the methods proposed for data privacy is that each data provider must separately train a network for this purpose. That puts a lot of burden on data providers. In our proposed method, the utility provider or a trusted authority can design and train an obfuscator according to the feature he wants to remain useful and share it among all data providers. Each data provider can then simply adjust the model’s utility-privacy trade-off by adjusting the noise level and the value of λ based on their sensitivity to their data.

3.2. Categorical Dataset

UCI Adult dataset includes census data of 48842 people and 14 categorical and continuous features of them [37]. Putting [28] as a benchmark, we convert the 14 features into a vector of length 106 by removing missing value data, normalizing the variables, and using one-hot encoding for categorical features. The binary features of income (more or less than 50K per year) and gender are considered the demands of the utility provider and the adversary, respectively. In addition, the binary features of gender and income, which are adversary and utility provider preferences, are removed from the dataset. The resulting dataset with a feature vector of length 102 is used for training networks.

DNNs Structures

All DNNs comprise 3 fully-connected layers, except R_ϕ , which comprises 2 fully-connected layers. The details of the DNNs are given in Table 2.

Utility Provider and Adversary Structures

Utility provider and adversary are 4-layer fully-connected networks that take an input of size 102 and convert it into 2 outputs (gender (male/female) for the

Table 2: DNNs architecture details for categorical datasets.

Component	Num	Layer	Output Size	Specs	Activation Function
Input Data		Samples	102		
Encoder	1	Linear	128		ReLU
	2	Linear	128		ReLU
	4	Linear	64		ReLU
Classifier	1	Linear	32		ReLU
	2	Linear	8		ReLU
	4	Linear	2		LogSoftMax
R_ϕ	1	Linear	64		ReLU
	2	Linear	62		ReLU
Decoder	1	Linear	128		ReLU
	2	Linear	128		ReLU
	4	Linear	102		Sigmoid
Adversary and	1	Linear	256	Dropout (p=0.2)	ReLU
	2	Linear	256	Dropout (p=0.3)	ReLU
Utility Provider	3	Linear	128	Dropout (p=0.4)	ReLU
	4	Linear	2		LogSoftMax

adversary and income ($\geq 50K$) for the utility provider). Details are in Table 2.

4. Experiments

The proposed scheme is implemented in this section, and its performance is compared with other methods.

4.1. Experiments Settings

To compare and evaluate the performance, the settings (Input dataset, the size ratio of training, validation and test sets, etc.) used for image and categorical datasets are similar to [24] and [28], respectively. For the image dataset, the desirable features of the utility provider are open mouth, smiling, and high cheekbone, and the private feature is gender. For the categorical dataset, the utility provider is interested in inferring income (is it more or less than 50K per year?), and the adversary is interested in understanding gender.

Dataset

We choose CelebA as our image dataset. CelebA is a large-scale facial feature dataset with 202,599 face images from 10,177 identities and 40 binary feature labels (such as gender, age, smile, etc.) in each image. 162,752 samples were used as a training set, and the rest were used as a test and validation set. The obfuscator is trained on the training set and is validated in each epoch using the validation set. Finally, the performance of the network is evaluated on the test set. To learn the adversary, again, the CelebA dataset is used considering the gender feature as the data to be inferred while the utility provider training is done on the obfuscated dataset, and then the utility provider performance is tested on the non-obfuscated dataset.

UCI Adult has been selected as a categorical dataset containing the census data of 48842 people. As explained in the previous section, with pre-processing corresponding to each record, we have a feature vector with length 106. The desired features of the utility provider and the adversary are income and gender. 80% of the dataset is used for training and the rest for testing.

Implementation Details

Our experiment is implemented by PyTorch [38] on Google Colab GPUs. The size of the images is $64 \times 64 \times 3$, and the mini-batch technique is used with a batch size of $n_b = 64$. We use Adam optimizer [39] to train all our networks and set the learning rate of all optimizers to 0.001. Obfuscator, adversary, and utility provider training have been conducted with similar setups. The initial value of the weights for all networks is randomly generated with a Gaussian random variable with a variance of 0.02 and an average of 0 for the convolutional layers and an average of 1 for the batch normalization layers. In addition, when the validation loss starts to increase, we stop the learning algorithm and use dropout with probability 0.5 in some layers to protect the training from overfitting.

For the categorical dataset, we have 48842 records, and corresponding to each record after one-hot encoding, there is a feature vector with length 106. We use 39074 records for training and the rest for testing. Also, we remove the features of income and gender, which are the demands of the utility provider

and the adversary, respectively, from the dataset and use the resulting dataset for obfuscator training. Weights are initialized by Gaussian distribution with mean 0 and variance 0.02.

The implementation codes of the proposed scheme are available in <https://github.com/bozorgmehr77/adjustable-privacy>.

4.2. Experiments Results

Let’s assume smiling is the desire of the utility provider, and gender is the adversary’s desire. By training the obfuscator on the CelebA dataset for several epochs, the validation and training errors of the autoencoder are plotted in Fig. 4. As can be seen from Fig. 4, with the increase in the number of epochs, the training and validation errors related to the autoencoder are both reduced and as a result, the obfuscated image \mathbf{X}' is getting closer to the original image \mathbf{X} . Therefore, according to Fig. 4, the PSNR between \mathbf{X} and \mathbf{X}' increases. Reducing the autoencoder error makes the obfuscated image remain as useful as possible despite the compression, while the classifier helps to preserve the desired feature (smiling here) well in the final dataset. Therefore, controlling the errors of the autoencoder preserves the desired feature and makes it well uncorrelated from other features. Later, in subsection 4.2.3, we will discuss the validation of the obfuscator training process.

After the training phase, we activate the functions $f(\cdot)$ and $g(\cdot)$. The function $g(\cdot)$ converts the output of the classifier to $\log(0) \approx \lambda$ and $\log(1) = 0$ depending on the smiling feature in the original dataset, and the function $f(\cdot)$ is an additive white Gaussian noise (AWGN) with $\mu = 0$ and variance proportional to ν . Here, ν is the average of the R output nodes for each record.

Suppose the model corresponding to the epoch number 183 is selected, the noise variance is set to 60ν , and $\lambda = -3000$. Then we convert the whole dataset \mathcal{D} to dataset \mathcal{D}' by passing through obfuscator. Now the adversary infers the gender from the obfuscated dataset, and its accuracy in recovering the gender, which is our privacy criterion, in weak and strong cases are %53.55 and %63.24, respectively, while the accuracy of gender inference from the main dataset is

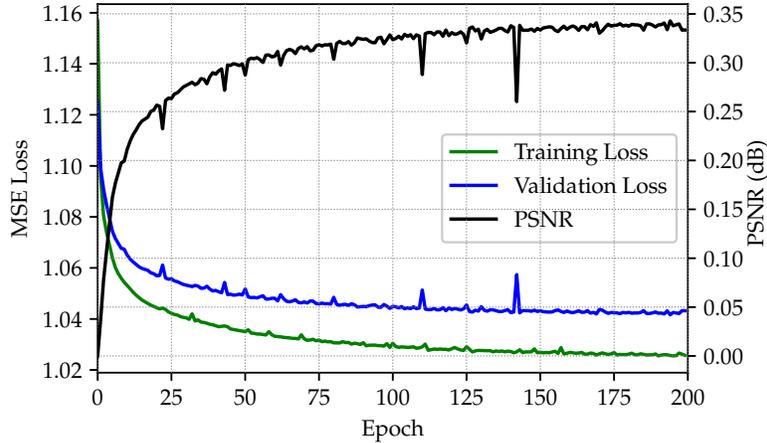


Figure 4: Training and validation losses related to autoencoder.

equal to %97.30. Therefore, the proposed algorithm has strengthened privacy in weak and strong adversarial cases by %43.75 and %34.06, respectively.

In terms of utility, assuming that the utility provider is trained with the obfuscated dataset, its inference accuracy for the smiling feature from the original dataset is equal to %85.43. But if it is trained on the main dataset, this value will equal %91.92. This slight decrease of %6.49 in usefulness shows the strength of the proposed method.

4.2.1. Effect of Noise Variance on Utility-Privacy Trade-off

The noise variance is a parameter the data provider can use to adjust privacy. Considering the CelebA-G-S case, we train the obfuscator for 200 epochs. We set λ equal to -3000 and change the noise variance in the interval $[0, 200\nu]$. Fig. 5 shows the changes in privacy and utility with noise variance. This figure gives the results for weak and strong adversaries and the case where the function $g(\cdot)$ is not applied. The following points are evident from Fig. 5.

- By increasing the noise, the utility increases slightly first because adding a small amount of noise makes the dataset more diverse. As a result, the

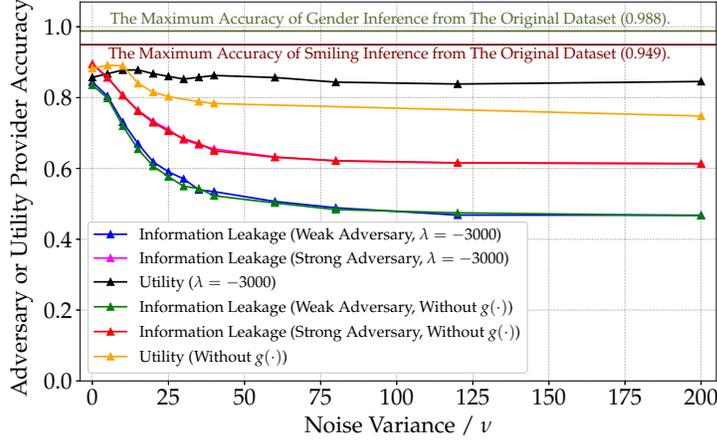


Figure 5: Utility and privacy in terms of the increase of noise variance for the CelebA-G-S and considering the cases where the function $g(\cdot)$ is active or inactive.

utility provider is better trained on it. With the further increase of the noise and in the presence of the function $g(\cdot)$, the utility decreases slightly, while without applying $g(\cdot)$, the decrease is more significant. Therefore, the presence of the function $g(\cdot)$ preserves the dataset’s utility regarding the smiling feature.

- Comparing the utility and privacy curves for the case where the $g(\cdot)$ function is disabled, and we don’t have noise (the starting point of the curves) with the maximum utility and privacy values shows that the privacy and utility have dropped by 0.066 and 0.093, respectively. The lower amount of utility loss is because the presence of the classifier makes the feature of smiling (non-confidential) better preserved.
- For both strong and weak adversaries, privacy is strengthened by increasing the noise variance until it reaches the saturation limit. The saturation value for a weak adversary is about 50%, which is the same as a random guess, and for a strong adversary, it is 60%.
- The comparison of privacy for two cases of active and inactive $g(\cdot)$ function

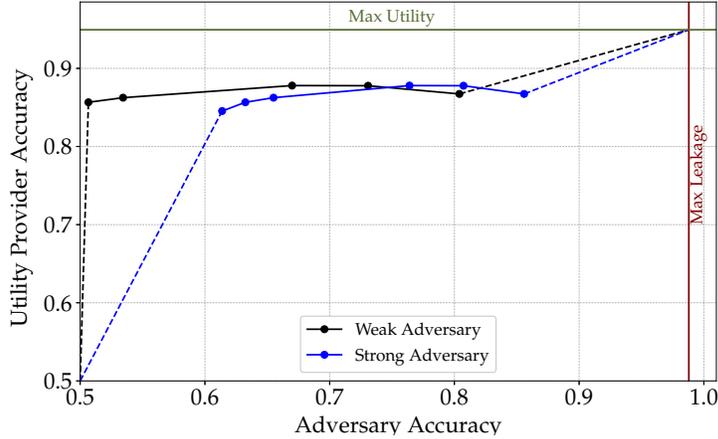


Figure 6: Utility-privacy trade-off for case CelebA-G-S with strong and weak adversaries.

shows that $g(\cdot)$ does not have much effect on privacy.

Moreover, the utility-privacy trade-off curves for strong and weak adversaries are depicted in Fig. 6. The horizontal axis represents information leakage as a measure of privacy, and the vertical axis represents utility. The ideal point is $(0.5, 1)$; the closer the curve is to this point; the better the algorithm is from the point of view of the utility-privacy trade-off. It can be seen that the curve has shifted to the right for the strong adversary.

4.2.2. Effect of λ on Utility-Privacy Trade-off

λ is a parameter that affects the utility of the dataset. Considering the case CelebA-G-S, the utility-privacy trade-off for $\lambda \in \{-5000, -3000, -1000\}$, the case where the function $g(\cdot)$ is not applied, and the noise variances of $\{0, 5, 10, 15, 20\}\nu$ are plotted in Fig. 7. The adversarial type is weak, and points related to not adding noise are marked with a cross. Further, the maximum utility and information leakage shown in the figure corresponds to the accuracy of inferring smiling and gender features by the utility provider and the adversary, respectively. The following points can be deduced from Fig. 7.

- By increasing the value of $|\lambda|$, the effect of the output value of the classifier

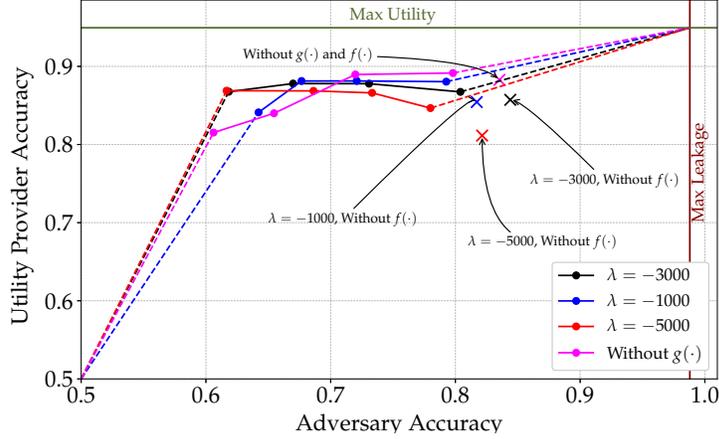


Figure 7: Comparison of the utility-privacy trade-off for case CelebA-G-S and different values of λ .

increases. As a result, the power of inferring the desired feature of the utility provider increases in high noise. On the other hand, the diversity of the dataset is reduced, and the effect of other features fades, which leads to a small increase in privacy. By decreasing the value of $|\lambda|$, the inference power of the utility provider decreases and, consequently, the utility. Therefore, to maintain the dataset's diversity and the accuracy of the utility provider's inference, the value of λ should be well adjusted.

- Increasing the value of $|\lambda|$ reduces the diversity of the dataset, and in low noise variances, it leads to a slight decrease in utility.
- With the addition of noise, the utility generally decreases. However, for $\lambda = -3000$, the utility is almost constant for a relatively wide range of noise variances, which makes it a suitable candidate.

Considering the case CelebA-G-S, the effect of different amounts of noise and λ is given visually in Fig. 8.



Figure 8: Obfuscated images using different values of noise variance and λ .

4.2.3. Effect of Epoch Number on Utility-Privacy Trade-off

As seen in Fig. 4, the autoencoder error continuously decreases with increasing epochs. However, to validate the training process, the overall performance of the obfuscator should be evaluated, and the epoch number should be determined based on this overall performance to have a satisfactory model. For this purpose, ignoring the functions $f(\cdot)$ and $g(\cdot)$, the inference accuracy of smiling and gender features by the well-trained networks is plotted in Fig. 9. As you can see, with the increase in the number of epochs, both curves climb to a high value, which shows that both the utility provider has a good performance and the dataset remains diverse in terms of features other than non-private ones. Therefore, by increasing the number of epochs, the overall performance of the obfuscator will be better.

Consider the case CelebA-G-S and choose the models corresponding to the epoch numbers 13 and 183. According to Fig. 9, in epoch 13, the accuracy of inferring smiling is high, while the accuracy of gender inference is low. Therefore, the feature of smiling is well preserved while other features are somewhat degraded. In epoch 183, where the autoencoder is well trained (based on Fig. 4), the accuracy of both smiling and gender inference is high, indicating that both

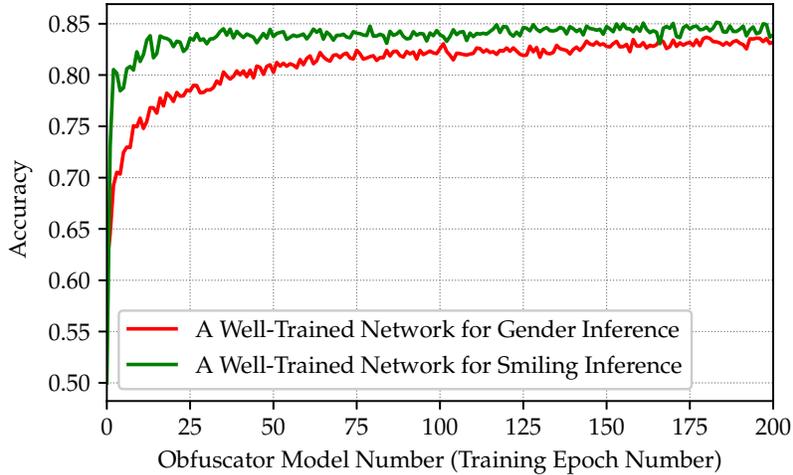


Figure 9: Accuracy of well-trained networks for extracting gender and smiling features by the number of epochs.

the classifier is well trained, and diversity of the obfuscated dataset is preserved.

Suppose λ is -3000 and -2000 for epochs 13 and 183, respectively, and we have a weak adversary; the utility-privacy trade-off is plotted in Fig. 10 for epochs 13 and 183. As you can see in the figure, more training has led to an increase in utility and a slight decrease in privacy. The reason for the reduction of privacy is that by reducing the autoencoder error, the output image preserves as many of the features of the input image as possible. As a result, the private feature is also present in the obfuscated dataset with higher quality. Regarding the utility, with less training of the autoencoder, the resulting dataset has less diversity, which will lead to less utility. Therefore, the model corresponding to a higher epoch performs better regarding the utility-privacy trade-off, and the obfuscator is not overfitted until at least epoch 200.

It should be noted that the values of $\lambda = -2000, -3000$ have been chosen so that the models related to epochs 13 and 183 perform well, as shown in Figs. 11 and 12. To depict these two figures, assuming that the function $f(\cdot)$ is inactive, the accuracy of two well-trained networks for inferring smiling and gender features has been evaluated in terms of the λ parameter on the obfuscated

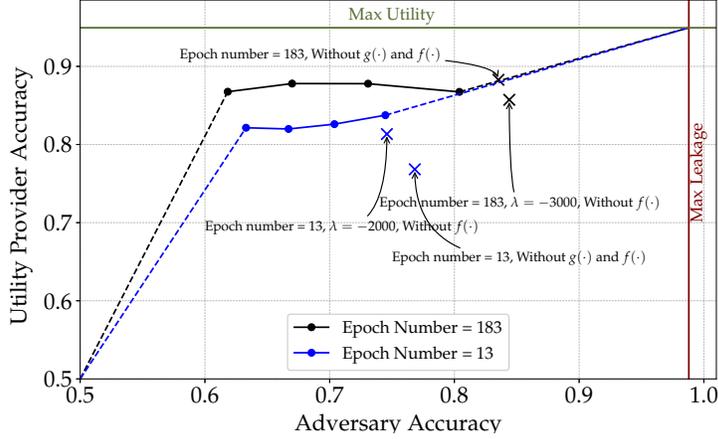


Figure 10: Comparison of utility-privacy trade-off for obfuscator corresponding to models No. 13 and 183.

test set. The optimal λ point, as shown in the figure, is where both networks have reached high accuracy, and the accuracy of the gender-inferring network has not dropped too much to have diversity in the constructed data set. Here, gender represents all the features except non-private ones.

4.2.4. Decorrelation

In this section, we show that the proposed obfuscator makes the utility provider’s desired feature well uncorrelated from the rest of the features. For this purpose, consider the case CelebA-G-S and set λ to -3000 . We also assume that the function $f(\cdot)$ is not applied. For each image, we randomly change the output of the classifier to 0 or -3000 (smiling or not) and give it to the decoder along with the R_ϕ output vector without applying noise to generate a new image. The output image is labeled smiling or not smiling by the utility provider. The histogram of inferring the smiling feature by the utility provider is given in Fig. 13. As you can see, the images are well separated regarding smiling or not smiling. This shows that the classifier’s output has almost complete control over the smiling feature, and the R_ϕ output has almost no effect on the inference of

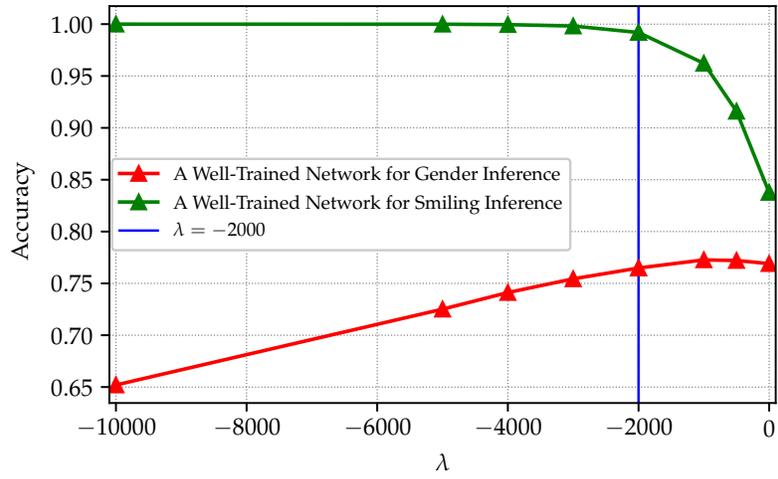


Figure 11: Choosing a suitable λ for obfuscator corresponding to the model number 13 and case CelebA-G-S.

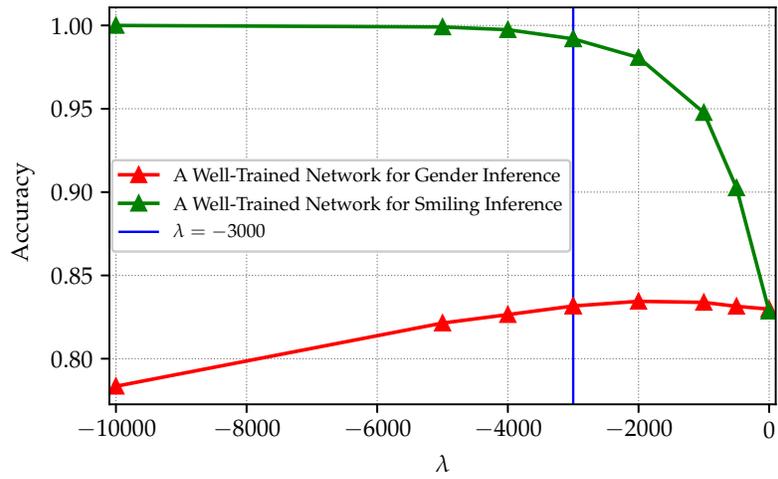


Figure 12: Choosing a suitable λ for obfuscator corresponding to the model number 183 and case CelebA-G-S.

this feature.

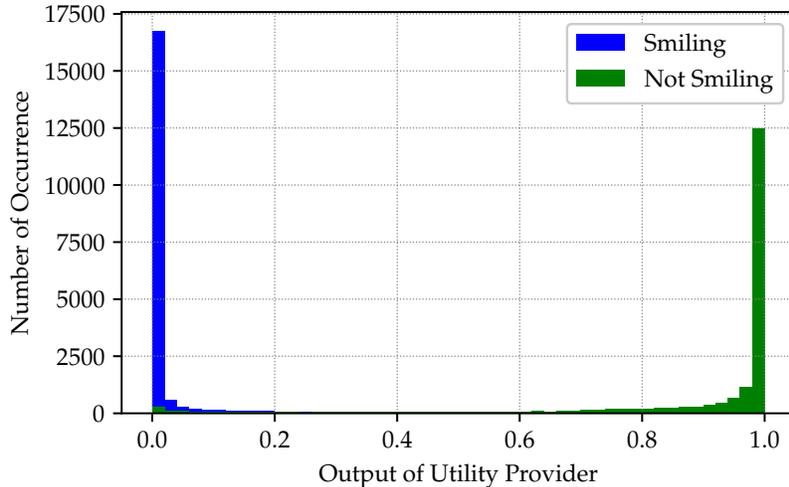


Figure 13: Histogram of the utility provider prediction when the classifier’s output is randomly labeled smiling or not smiling.

4.2.5. Comparison with Similar Works

In this section, we have compared the proposed scheme with the methods of different references regarding the utility-privacy trade-off. For this purpose, the results obtained in [24] have been used. This comparison has been made for all three cases of CelebA-G-M, CelebA-G-S, and CelebA-G-C and for both weak and strong adversaries. As is evident from Figs. 14, 15 and 16, the proposed model outperforms other methods. The references of the algorithms used for comparison are mentioned in the figure. “Noise” refers to adding Gaussian noise with zero mean and variance 40 as done in [22]. In the learned noise method, first, a noise gets into a DNN, and the output is added to the dataset [20, 21]. To make a fair comparison, the adversary and the utility provider are designed to be almost identical to the previous works regarding maximum privacy and utility. Another positive point about the proposed scheme is that utility provider converges in low epochs, generally less than 10. Strong adversaries are considered in the case of curves related to previous work. In addition, the utility provider is trained on the obfuscated dataset, and its accuracy is checked on the

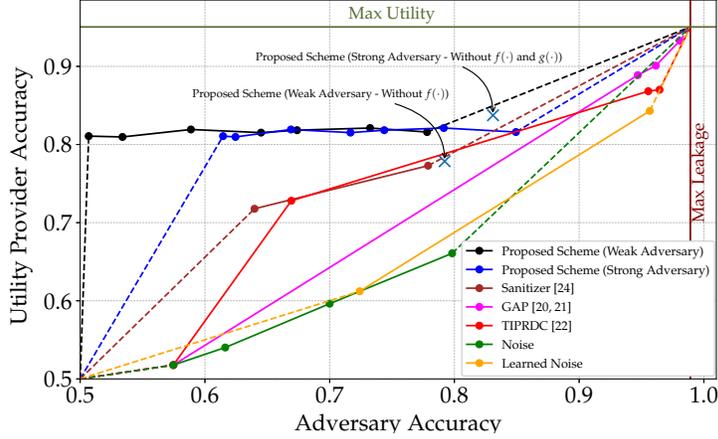


Figure 14: Comparison of utility-privacy trade-off for different methods for case CelebA-G-M.

obfuscated dataset, while the logical assumption is that the accuracy should be calculated on the original dataset (similar to what we did in this paper). As you can see in the figure, the performance curve will be better assuming the same utility provider as the previous works.

4.2.6. Categorical Dataset

For the categorical dataset, the performance of the proposed algorithm is compared with the previous methods in terms of utility-privacy trade-off in Fig. 17. All curve values of previous works are taken from [28]. AE-PUPET, UAE-PUPET, VAE-PUPET, and b-VAE-PUPET methods are related to [28], and VFAE, LMFIR, and emb-g-filter methods are introduced in [29], [30], and [31], respectively.

4.2.7. Area Under Curve

Calculating the area under the curve (AUC) can be a good measure to compare different algorithms regarding the utility-privacy trade-off. For this purpose, we consider the convex hull of the curve and then calculate AUC [28]. AUC values for different schemes are compared in Table 3. The AUC of

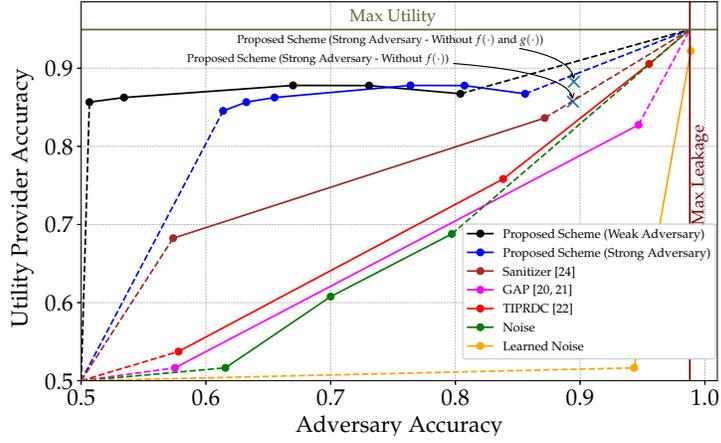


Figure 15: Comparison of utility-privacy trade-off for different methods for case CelebA-G-S.

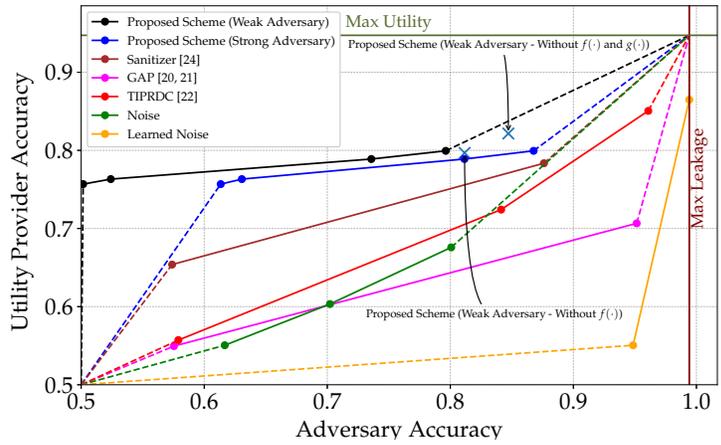


Figure 16: Comparison of utility-privacy trade-off for different methods for case CelebA-G-C.

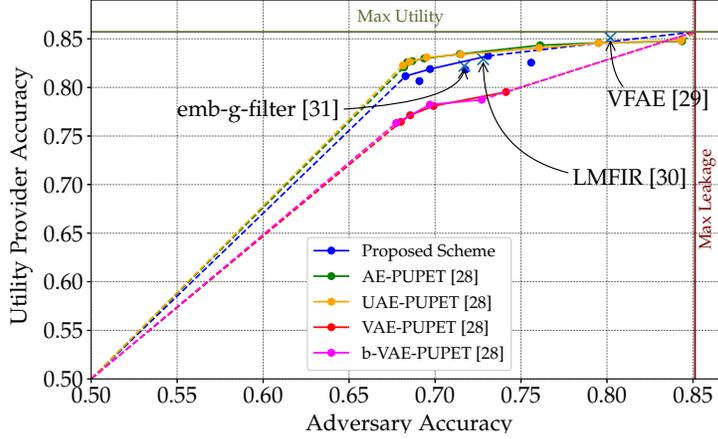


Figure 17: Comparison of utility-privacy trade-off for different methods on UCI Adult dataset. The desirable features of the adversary and the utility provider are gender and income, respectively.

the proposed scheme for the image dataset is more than other methods, and it competes with other methods in the case of the categorical dataset. For comparison, to calculate AUC similar to [24], the origin point is considered $(0, 0)$, while it is logical to take the point $(0.5, 0.5)$ as the origin.

4.3. Complexity Analysis

As you have seen in the previous sections, the proposed scheme for image datasets performs better than the earlier works. In the case of categorical datasets, it is almost close to the performance of the best available algorithm. In addition to performance, the proposed algorithm has advantages over other algorithms in terms of complexity and convergence time. The absence of a GAN structure in the proposed scheme prevents convergence and stability problems [40, 41, 42], which makes convergence faster than GAN-based algorithms. For the categorical dataset, although our performance is slightly worse than the performance of the algorithm in [28], as shown in Table 4, our obfuscator parameters are about 25% of the parameters in the proposed network in [28]. Therefore, our network structure is simpler and converges faster. Since [28] has used the

Image Dataset	(CelebA-G-M) Convex Hull AUC	(CelebA-G-S) Convex Hull AUC	(CelebA-G-C) Convex Hull AUC
Sanitizer [24]	0.5210	0.5337	0.5241
TIPRDC [22]	0.5121	0.4690	0.4736
GAP [20, 21]	0.4701	0.4690	0.4715
Gaussian Noise	0.4701	0.4690	0.4709
Learned Noise	0.4701	0.4693	0.4709
Proposed Scheme	0.5790	0.5963	0.5566
Categorical Dataset	Convex Hull AUC		
AE-PUPET [28]	0.4236		
UAE-PUPET [28]	0.4234		
VAE-PUPET [28]	0.3995		
b-VAE-PUPET [28]	0.4001		
Proposed Scheme	0.4183		

Table 3: Comparison of AUC convex hulls for utility-privacy trade-off curves.

Method	Components	Number of Parameters
Mandal et al. [28]	Obfuscator, Utility Provider, and Adversary	338, 214
Proposed Scheme	Obfuscator	88, 494

Table 4: Comparison of the complexity of the proposed scheme with the method presented in [28] in terms of the number of parameters.

GAN structure, it should also train the adversary and the utility provider in addition to the obfuscator, which increases the number of parameters.

5. Conclusion

This paper introduced a private method of data publishing using a structure based on an autoencoder. The simulation results show that the proposed method establishes a good balance in trade-off between utility and privacy compared to previous methods. Moreover, the presented method has two main advantages over the earlier methods. First, each data provider can adjust the privacy and utility level; secondly, there is no need to specify the private feature for the obfuscation design.

Acknowledgment

This work was partially supported by Iran National Science Foundation (INSF) under contract No. 97011231.

References

- [1] L. Sankar, S. R. Rajagopalan, H. V. Poor, Utility-privacy tradeoffs in databases: An information-theoretic approach, *IEEE Transactions on Information Forensics and Security* 8 (6) (2013) 838–852.
- [2] K. Kalantari, L. Sankar, O. Kosut, On information-theoretic privacy with general distortion cost functions, in: 2017 IEEE International Symposium on Information Theory (ISIT), IEEE, 2017, pp. 2865–2869.
- [3] L. Sweeney, k -anonymity: A model for protecting privacy, *International journal of uncertainty, fuzziness and knowledge-based systems* 10 (05) (2002) 557–570.
- [4] A. Machanavajjhala, D. Kifer, J. Gehrke, M. Venkatasubramanian, ℓ -diversity: Privacy beyond k -anonymity, *ACM Transactions on Knowledge Discovery from Data (TKDD)* 1 (1) (2007) 3–es.
- [5] N. Li, T. Li, S. Venkatasubramanian, t -closeness: Privacy beyond k -anonymity and ℓ -diversity, in: 2007 IEEE 23rd international conference on data engineering, IEEE, 2006, pp. 106–115.
- [6] C. Dwork, A. Roth, et al., The algorithmic foundations of differential privacy, *Foundations and Trends® in Theoretical Computer Science* 9 (3–4) (2014) 211–407.
- [7] M. S. Alvim, M. E. Andrés, K. Chatzikokolakis, P. Degano, C. Palamidessi, Differential privacy: On the trade-off between utility and information leakage., *Formal Aspects in Security and Trust* 7140 (2011) 39–54.

- [8] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, L. Zhang, Deep learning with differential privacy, in: Proceedings of the 2016 ACM SIGSAC conference on computer and communications security, 2016, pp. 308–318.
- [9] N. Phan, Y. Wang, X. Wu, D. Dou, Differential privacy preservation for deep auto-encoders: an application of human behavior prediction, in: Thirtieth AAAI Conference on Artificial Intelligence, 2016.
- [10] N. Papernot, M. Abadi, U. Erlingsson, I. Goodfellow, K. Talwar, Semi-supervised knowledge transfer for deep learning from private training data, arXiv preprint arXiv:1610.05755.
- [11] E. Hesamifard, H. Takabi, M. Ghasemi, Cryptodl: Deep neural networks over encrypted data, arXiv preprint arXiv:1711.05189.
- [12] F. Emekçi, O. D. Sahin, D. Agrawal, A. El Abbadi, Privacy preserving decision tree learning over multiple parties, *Data & Knowledge Engineering* 63 (2) (2007) 348–361.
- [13] P. Mohassel, Y. Zhang, Secureml: A system for scalable privacy-preserving machine learning, in: 2017 IEEE symposium on security and privacy (SP), IEEE, 2017, pp. 19–38.
- [14] A.-T. Tran, T.-D. Luong, J. Karnjana, V.-N. Huynh, An efficient approach for privacy preserving decentralized deep learning models based on secure multi-party computation, *Neurocomputing* 422 (2021) 245–262.
- [15] R. Shokri, V. Shmatikov, Privacy-preserving deep learning, in: Proceedings of the 22nd ACM SIGSAC conference on computer and communications security, 2015, pp. 1310–1321.
- [16] H. Zhu, R. Wang, Y. Jin, K. Liang, J. Ning, Distributed additive encryption and quantization for privacy preserving federated deep learning, *Neurocomputing* 463 (2021) 309–327.

- [17] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, H. V. Poor, Federated learning with differential privacy: Algorithms and performance analysis, *IEEE Transactions on Information Forensics and Security* 15 (2020) 3454–3469.
- [18] Z. Kuang, Z. Guo, J. Fang, J. Yu, N. Babaguchi, J. Fan, Unnoticeable synthetic face replacement for image privacy protection, *Neurocomputing* 457 (2021) 322–333.
- [19] N. Raval, A. Machanavajjhala, J. Pan, Olympus: Sensor privacy through utility aware obfuscation., *Proc. Priv. Enhancing Technol.* 2019 (1) (2019) 5–25.
- [20] C. Huang, P. Kairouz, X. Chen, L. Sankar, R. Rajagopal, Context-aware generative adversarial privacy, *Entropy* 19 (12) (2017) 656.
- [21] P. Kairouz, J. Liao, C. Huang, M. Vyas, M. Welfert, L. Sankar, Generating fair universal representations using adversarial models, *IEEE Transactions on Information Forensics and Security* 17 (2022) 1970–1985.
- [22] A. Li, Y. Duan, H. Yang, Y. Chen, J. Yang, Tiprdc: task-independent privacy-respecting data crowdsourcing framework for deep learning with anonymized intermediate representations, in: *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 824–832.
- [23] A. Li, J. Guo, H. Yang, Y. Chen, Deepobfuscator: Adversarial training framework for privacy-preserving image classification, *arXiv preprint arXiv:1909.04126* 2 (3).
- [24] A. Singh, E. Garza, A. Chopra, P. Vepakomma, V. Sharma, R. Raskar, Decouple-and-sample: Protecting sensitive information in task agnostic data release, in: *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, Springer, 2022, pp. 499–517.

- [25] S. A. Osia, A. Taheri, A. S. Shamsabadi, K. Katevas, H. Haddadi, H. R. Rabiee, Deep private-feature extraction, *IEEE Transactions on Knowledge and Data Engineering* 32 (1) (2018) 54–66.
- [26] S. A. Osia, A. S. Shamsabadi, S. Sajadmanesh, A. Taheri, K. Katevas, H. R. Rabiee, N. D. Lane, H. Haddadi, A hybrid deep learning architecture for privacy-preserving mobile analytics, *IEEE Internet of Things Journal* 7 (5) (2020) 4505–4518.
- [27] H. Nguyen, D. Zhuang, P.-Y. Wu, M. Chang, Autogan-based dimension reduction for privacy preservation, *Neurocomputing* 384 (2020) 94–103.
- [28] B. Mandal, G. Amariuca, S. Wei, Uncertainty-autoencoder-based privacy and utility preserving data type conscious transformation, in: *2022 International Joint Conference on Neural Networks (IJCNN)*, IEEE, 2022, pp. 1–8.
- [29] C. Louizos, K. Swersky, Y. Li, M. Welling, R. Zemel, The variational fair autoencoder, *arXiv preprint arXiv:1511.00830*.
- [30] J. Song, P. Kalluri, A. Grover, S. Zhao, S. Ermon, Learning controllable fair representations, in: *The 22nd International Conference on Artificial Intelligence and Statistics*, PMLR, 2019, pp. 2164–2173.
- [31] X. Chen, T. Navidi, S. Ermon, R. Rajagopal, Distributed generation of privacy preserving data with user customization, *arXiv preprint arXiv:1904.09415*.
- [32] F. Tramer, D. Boneh, Differentially private learning needs better features (or much more data), *arXiv preprint arXiv:2011.11660*.
- [33] N. Raval, A. Machanavajjhala, L. P. Cox, Protecting visual secrets using adversarial nets, in: *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, IEEE, 2017, pp. 1329–1332.

- [34] Z. Liu, P. Luo, X. Wang, X. Tang, Deep learning face attributes in the wild, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 3730–3738.
- [35] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.
- [36] P. Contributors, Softmax - pytorch 2.0 documentation, accessed on March 19, 2023 (2023).
URL <https://pytorch.org/docs/stable/generated/torch.nn.Softmax.html>
- [37] D. Dua, C. Graff, et al., Uci machine learning repository.
- [38] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, Advances in neural information processing systems 32.
- [39] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980.
- [40] I. Goodfellow, Nips 2016 tutorial: Generative adversarial networks, arXiv preprint arXiv:1701.00160.
- [41] N. Kodali, J. Abernethy, J. Hays, Z. Kira, On convergence and stability of gans, arXiv preprint arXiv:1705.07215.
- [42] S. A. Barnett, Convergence problems with generative adversarial networks (gans), arXiv preprint arXiv:1806.11382.