

## Privacy-Preserving Adversarial Networks

Wang, Y.; Koike-Akino, T.; Erdogmus, D.

TR2018-062 July 10, 2018

### Abstract

We propose a data-driven framework for optimizing privacy-preserving data release mechanisms to attain the information-theoretically optimal tradeoff between minimizing distortion of useful data and concealing sensitive information. Our approach employs adversarially-trained neural networks to implement randomized mechanisms and to perform a variational approximation of mutual information privacy. We validate our Privacy-Preserving Adversarial Networks (PPAN) framework via experiments on discrete and continuous synthetic data, as well as the MNIST handwritten digits dataset. For synthetic data, we find that our model-agnostic PPAN approach achieves tradeoff points very close to the optimal tradeoffs that are analytically-derived from model knowledge. In experiments with the MNIST data, we visually demonstrate a learned tradeoff between minimizing the pixel-level distortion versus concealing the written digit.

*arXiv*

This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of Mitsubishi Electric Research Laboratories, Inc.; an acknowledgment of the authors and individual contributions to the work; and all applicable portions of the copyright notice. Copying, reproduction, or republishing for any other purpose shall require a license with payment of fee to Mitsubishi Electric Research Laboratories, Inc. All rights reserved.



---

# Privacy-Preserving Adversarial Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose a data-driven framework for optimizing privacy-preserving data re-  
2 lease mechanisms to attain the information-theoretically optimal tradeoff between  
3 minimizing distortion of useful data and concealing sensitive information. Our  
4 approach employs adversarially-trained neural networks to implement randomized  
5 mechanisms and to perform a variational approximation of mutual information pri-  
6 vacy. We validate our Privacy-Preserving Adversarial Networks (PPAN) framework  
7 via experiments on discrete and continuous synthetic data, as well as the MNIST  
8 handwritten digits dataset. For synthetic data, we find that our model-agnostic  
9 PPAN approach achieves tradeoff points very close to the optimal tradeoffs that are  
10 analytically-derived from model knowledge. In experiments with the MNIST data,  
11 we visually demonstrate a learned tradeoff between minimizing the pixel-level  
12 distortion versus concealing the written digit.

## 13 1 Introduction

14 Our work addresses the problem of privacy-preserving data release, where the goal is to release  
15 useful data while also limiting the exposure of associated sensitive information. Approaches that  
16 involve data modification must consider the tradeoff between concealing sensitive information and  
17 minimizing distortion to preserve data utility. However, practical optimization of this tradeoff can be  
18 challenging when we wish to quantify privacy via statistical measures (such as mutual information)  
19 and the actual statistical distributions of data are unknown. In this paper, we propose a data-driven  
20 framework involving adversarially trained neural networks to design privacy-preserving data release  
21 mechanisms that approach the theoretically optimal privacy-utility tradeoffs.

22 Privacy-preserving data release is a broad and widely explored field, where the study of principled  
23 methods have been well motivated by highly publicized leaks stemming from the inadequacy of simple  
24 anonymization techniques, such as reported in [26, 22]. A wide variety of methods to statistically  
25 quantify and address privacy have been proposed, such as  $k$ -anonymity [27],  $L$ -diversity [17],  $t$ -  
26 closeness [15], and differential privacy [5]. In our work, we focus on an information-theoretic  
27 approach where privacy is quantified by the mutual information between the data release and the  
28 sensitive information [32, 24, 3, 25, 2]. Unlike the methods mentioned earlier, measuring privacy  
29 via mutual information implicitly requires consideration of the statistical distribution of the data.  
30 Ignoring the data distribution can weaken the scope of privacy guarantees. For example, an adversary  
31 armed with only mild knowledge about the correlation of the data<sup>1</sup> can undermine the practical  
32 privacy protection of differential privacy, as noted in examples given by [12, 3, 16, 31]. While model  
33 assumptions are avoided in the definition of differential privacy, independence across individuals  
34 in the dataset is implicitly required to avoid undermining privacy guarantees [12]. The example  
35 in [3, Sec. V] demonstrates that an  $\epsilon$ -differentially private mechanism can leak an unbounded amount

---

<sup>1</sup>Note that even when data samples are inherently independent, the prior knowledge of an adversary could become correlated when conditioned on particular side information.

36 of sensitive information (on the order of  $O(\epsilon^2 \log n)$  where  $n$  is size of the dataset). The *linkage*  
37 *inequality* [31, Def. 2] formulates certain properties that one may reasonably require of privacy  
38 measures; however, it is not satisfied by differential privacy.

39 We build upon the non-asymptotic, information-theoretic framework introduced by [24, 3], where  
40 the sensitive and useful data are respectively modeled as random variables  $X$  and  $Y$ . We also  
41 adopt the extension considered in [2], where only a (potentially partial and/or noisy) observation  
42  $W$  of the data is available. In this framework, the design of the privacy-preserving mechanism to  
43 release  $Z$  is formulated as the optimization of the tradeoff between minimizing privacy-leakage  
44 quantified by the mutual information  $I(X; Z)$  and minimizing an expected distortion  $\mathbb{E}[d(Y, Z)]$ .  
45 This non-asymptotic framework has strong connections to generalized rate-distortion problems (see  
46 discussion in [24, 3, 31]), as well as related asymptotic privacy frameworks where communication  
47 efficiency is also considered in a rate-distortion-privacy tradeoff [32, 25].

48 In principle, when the data distribution is known, the optimal design of the privacy-preserving  
49 mechanism can be tackled as a convex optimization problem [24, 3]. However, in practice, model  
50 knowledge is often missing or inaccurate for realistic data sets, and the optimization becomes  
51 intractable for high-dimensional and continuous data. Addressing these challenges, we propose a  
52 data-driven approach that optimizes the privacy-preserving mechanism to attain the theoretically  
53 optimal privacy-utility tradeoffs, by learning from a set of training data rather than requiring model  
54 knowledge. We call this approach *Privacy-Preserving Adversarial Networks* (PPAN) since the  
55 mechanism, realized as a randomized neural network, is trained along with an adversarial network  
56 that attempts to recover the sensitive information from the released data. The key to attaining  
57 information-theoretic privacy is that the adversarial network specifically estimates the posterior  
58 distribution (rather than only the value) of the sensitive variable given the released data to enable  
59 a variational approximation of mutual information [1]. While the adversary is trained to minimize  
60 the log-loss with respect to this posterior estimate, the mechanism network is trained to attain the  
61 dual objectives of minimizing distortion and concealing sensitive information (by maximizing the  
62 adversarial loss).

## 63 1.1 Related Work

64 The general concept of adversarial training of neural networks was introduced by [7], which proposed  
65 *Generative Adversarial Networks* (GAN) for learning generative models that can synthesize new data  
66 samples. Since their introduction, GANs have inspired a large and growing number of adversarially  
67 trained neural network architectures for a wide variety of purposes [9].

68 The earlier works of [6, 8] have also proposed adversarial training frameworks for optimizing privacy-  
69 preserving mechanisms, where the adversarial network is realized as a classifier that attempts to  
70 recover a discrete sensitive variable. In [6], the mechanism is realized as an autoencoder, and  
71 the adversary attempts to predict a binary sensitive variable from the latent representation. In the  
72 framework of [8], a deterministic mechanism is trained with the adversarial network realized as a  
73 classifier attempting to predict the sensitive variable from the output of the mechanism. Both of  
74 these frameworks additionally propose using an optional predictor network that attempts to predict  
75 a useful variable from the output of the mechanism network. Thus, while the adversarial network  
76 is trained to recover the sensitive variable, the mechanism and predictor (if present) networks are  
77 trained to realize multiple objectives: maximizing the loss of the adversary as well as minimizing the  
78 reconstruction loss of the mechanism network and/or the prediction loss of the predictor network.  
79 However, a significant limitation of both of these approaches is that they consider only deterministic<sup>2</sup>  
80 mechanisms, which generally do not achieve the optimal privacy-utility tradeoffs, although neither  
81 attempts to address information-theoretic privacy.

82 The recent, independent work of [10] proposes a similar adversarial training framework, which also  
83 realizes the necessity of and proposes randomized mechanism networks, in order to address the  
84 information-theoretically optimal privacy-utility tradeoffs. They also rediscover the earlier realization  
85 of [3] that mutual information privacy arises from an adversary (which outputs a distribution) that  
86 is optimized with respect to log-loss. However, their framework does not make the connections  
87 to a general variational approximation of mutual information applicable to arbitrary (i.e., discrete,

---

<sup>2</sup>While [8] does also consider a “noisy” version of their mechanism, the randomization is limited to only independent, additive noise before or after deterministic filtering.

88 continuous, and/or multivariate) sensitive variable alphabets, and hence their data-driven formulation  
 89 and empirical evaluation is limited to only binary sensitive variables.

## 90 1.2 Contributions and Paper Outline

91 Our framework, presented in Section 2, provides the first data-driven approach for optimizing  
 92 privacy-preserving data release mechanisms that approaches the information-theoretically optimal  
 93 privacy-utility tradeoffs. A key novelty of our approach is the use of adversarial training to perform a  
 94 variational approximation of mutual information privacy. Unlike previous work, our approach can  
 95 handle randomized data release mechanisms where the input to the mechanism can be a general  
 96 observation of the data, e.g., a full or potentially noisy/partial view of the sensitive and useful  
 97 variables.

98 In our proposed framework all of the variables that are involved can be discrete, continuous, and/or  
 99 high-dimensional vectors. We develop specific network architectures and sampling methods ap-  
 100 propriate for various scenarios in Section 2.3. In particular, when all of the variables have finite  
 101 alphabets, we demonstrate that the network architectures can be efficiently minimalized to essentially  
 102 just the matrices describing the conditional distributions, and that replacing sampling with a directly  
 103 computed expectation improves training performance.

104 We evaluate our PPA approach in Section 3 with experiments on multivariate Gaussian synthetic  
 105 data and the MNIST handwritten digit dataset. For the synthetic data experiment, we demonstrate  
 106 that PPA closely approaches the theoretically optimal privacy-utility tradeoff.

107 In the supplementary document, we present a number of additional results, detailed technical deriva-  
 108 tions, and extensions which are of independent interest which could not be accommodated within the  
 109 main paper due to space limitations. In Section A, we discuss how to handle mutual information as a  
 110 utility function in the PPA framework as opposed to expected distortion that we focus on in this  
 111 work. In Section B, we discuss how the technique of sampling from a multivariate Gaussian described  
 112 in Section 2.3.2 can be extended to GMMs. In Section C, we consider synthetic discrete-valued data  
 113 following a *symmetric pair* distribution and compare the privacy-utility tradeoff results with an ap-  
 114 proach addressing the same problem in [19]. The analytical expressions for the theoretically-optimal  
 115 privacy-utility tradeoffs for symmetric pair distribution are presented in Section D. In Section E, we  
 116 consider scalar jointly Gaussian sensitive and useful attributes and benchmark the performance of  
 117 PPA against the theoretically optimal privacy-utility tradeoff. In Section G, we demonstrate how  
 118 the PPA framework can be used to generate rate-distortion curves studied in information theory,  
 119 purely from samples. Finally in Sections H and I, we provide and derive analytical expressions for  
 120 the optimal privacy-utility tradeoffs for Gaussian distributed data and mean square error distortion.

## 121 2 Problem Formulation and PPA Methods

### 122 2.1 Privacy-Utility Tradeoff Optimization

123 We consider the privacy-utility tradeoff optimization problem described in [2], which extends the  
 124 frameworks initiated by [24, 3]. Figure 1 depicts the problem setting where observed data  $W$ ,  
 125 sensitive attributes  $X$ , and useful attributes  $Y$  are modeled as random variables that are jointly  
 126 distributed according to a data model  $P_{W,X,Y}$  over the space  $\mathcal{W} \times \mathcal{X} \times \mathcal{Y}$ . The observed data  $W$  is a  
 127 potentially noisy/partial observation of the sensitive and useful data attributes  $(X, Y)$ . The goal is to  
 128 design and optimize the data release mechanism, i.e., a system that processes the observed data  $W$  to  
 129 produce a release  $Z \in \mathcal{Z}$  that minimizes the privacy-leakage of the sensitive attributes  $X$ , while also  
 130 maximizing the utility gained from revealing information about  $Y$ . This system is specified by the  
 131 *release mechanism*  $P_{Z|W}$ , with  $(W, X, Y, Z) \sim P_{W,X,Y}P_{Z|W}$ , and thus  $(X, Y) \leftrightarrow W \leftrightarrow Z$  forms a  
 132 Markov chain. Privacy-leakage is quantified by the mutual information  $I(X; Z)$  between the sensitive  
 133 attributes  $X$  and the release  $Z$ . Utility is inversely quantified by the expected distortion  $\mathbb{E}[d(Y, Z)]$   
 134 between the useful attributes  $Y$  and the release  $Z$ , where the distortion function  $d : \mathcal{Y} \times \mathcal{Z} \rightarrow [0, \infty)$   
 135 is given by the application. The design of the release mechanism  $P_{Z|W}$  is formulated as the following  
 136 privacy-utility tradeoff optimization problem,

$$136 \min_{P_{Z|W}: (X,Y) \leftrightarrow W \leftrightarrow Z} I(X; Z), \quad \text{s.t.} \quad \mathbb{E}[d(Y, Z)] \leq \delta, \quad (1)$$

137 where the parameter  $\delta$  indicates the distortion (or *disutility*) budget allowed for the sake of preserving  
 138 privacy.

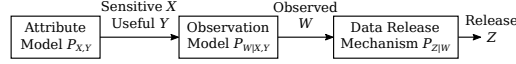


Figure 1: Setting for privacy-utility tradeoff optimization.

139 As noted in [2], given a fixed data model  $P_{W,X,Y}$  and distortion function  $d$ , the problem in (1) is a  
 140 convex optimization problem, since the mutual information objective  $I(X; Z)$  is a convex functional  
 141 of  $P_{Z|X}$ , which is in turn a linear functional of  $P_{Z|W}$ , and the expected distortion  $\mathbb{E}[d(Y, Z)]$  is a  
 142 linear functional of  $P_{Y,Z}$  and hence also of  $P_{Z|W}$ . While the treatment in [2] considers discrete  
 143 variables over finite alphabets, the formulation of (1) need not be limited those assumptions. Thus, in  
 144 this work, we seek to also address this problem with high-dimensional, continuous variables.

## 145 2.2 Adversarial Training for an Unknown Data Model

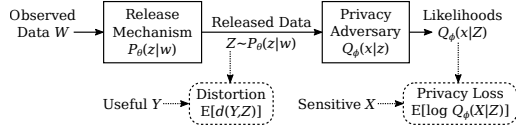


Figure 2: Adversarial training framework.

146 Our aim is to solve the privacy-utility tradeoff optimization problem when the data model  $P_{W,X,Y}$  is  
 147 unknown but instead a set of training samples is available:  $\{(w_i, x_i, y_i)\}_{i=1}^n \sim \text{i.i.d. } P_{W,X,Y}$ . A key  
 148 to our approach is approximating  $I(X; Z)$  via a variational lower bound given by [1] and also used  
 149 in [4]. This bound is based on the following identity which holds for any distribution  $Q_{X|Z}$  over  $\mathcal{X}$   
 150 given values in  $\mathcal{Z}$

$$-h(X|Z) = \text{KL}(P_{X|Z} \| Q_{X|Z}) + \mathbb{E}[\log Q_{X|Z}(X|Z)],$$

151 where  $\text{KL}(\cdot \| \cdot)$  denotes the Kullback-Leibler (KL) divergence. Therefore, since  $I(X; Z) = h(X) -$   
 152  $h(X|Z)$  and KL divergence is nonnegative,

$$h(X) + \max_{Q_{X|Z}} \mathbb{E}[\log Q_{X|Z}(X|Z)] = I(X; Z), \quad (2)$$

153 where the maximum is attained when the variational posterior  $Q_{X|Z} = P_{X|Z}$ . Using (2) with  
 154 the constant  $h(X)$  term dropped, we convert the formulation of (1) to an unconstrained minimax  
 155 optimization problem,

$$\min_{P_{Z|W}} \max_{Q_{X|Z}} \mathbb{E}[\log Q_{X|Z}(X|Z)] + \lambda \mathbb{E}[d(Y, Z)], \quad (3)$$

156 where the expectations are with respect to  $(W, X, Y, Z) \sim P_{W,X,Y} P_{Z|W}$ , and the parameter  $\lambda > 0$   
 157 can be adjusted to obtain various points on the optimal privacy-utility tradeoff curve. Alternatively,  
 158 to target a specific distortion budget  $\delta$ , the second term in (3) could be replaced with a penalty term  
 159  $\lambda(\max(0, \mathbb{E}[d(Y, Z)] - \delta))^2$ , where  $\lambda > 0$  is made relatively large to penalize exceeding the budget.  
 160 The expectations in (3) can be conveniently approximated by Monte Carlo sampling over training set  
 161 batches.

162 The minimax formulation of (3) can be interpreted and realized in an adversarial training framework  
 163 (as illustrated by Figure 2), where the variational posterior  $Q_{X|Z}$  is viewed as the posterior likelihood  
 164 estimates of the sensitive attributes  $X$  made by an adversary observing the release  $Z$ . The data  
 165 release mechanism is trained to minimize both the distortion and privacy loss terms, while the  
 166 adversary is trained to maximize the privacy loss. Specifically, the adversary attempts to maximize  
 167 the negative log-loss  $\mathbb{E}[\log Q_{X|Z}(X|Z)]$ , which the release mechanism  $P_{Z|W}$  attempts to minimize.  
 168 The release mechanism and adversary are realized as neural networks, which take as inputs  $W$  and  $Z$ ,  
 169 respectively, and produce the parameters that specify their respective distributions  $P_{Z|W}$  and  $Q_{X|Z}$   
 170 within parametric families that are appropriate for the given application. For example, a release  
 171 mechanism suitable for the release space  $\mathcal{Z} = \mathbb{R}^d$  could be the multivariate Gaussian

$$P_{Z|W}(z|w) = \mathcal{N}(z; (\boldsymbol{\mu}, \boldsymbol{\Sigma}) = f_{\theta}(w)),$$

172 where the mean  $\mu$  and covariance  $\Sigma$  are determined by a neural network  $f_\theta$  as a function of  $w$  and  
 173 controlled by the parameters  $\theta$ . For brevity of notation, we will use  $P_\theta(z|w)$  to denote the distribution  
 174 defined by the release mechanism network  $f_\theta$ . Similarly, we will let  $Q_\phi(x|z)$  denote the parametric  
 175 distribution defined by the adversary network that is controlled by the parameters  $\phi$ . For each training  
 176 sample tuple  $(w_i, x_i, y_i)$ , we sample  $k$  independent releases  $\{z_{i,j}\}_{j=1}^k \stackrel{\text{iid}}{\sim} P_\theta(z|w_i)$  to approximate  
 177 the loss term with

$$\mathcal{L}^i(\theta, \phi) := \frac{1}{k} \sum_{j=1}^k [\log Q_\phi(x_i|z_{i,j}) + \lambda d(y_i, z_{i,j})]. \quad (4)$$

178 The networks are optimized with respect to these loss terms averaged over the training data (or  
 179 mini-batches)

$$\min_{\theta} \max_{\phi} \frac{1}{n} \sum_{i=1}^n \mathcal{L}^i(\theta, \phi), \quad (5)$$

180 which approximates the theoretical privacy-utility tradeoff optimization problem as given in (3), since  
 181 by the law of large numbers, as  $n \rightarrow \infty$ ,

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}^i(\theta, \phi) \xrightarrow{\text{a.s.}} \mathbb{E}[\log Q_\phi(X|Z) + \lambda d(Y, Z)],$$

182 where the expectation is with respect to  $(W, X, Y, Z) \sim P_{W,X,Y} P_\theta(z|w)$ . Similarly, the second  
 183 term in (4) could be replaced with a penalty term  $\lambda(\max(0, d(y_i, z_{i,j}) - \delta))^2$  to target a specific  
 184 distortion budget  $\delta$ . Similar to GANs [7], the minimax optimization in (5) can be more practically  
 185 handled by alternating gradient descent/ascent between the two networks (possibly with multiple  
 186 inner maximization updates per outer minimization update) rather than optimizing the adversary  
 187 network until convergence for each release mechanism network update.

### 188 2.3 Sampling the Release Mechanism

189 To allow optimization of the networks via gradient methods, the release samples need to be generated  
 190 such that the gradients of the loss terms can be readily calculated. Various forms of the release  
 191 mechanism distribution  $P_\theta(z|w)$  are appropriate for different applications, and each require their  
 192 own specific sampling methods. In this section, we outline some of these forms and their associated  
 193 sampling methods.

#### 194 2.3.1 Finite Alphabets

195 When the release space  $\mathcal{Z}$  is a finite discrete set, we can forgo sampling altogether and calculate the  
 196 loss terms via

$$\mathcal{L}_{\text{disc}}^i(\theta, \phi) := \sum_{z \in \mathcal{Z}} P_\theta(z|w_i) (\log Q_\phi(x_i|z) + \lambda d(y_i, z)), \quad (6)$$

198 which replaces the empirical average over  $k$  samples with the direct expectation over  $Z$ . We found  
 199 that this direct expectation produced better results than estimation via sampling, such as by applying  
 200 the Gumbel-softmax categorical reparameterization trick (see [18, 11]).

201 Further, if  $\mathcal{W}$  and  $\mathcal{X}$  are also finite alphabets, then  $P_\theta(z|w)$  and  $Q_\phi(x|z)$  can be exactly parameterized  
 202 by matrices of size  $|\mathcal{Z}| \times |\mathcal{W}|$  and  $|\mathcal{X}| \times |\mathcal{Z}|$ , respectively. Thus, in the purely finite alphabet case,  
 203 with the variables represented as one-hot vectors, the mechanism and adversary are most efficiently  
 204 realized as networks with no hidden layers and softmax applied to the output (to yield stochastic  
 205 vectors).

#### 206 2.3.2 Gaussian Approximations for Reals

207 A multivariate Gaussian release mechanism can be sampled by employing the reparameterization trick  
 208 of [14], which first samples a vector of independent standard normal variables  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and then  
 209 generates  $z = \mathbf{A}\mathbf{u} + \mu$ , where the parameters  $(\mu, \mathbf{A}) = f_\theta(w)$  are produced by the release mechanism  
 210 network to specify a conditional Gaussian with mean  $\mu$  and covariance  $\Sigma = \mathbf{A}\mathbf{A}^T$ . This approach  
 211 can be extended to Gaussian Mixture Models as explained in Section B of the supplementary.

212 **2.3.3 Universal Approximators**

213 Another approach, as seen in [21], is to directly produce the release sample as  $z = f_\theta(w, u)$  using  
 214 a neural network that takes random seed noise  $u$  as an additional input. The seed noise  $u$  can be  
 215 sampled from a simple distribution (e.g., uniform, Gaussian, etc.) and provides the randomization  
 216 of  $z$  with respect to  $w$ . Since the transformations applying the seed noise can, in principle, be  
 217 learned, this approach could potentially approximate any “nice” distribution due to the universal  
 218 approximation properties of neural networks. However, although it is not needed for training, it is  
 219 generally intractable to produce an explicit expression for  $P_\theta(z|w)$  as implied by the network.

220 **3 Experimental Results**

221 In this section, we present the privacy-utility tradeoffs that are achieved by our PPAN framework in  
 222 experiments with synthetic and real data. For the synthetic data experiment, we consider Gaussian  
 223 joint distribution over the sensitive, useful, and observed data, for which we can compare the  
 224 results obtained by PPAN against the theoretically optimal tradeoffs (derived in Section H of the  
 225 supplementary material). We use the MNIST handwritten digits dataset to illustrate the application  
 226 of the PPAN framework to real data in Section 3.2. We demonstrate optimized networks that can  
 227 trace the tradeoff between concealing the digit and reducing image distortion. Our experiments were  
 228 implemented using the Chainer deep learning framework [30], with optimization performed by their  
 229 implementation of Adam [13]. Experiments on synthetic discrete data are presented in Section C.

230 **3.1 Gaussian Synthetic Data**

231 Consider multivariate jointly Gaussian sensitive and useful attributes  $\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} I_5 & \text{diag}(\rho) \\ \text{diag}(\rho) & I_5 \end{bmatrix}\right)$   
 232 where both  $X, Y \in \mathbb{R}^5$  and  $\rho = [0.47, 0.24, 0.85, 0.07, 0.66]$ . The observed data is  $W = Y$ . As we  
 233 note in Section H of supplementary material, the optimum release in this case is jointly Gaussian with  
 234 the attributes. Thus we could use a mechanism network architecture that can realize the procedure  
 235 described in Section 2.3.2 to generate the release. However, since the optimal release distribution is  
 236 not known for general attribute models, we use the universal approximator technique described in  
 237 Section 2.3.3.

238 The mechanism implemented in these experiments consists of three fully connected layers, with the  
 239 ReLU activation function applied at the outputs of the two hidden layers, and no activation function  
 240 is used at the output layer. The mechanism takes as input observation  $W$  and seed noise  $U$ , and  
 241 generates the release  $Z = f_\theta(W, U)$ , where  $\theta$  denotes the parameters of the mechanism network. The  
 242 seed noise vector has 8 components, each of which are i.i.d. Uniform $[-1, 1]$ . The adversary network,  
 243 with parameters denoted by  $\phi$ , models the posterior probability  $Q_\phi(X|Z)$  of the sensitive attribute  
 244 given the release. We assume that  $Q_\phi(\cdot|z)$  is a normal distribution with mean vector  $\mu_\phi(z)$  and  
 245 covariance matrix  $\text{diag}(\sigma_\phi^2(z))$ , i.e., they are functions of the release  $z$ . The adversary network has  
 246 three fully connected layers to learn the mean and variances. The network takes as input the release  
 247  $z$  and outputs the pair  $(\mu_\phi(z), \log \sigma_\phi^2(z))$ , where the log is applied componentwise on the variance  
 248 vector. We use the adversarial networks to solve the min-max optimization problem described in (5).  
 249 We choose  $k = 1$  in (4), and similar to the previous section, we use the penalty modification of the  
 250 distortion term, i.e.,

$$\mathcal{L}_{\text{gauss}}^i(\theta, \phi) = \log Q_\phi(x_i|z_i) + \lambda(\max(0, \|y_i - z_i\|^2 - \delta))^2.$$

251 We choose the multiplier  $\lambda = 10$  and 20 linearly spaced values for  $\delta$  in the range  $[0, 4.5]$ . For each  
 252 value of  $\delta$ , we sample the data model to obtain an independent dataset realization and use it to train  
 253 and test the adversarial networks. We use 8000 training samples and evaluate the performance of  
 254 PPAN on 4000 test samples. Each hidden layer has 20 nodes. The adversarial networks were trained  
 255 for 250 epochs with a minibatch size of 200. In each iteration we do 5 gradient descent steps to  
 256 update the parameters of the adversary network before updating the mechanism network.

257 We plot the privacy-leakage and distortion values returned by the PPAN mechanism on the test set  
 258 along with the optimal tradeoff curve (from Proposition 2 of the supplementary material) in Figure 3.  
 259 The privacy-leakage values were estimated following the procedure in Section F. The performance of  
 260 the PPAN mechanism is very close to the theoretically optimum tradeoff curve over a wide range of  
 261 target distortion values.



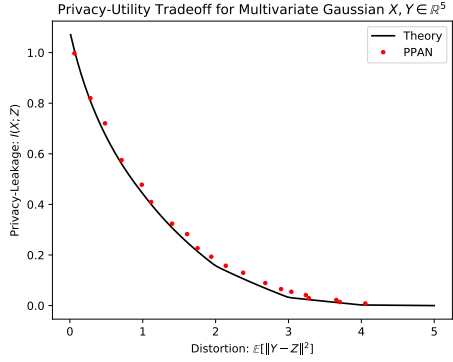


Figure 3: Comparison of the results achieved by PPAN versus the theoretical optimum tradeoff curve for the useful data only observation model (i.e.,  $W = Y$ ) for multivariate Gaussian  $(X, Y)$ .

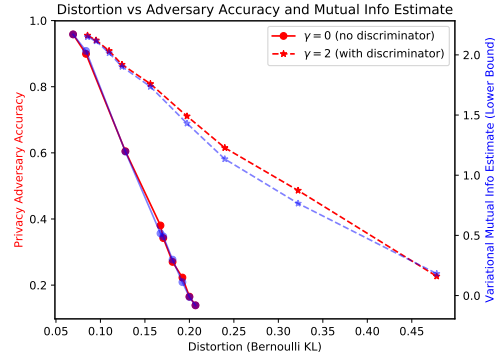
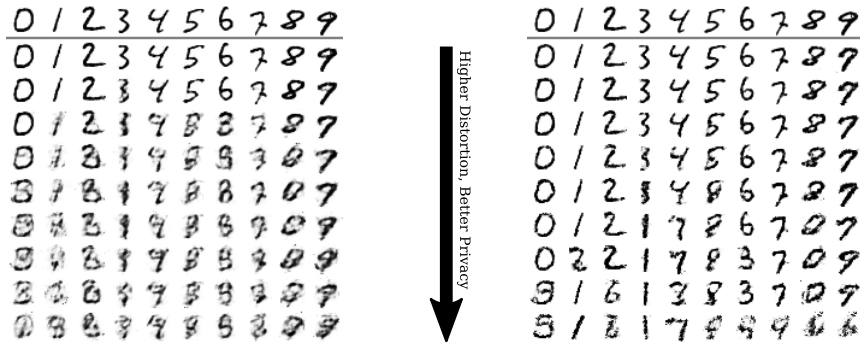


Figure 4: Evaluation of the distortion vs privacy tradeoffs for PPAN applied to the MNIST test set, with privacy measured by adversary accuracy (in red) and estimated mutual information (in blue).



(a) Without additional discriminator (i.e.,  $\gamma = 0$ ).

(b) With additional discriminator ( $\gamma = 2$ ).

Figure 5: Examples from applying PPAN to conceal MNIST handwritten digits. Top row consists of the original test set examples input to the mechanism, while other rows show corresponding mechanism outputs at different tradeoff points.

### 262 3.2 MNIST Handwritten Digits

263 The MNIST dataset consists of 70K labeled images of handwritten digits split into training and test  
 264 sets of 60K and 10K images, respectively. Each image consists of  $28 \times 28$  grayscale pixels, which  
 265 we handle as vectors in  $[0, 1]^{784}$ . In the first set of experiments, we consider the image to be both  
 266 the useful and the observed data, i.e.,  $W = Y$ , the digit label to be the sensitive attribute  $X$ , and the  
 267 mechanism release as an image  $Z \in [0, 1]^{784}$ . We measure the distortion between the original and  
 268 released images  $Y, Z$  as

$$d(Y, Z) := \frac{-1}{784} \sum_{i=1}^{784} Y[i] \log(Z[i]) + (1 - Y[i]) \log(1 - Z[i]),$$

269 which, for a fixed  $Y$ , corresponds to minimizing the average KL-divergence between corresponding  
 270 pixels that are each treated as a Bernoulli distribution. Thus, the privacy objective is to conceal the  
 271 digit, while the utility objective is to minimize (average pixel-level) image distortion.

272 The mechanism and adversary networks both use two hidden layers with 1000 nodes each and  
 273 fully-connected links between all layers. The hidden layers use tanh as the activation function.  
 274 The mechanism input layer uses  $784 + 20$  nodes for the image concatenated with 20 random  
 275 Uniform $[-1, 1]$  seed noise values. The mechanism output layer uses 784 nodes with the sigmoid  
 276 activation function to directly produce an image in  $[0, 1]^{784}$ . Note that the mechanism network is an  
 277 example of the universal approximator architecture mentioned in Section 2.3.3. The attacker input

278 layer uses 784 nodes to receive the image produced by the mechanism. The attacker output layer  
 279 uses 10 nodes normalized with a softmax activation function to produce a distribution over the digit  
 280 labels  $\{0, \dots, 9\}$ .

281 In a second set of experiments, we employ the standard GAN approach of adding a discriminator  
 282 network to further encourage the mechanism to produce output images that resemble realistic digits.  
 283 The discriminator network architecture uses a single hidden layer with 500 nodes, and has an output  
 284 layer with one node that uses the sigmoid activation function. The discriminator network, denoted  
 285 by  $D_\psi$  with parameters  $\psi$ , attempts to distinguish the outputs of the mechanism network from the  
 286 original training images. Its contribution to the overall loss is controlled by a parameter  $\gamma \geq 0$  (with  
 287 zero indicating its absence). Incorporating this additional network, the training loss terms are given  
 288 by

$$\mathcal{L}_{\text{mnist}}^i(\theta, \phi, \psi) := \log Q_\phi(x_i|z_i) + \lambda d(y_i, z_i) + \gamma \log D_\psi(z_i) + \gamma \log(1 - D_\psi(y_i)), \quad (7)$$

289 where  $z_i$  is generated from the input image  $w_i = y_i$  by the mechanism network controlled by the  
 290 parameters  $\theta$ . The overall adversarial optimization objective with both the privacy adversary and the  
 291 discriminator is given by

$$\min_{\theta} \max_{\phi, \psi} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{mnist}}^i(\theta, \phi, \psi).$$

292 We used the 10K test images to objectively evaluate the performance of the trained mechanisms for  
 293 Figure 4, which depicts image distortion versus privacy measured by the accuracy of the adversary in  
 294 recognizing the original digit and the variational lower bound for mutual information obtained by  
 295 using the posterior distribution of the sensitive attribute learnt by the adversary in (2).

296 Figure 5 shows example results from applying trained privacy mechanisms to MNIST test set  
 297 examples. The first row depicts the original test set examples input to the mechanism, while the  
 298 remaining rows each depict the corresponding outputs from a mechanism trained with different  
 299 values for  $\lambda$ . From the second to last rows, the value of  $\lambda$  is decreased (from 35 to 8), reducing the  
 300 emphasis on minimizing distortion. We see that the outputs start from accurate reconstructions and  
 301 become progressively more distorted while the digit becomes more difficult to correctly recognize  
 302 as  $\lambda$  decreases. The left side of Figure 5 shows the results with the standard PPAN formulation,  
 303 trained via (7) with  $\gamma = 0$ , where we see that the mechanism seems to learn to minimize distortion  
 304 while rendering the digit unrecognizable, which in some cases results in an output that resembles a  
 305 different digit. The right side of Figure 5 shows the results for the second set of experiments when the  
 306 additional discriminator network is introduced, which is jointly trained via (7) with  $\gamma = 2$ . There we  
 307 see that the additional discriminator network encourages outputs that more cleanly resemble actual  
 308 digits, which required lower values for  $\lambda$  (ranging from 15 to 2) to generate distorted images and  
 309 also led to a more abrupt shift toward rendering a different digit. For both sets of experiments, the  
 310 networks were each alternately updated once per batch (of 100 images) over 50 epochs of the 60K  
 311 MNIST training set images.

## 312 4 Conclusion

313 In this work we introduced and developed a practical, data-driven method for optimizing privacy-  
 314 preserving data release mechanisms within the well-established information-theoretic framework. The  
 315 key to this approach is the application of adversarially-trained neural networks, where the mechanism  
 316 is realized as a randomized network, and a second network acts as a privacy adversary that attempts to  
 317 recover sensitive information. By estimating the posterior distribution of the sensitive variable given  
 318 the released data, the adversarial network enables a variational approximation of mutual information.  
 319 This allows our method to approach the information-theoretically optimal privacy-utility tradeoffs,  
 320 which we demonstrate in experiments with discrete and continuous synthetic data. We also conducted  
 321 experiments with the MNIST handwritten digits dataset, where we trained a mechanism that trades  
 322 off between minimizing the pixel-level image distortion and concealing the digit.

## 323 References

- 324 [1] D. Barber and F. Agakov. The IM algorithm: A variational approach to information maximiza-  
325 tion. In *Proceedings of the 16th International Conference on Neural Information Processing*  
326 *Systems*, NIPS'03, pages 201–208, Cambridge, MA, USA, 2003. MIT Press.
- 327 [2] Y. O. Basciftci, Y. Wang, and P. Ishwar. On privacy-utility tradeoffs for constrained data release  
328 mechanisms. In *Information Theory and Applications Workshop*, Feb. 2016.
- 329 [3] F. P. Calmon and N. Fawaz. Privacy against statistical inference. In *Allerton Conf. on Comm.,*  
330 *Ctrl., and Comp.*, pages 1401–1408, 2012.
- 331 [4] X. Chen, X. Chen, Y. Duan, R. Houthoof, J. Schulman, I. Sutskever, and P. Abbeel. Infogan:  
332 Interpretable representation learning by information maximizing generative adversarial nets. In  
333 D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural*  
334 *Information Processing Systems 29*, pages 2172–2180. Curran Associates, Inc., 2016.
- 335 [5] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private  
336 data analysis. In *Theory of Cryptography*, pages 265–284. Springer, 2006.
- 337 [6] H. Edwards and A. J. Storkey. Censoring representations with an adversary. In *Proceedings of*  
338 *the International Conference on Learning Representations (ICLR)*, 2015.
- 339 [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and  
340 Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*,  
341 pages 2672–2680, 2014.
- 342 [8] J. Hamm. Enhancing utility and privacy with noisy minimax filters. In *2017 IEEE International*  
343 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6389–6393, March  
344 2017.
- 345 [9] A. Hindupur. The GAN zoo. <https://deephunt.in/the-gan-zoo-79597dc8c347>, 2017.
- 346 [10] C. Huang, P. Kairouz, X. Chen, L. Sankar, and R. Rajagopal. Context-aware generative  
347 adversarial privacy. *Entropy*, 19(12), 2017.
- 348 [11] E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. In  
349 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- 350 [12] D. Kifer and A. Machanavajjhala. No free lunch in data privacy. In *Proceedings of the 2011*  
351 *ACM SIGMOD International Conference on Management of data*, pages 193–204. ACM, 2011.
- 352 [13] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proceedings of the*  
353 *International Conference on Learning Representations (ICLR)*, 2015.
- 354 [14] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *Proceedings of the*  
355 *International Conference on Learning Representations (ICLR)*, 2014.
- 356 [15] N. Li, T. Li, and S. Venkatasubramanian. t-closeness: Privacy beyond k-anonymity and  
357 l-diversity. In *IEEE Intl. Conf. on Data Eng.*, pages 106–115. IEEE, 2007.
- 358 [16] C. Liu, S. Chakraborty, and P. Mittal. Dependence makes you vulnerable: Differential privacy  
359 under dependent tuples. In *Network and Distributed System Security Symposium*, pages 21–24,  
360 2016.
- 361 [17] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian. L-diversity: Privacy  
362 beyond k-anonymity. *ACM Trans. Knowl. Discov. Data*, 1(1), Mar. 2007.
- 363 [18] C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation  
364 of discrete random variables. In *Proceedings of the International Conference on Learning*  
365 *Representations (ICLR)*, 2017.
- 366 [19] A. Makhdoumi and N. Fawaz. Privacy-utility tradeoff under statistical uncertainty. In *Allerton*  
367 *Conf. on Comm., Ctrl., and Comp.*, pages 1627–1634, 2013.

- 368 [20] A. Makhdoumi, S. Salamatian, N. Fawaz, and M. Médard. From the information bottleneck to  
369 the privacy funnel. In *IEEE Information Theory Workshop*, pages 501–505, 2014.
- 370 [21] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *arXiv*  
371 *preprint arXiv:1511.05644*, 2015.
- 372 [22] A. Narayanan and V. Shmatikov. Robust de-anonymization of large sparse datasets. In *IEEE*  
373 *Symp. on Security and Privacy*, pages 111–125. IEEE, 2008.
- 374 [23] B. Póczos and J. Schneider. Nonparametric estimation of conditional information and diver-  
375 gences. In *Proceedings of the Fifteenth International Conference on Artificial Intelligence and*  
376 *Statistics*, volume 22 of *Proceedings of Machine Learning Research*, pages 914–923, 21–23  
377 Apr 2012.
- 378 [24] D. Rebollo-Monedero, J. Forné, and J. Domingo-Ferrer. From t-closeness-like privacy to  
379 postrandomization via information theory. *IEEE Trans. Knowl. Data Eng.*, 22(11):1623–1636,  
380 2010.
- 381 [25] L. Sankar, S. R. Rajagopalan, and H. V. Poor. Utility-privacy tradeoffs in databases: An  
382 information-theoretic approach. *IEEE Trans. on Information Forensics and Security*, 8(6):838–  
383 852, 2013.
- 384 [26] L. Sweeney. Simple demographics often identify people uniquely. *Carnegie Mellon University,*  
385 *Data Privacy Working Paper*, 2000.
- 386 [27] L. Sweeney. k-anonymity: A model for protecting privacy. *Intl. Journal of Uncertainty,*  
387 *Fuzziness and Knowledge-Based Systems*, 10(5):557–570, 2002.
- 388 [28] T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2 edition,  
389 2012.
- 390 [29] N. Tishby, F. C. Pereira, and W. Bialek. The information bottleneck method. In *Allerton Conf.*  
391 *on Comm., Ctrl., and Comp.*, pages 368—377, 1999.
- 392 [30] S. Tokui, K. Oono, S. Hido, and J. Clayton. Chainer: a next-generation open source framework  
393 for deep learning. In *Proceedings of Workshop on Machine Learning Systems (LearningSys) in*  
394 *The Twenty-ninth Annual Conference on Neural Information Processing Systems (NIPS)*, 2015.
- 395 [31] Y. Wang, Y. O. Basciftci, and P. Ishwar. Privacy-utility tradeoffs under constrained data release  
396 mechanisms. *arXiv preprint arXiv:1710.09295*, 2017.
- 397 [32] H. Yamamoto. A source coding problem for sources with additional outputs to keep secret from  
398 the receiver or wiretappers. *IEEE Trans. on Information Theory*, 29(6):918–923, 1983.

Table 1: The models used for obtaining synthetic training and test datasets in our experiments.

Case	Attribute Model	Observation	Distortion Metric
Discrete, Sec. C	$(X, Y)$ symmetric pair for $m = 10, p = 0.4$ , see (8)	$W = Y$ and $W = (X, Y)$	$\Pr[Y \neq Z]$
Continuous, Sec. 3.1	$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} I_5 & \text{diag}(\rho) \\ \text{diag}(\rho) & I_5 \end{bmatrix}\right)$ , $\rho = [0.47, 0.24, 0.85, 0.07, 0.66]$	$W = Y$	$\mathbb{E}[\ Y - Z\ ^2]$
Continuous, Sec. E	$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} 1 & 0.85 \\ 0.85 & 1 \end{bmatrix}\right)$	$W = Y$ and $W = (X, Y)$	$\mathbb{E}[(Y - Z)^2]$
Continuous, Sec. G	$X = Y \sim \mathcal{N}\left(\mathbf{0}, \text{diag}(\sigma^2)\right)$ , $\sigma^2 = [0.47, 0.24, 0.85, 0.07, 0.66]$	$W = X = Y$	$\mathbb{E}[\ Y - Z\ ^2]$

## 399 A Mutual Information Utility

400 While we focused on expected distortion to measure (dis)utility, our framework can be adapted to  
 401 other general utility measures, for example, mutual information between the useful information  
 402 and the released data. The conditional entropy  $h(Y|Z)$  is an alternative measure for distortion,  
 403 which corresponds to the utility objective of maximizing the mutual information  $I(Y; Z)$ , since  
 404  $h(Y)$  is fixed. When  $h(Y|Z)$  is used as the distortion measure in a scenario where the observation  
 405  $W = X$ , the privacy-utility tradeoff optimization problem, as described in Section 2.1, becomes  
 406 equivalent to the *Information Bottleneck* problem considered in [29]. In other scenarios where the  
 407 observation  $W = Y$ , this problem becomes the *Privacy Funnel* problem introduced by [20]. The  
 408 formulation of (3) can be modified to address conditional entropy distortion by introducing another  
 409 variational posterior  $Q_{Y|Z}$  and using the following optimization, which applies a second variational  
 410 approximation of mutual information,

$$\min_{P_{Z|W}, Q_{Y|Z}} \max_{Q_{X|Z}} \mathbb{E}[\log Q_{X|Z}(X|Z)] - \lambda \mathbb{E}[\log Q_{Y|Z}(Y|Z)],$$

411 where the expectations are with respect to  $(W, X, Y, Z) \sim P_{W,X,Y}P_{Z|W}$ , and the parameter  $\lambda > 0$   
 412 can be adjusted to obtain various points along the optimal tradeoff curve. In a similar fashion to the  
 413 approach in Section 2.2, this optimization problem can be practically addressed via the training of  
 414 three neural networks, which respectively parameterize the mechanism  $P_{Z|W}$  and the two variational  
 415 posteriors  $Q_{X|Z}$  and  $Q_{Y|Z}$ .

## 416 B Sampling from a Gaussian Mixture Model (GMM)

417 The technique of sampling from a multivariate Gaussian described in Section 2.3.2 can be extended  
 418 to GMMs as follows. The release mechanism can be realized with a neural network  $f_\theta$  that produces  
 419 the set of parameters  $\{(\boldsymbol{\mu}_{l,i}, \mathbf{A}_{l,i}, \pi_{l,i})\}_{l=1}^m = f_\theta(w_i)$ , where  $\pi_{l,i}$  are the mixture weights. We then  
 420 sample  $z_{l,i} = \mathbf{A}_{l,i}\mathbf{u}_{l,i} + \boldsymbol{\mu}_{l,i}$  for each component distribution of the GMM, and compute the loss  
 421 terms via

$$\mathcal{L}_{\text{GMM}}^i(\theta, \phi) := \sum_{l=1}^m \pi_{l,i} (\log Q_\phi(x_i|z_{l,i}) + \lambda d(y_i, z_{l,i})),$$

422 which combines the Gaussian sampling reparameterization trick with a direct expectation over the  
 423 mixture component selection.

## 424 C Discrete Synthetic Data

425 In our experiments with discrete data, we will consider two observation models, full data (where  
 426  $W = (X, Y)$ ) and useful data only (where  $W = Y$ ). We use a toy distribution for the attributes  
 427 for which the theoretically optimal privacy-utility tradeoffs have been analytically derived in [31],  
 428 using probability of error as the distortion metric, i.e.,  $\mathbb{E}[\mathbf{1}(Y \neq Z)] = \Pr[Y \neq Z]$ . Specifically,  
 429 we consider sensitive and useful attributes that are distributed over the finite alphabets  $\mathcal{X} = \mathcal{Y} =$   
 430  $\{0, \dots, m-1\}$ , with  $m \geq 2$  and parameter  $p \in [0, 1]$ , according to the *symmetric pair* distribution

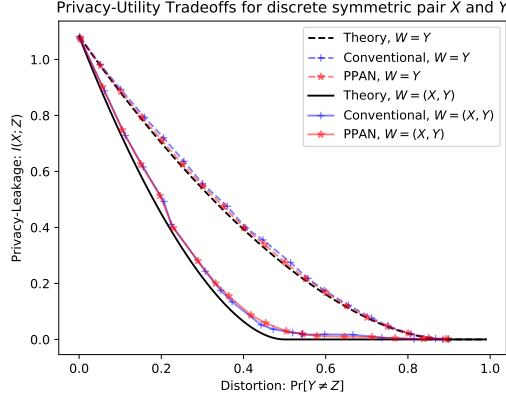


Figure 6: Comparison of PPAN performance against the conventional model estimation approach of [19] and the theoretical optimum, for two observation scenarios: full data observed, i.e.  $W = (X, Y)$ , shown by solid lines, and only useful attribute observed, i.e.  $W = Y$ , shown by dashed lines.

431 given by

$$P_{X,Y}(x, y) = \begin{cases} \frac{1-p}{m}, & \text{if } x = y, \\ \frac{p}{m(m-1)}, & \text{otherwise.} \end{cases} \quad (8)$$

### 432 C.1 Network Architecture and Evaluation

433 As mentioned in Section 2.3.1, the network architecture for the release mechanism and adversary can  
 434 be reduced to a bare minimum when all of the variables are finite-alphabet. Each network simply  
 435 applies a single linear transformation (with no bias term) on the one-hot encoded input, followed by  
 436 the softmax operation to yield a stochastic vector. The mechanism network takes as input  $w$  encoded  
 437 as a one-hot column vector  $\mathbf{w}$  and outputs  $P_\theta(\cdot|w) = \text{softmax}(\mathbf{G}\mathbf{w})$ , where the network parameters  
 438  $\theta = \mathbf{G}$  are the entries of a  $|\mathcal{Z}| \times |\mathcal{W}|$  real matrix. Note that applying the softmax operation to each  
 439 column of  $\mathbf{G}$  produces the conditional distribution  $P_{Z|W}$  describing the mechanism. Similarly, the  
 440 attacker network is realized as  $Q_\phi(\cdot|z) = \text{softmax}(\mathbf{A}\mathbf{z})$ , where  $\mathbf{z}$  is the one-hot encoding of  $z$ , and  
 441 the network parameters  $\phi = \mathbf{A}$  are entries of a  $|\mathcal{X}| \times |\mathcal{Z}|$  real matrix. We optimize these networks  
 442 according to (5), using the penalty term modification of the loss terms in (6) as given by

$$\mathcal{L}_{\text{disc}}^i(\theta, \phi) := \sum_{z \in \mathcal{Z}} P_\theta(z|w_i) (\log Q_\phi(x_i|z) + \lambda \max(0, d(y_i, z) - \delta)^2).$$

443 We use  $\lambda = 500$  in these experiments.

444 In Figure 6, we compare the results of PPAN against the theoretical baselines<sup>3</sup> given by [31], as  
 445 well as against a conventional approach suggested by [19], where the joint distribution of  $P_{W,X,Y}$   
 446 is estimated from the training data and then used in the convex optimization of (1). We can see that  
 447 the PPAN mechanism learns a data release distribution that has close to optimal privacy leakage  
 448 for a wide range of distortion values. We used 1000 training samples generated according to the  
 449 symmetric pair distribution in (8) with  $m = 10$  and  $p = 0.4$ . The PPAN networks were trained for  
 450 2500 epochs (for the full data observation case) with a minibatch size of 100, with each network  
 451 alternatingly updated once per iteration. For the useful data only observation case, 2000 epochs  
 452 were used. For evaluating both the PPAN and conventional approaches, we computed the mutual  
 453 information and probability of error from the joint distribution that combines the optimized  $P_{Z|W}$   
 454 with the true  $P_{X,Y,W}$ .

<sup>3</sup>For convenience, the specific expressions are reproduced in Section D of the supplementary material.

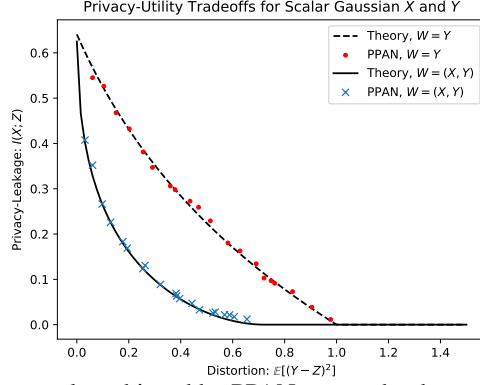


Figure 7: Comparison of the results achieved by PPAN versus the theoretical optimum tradeoff curve, with jointly Gaussian scalar  $(X, Y)$ , for the useful data only (i.e.,  $W = Y$ ) and the full data (i.e.,  $W = (X, Y)$ ) observation models.

## 455 D Theoretically Optimal Privacy-Utility Tradeoffs for Symmetric Pair 456 Distribution

457 The mutual information of the symmetric pair distribution (see (8)) is given by [31] as

$$I(X; Y) = \log m - p \log(m-1) - h_2(p) =: r_m(p),$$

458 where  $h_2(p) := -p \log p - (1-p) \log(1-p)$  is the binary entropy function, and for convenience in  
459 later discussion, we define  $r_m(p)$  as a function of the distribution parameters  $m$  and  $p$ .

460 For sensitive and useful attributes jointly distributed according to the symmetric pair distribution,  
461 the theoretically optimal privacy-utility tradeoffs, as defined by (1), are analytically derived in [31]  
462 for several data observation models, while using probability of error as the distortion metric, i.e.,  
463  $\mathbb{E}[\mathbf{1}(Y \neq Z)] = \Pr[Y \neq Z]$ . In one case, when the observation is the full data, i.e.,  $W = (X, Y)$ ,  
464 the optimal mutual information privacy-leakage as a function of the distortion (probability of error)  
465 limit  $\delta \in [0, 1]$  is given by

$$I_{W=(X,Y)}^*(\delta) = \begin{cases} r_m(p + \delta), & \text{if } \delta \leq 1 - \frac{1}{m} - p, \\ r_m(p - \delta), & \text{if } \delta \leq p - (1 - \frac{1}{m}), \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

466 In another case, when the observation is only the useful attribute, i.e.,  $W = Y$ , the optimal privacy-  
467 leakage as a function  $\delta \in [0, 1]$  is given by

$$I_{W=Y}^*(\delta) = \begin{cases} r_m\left(p + \delta \left(1 - \frac{pm}{m-1}\right)\right), & \text{if } \delta < 1 - \frac{1}{m}, \\ 0, & \text{otherwise.} \end{cases} \quad (10)$$

## 468 E Scalar Gaussian Attributes

469 Consider jointly Gaussian sensitive and useful attributes such that  $\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.85 \\ 0.85 & 1 \end{bmatrix}\right)$ . We  
470 analyze two different observation models here:  $W = Y$ , called useful data only (UD) and  $W =$   
471  $(X, Y)$ , called full data (FD). The distortion metric is the mean squared error between the release and  
472 the useful attribute. The values of the multipliers chosen are:  $\lambda^{\text{UD}} = 10$  and  $\lambda^{\text{FD}} = 50$ . In each case,  
473 we run experiments for 20 different values of the target distortion with  $\delta^{\text{UD}} \in [0, 1]$  and  $\delta^{\text{FD}} \in [0, 0.8]$ .  
474 The privacy-leakage and distortion values returned by the PPAN mechanism on the test set are plotted  
475 along with the optimal tradeoff curves (from Propositions 1 and 3 of the supplementary material) in  
476 Figure 7. In both the observation models, we observe that the PPAN mechanism generates releases  
477 that have nearly optimal privacy-leakage over a range of distortion values.

## 478 F Estimating Mutual Information Leakage

479 Distortion caused by a release is estimated by the empirical mean squared error with respect to  
480 the testing samples. However, estimating mutual information to evaluate privacy leakage is less

481 straightforward since the joint distribution  $P_{X,Z}$  as realized by the optimized mechanism is not  
482 available explicitly. Since for these experiments, the optimal release  $Z$  is jointly Gaussian with  
483  $X$  (as we show in Section H of supplementary material), we estimate  $I(X; Z)$  via a Gaussian  
484 approximation. Specifically, we use the expression for the mutual information of jointly Gaussian  
485 random vectors and replace all covariance matrices that appear there by their empirical counterparts,  
486 i.e.,  $\hat{I}(X; Z) = 0.5 \log(\det(\hat{\Sigma}_X) / \det(\hat{\Sigma}_{X|Z}))$ , where  $\hat{\Sigma}_{X|Z} := \hat{\Sigma}_X - \hat{\Sigma}_{X,Z} \hat{\Sigma}_Z^+ \hat{\Sigma}_{X,Z}^T$  and  $\hat{\Sigma}_X$   
487 denotes the empirical self covariance matrix of  $X$ ,  $\hat{\Sigma}_Z^+$  denotes the pseudoinverse of the empirical  
488 self covariance matrix of  $Z$ , and  $\hat{\Sigma}_{X,Z}$  denotes their empirical cross covariance matrix. This  
489 underestimates the true mutual information leakage since

$$\begin{aligned} I(X; Z) &= h(X) - h(X - \hat{\mathbb{E}}[X|Z]|Z) \\ &\geq h(X) - h(X - \hat{\mathbb{E}}[X|Z]) = \hat{I}(X; Z), \end{aligned}$$

490 where  $\hat{\mathbb{E}}[X|Z]$  is the linear MMSE estimate of  $X$  as a function of  $Z$ . We use this estimate only for  
491 its simplicity, and one could use other non-parametric estimates of mutual information [23].

## 492 G Rate Distortion

493 We can apply the PPAN framework to the problem of computing the minimum required rate of a  
494 code that describes a multivariate source  $X$  to within a target value of expected distortion. This is  
495 a standard problem in information theory when the source distribution is known, for example, see  
496 Chapter 10 of [28]. However, the PPAN framework can be used to empirically approximate the  
497 rate-distortion curve from i.i.d. samples of the source without knowledge of the source distribution.  
498 The computation of the rate-distortion function can be viewed as a degenerate case of the PPAN  
499 framework with  $W = X = Y$ , i.e., the sensitive and useful attributes are the same and the observed  
500 dataset is the attribute. The release  $Z$  corresponds to an estimate  $\hat{X}$  with expected distortion less than  
501 a target level while retaining as much expected uncertainty about  $X$  as possible.

502 We illustrate the PPAN approach using a Gaussian source  $X \in \mathbb{R}^5$  and mean squared error distortion.  
503 For the experiment, we choose the attribute model  $X \sim \mathcal{N}(\mathbf{0}, \text{diag}(0.47, 0.24, 0.85, 0.07, 0.66))$  and  
504 the value  $\lambda = 500$ . We run the experiment for 20 different values of the target distortion, linearly  
505 spaced between 0 to 2.5. The inputs to the adversarial network are realizations of the attributes  
506 and seed noise. The seed noise is chosen to be a random vector of length 8 with each component  
507 i.i.d. Uniform $[-1, 1]$ . The network architecture and values of other hyperparameters are the same as  
508 those used for multivariate Gaussian attributes in Section 3.1. Using the learned parameters  $\theta^*$ , the  
509 mechanism network generates a release as  $Z = f_{\theta^*}(W, U) = f_{\theta^*}(X, U)$ . The distortion is estimated  
510 by the empirical mean squared error of the release with respect to the training samples. The privacy  
511 loss is quantified by the estimate  $\hat{I}(X; Z)$  as described in Section F.

512 The optimal privacy-utility tradeoff (or, rate-distortion) curve is  $R(D) =$   
513  $\min_{P(Z|X) : \mathbb{E}\|X-Z\|^2 \leq D} I(X; Z) = \sum_{j=1}^5 \max\{0, 0.5 \log((\sigma[j])^2 / D_j)\}$  [28], where  $\sigma^2$   
514 are the true variance parameters of the attribute distribution and  $\sum_{j=1}^5 D_j = D$ . The values of  $D_j$   
515 for each component is obtained using the Karush-Kuhn-Tucker (KKT) conditions for the constrained  
516 optimization problem, the solution of which is a standard waterfilling procedure.

517 We plot the (privacy-leakage, utility loss) pairs returned by the PPAN mechanism along with the  
518 optimal tradeoff curve in Figure 8. One can see that the operating points attained by the PPAN  
519 mechanism are very close to the theoretical optimum tradeoff for a wide range of target distortion  
520 values.

## 521 H Optimum Privacy Utility Tradeoff for Gaussian Attributes

522 In Section 3 we compare the (privacy, distortion) pairs achieved by the model-agnostic PPAN  
523 mechanism with the optimal model-aware privacy-utility tradeoff curve. For jointly Gaussian  
524 attributes and mean squared error distortion, we can obtain, in some cases, analytical expressions  
525 for the optimal tradeoff curve as described below. Some of the steps in the proofs use bounding  
526 techniques from rate-distortion theory, which is to be expected given the tractability of the Gaussian  
527 model and the choice of mutual information and mean squared error as the privacy and utility metrics  
528 respectively.



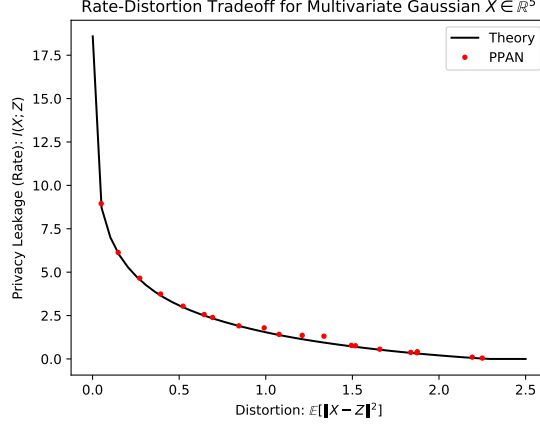


Figure 8: Comparison of results obtained by PPAN versus the the optimal rate-distortion curve, for the rate-distortion problem where  $X = Y = W$  is multivariate Gaussian.

529 **Proposition 1.** (*Useful Data only: Scalar Gaussian with mean squared error*) In problem (1), let  
 530  $X, Y$  be jointly Gaussian scalars with zero means  $\mu_X = \mu_Y = 0$ , variances  $\sigma_X^2, \sigma_Y^2$  respectively,  
 531 and correlation coefficient  $\rho \in [-1, 1]$ . Let mean squared error be the distortion measure. If the  
 532 observation  $W = Y$  (Useful Data only observation model), then the optimal release  $Z$  corresponding  
 533 to

$$\min_{P_{Z|Y}} I(X; Z), \quad \text{such that } \mathbb{E}(Y - Z)^2 \leq \delta \quad \text{and} \quad X \leftrightarrow Y \leftrightarrow Z \quad (11)$$

534 is given by

$$Z = \begin{cases} 0, & \text{if } \delta \geq \sigma_Y^2 \\ (1 - \delta/\sigma_Y^2)Y + U, & \text{if } \delta < \sigma_Y^2 \end{cases}$$

535 where  $U \perp (X, Y)$  and  $U \sim \mathcal{N}(0, \delta(1 - \delta/\sigma_Y^2))$ . The mutual information leakage caused by  
 536 releasing  $Z$  is

$$I(X; Z) = \max \left\{ 0, \frac{1}{2} \log \left( \frac{1}{1 - \rho^2 + \rho^2 \delta / \sigma_Y^2} \right) \right\}.$$

537 The result of Proposition 1 is known in the existing literature, e.g., [24] (see Eq. 8) and [25] (see  
 538 Example 2). For completeness, we present the proof of this result in Section I.1. The theoretical  
 539 tradeoff curve in Figure 7 was obtained using the expressions in Proposition 1.

540 The case of Useful Data only observation model for jointly Gaussian *vector* attributes and mean  
 541 squared error is also considered in [24], where they provide a numerical procedure to evaluate the  
 542 tradeoff curve. Here, we focus on a special case where we can compute the solution analytically.

543 Consider the generalization to vector variables of problem (11)

$$\min_{P_{Z|Y}} I(X; Z) \text{ such that } \mathbb{E}(Y - Z)^T(Y - Z) \leq \delta \text{ and } X \leftrightarrow Y \leftrightarrow Z. \quad (12)$$

544 Let  $X, Y$  be jointly Gaussian vectors of dimensions  $m$  and  $n$  respectively. We assume that  $X, Y$  have  
 545 zero means  $\mu_X = \mu_Y = 0$  and non-singular covariance matrices  $\Sigma_X, \Sigma_Y \succ 0$ . Let  $\Sigma_{XY}$  denote  
 546 the cross-covariance matrix and  $P := \Sigma_X^{-\frac{1}{2}} \Sigma_{XY} \Sigma_Y^{-\frac{1}{2}}$  the normalized cross-covariance matrix with  
 547 singular value decomposition  $P = U_P \Lambda_P V_P^T$ . We assume that all singular values of  $P$ , denoted by  
 548  $\rho'_i, i = 1, \dots, \min\{m, n\}$ , are strictly positive. If

$$X' := U_P^T \Sigma_X^{-\frac{1}{2}} X, \quad Y' := V_P^T \Sigma_Y^{-\frac{1}{2}} Y, \quad \text{and } Z' := V_P^T \Sigma_Y^{-\frac{1}{2}} Z$$

549 denote reparameterized variables, then  $X', Y'$  are zero-mean, jointly Gaussian, with identity co-  
 550 variance matrices  $I_m, I_n$  respectively and  $m \times n$  diagonal cross-covariance matrix  $\Lambda_P$ . Since the

551 transformation from  $(X, Z)$  to  $(X', Z')$  is invertible,  $I(X'; Z') = I(X; Z)$ . The mean squared error  
 552 between  $Y$  and  $Z$ :

$$\mathbb{E}[(Y - Z)^\top(Y - Z)] = \mathbb{E}[(Y' - Z')^\top(V_P^\top \Sigma_Y V_P)(Y' - Z')].$$

553 For the special case when  $V_P^\top \Sigma_Y V_P = cI_n$  for some  $c > 0$ , the vector problem (12) reduces to the  
 554 following problem

$$\min_{P_{Z'|Y'}} I(X'; Z') \text{ such that } \mathbb{E}(Y' - Z')^\top(Y' - Z') \leq \delta/c \text{ and } X' \leftrightarrow Y' \leftrightarrow Z'. \quad (13)$$

555 **Proposition 2.** *If  $\begin{bmatrix} X' \\ Y' \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0}_m \\ \mathbf{0}_n \end{bmatrix}, \begin{bmatrix} I_m & \Lambda_P \\ \Lambda_P^\top & I_n \end{bmatrix}\right)$ , then the minimizer of (13) is given by*

$$Z'_i = (1 - \delta'_i)Y'_i + U_i, i = 1, \dots, \min\{m, n\},$$

556 where  $(U_1, \dots, U_{\min\{m, n\}}) \perp (X', Y')$  and for all  $i$ ,  $U_i \sim \mathcal{N}(0, \delta'_i(1 - \delta'_i))$ ,  $\delta'_i := \min\{1, t -$   
 557  $(\rho'_i)^{-2} - 1\}$ , where  $\rho'_i > 0$  denotes the  $i$ -th main diagonal entry of  $\Lambda_P$ , and the value of parameter  
 558  $t$  can be found by the equation  $\sum_i \delta'_i = \delta/c$ . The mutual information between the release and the  
 559 sensitive attribute is  $I(X', Z') = \sum_{i=1}^{\min\{m, n\}} \max\{0, -0.5 \log(1 - \rho'_i{}^2 + \rho'_i{}^2 \delta'_i)\}$ .

560 The proof of the above proposition is given in Section I.2. We evaluate the above parametric  
 561 expression for various values of  $\delta$  in order to obtain the theoretical tradeoff curves in Figure 3.

562 For the case of full data observation, we have the following result.

563 **Proposition 3.** *(Full Data: Scalar Gaussian with mean squared error) In problem (1), let  $X, Y$  be*  
 564 *jointly Gaussian scalars with zero means, unit variances, and correlation coefficient  $\rho \in [0, 1]$ . Let*  
 565 *mean squared error be the distortion measure. If the observation  $W = (X, Y)$  (full data observation*  
 566 *model), then the optimal release  $Z$  corresponding to*

$$\min_{P_{Z|X, Y}} I(X; Z), \quad \text{such that } \mathbb{E}(Y - Z)^2 \leq \delta \quad (14)$$

567 is given by

$$Z = (1 - \delta)Y - (X - \rho Y) \sqrt{\frac{\delta(1 - \delta)}{1 - \rho^2}}.$$

568 The mutual information leakage caused by this release is

$$I(X; Z) = \begin{cases} 0, & \text{if } \delta \geq \rho^2 \\ \frac{1}{2} \log \left( \frac{1}{1 - \left( \sqrt{\rho^2(1 - \delta)} - \sqrt{(1 - \rho^2)\delta} \right)^2} \right), & \text{if } \delta < \rho^2. \end{cases}$$

569 The proof of the above proposition is presented in Section I.3. The theoretical tradeoff curve in  
 570 Figure 7 was obtained using the above expression.

## 571 I Proofs of Propositions

### 572 I.1 Proof of Proposition 1

573 *Proof.* We can expand the mutual information term as follows,

$$\begin{aligned} I(X; Z) &= h(X) - h(X|Z), \\ &= h(X) - h(X - \rho\sigma_X Z / \sigma_Y | Z), \\ &\geq h(X) - h(X - \rho\sigma_X Z / \sigma_Y), \end{aligned} \quad (15)$$

$$\geq 0.5 \log 2\pi e \sigma_X^2 - h(\mathcal{N}(0, \mathbb{E}[(X - \rho\sigma_X Z / \sigma_Y)^2])), \quad (16)$$

$$= \frac{1}{2} \log \left( \frac{\sigma_X^2}{\mathbb{E}[(X - \rho\sigma_X Z / \sigma_Y)^2]} \right). \quad (17)$$

574 Inequality (15) is true because conditioning can only reduce entropy and inequality (16) is true since  
575 the zero-mean normal distribution has the maximum entropy for a given value of the second moment.  
576 Let  $T := X - \rho\sigma_X Y/\sigma_Y$ , then  $T$  is jointly Gaussian and we have that

$$\text{Cov} \begin{pmatrix} T \\ Y \end{pmatrix} = \begin{bmatrix} 1 & -\rho\sigma_X/\sigma_Y \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -\rho\sigma_X/\sigma_Y & 1 \end{bmatrix} = \begin{bmatrix} \sigma_X^2(1-\rho^2) & 0 \\ 0 & \sigma_Y^2 \end{bmatrix}.$$

577 Hence,  $T$  is independent of  $Y$ . Since  $X \leftrightarrow Y \leftrightarrow Z$  also forms a Markov chain, we have that  $T$   
578 is conditionally independent of  $Z$  given  $Y$ . Due to the distortion constraint, we can upper bound  
579  $\mathbb{E}[(X - \rho\sigma_X Z/\sigma_Y)^2]$  in the following manner.

$$\begin{aligned} \mathbb{E} \left[ \left( X - \frac{\rho\sigma_X}{\sigma_Y} Z \right)^2 \right] &= \sigma_X^2 + \frac{\rho^2\sigma_X^2}{\sigma_Y^2} \mathbb{E}[Z^2] - 2\frac{\rho\sigma_X}{\sigma_Y} \mathbb{E} \left[ \left( T + \frac{\rho\sigma_X}{\sigma_Y} Y \right) Z \right], \\ &\leq \sigma_X^2 + \frac{\rho^2\sigma_X^2}{\sigma_Y^2} (\delta - \sigma_Y^2 + 2\mathbb{E}[YZ]) - 2\frac{\rho\sigma_X}{\sigma_Y} \left( \frac{\rho\sigma_X}{\sigma_Y} \mathbb{E}[YZ] + \mathbb{E}[TZ] \right), \end{aligned} \quad (18)$$

$$= \sigma_X^2(1-\rho^2) + \rho^2\delta\sigma_X^2/\sigma_Y^2. \quad (19)$$

580 Inequality (18) is true because  $\mathbb{E}[(Y - Z)^2] \leq \delta$ , and equation (19) is true because

$$\mathbb{E}[TZ] = \mathbb{E}_Y[\mathbb{E}_{T,Z|Y}[TZ]] \stackrel{(i)}{=} \mathbb{E}_Y[\mathbb{E}_{T|Y}[T]\mathbb{E}_{Z|Y}[Z]] \stackrel{(ii)}{=} \mathbb{E}_Y[\mathbb{E}_T[T]\mathbb{E}_{Z|Y}[Z]] \stackrel{(iii)}{=} 0,$$

581 where (i) is true because  $T \perp\!\!\!\perp Z | Y$ , (ii) is true because  $T \perp\!\!\!\perp Y$  and (iii) is true because  $T$  has zero  
582 mean. Thus by equations (17) and (19), we get that

$$\min_{X \leftrightarrow Y \leftrightarrow Z, \mathbb{E}[(Y-Z)^2] \leq \delta} I(X; Z) \geq \max \left\{ 0, \frac{1}{2} \log \left( \frac{1}{1-\rho^2 + \rho^2\delta/\sigma_Y^2} \right) \right\}. \quad (20)$$

583 For the choice of  $Z$  as stated in the proposition, we can check that  $X$  and  $Z$  are jointly Gaussian with  
584  $I(X; Z) = 0.5 \log_2(\sigma_X^2/\text{Var}(X|Z))$  and  $\text{Var}(X|Z) = \sigma_X^2(1-\rho^2 + \rho^2\delta/\sigma_Y^2)$ . Thus  $Z$  attains the  
585 lower bound for the privacy-leakage in (20) when  $\delta < \sigma_Y^2$ . Otherwise, the lower bound on mutual  
586 information is 0 and can be attained by  $Z = 0$ .  $\square$

## 587 I.2 Proof of Proposition 2

588 *Proof.* (a) If  $m \leq n$ , then  $(X'_1, Y'_1), \dots, (X'_m, Y'_m), Y'_{m+1}, \dots, Y'_n$  are independent because they are  
589 jointly Gaussian and for all  $i \neq j$ ,  $\text{Cov}(X'_i, X'_j) = \text{Cov}(X'_i, Y'_j) = \text{Cov}(Y'_i, Y'_j) = 0$ . Similarly, if  
590  $m \geq n$ , then  $(X'_1, Y'_1), \dots, (X'_n, Y'_n), X'_{n+1}, \dots, X'_m$  are independent.

591 In the following, we use the following well-known properties of mutual information and conditional  
592 mutual information. For any three random variables  $A, B, C$ , (b)  $I(A; B) = I(B; A) \geq 0$ , (c)  
593  $I(A; B) = 0 \Leftrightarrow A \perp\!\!\!\perp B$ , (d)  $I(A; C|B) \geq 0$ , (e)  $I(A; C|B) = 0 \Leftrightarrow (A \perp\!\!\!\perp C)|B$ , (f)  $I(A; B, C) =$   
594  $I(A; B) + I(A; C|B)$  so that  $I(A; B, C) \geq I(A; B)$ .

595 If  $m \leq n$ , then  $I(X'; Z') = I(X'_1, \dots, X'_m; Z') \stackrel{(f)}{=} \sum_{i=1}^m I(X_i; Z'|X_1, \dots, X_{i-1}) \stackrel{(f,a,c)}{=} \sum_{i=1}^m I(X_i; Z', X_1, \dots, X_{i-1}) \geq \sum_{i=1}^m I(X_i; Z') \geq \sum_{i=1}^m I(X_i; Z_i)$ . Similarly, if  $m \geq n$ ,  
597  $I(X'; Z') \geq \sum_{i=1}^m I(X_i; Z') \geq \sum_{i=1}^n I(X_i; Z_i)$ . Thus in general,

$$I(X'; Z') \geq \sum_{i=1}^{\min\{m,n\}} I(X_i; Z_i).$$

598 The distortion constraint in (13) implies that  $\sum_{i=1}^{\min\{m,n\}} \mathbb{E}[(Y'_i - Z'_i)^2] \leq \delta/c$ . Thus the optimal  
599 function value in (13) is lower bounded by the optimum value of the following problem.

$$\begin{aligned} &\min_{P_{Z'|Y'}} \sum_{i=1}^{\min\{m,n\}} I(X'_i; Z'_i) \quad \text{s.t.} \quad \sum_{i=1}^{\min\{m,n\}} \mathbb{E}[(Y'_i - Z'_i)^2] \leq \delta/c \quad \text{and} \quad X' \leftrightarrow Y' \leftrightarrow Z'. \\ &\equiv \min_{\sum \delta'_i \leq \delta/c, P_{Z'|Y'}} \sum_{i=1}^{\min\{m,n\}} I(X'_i; Z'_i) \quad \text{s.t.} \quad \forall i, \mathbb{E}[(Y'_i - Z'_i)^2] \leq \delta'_i \quad \text{and} \quad X' \leftrightarrow Y' \leftrightarrow Z'. \end{aligned} \quad (21)$$

600 Let  $Y_{\sim i} := \{Y_1, \dots, Y_n\} \setminus Y_i$ . Since  $X' - Y' - Z'$  forms a Markov chain, if  $m \leq n$ , we have 0  $\stackrel{(e)}{=} I(X'; Z'|Y') = I(X_1, \dots, X_m; Z'|Y_1, \dots, Y_n) \stackrel{(f)}{=} \sum_{i=1}^m I(X_i; Z'|Y_1, \dots, Y_n, X_1, \dots, X_{i-1}) \stackrel{(f,a,c)}{=} \sum_{i=1}^m I(X_i; Z', X_1, \dots, X_{i-1}, Y_{\sim i}|Y_i) \stackrel{(f)}{\geq} \sum_{i=1}^m I(X_i; Z_i|Y_i) \stackrel{(d)}{\geq} 0$ . Thus,

$$0 \geq \sum_{i=1}^m I(X_i; Z_i|Y_i) \geq 0.$$

603 A similar expression can be derived for the case  $m \geq n$ . In general, for all  $i = 1, \dots, \min\{m, n\}$ ,  
604  $X_i \leftrightarrow Y_i \leftrightarrow Z_i$  forms a Markov chain. Thus for output perturbation, the Markov constraint on the  
605 vectors passes through as a Markov constraint on the individual components of the variables. We can  
606 therefore rewrite problem (21) as follows,

$$\min_{\sum \delta'_i \leq \delta/c} \sum_{i=1}^{\min\{m,n\}} \min_{P_{Z'_i|Y'_i}} I(X'_i; Z'_i), \quad \text{s.t.} \quad \forall i, \mathbb{E}[(Y'_i - Z'_i)^2] \leq \delta'_i \quad \text{and} \quad X'_i \leftrightarrow Y'_i \leftrightarrow Z'_i.$$

607 For each  $i$ , the solution to the inner constrained minimization problem is given by Proposition 1.  
608 Plugging in the solution we arrive at the following constrained convex minimization problem

$$\min_{\sum \delta'_i \leq \delta/c} \sum_{i=1}^{\min\{m,n\}} \max \left\{ 0, \frac{1}{2} \log \left( \frac{1}{1 - \rho_i'^2 + \rho_i'^2 \delta'_i} \right) \right\}$$

609 where  $\rho'_i = \mathbb{E}[X'_i Y'_i]$  and we have used the expression for the optimal privacy-leakage in the scalar  
610 case, i.e., Eq. (20) with  $\sigma_Y^2 = 1$ . The Lagrangian of the above convex program has the following  
611 form

$$\mathcal{L}(\delta', \eta, \zeta) := \sum_{i=1}^{\min\{m,n\}} \frac{1}{2} \log \left( \frac{1}{1 - \rho_i'^2 + \rho_i'^2 \delta'_i} \right) + \zeta \left( \sum_{i=1}^{\min\{m,n\}} \delta'_i - \frac{\delta}{c} \right) + \sum_{i=1}^{\min\{m,n\}} \eta_i (\delta'_i - 1),$$

612 where  $\delta' = (\delta'_1, \dots, \delta'_{\min\{m,n\}})$ ,  $\eta = (\eta_1, \dots, \eta_{\min\{m,n\}})$ , and  $\zeta$  and all the  $\eta_i$ 's are non-negative  
613 Lagrange multipliers. Here, the non-negativity condition associated with  $\max\{0, \cdot\}$  has been sub-  
614 summed by requirement that  $\delta'_i \leq 1$  for all  $i$ . The Karush-Kuhn-Tucker (KKT) conditions for optimality  
615 are as follows,

$$\sum_{i=1}^{\min\{m,n\}} \delta'_i = \frac{\delta}{c}, \quad \forall i, 0 \leq \delta'_i \leq 1, \eta_i \geq 0, \eta_i (\delta'_i - 1) = 0, \frac{\partial \mathcal{L}}{\partial \delta'_i} = 0 \Rightarrow \eta_i = \frac{1}{2(\delta'_i - 1 + \rho_i'^{-2})} - \zeta.$$

616 This implies that if for any  $i$ ,

$$(2\zeta)^{-1} > \rho_i'^{-2} \Leftrightarrow \eta_i > 0, \text{ then } \delta'_i = 1, \text{ otherwise } \delta'_i = (2\zeta)^{-1} - (\rho_i'^{-2} - 1).$$

617 The value of  $(2\zeta)^{-1}$  can be found by the equation

$$\sum_{i=1}^{\min\{m,n\}} \max \left\{ 0, \min\{1, (2\zeta)^{-1} - (\rho_i'^{-2} - 1)\} \right\} = \frac{\delta}{c},$$

618 which is a modified water-filling solution. Based on the value of  $\delta'_i$ , we can construct a  $Z'_i$  that attains  
619 the lower bound on the mutual information by setting  $\sigma_Y^2 = 1$  in the results of Proposition (1).  $\square$

### 620 I.3 Proof of Proposition 3

621 *Proof.* In this proposition,  $X$  and  $Y$  are assumed to be jointly Gaussian with zero means, unit  
622 variances, and correlation coefficient  $\rho \in [0, 1]$ . Consider the Linear Minimum Mean Squared Error  
623 (LMMSE) estimate of  $X$  given  $Z$  denoted as  $\widehat{\mathbb{E}}[X|Z] = \mathbb{E}[XZ|Z]/\mathbb{E}[Z^2]$ . Then, similar to the proof  
624 of Proposition 1, we can expand the mutual information in the following manner.

$$\begin{aligned} I(X; Z) &= h(X) - h(X|Z) = h(X) - h(X - \widehat{\mathbb{E}}[X|Z]|Z) \\ &\geq h(X) - h(X - \widehat{\mathbb{E}}[X|Z]) \geq h(X) - h\left(\mathcal{N}\left(0, \mathbb{E}\left[\left(X - \widehat{\mathbb{E}}[X|Z]\right)^2\right]\right)\right) \\ &= -\frac{1}{2} \log \left( 1 - \frac{(\mathbb{E}[XZ])^2}{\mathbb{E}[Z^2]} \right), \end{aligned}$$

625 where in writing the last equality we have used the fact that  $\mathbb{E} \left[ \widehat{\mathbb{E}}[X|Z](X - \widehat{\mathbb{E}}[X|Z]) \right] = 0$  by the  
 626 orthogonality principle of least squares estimation. Thus, we have that

$$\min_{\mathbb{E}[(Y-Z)^2] \leq \delta} I(X; Z) \geq -\frac{1}{2} \log \left[ 1 - \min_{\mathbb{E}[(Y-Z)^2] \leq \delta} \frac{(\mathbb{E}[XZ])^2}{\mathbb{E}[Z^2]} \right] \quad (22)$$

627 Below, we focus on the minimization problem on the right side of Eq. (22). It will turn out that for  
 628 the minimizing  $Z^*$ , we will have equality in Eq. (22). In what follows, it is helpful to think of the  
 629 random variables  $X, Y, Z$  as vectors in the vector space  $\mathcal{L}_2$  of all random variables with finite second  
 630 moments over the underlying probability space. We will emphasize the vector nature by denoting  
 631 the random variables  $X, Y, Z$  by their corresponding bold lowercase letters  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  respectively. The  
 632 expectation operator on the product of two random variables in  $\mathcal{L}_2$  is an inner product, and hence we  
 633 can write the optimization problem of interest as follows,

$$\mathbf{z}^* := \arg \min_{\mathbf{z}: \|\mathbf{z}-\mathbf{y}\|^2 \leq \delta} \frac{|\langle \mathbf{x}, \mathbf{z} \rangle|^2}{\|\mathbf{z}\|^2} = \arg \min_{\mathbf{z}: \|\mathbf{z}-\mathbf{y}\|^2 \leq \delta} \left| \langle \mathbf{x}, \frac{\mathbf{z}}{\|\mathbf{z}\|} \rangle \right|^2, \quad (23)$$

634 where,  $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$ , and  $\langle \mathbf{x}, \mathbf{y} \rangle = \rho$ . Let  $\hat{i} := \mathbf{x}$ ,  $\hat{j} := \frac{1}{\sqrt{1-\rho^2}}(\mathbf{y} - \rho\mathbf{x}) = \frac{\mathbf{y} - \text{Proj}_{\text{Span}(\mathbf{x})}(\mathbf{y})}{\|\mathbf{y} - \text{Proj}_{\text{Span}(\mathbf{x})}(\mathbf{y})\|}$ ,

635 and  $\hat{k} := \frac{\mathbf{z} - \text{Proj}_{\text{Span}(\mathbf{x}, \mathbf{y})}(\mathbf{z})}{\|\mathbf{z} - \text{Proj}_{\text{Span}(\mathbf{x}, \mathbf{y})}(\mathbf{z})\|}$ . Then  $\hat{i}, \hat{j}, \hat{k}$  are unit vectors along three orthogonal coordinate axes and

636  $\mathbf{y} = \rho\hat{i} + \sqrt{1-\rho^2}\hat{j}$ . Let  $\mathbf{t} := \mathbf{z} - \mathbf{y} = t_1\hat{i} + t_2\hat{j} + t_3\hat{k}$  so that  $\mathbf{z} = (t_1 + \rho)\hat{i} + (t_2 + \sqrt{1-\rho^2})\hat{j} + t_3\hat{k}$ .  
 637 Then the problem in Eq. (23) is equivalent to the following one

$$(t_1^*, t_2^*, t_3^*) := \arg \min_{t_1, t_2, t_3: t_1^2 + t_2^2 + t_3^2 \leq \delta} \left[ \frac{(t_1 + \rho)^2}{t_1^2 + t_2^2 + t_3^2 + 2t_1\rho + 2t_2\sqrt{1-\rho^2} + 1} \right]. \quad (24)$$

638 **Case  $\rho^2 \leq \delta$ :** If  $\rho^2 \leq \delta$ , then  $t_1^* = -\rho, t_2^* = t_3^* = 0$  is a minimizer of the problem in (24) and  
 639  $\mathbf{z}^* = \sqrt{1-\rho^2}\hat{j} = \mathbf{y} - \rho\mathbf{x}$ . This solution is displayed along with  $\mathbf{x}$  and  $\mathbf{y}$  in Figure 9 and has an  
 640 immediate geometric interpretation. One can see that  $\langle \mathbf{x}, \mathbf{z} \rangle = 0$  which implies that  $X \perp Z$  because  
 641 then  $Z$  and  $X$  are uncorrelated and  $Z$ , being a linear combination of  $X$  and  $Y$ , is jointly Gaussian  
 642 with them. Also then,  $\|\mathbf{z} - \mathbf{y}\|^2 \leq \delta$ , or equivalently,  $\mathbb{E}[(Y - Z)^2] \leq \delta$ . Thus, the lower bound of 0  
 for  $I(X; Z)$  is attained in (14) by this solution.

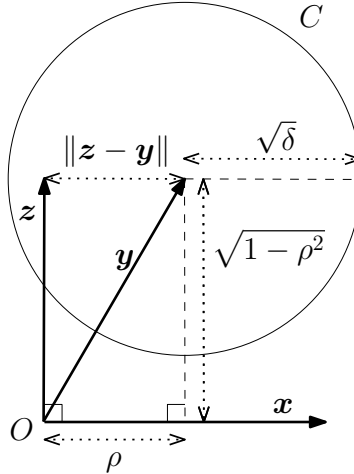


Figure 9: Solution to problem (14) for the case when  $0 \leq \rho \leq +\sqrt{\delta}$ . In the figure,  $\mathbf{x}, \mathbf{y}$  represent unit length vectors with inner product equal to  $\rho \in [0, 1]$ . The vector  $\mathbf{z}$  is perpendicular to  $\mathbf{x}$  and lies within a distortion sphere of radius  $\sqrt{\delta}$  around  $\mathbf{y}$ . The circle  $C$  is the projection of the distortion sphere on the  $\mathbf{x}$ - $\mathbf{y}$  plane and point  $O$  is the origin. The dotted lines with hollow arrowheads denote the lengths of various quantities.

643

644 **Case  $\rho^2 > \delta$ :** If  $(t_1, t_2, t_3)$  is feasible in (24), i.e.,  $t_1^2 + t_2^2 + t_3^2 \leq \delta$  then so is  $(t'_1, t'_2, t'_3) :=$   
 645  $(t_1, +\sqrt{t_2^2 + t_3^2}, 0)$ . If  $t_3 \neq 0$  then  $(t'_1, t'_2, t'_3)$  strictly dominates  $(t_1, t_2, t_3)$  because the denominator

646 of the objective function in (24) is strictly larger for  $(t'_1, t'_2, t'_3)$  than for  $(t_1, t_2, t_3)$ . Thus, we must  
 647 have

$$t_3^* = 0, \quad (25)$$

648 otherwise we can strictly improve (i.e., strictly decrease) the objective function value contradicting  
 649 the optimality of  $(t_1^*, t_2^*, t_3^*)$ . Geometrically, this means that  $\mathbf{z}^*$  must lie in the two dimensional  
 650 subspace spanned by  $\mathbf{x}$  and  $\mathbf{y}$ . Consequently, the minimization problem in (24) reduces to

$$(t_1^*, t_2^*) = \arg \min_{t_1, t_2: t_1^2 + t_2^2 \leq \delta} \left[ \frac{(t_1 + \rho)^2}{t_1^2 + t_2^2 + 2t_1\rho + 2t_2\sqrt{1 - \rho^2} + 1} \right]. \quad (26)$$

651 If  $(t_1, t_2)$  is feasible in (26), i.e.,  $t_1^2 + t_2^2 \leq \delta$  then so is  $(t'_1, t'_2) := (t_1, +\sqrt{t_2^2 + (\delta - t_1^2 - t_2^2)})$ . If  
 652  $t_1^2 + t_2^2 < \delta$ , then  $(t'_1, t'_2)$  strictly dominates  $(t_1, t_2)$  because the denominator of the objective function  
 653 in (26) is strictly larger for  $(t'_1, t'_2)$  than for  $(t_1, t_2)$ . Thus, we must have

$$(t_1^*)^2 + (t_2^*)^2 = \delta, \quad (27)$$

654 otherwise we can strictly improve (i.e., strictly decrease) the objective function value contradicting  
 655 the optimality of  $(t_1^*, t_2^*)$ . Geometrically, this means that  $\mathbf{z}^*$  must lie on the circle of radius  $+\sqrt{\delta}$   
 656 centered at  $\mathbf{y}$ .

657 Finally, we observe that if  $\mathbf{z}$  is feasible in (23), i.e.,  $\|\mathbf{y} - \mathbf{z}\|^2 \leq \delta$ , then so is  $\mathbf{z}' := \text{Proj}_{\text{Span}(\mathbf{z})}(\mathbf{y}) =$   
 658  $\gamma\mathbf{z}$ , where  $\gamma = \frac{\langle \mathbf{y}, \mathbf{z} \rangle}{\|\mathbf{z}\|^2}$ . This is because the orthogonal projection of a vector onto a subspace is the  
 659 vector in the subspace closest to it so that  $\|\mathbf{y} - \mathbf{z}'\|^2 \leq \|\mathbf{y} - \mathbf{z}\|^2$ . Also observe that the value of the  
 660 objective function in (23) is the same for both  $\mathbf{z}$  and  $\gamma\mathbf{z}$  and that  $(\mathbf{y} - \mathbf{z}') \perp \mathbf{z}'$ . Thus, we may assume  
 661 that there is an optimal solution  $\mathbf{z}^*$  such that  $(\mathbf{y} - \mathbf{z}^*) \perp \mathbf{z}^*$  for if not, we can rescale  $\mathbf{z}^*$  suitably  
 662 to ensure this property without affecting the objective function or violating the distortion constraint.  
 663 Since  $(\mathbf{z}^* - \mathbf{y}) = t_1^*\hat{i} + t_2^*\hat{j}$  and  $\mathbf{y} = \rho\hat{i} + \sqrt{1 - \rho^2}\hat{j}$ , the orthogonality condition  $(\mathbf{y} - \mathbf{z}^*) \perp \mathbf{z}^*$  can  
 664 be restated as

$$t_1^*(t_1^* + \rho) + t_2^*(t_2^* + \sqrt{1 - \rho^2}) = 0$$

665 which simplifies to

$$(t_1^*)^2 + (t_2^*)^2 + t_1^*\rho + t_2^*\sqrt{1 - \rho^2} = 0 \quad (28)$$

666 Combining (28) with (27) we get

$$\delta + t_1^*\rho + \sqrt{(\delta - (t_1^*)^2)(1 - \rho^2)} = 0.$$

667 This reduces to the following quadratic equation for  $t_1^*$  with two real roots

$$(t_1^*)^2 + 2\delta\rho t_1^* + \delta^2 - \delta(1 - \rho^2) = 0 \Rightarrow t_1^* = -\delta\rho \pm \sqrt{\delta(1 - \delta)(1 - \rho^2)}.$$

668 We note that  $\delta < 1$  since we are considering the case  $\delta < \rho^2 \leq 1$ . Of the two real roots,  $t_1^* =$   
 669  $-\delta\rho - \sqrt{\delta(1 - \delta)(1 - \rho^2)}$  has a lower objective value in (26). Using this value for  $t_1^*$  and setting  
 670  $t_2^* = \sqrt{\delta - (t_1^*)^2}$ ,  $t_3^* = 0$ , we can conclude that for the case when  $\delta < \rho^2$ , the random variable

$$Z^* := (1 - \delta)Y - (X - \rho Y) \sqrt{\frac{\delta(1 - \delta)}{1 - \rho^2}} \quad (29)$$

671 attains the lower bound on the mutual information, which equals

$$I(X; Z) = \frac{1}{2} \log \left( \frac{1}{1 - \left( \sqrt{\rho^2(1 - \delta)} - \sqrt{(1 - \rho^2)\delta} \right)^2} \right). \quad (30)$$

672 We can interpret the solution geometrically as shown in Figure 10. Unlike the previous case  
 673 ( $0 \leq \rho \leq +\sqrt{\delta}$ ), here the feasible distortion sphere does not allow  $\mathbf{z}$  to be perpendicular to  $\mathbf{x}$ .  
 674 However, one can see that the solution must lie on a tangent from the origin to the distortion sphere.  
 675 The optimum  $\mathbf{z}$  in this case (Eq. (29)) is the addition of two vectors, one along  $\mathbf{y}$  and the other  
 676 perpendicular to  $\mathbf{y}$  (the unit vector along which is  $-(\mathbf{x} - \rho\mathbf{y})/\sqrt{1 - \rho^2}$ ). The coefficients for the  
 677 linear combination can be inferred from the geometry of the figure.  $\square$

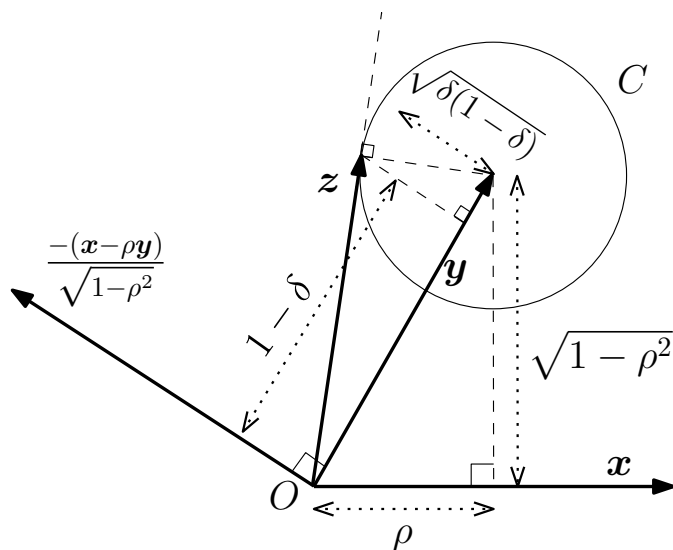


Figure 10: Geometric interpretation of the solution to the problem in (26) for the case when  $\rho > +\sqrt{\delta}$ . Here,  $x, y$  are unit length vectors with inner product equal to  $\rho$ . The unit vector perpendicular  $y$  is given by  $-(x - \rho y) / \sqrt{1 - \rho^2}$ . The circle  $C$  is the projection of the feasible distortion sphere onto the  $x$ - $y$  plane. The problem in (26) can be stated as finding  $z$  within the feasible distortion sphere which minimizes the cosine of the angle between  $z$  and  $x$ . The optimum  $z$  lies on the tangent from the origin  $O$  to the circle  $C$  on the far side of  $x$ . The dotted lines with hollow arrowheads denote the lengths of various quantities.