

Obfuscation via Information Density Estimation

Hsiang Hsu, Shahab Asoodeh, and Flavio du Pin Calmon*

Abstract

Identifying features that leak information about sensitive attributes is a key challenge in the design of information obfuscation mechanisms. In this paper, we propose a framework to identify information-leaking features via information density estimation. Here, features whose information densities exceed a pre-defined threshold are deemed information-leaking features. Once these features are identified, we sequentially pass them through a targeted obfuscation mechanism with a provable leakage guarantee in terms of E_γ -divergence. The core of this mechanism relies on a data-driven estimate of the trimmed information density for which we propose a novel estimator, named the *trimmed information density estimator* (TIDE). We then use TIDE to implement our mechanism on three real-world datasets. Our approach can be used as a data-driven pipeline for designing obfuscation mechanisms targeting specific features.

1 Introduction

A challenging problem in dataset and information sharing platforms is limiting the leakage of sensitive or private information. Sensitive information leakage can be controlled by *obfuscating* samples in a dataset prior to disclosure; i.e., perturbing the sample in a way that sensitive information cannot be effectively inferred [1–4]. Samples may contain several *features*, only some of which might leak information about sensitive attributes. For example, not all areas in a facial image equally disclose emotion (as a sensitive attribute), and not all terms used in Tweets equally reveal a user’s political preference. Given a set of sensitive attributes, an information obfuscation mechanism should ideally target only those features of the data that leak excessive amount of sensitive information. Such mechanisms usually achieve higher utility (e.g., the quality of the image) by incorporating either complete (cf. information-theoretic privacy [5–10]) or partial (cf. generative adversarial privacy [3, 11]) knowledge of the underlying data distribution.

In this paper, we propose a data-driven information-obfuscation mechanism. As a natural first step, we *identify* the information-leaking features in the data via an information-theoretic quantity called the *information density* [12, 13]. This quantity is at the heart of most information-theoretic measures of privacy [6–8] as well as differential privacy (DP) [14–19]. Intuitively, the information density captures the change of the belief about a sensitive attribute upon an observation of a sample in a disclosed dataset.

Features whose information density are above a certain threshold (which we call information-leaking features) can be randomized (e.g., perturbed) via an obfuscation mechanism. The goal of the obfuscation mechanism is to limit unwanted inferences about a sensitive attribute from disclosed data. We argue that this objective can be mathematically formulated in terms of a specific type of f -divergence [20], called the E_γ -divergence, which captures the tail distribution of the information

*H. Hsu, S. Asoodeh and F. P. Calmon are with the John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA, 02138. E-mails: hsianghsu@g.harvard.edu, shahab@seas.harvard.edu, flavio@seas.harvard.edu.

density. We propose a feature-dependent Gaussian mechanism that ensures obfuscation in terms of E_γ -divergence by targeting only the information-leaking features.

The methodology proposed here aims to develop a theoretical foundation for expounding existing approaches that completely rely on neural networks to identify and obfuscate the information-leaking features [1–3]. Despite its theoretical nature, our approach has a comparable performance in terms of sensitive information leakage as [1], without a specific “utility” target having to be pre-determined by a user. Furthermore, it adds a layer of interpretability, enabling features that pose an excessive leakage risk to be identified and communicated to the data owner.

In practice, we need to estimate the information density from samples. This estimation problem is inherently connected to mutual information estimation (since the expected value of information density is equal to the mutual information) which is known to be challenging [21–23] unless an adequate parametric model is assumed [24]. The main difficulty lies in the unboundedness of the information density, which leads to high sample complexity for reliable estimation. However, since our mechanism perturbs only information-leaking features, it requires the *trimmed information density* whose estimation is a much easier task than the original information density estimation problem. Inspired by [25,26], we develop the trimmed information density estimator (TIDE), based on the variational representations of f -divergences [25,27].

The contributions of this paper, from theoretical results to practice, are listed as follows:

1. We propose a framework for identifying information-leaking features in terms of the trimmed information density, and use the E_γ -divergence between the distributions over a sensitive attribute prior and posterior to observing a disclosed sample to measure the information leakage. Moreover, we demonstrate that obfuscation mechanisms that aim to minimize the E_γ -divergence satisfies several desirable properties in terms of information leakage guarantees (cf. Section 2).
2. We propose an estimator for the trimmed (thresholded) information density, named TIDE, and derive accompanying consistency and sample complexity guarantees. On the practical side, we present a neural network-based implementation for the TIDE (cf. Section 3).
3. We apply the obfuscation mechanism in Section 2 for image obfuscation [28–30] with GENKI-4k [31] and Celebrity Attributes (CelebA) [32] datasets, and for identifying politically-charged terms in Tweets collected from online media [33] (cf. Section 4). These experiments provide evidence that the TIDE can potentially serve as a building block in the design of obfuscation mechanisms.

It is worth mentioning that information obfuscation, being inherently prior-dependent, has several limitations [11]. In Section 5, we list some of these limitations and together with our final remarks. Proofs, experimental details, and additional experiments on synthetic data are provided in the Supplementary Material. Source code for the experiments will be made publicly available after review.

1.1 Related Work

The problem of balancing the competing objectives of providing meaningful information and inference from disclosed data, on the one hand, and obfuscating sensitive information, on the other hand, has been widely studied in information-theoretic privacy [5–9, 34, 35]. Following the information-theoretic trend, these works exploit average measures (in particular mutual information and its variants) to obfuscate data. Recently, information obfuscation has been achieved using neural networks. For example, in [1], an optimization problem similar to the *privacy funnel* [36] is formulated

to train a neural network to automatically obfuscate sensitive information while maintaining utility. In [2, 3, 11], neural generative models are introduced to generate “privatized” data that resemble the original data. These works rely on neural networks to select and perturb features, while our approach is different in the sense that we first identify the information-leaking features using the information density and apply local obfuscation only on these features.

Our approach of first identifying the information-leaking features and then perturbing those features is inspired by the instance-based additive mechanism of [37] in the DP setting. In fact, the information density appears in DP under the name of *privacy loss* variable [15–19], thereby connecting DP and information-theoretic quantities, e.g. mutual information DP [38] and Rényi DP [39]. Despite this connection, we emphasize that our approach is fundamentally different from DP, in that we consider prior distribution on sensitive attributes and also we allow correlation among features (see, e.g., [40] for the limitations of DP for correlated data).

Estimating information density from samples is connected to density ratio estimation [26, 27, 41] — a fundamental task in various applications of machine learning and statistics, including outlier detection [42], transfer learning [43], and generative adversarial networks [44]. A naïve approach to determine the density ratio is to use the plug-in estimator, which is known to perform poorly [24] unless adequate parametric models (e.g., linear [41], kernel [45], or exponential family [26] models) are assumed. The two closest approaches to the trimmed information density estimation in this paper are (i) [27], which proposed using the variational representation of f -divergences to convert information density estimation into an optimization problem over finite-complexity set of functions and (ii) [26], which estimated the trimmed density ratio of variables from exponential family distributions. We enforce a threshold on the information density when solving the optimization problem in the variational representation of f -divergences (see Section 3).

1.2 Notation

Capital letters (e.g., X) denote random variables, and calligraphic letters (e.g., \mathcal{X}) denote sets. We denote the probability measure of $X \times S$ by $P_{X,S}$, the conditional probability measure of S given X by $P_{S|X}$, and the marginal probability measure of X and S by P_X and P_S respectively. We use $P_{S|X}(\cdot|x)$ and $P_{S|x}$ interchangeably. We represent the fact that X is distributed according to P_X by $X \sim P_X$. KL-divergence is given by $D_{\text{KL}}(P_{S,X}||P_S P_X) = \mathbb{E}_{P_{S,X}}[\log(P_{S,X}/P_S P_X)]$. We denote the realization (i.e., sample) drawn from a probability distribution by $x = (x_1, \dots, x_j, \dots, x_m)$, where x_j is the j^{th} feature for $j = 1, \dots, m$. Similarly, X_j is the j^{th} feature of the data variable. We denote $[k] = [1, \dots, k]$, $x^k = [x_1, \dots, x_k]$, and $(z)_+ = \max\{z, 0\}$ for a scalar z . Finally, $I_{d \times d}$ is the identity matrix of dimension d , and $\mathbf{1}_{\{\cdot\}}$ is the indicator function.

2 Problem Formulation

We consider the setting where a user wishes to disclose data X (e.g., image, tweet) while controlling the information revealed about a (correlated) sensitive attribute S (e.g., emotion, political preference). The goal is to produce an obfuscated representation Y of X that discloses only negligible information about S . We assume that X consists of m features, i.e., $X = (X_1, \dots, X_m)$, where each feature takes values in a compact set \mathcal{X} . Throughout this section, we assume that $(S, X) \sim P_{S,X}$ and the underlying distribution $P_{S,X}$ is given. This restrictive assumption will be dropped in the subsequent section.

One possible approach to obfuscate X is to independently perturb each feature (e.g., by adding noise to each pixel of an image). However, in many applications, only a few features of the data are correlated with the sensitive attribute, rendering adding independent noise highly sub-optimal. In

this section, we propose an information-theoretic framework for data obfuscation which consists of two parts: First, we *identify information-leaking features*, and then obfuscate *only* those features. This way, many features need not be perturbed, leading to an improvement in the utility of the disclosed data.

Our framework relies on an information-theoretic quantity called the *information density*, a term coined in [12] and has since been used in numerous applications in information theory and statistics, particularly in binary hypothesis testing (see, e.g., Neyman-Pearson Lemma [46]).

Definition 1 (Information Density). Given a pair of realization (s, x) of $(S, X) \sim P_{S, X}$, the information density between s and x is defined as

$$i(s; x) \triangleq \log \frac{P_{S, X}(s, x)}{P_S(s)P_X(x)} = \log \frac{P_{X|S}(x|s)}{P_X(x)}. \quad (1)$$

Similarly, information density can be defined for each feature x_j as

$$i(s; x_j) \triangleq \log \frac{P_{S, X_j}(s, x_j)}{P_S(s)P_{X_j}(x_j)}, \quad (2)$$

and the conditional information density between s and x_j given another feature x_r as

$$i(s; x_j|x_r) \triangleq \log \frac{P_{S, X_j|X_r}(s, x_j|x_r)}{P_{S|X_r}(s|x_r)P_{X_j|X_r}(x_j|x_r)}. \quad (3)$$

Intuitively, $i(s; x_j)$ evaluates the change of belief about s upon observing x_j . In particular, if $|i(s; x_j)|$ is small, then x_j does not significantly contribute in increasing the belief of an adversary about s , since $P_{S|X}(s|x_j) \approx P_S(s)$. This, however, does not mean that x_j can be disclosed “as is” without incurring an information leakage risk. To see why, consider, for example, that $m = 2$, X_1 and X_2 are independent and uniform binary random variables, and $S = X_1 + X_2$ (modulo 2). Although $i(s; x_1) = i(s; x_2) = 0$ for any realization (s, x_1, x_2) of (S, X_1, X_2) , the release of both x_1 and x_2 would allow perfect reconstruction of s . To account for such inferences of sensitive attributes, we consider the conditional information density as a yardstick for identifying information-leaking features.

Definition 2 (Information-Leaking Feature). Given an observed sample $x = (x_1, \dots, x_m)$, $j \in [m]$, and $\varepsilon \geq 0$, the feature x_j is said to be an ε -*information-leaking feature* if there exists a sensitive attribute s such that $|i(s; x_j|x^{j-1})| > \varepsilon$.

The threshold ε is a tradeoff parameter between information leakage risk and the utility of the disclosed data (e.g., the quality of an image). Notice that if the data is not equipped with a natural ordering (e.g., time series), we can choose an arbitrary ordering for the conditioning features x^{j-1} (cf. Section 4.1 for an example in images).

2.1 A Naïve Obfuscation Mechanism

Given any $j \in [m]$, $\varepsilon \geq 0$, and all features x^{j-1} , define

$$B_j^\varepsilon(x^{j-1}) \triangleq \{x \in \mathcal{X} : |i(s; x_j|x^{j-1})| > \varepsilon \text{ for some } s \in \mathcal{S}\}. \quad (4)$$

If $x_j \notin B_j^\varepsilon(x^{j-1})$, then it can be disclosed “as is” because it cannot be used to infer sensitive attributes given all the previous features. On the other hand, each feature $x_j \in B_j^\varepsilon(x^{j-1})$ is required

to be obfuscated. To do so, we shall pass all such features sequentially through an *obfuscation mechanisms* to ensure that they no longer belong to $B_j^\varepsilon(x^{j-1})$.

Consider the mechanisms $\mathcal{M}_j : \mathcal{X} \rightarrow \mathcal{X}$ such that if $x_j \notin B_j^\varepsilon(x^{j-1})$ then $\mathcal{M}_j(x_j) = x_j$ (deterministic) and if $x_j \in B_j^\varepsilon(x^{j-1})$ then $\mathcal{M}_j(x_j)$ generates Y_j a random variable from a distribution to be designed. A natural question raised here is: how should information obfuscation be measured? To answer this question, we introduce the *obfuscation metric* $\Pr(|i(s; Y_j|y^{j-1})| > \varepsilon)$ and require

$$\Pr(|i(s; Y_j|y^{j-1})| > \varepsilon) \leq \frac{\delta}{m}, \quad (5)$$

for all $s \in \mathcal{S}$, where y^{j-1} is any output of the $\mathcal{M}_1(x_1), \dots, \mathcal{M}_{j-1}(x_{j-1})$. Although this metric is intuitive, it presents a serious drawback for use in practice. Any reasonable mechanism must be immune to post-processing: any processing of the mechanism's output should only decrease the information leakage risk or equivalently the obfuscation metric. However, the obfuscation metric in (5) may violate this property. To see this, let $m = 1$ and \tilde{Y} be obtained by applying an arbitrary post-processing to Y the output of the mechanism \mathcal{M}_1 satisfying the obfuscation metric $\Pr(i(Y; s) > \varepsilon) \leq \delta$ for all s . Immunity to post-processing is then equivalent to requiring

$$\Pr(i(s; \tilde{Y}) > \varepsilon) \leq \Pr(i(s; Y) > \varepsilon), \quad (6)$$

for all $s, \varepsilon \geq 0$ and $\delta \in [0, 1]$. However, we show in the following that there must exist some ε for which (6) is violated. To see this, notice that $\mathbb{E} \left[\frac{P_{\tilde{Y}|s}(\tilde{Y}|s)}{P_{\tilde{Y}}(\tilde{Y})} \right] = \mathbb{E} \left[\frac{P_{Y|s}(Y|s)}{P_Y(Y)} \right] = 1$ and hence we have

$$\int_0^\infty \Pr(e^{i(s; \tilde{Y})} \geq t) dt = \int_0^\infty \Pr(e^{i(s; Y)} \geq t) dt. \quad (7)$$

Therefore, Eq. (6) must hold with equality for all $\varepsilon \geq 0$ which in turn implies

$$D_{\text{KL}}(P_{\tilde{Y}|s} \| P_{\tilde{Y}}) = D_{\text{KL}}(P_{Y|s} \| P_Y). \quad (8)$$

However, according to data processing inequality for KL divergence, Eq. (8) cannot hold true in general. Therefore, there must exist some ε for which (6) does not hold. For more details about this construction, see [47].

Next, we propose another metric in terms of a certain f -divergence, the so-called E_γ -divergence, and show that it implies (5) while being immune to post-processing.

2.2 E_γ -Divergence

To address the issue raised above, we resort to a particular divergence metric between two probability distributions called E_γ -divergence, and show that this divergence bounds an appropriately weighted tail distributions of $i(s; Y)$.

Definition 3 (E_γ -Divergence [48]). Given two probability distributions P and Q defined on the same support set \mathcal{A} and $\gamma \geq 1$, we define E_γ -divergence as

$$E_\gamma(P \| Q) \triangleq \sup_{A \subset \mathcal{A}} P(A) - \gamma Q(A) \quad (9)$$

$$= \int_{a \in \mathcal{A}} (dP(a) - \gamma dQ(a))_+, \quad (10)$$

where the equality comes from the fact that the optimizer in (9) is $\mathcal{A}^* = \{a \in \mathcal{A} | P(a) - \gamma Q(a) \geq 0\}$.

E_γ -divergence has been considered in various fields; for example, it appears in DP literature as an equivalent definition for differentially private mechanisms (see e.g. [49, 50]), in statistics as the probability of correct decision in Bayesian binary hypothesis testing [48], and in information theory for proving general channel coding converse results [48, 51].

Notice that $E_\gamma(P\|Q) \leq 1$ for all $\gamma \geq 1$ and any pair of distributions (P, Q) . It is clear that the constraint $E_\gamma(P_Y\|P_{Y|s}) \leq \delta$ for some $\delta \in (0, 1)$ ensures that $P_Y(A) - \gamma P_{Y|s}(A) \leq \delta$ for *all* subsets $A \subset \mathcal{X}$ and in particular $P_Y(\mathcal{A}^*) \leq \delta$. Since for $\gamma = e^\varepsilon$, the set \mathcal{A}^* corresponds to the tail events of the random variable $i(Y; s)$, we henceforth assume $\gamma = e^\varepsilon$. Note also that to have control on both tail events $\{i(Y; s) < -\varepsilon\}$ and $\{i(Y; s) > \varepsilon\}$, we need to consider both $E_{e^\varepsilon}(P_Y\|P_{Y|s}) \leq \delta$ and $E_{e^\varepsilon}(P_{Y|s}\|P_Y) \leq \delta$. In the sequel, we present our results only for $E_{e^\varepsilon}(P_{Y|s}\|P_Y) \leq \delta$. The results for the reversed divergence can be derived *mutatis mutandis*.

Having this divergence at our disposal, we can now propose obfuscation criteria for the mechanisms $\{\mathcal{M}_j\}$. As before, if $x_j \notin B_j^\varepsilon(x^{j-1})$, we set $\mathcal{M}_j(x_j) = x_j$; otherwise, we shall construct randomized mechanism $\mathcal{M}_j : \mathcal{X} \rightarrow \mathcal{X}$ such that $\mathcal{M}_j(x_j) = Y_j$ satisfies

$$E_{e^\varepsilon}(P_{Y_j|s, y^{j-1}}\|P_{Y_j|y^{j-1}}) \leq \frac{\delta}{m}, \quad (11)$$

where y^{j-1} is a realization of all previous mechanisms $\mathcal{M}_1(x_1), \dots, \mathcal{M}_{j-1}(x_{j-1})$. The factor $\frac{1}{m}$ in the right-hand side of (11) is only for the sake of normalization (to be clarified in Theorem 2).

It is clear from (9) that upper bounds on $E_\gamma(P_{Y_j|s, y^{j-1}}\|P_{Y_j|y^{j-1}})$ directly translate into low-leakage guarantee (5). Furthermore, since E_γ -divergence belongs to the family of f -divergences [52], it satisfies the data processing inequality which in turn implies that mechanisms satisfying (11) are immune to post-processing.

To even further justify the choice of E_γ -divergence as a “proxy” for the obfuscation metric in (5), we prove in the following theorem an equivalent formula for $E_{e^\varepsilon}(P_{Y_j|s, y^{j-1}}\|P_{Y_j|y^{j-1}})$ in terms of the tail distribution $\Pr(i(s; Y_j|y^{j-1}) > t)$ for $t \geq 0$.

Theorem 1 (Tail Distribution Formula). Given distributions $P_{Y_j|s, y^{j-1}}$ and $P_{Y_j|y^{j-1}}$, we have

$$E_{e^\varepsilon}(P_{Y_j|s, y^{j-1}}\|P_{Y_j|y^{j-1}}) = e^\varepsilon \int_\varepsilon^\infty e^{-t} \Pr(i(s; Y_j|y^{j-1}) > t) dt. \quad (12)$$

This result provides an operational interpretation for E_γ -divergence for our obfuscation setting. More precisely, $E_{e^\varepsilon}(P_{Y_j|s, y^{j-1}}\|P_{Y_j|y^{j-1}}) \leq \delta$ enforces the events $\{i(Y; s) > t\}$ to have small aggregate (weighted) probability for *all* $t \geq \varepsilon$.

Next, we address the composition property of the above mechanisms: If each mechanism \mathcal{M}_j satisfies (11), then so does the composed mechanism $\mathcal{M} = (\mathcal{M}_1, \dots, \mathcal{M}_m)$ with parameters $m\varepsilon$ and δ . Recall that $Y = (Y_1, \dots, Y_m)$ is the output of the mechanism \mathcal{M} .

Theorem 2 (Composition). For all mechanisms $\mathcal{M}_j, j \in [m]$ satisfying (11), we have for all $s \in \mathcal{S}$

$$E_{e^{m\varepsilon}}(P_{Y|s}\|P_Y) \leq \delta. \quad (13)$$

This theorem states that a guarantee for each feature, given by (11), will result in a meaningful guarantee for the whole sample. This, in particular, demonstrates the need for conditional information density in Definition 2, as opposed to the unconditional one.

2.3 A Gaussian Obfuscation Mechanism

We next give an explicit construction of mechanisms $\{\mathcal{M}_j\}$ satisfying (11). Here, we assume that each feature $x_j \in C$ where C is a compact subset of \mathbb{R}^r . Recall that each mechanism \mathcal{M}_j is required to generate Y_j satisfying (11). As a simple approach to enforce this guarantee, we propose the additive Gaussian mechanism; that is, for each given $j \in [m]$, ε , and $x^{j-1} \in \mathcal{X}^{j-1}$, we consider the following mechanism

$$Y_j = x_j + \lambda \mathbf{1}_{\{x_j \in B_\varepsilon^r(x^{j-1})\}} N, \quad (14)$$

where N is an independent standard Gaussian noise $\mathcal{N}(0, \mathbf{I}_{r \times r})$ and $\lambda > 0$ is determined according to the following theorem.

Theorem 3 (Gaussian Obfuscation). The Gaussian obfuscation mechanism (14) satisfies (11) if λ satisfies

$$\theta_{e^\varepsilon}(K, \lambda) \leq \frac{\delta}{m}, \quad (15)$$

where K is the radius of C , i.e., $K = \max_{w \in C} \|w\|$, and for any $a > 0$

$$\theta_{e^\varepsilon}(a, \lambda) \triangleq \mathbf{Q}\left(\frac{\lambda\varepsilon}{a} - \frac{a}{2\lambda}\right) - e^\varepsilon \mathbf{Q}\left(\frac{\lambda\varepsilon}{a} + \frac{a}{2\lambda}\right), \quad (16)$$

where $\mathbf{Q}(v) = \Pr(\mathcal{N}(0, 1) \geq v) = \int_v^\infty \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$.

In light of this theorem, if $\varepsilon \approx 0$, then the noise variance λ must be of order $O\left(\frac{K}{-\log(1-\frac{\delta}{m})}\right)$. The exact value of noise variance, however, cannot be derived as there is no analytic expression for the \mathbf{Q} function.

We have thus far made the information-theoretic assumption that the underlying distribution $P_{S,X}$ is given and, consequently, the information density is known exactly. In the following section, we propose a *data-driven* estimator for information density which renders our proposed mechanism applicable to real-world datasets.

3 Trimmed Information Density

The obfuscation mechanism in Section 2 relies on the conditional information density $i(s; x_j | x^{j-1})$ to identify the set of information-leaking. Notice that, since information density satisfies the chain rule, i.e.

$$i(s; x_j | x^{j-1}) = i(s; x^j) - i(s; x^{j-1}), \quad (17)$$

an estimate of $i(s; x_j | x^{j-1})$ can be constructed by estimates of $i(s; x^j)$ and $i(s; x^{j-1})$.

In general, exact estimation of the information density is hard due to its unboundeness. However, we do not need the exact estimation; instead, we only need to know if the absolute value of the conditional information density is larger than the threshold ε (Definition 2). In other words, estimating the *trimmed* information density is sufficient for obfuscation purposes. Moreover, the tail of the information density satisfies [48]

$$\Pr \{i(s; X_j) > t\} \leq e^{-t}, \quad \forall s, \quad (18)$$

indicating that the estimation error caused by trimming the information densities can be controlled. In this section, we propose a consistent and scalable estimator for the *trimmed information density*, called the TIDE, and show that estimating the trimmed information density can be easier than estimating the exact information density in terms of sample complexity.

3.1 Trimmed Information Density Estimator

TIDE is based on a variational representation of KL divergence¹ known as the Donsker-Varadhan (DV) representation [53], given by

$$D_{\text{KL}}(P_{S,X} \| P_S P_X) = \sup_{g: \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}} \{ \mathbb{E}_{P_{S,X}}[g(S, X)] - \log \mathbb{E}_{P_S P_X}[e^{g(S, X)}] \}. \quad (19)$$

Recall that $D_{\text{KL}}(P_{S,X} \| P_S P_X)$ is equal to the mutual information $I(S; X)$ between S and X , which is in fact the expected information density $\mathbb{E}_{P_{S,X}}[i(S, X)]$. It can be shown that the maximizer g^* of (19) is exactly the information density, i.e., $g^*(s, x) = i(s; x)$. Hence, the problem of estimating information density is equivalent to solving the functional optimization problem (19) given access to samples drawn from $P_{S,X}$.

Since the search space in (19) is unconstrained, directly solving the optimization by computing the empirical expectations would fail in general. One practical approach is to restrict the search space to a family $\mathcal{G}(\Theta)$ of continuous and bounded functions g_θ parameterized by θ in a compact domain $\Theta \subset \mathbb{R}^d$, where d is the number of parameters. The new constrained optimization problem corresponds to approximating the information density by a bounded function, thus the name *trimmed* information density.

The TIDE is then given by

$$\hat{g}_n \triangleq \operatorname{argmax}_{g_\theta \in \mathcal{G}(\Theta)} \{ \mathbb{E}_{P_{S_n, X_n}}[g_\theta(S, X)] - \log \mathbb{E}_{P_{S_n} P_{X_n}}[e^{g_\theta(S, X)}] \}, \quad (20)$$

where P_{S_n, X_n} and $P_{S_n} P_{X_n}$ denote the empirical distributions of $P_{S,X}$ and $P_S P_X$ by n samples, respectively.

3.2 Consistency and Sample Complexity

The TIDE obtained by solving (20) belongs to a broader class of *extremum estimators* [54] of the form $\hat{a} = \operatorname{argmax}_{a \in \mathcal{A}} \Lambda_n(a)$, where $\Lambda_n(a)$ is an objective function and \mathcal{A} is a parameter space. The consistency of extremum estimators is guaranteed by the Newey-McFadden Lemma [55] (cf. Supplementary Material), which in turn implies the consistency of the TIDE, as stated in the following theorem.

Theorem 4 (Consistency). If $\mathcal{G}(\Theta)$ is the family of continuous and bounded functions parameterized by θ taking values in a compact domain Θ , then the TIDE (20) is consistent, i.e., for any $\eta > 0$, there exist $N > 0$ such that for all $n > N$, we have $|\hat{g}_n(s, x) - g^*(s, x)| \leq \eta$ with probability one for all $s \in \mathcal{S}$ and $x \in \mathcal{X}$.

We turn our attention to deriving the sample complexity of the TIDE. We make further assumption that functions in $\mathcal{G}(\Theta)$ are Lipschitz, and use (18) to prove the following theorem. To avoid technical complications, we assume that $\mathbb{E}_{P_{S,X}}[g(S, X)]$ and $\mathbb{E}_{P_S P_X}[e^{g(S, X)}]$ are finite for all functions g in $\mathcal{G}(\Theta)$.

Theorem 5 (Sample Complexity). Assume that functions in $\mathcal{G}(\Theta)$ are bounded by M and Lipschitz with respect to θ , and $\Theta \subset \mathbb{R}^d$ is compact. Then we have $|\hat{g}_n(s, x) - g^*(s, x)| \leq \eta$ with probability at least $1 - e^{-M}$, for all $s \in \mathcal{S}$ and $x \in \mathcal{X}$, where $n = O(\frac{M^3 d}{\eta^2})$.

¹Other f -divergences [20, 52] could also be used, see the Supplementary Material for more details.

Observe that trimming the information density is crucial for the bound in the previous theorem to hold: if $M \rightarrow \infty$ (i.e., estimating the exact information density), the sample complexity of the TIDE grows to infinity and the result is vacuous. In fact, we need to restrict the search space to all continuous and bounded functions \mathcal{G} to exactly approximate the trimmed information density. However, for computational reason, we assume that these functions can be parameterized by a compact domain Θ , and the complexity of the family $\mathcal{G}(\Theta)$ is characterized by its number of parameters d . As the complexity of the functions $d \rightarrow \infty$, meaning the search space is too large, the sample complexity goes to infinity. This assumption allows us to approximate the functions in $\mathcal{G}(\Theta)$ by neural networks, where Θ is the weights in all layers, as we will see next.

3.3 Implementation

In practice, we use the set of functions representable by a neural network with output clipped to $[-M, M]$ to approximate the set of continuous and bounded functions $g(s, x)$ in \mathcal{G} . By sampling (s, x) from $P_{S, X}$ and from $P_S \times P_X$ for the first and second expectations in (20), we can fit the neural network. After training, the $g(s, x)$ outputs the estimate of the trimmed information density of samples $|i(s; x)| \leq M$. In order to reconstruct the conditional information density by the chain rule (17), we compute $g(s, x^j)$ for $i(s, x^j)$ and $g(s, x^{j-1})$ for $i(s, x^{j-1})$; then the $i(s, x^j) - i(s, x^{j-1})$ gives the desired conditional information density $|i(s; x_j | x^{j-1})| \leq 2M$.

4 Experiments

The experiments contain two parts. First, we investigate image obfuscation [28–30] as a common use case of our approach with the GENKI-4k [31] and Celebrity Attributes [32] datasets. Second, we demonstrate how TIDE can be possibly used to discover politically-charged terms in the Tweets of online media [33]. Detailed experimental setups (e.g., architecture of the function g in TIDE, training details) and additional experiments on Gaussian synthetic data are provided in the Supplementary Material.

4.1 Image Obfuscation

A common application of information obfuscation is image obfuscation [28–30], where we aim to hide information related to a given sensitive attribute in an image. Unlike existing works which rely on neural networks to select and perturb features [1, 28], we apply the TIDE to identify information-leaking features for the Gaussian obfuscation mechanism (Section 2). We split x into a grid, where each “patch” of size $p \times p$ pixels in the grid represents the low-level features x_j of the image x . It is a common method to extract low-level features in an image [28]. We number each x_j in an image from the upper-left corner to the lower-right, and use the TIDE (with $M = 3$) to determine the information-leaking features by (4), and demonstrate our obfuscation approach on two datasets: the GENKI-4k and Celebrity Attributes datasets.

4.1.1 GENKI-4K Dataset

This dataset contains 2400 images for training and 600 for testing, where each image x is a 64×64 pixels face that has emotion smiling ($s = 1$) or not ($s = 0$). We select 10 faces for illustration in Figure 1. When the patch size is 32×32 (4 patches), the TIDE simply flags the lower two patches to be information-leaking. As the patch becomes finer, the information-leaking patches concentrate to the mouse area; thus when applying the Gaussian obfuscation mechanism, it is visually possible

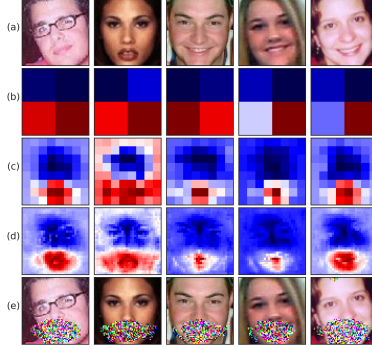


Figure 1: Row (a) shows original images. Rows (b), (c) and (d) show the information-leaking patches found by the TIDE (20) with patch sizes 32×32 , 8×8 and 2×2 pixels respectively (color red indicates higher value). Row (e) shows the Gaussian obfuscation mechanism (14) on row (d) with $\varepsilon = 0.5$ and $\lambda = 1.0$, which successfully hide the sensitive attribute of emotion. The information-leaking patches is easy to interpret: the TIDE focuses more on the mouth area as the patches become finer.

Table 1: Classification accuracy of emotion obfuscation for the GENKI-4k dataset with different patch sizes $p \times p$ and threshold ε . Results on obfuscating the lower-half image by Gaussian noise (LHI) and on random guessing are shown as comparison.

		Classification Accuracy %				
		ε	0.5	0.6	0.7	0.8
$p \times p$	32×32	50.54	50.54	92.04	92.04	92.04
	16×16	50.72	51.46	79.14	89.52	92.04
	8×8	50.93	68.94	78.71	87.33	92.04
	4×4	50.60	65.06	75.23	83.89	92.04
	2×2	50.64	62.25	68.59	80.26	92.04
	LHI	50.58	-	-	-	-
Guess	50.41	-	-	-	-	

to identify the gender of the subject but with their emotion obfuscated. The leakage guarantee in Theorem 3, $\delta/m \approx 0.24$, can be computed by (16) with $\varepsilon = 0.5$ and $K = 1$ since the images are normalized. Note that the TIDE can not only reveal the patches informative of emotion, but also captures the contour of faces.

We train an adversary that can classify the emotion of the subject with accuracy 92.04%, and report the classification accuracy of the Gaussian obfuscation mechanism ($\lambda = 1$ in (14)) under different patch sizes and threshold ε in Table 1. When $\varepsilon = \infty$ (i.e. $B_j^\varepsilon(x^{j-1}) = \phi$ for all j), no patch is identified by the TIDE, and therefore the performances are the same as the adversary. A simple mechanism to hide the emotion in images is adding Gaussian noise onto the Lower Half of the Image (LHI). As a comparison, the results of LHI and random guessing are also included in Table 1. The LHI gives similar performance when the patch size is 32×32 since when $\varepsilon = 0.5$, the lower two patches of the image will be identified as information-leaking for the mechanism (Figure 1 row (b)), but LHI will erase too much information that is not related to the emotion. The random guessing values correspond to the prior distribution of the emotion labels in the training set.

4.1.2 Celebrity Attributes (CelebA) Dataset

This more challenging dataset contains 202599 colorful high-resolution images, where each image is a 218×178 -pixel face image of a celebrity with 40 distinct binary labels, including **smiling**,

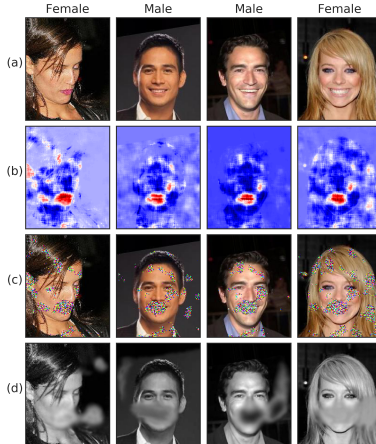


Figure 2: Row (a) shows original images. Row (b) shows the information-leaking patches with size 2×2 by the TIDE (color red indicates higher value). Row (c) shows the Gaussian obfuscation mechanism on row (b) with $\varepsilon = 0.74$ and $\lambda = 1$, and row (d) shows information obfuscation in [1] with the sensitive information budget equal to 0.72 bits.

Table 2: Comparison between our approach (with patch size 2×2) and [1] (ε here stands for the tolerance of sensitive information leak) on emotion and gender classification accuracy for the CelebA dataset.

Threshold	Classification Accuracy %			
	Our approach		Method in [1]	
	Emotion	Gender	Emotion	Gender
ε				
∞	92.04	94.29	92.04	94.29
0.8	85.97	91.48	85.59	92.53
0.7	75.15	90.39	76.40	91.20
0.6	71.33	87.61	70.88	89.77
0.5	69.01	86.97	68.60	89.47
LHI	53.91	69.35	53.91	69.35
Guess	51.79	58.32	51.79	58.32

gender, Arched Eyebrows, etc. We select 100k images as X and the sensitive attribute S to be emotion as well for training the TIDE. In Figure 2, we randomly pick 4 images for illustration. Given a small patch size, the Gaussian obfuscation mechanism ($\lambda = 1$ in (14)) perturbs selective patches to hide the sensitive attribute while keeping other useful information (e.g. gender) intact. The leakage guarantee (Theorem 3), $\delta/m \approx 0.18$, can be computed by (16) with $\varepsilon = 0.74$ and $K = 1$. In Figure 2 row (d), we reproduce the method by [1] since it is the state-of-the-art result in information obfuscation and its implementation is publicly available. The main difference between our approach and [1] is that [1] requires an additionally pre-specified utility (i.e. the labels of gender), while our approach does not require such labels. As we can see, both methods shown in Figure 2 rows (c) and (d) obfuscate the mouth and some other area. However, our approach tends to obfuscate less of the subject’s face.

We train two classifiers for emotion and gender, and report the accuracy of our approach and [1] in Table 2. Both methods block emotion recognition, effectively pushing the accuracy of the emotion classifier towards random guessing. More importantly, the gender classifier still performs well over the sanitized images.

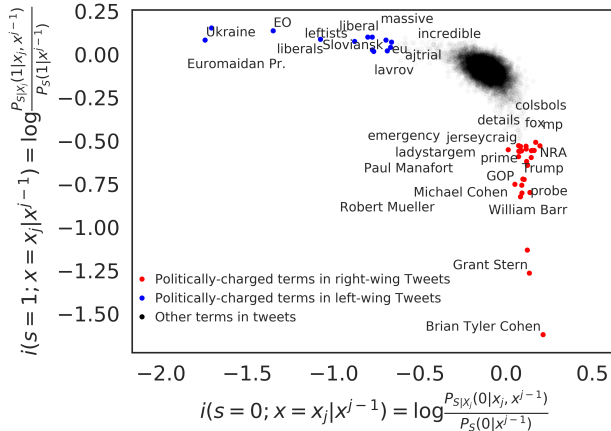


Figure 3: $i(s; x_j | x^{j-1})$ for terms in Tweets. GOP: Grand Old Party (i.e. the Republican Party), NRA: National Rifle Association, EO: Entrepreneurs’ Organization, Euromaidan Pr.: Euromaidan Press.

4.2 Information-Leaking Terms in Tweets

Finally, we showcase how the TIDE can be used in natural language to identify politically-charged terms in the Tweets from online media [33]. The information density is called the pointwise mutual information (PMI) in natural language processing to measure associations between words and labels [56]. Since perturbation on languages is not yet well-defined [57], we do not perform the mechanism in (14), but focus on identifying information-leaking terms.

We collect $N = 75946$ Tweets from more than 20 online publishers (e.g. CNN, Bloomberg, New York Times), and determine their private attribute S as the political preference of being right-wing ($s = 0$) and left-wing ($s = 1$) according to [33], where the numbers of samples with each political bias are equivalent. We pre-process the Tweets to keep only meaningful terms (i.e. pieces of words) and use bag-of-words representation [58] to tokenize all the pieces of words for each Tweet according to term frequency, ending up with 24657 terms (i.e. features $x_j, j \in [24657]$). We train the TIDE using the tokenized Tweets as x . In Figure 3, we show the estimate of trimmed conditional information density $i(s; x_j | x^{j-1})$ of each term. It is clear that some terms carry more information about the political bias. For instance, terms such as “Grand Old Party” and “National Rifle Association” associate with right-wing politics, and terms “Europe” and “liberal(s)” with the left. In this scenario, our approach could be eventually deployed as a plug-in to warn the users about potential political preference leaks before posting Tweets.

5 Final Remarks

We introduced a new information obfuscation framework that first identifies information-leaking features using the trimmed information density, and then tailors the obfuscation mechanism only on these features. To our knowledge, this framework is the first formal application of information density to quantify information-leaking features, and could potentially serve as a data-driven tool for designing obfuscation mechanism for high-dimensional data.

Limitations. In order to estimate the information density, we make two key assumptions: (i) we know *a priori* sensitive attributes that we wish to hide (e.g., political preference), and (ii) we have access to a reference dataset from which we can fit the TIDE (though this is difficult to avoid as discussed in [59]). Although these assumptions are restrictive in practice, they allow us to

develop systematic machinery to discover information-leaking samples and features in an entirely data-driven manner.

Here, we give proofs of theorems and other technical discussions omitted from Sections 2 and 3 and also provide further details about the experiment setups and also the training phase.

A Proofs and Theoretical Backgrounds

In this section, we provide proofs omitted in the main text, as well as some discussions on the relationship between the TIDE and variational representations of f divergences and the Newey-McFadden lemma.

A.1 Proof of Theorem 1

For notational brevity, we drop y^{j-1} from the conditioning part of $P_{Y_j|s,y^{j-1}}$ and $P_{Y_j|y^{j-1}}$ and also write P and Q for $P_{Y_j|s}$ and P_{Y_j} , respectively. To prove this theorem, note that according to Definition 3, we can write

$$\mathbb{E}_{e^\varepsilon}(P\|Q) = P(i(s; Y_j) > \varepsilon) - e^\varepsilon Q(i(s; Y_j) > \varepsilon).$$

Hence, letting \mathcal{C} denote the tail event $\{y : i(s; y) > \varepsilon\}$ for a given s , we have

$$\begin{aligned} \mathbb{E}_{e^\varepsilon}(P\|Q) &= P(\mathcal{C}) - e^\varepsilon Q(\mathcal{C}) \\ &= \mathbb{E}_P \left[\mathbf{1}_{\{Y_j \in \mathcal{C}\}} \right] - e^\varepsilon \mathbb{E}_Q \left[\mathbf{1}_{\{Y_j \in \mathcal{C}\}} \right] \\ &\stackrel{(a)}{=} \mathbb{E}_Q \left[e^{i(s; Y_j)} \mathbf{1}_{\{Y_j \in \mathcal{C}\}} \right] - e^\varepsilon \mathbb{E}_Q \left[\mathbf{1}_{\{Y_j \in \mathcal{C}\}} \right] \\ &= \mathbb{E}_Q \left[\left(e^{i(s; Y_j)} - e^\varepsilon \right) \mathbf{1}_{\{Y_j \in \mathcal{C}\}} \right] \\ &= \mathbb{E}_Q \left[\left(e^{i(s; Y_j)} - e^\varepsilon \right)_+ \right] \\ &= \mathbb{E}_Q \left[e^{i(s; Y_j)} e^{-i(s; Y_j)} \left(e^{i(s; Y_j)} - e^\varepsilon \right)_+ \right] \\ &\stackrel{(b)}{=} \mathbb{E}_P \left[\left(1 - e^\varepsilon e^{-i(s; Y_j)} \right)_+ \right] \\ &= \int_0^\infty \Pr \left(\left(1 - e^\varepsilon e^{-i(s; Y_j)} \right) \mathbf{1}_{\{Y_j \in \mathcal{C}\}} \geq t \right) dt, \end{aligned}$$

where both (a) and (b) follow from the simple change-of-variable argument $\mathbb{E}_P [f(Y)] = \mathbb{E}_Q [e^{i(s; Y_j)} f(Y)]$ for any function f .

Furthermore, since $(1 - e^{\varepsilon - i(s; Y_j)}) 1_{\{i(s; Y_j) > \varepsilon\}} < 1$ with probability one, we have

$$\begin{aligned}
\mathbf{E}_{e^\varepsilon}(P \| Q) &= \int_0^\infty \Pr \left((1 - e^{\varepsilon - i(s; Y_j)}) 1_{\{Y_j \in \mathcal{C}\}} \geq t \right) dt \\
&= \int_0^1 \Pr \left((1 - e^{\varepsilon - i(s; Y_j)}) 1_{\{Y_j \in \mathcal{C}\}} \geq t \right) dt \\
&= \int_0^1 \Pr \left(1 - e^{\varepsilon - i(s; Y_j)} \geq t \right) dt \\
&= \int_0^1 \Pr \left(e^{-i(s; Y_j)} \leq (1 - t)e^{-\varepsilon} \right) dt \\
&= e^\varepsilon \int_0^{e^{-\varepsilon}} \Pr \left(e^{-i(s; Y_j)} \leq b \right) db \\
&= e^\varepsilon \int_\varepsilon^\infty e^{-t} \Pr \left(i(s; Y_j) \geq t \right) dt.
\end{aligned}$$

A.2 Proof of Theorem 2

First assume that $m = 2$. For any set $A \subset \mathcal{X}^2$ and $s \in \mathcal{S}$, we have

$$\begin{aligned}
P_{Y_1 Y_2 | s}(A) &= \sum_{y_1 \in \mathcal{X}} P_{Y_1 | s}(y_1) \Pr((y_1, Y_2) \in A | s) \\
&\leq \sum_{y_1 \in \mathcal{X}} P_{Y_1 | s}(y_1) \min \{1, e^\varepsilon \Pr((y_1, Y_2) \in A) + \delta'\} \\
&\leq \sum_{y_1 \in \mathcal{X}} P_{Y_1 | s}(y_1) \min \{1, e^\varepsilon \Pr((y_1, Y_2) \in A)\} + \delta' \\
&\leq \sum_{y_1 \in \mathcal{X}} (e^\varepsilon P_{Y_1}(y_1) + \zeta(y_1)) \min \{1, e^\varepsilon \Pr((y_1, Y_2) \in A)\} + \delta' \\
&\leq \sum_{y_1 \in \mathcal{X}} e^\varepsilon P_{Y_1}(y_1) \min \{1, e^\varepsilon \Pr((y_1, Y_2) \in A)\} + \sum_{y_1 \in \mathcal{X}} \zeta(y_1) + \delta' \\
&\leq e^{2\varepsilon} \sum_{y_1 \in \mathcal{X}} P_{Y_1}(y_1) \Pr((y_1, Y_2) \in A) + \sum_{y_1 \in \mathcal{X}} \zeta(y_1) + \delta' \\
&\leq e^{2\varepsilon} P_{Y_1 Y_2}(A) + \sum_{y_1 \in \mathcal{X}} \zeta(y_1) + \delta' \\
&= e^{2\varepsilon} P_{Y_1 Y_2}(A) + \mathbf{E}_{e^\varepsilon}(P_{Y_1 | s} \| P_{Y_1}) + \delta' \\
&\leq e^{2\varepsilon} P_{Y_1 Y_2}(A) + 2\delta'
\end{aligned}$$

where $\delta' = \frac{\delta}{2}$ and $\zeta(a) := \left(P_{Y_1 | s}(a) - e^\varepsilon P_{Y_1}(a) \right)_+$ for any $a \in \mathcal{X}$. The last step follows from the fact that $\mathbf{E}_\gamma(P \| Q) = \sum_{a \in \mathcal{X}} (P(a) - \gamma Q(a))_+$. Consequently, we obtain that

$$P_{Y_1 Y_2 | s}(A) \leq e^{2\varepsilon} P_{Y_1 Y_2}(A) + \delta,$$

for any set $A \subset \mathcal{X}^2$ for $m = 2$. Repeating this argument $(m - 1)$ times, we can write

$$P_{Y | s}(A) \leq e^{m\varepsilon} P_Y(A) + \delta,$$

for any set $A \subset \mathcal{X}^m$ and $s \in \mathcal{S}$ from which we conclude

$$\mathbf{E}_{e^\varepsilon}(P_{Y | s} \| P_Y) \leq \delta.$$

A.3 Proof of Theorem 3

For any $\gamma \geq 1$ and $y^{j-1} \in \mathcal{X}^{j-1}$, we have

$$\begin{aligned} \mathbb{E}_\gamma(P_{Y_j|s,y^{j-1}} \| P_{Y_j|y^{j-1}}) &\leq \sup_{x^{j-1}} \mathbb{E}_\gamma(P_{Y_j|s,x^{j-1},y^{j-1}} \| P_{Y_j|x^{j-1},y^{j-1}}) \\ &= \sup_{x^{j-1}} \mathbb{E}_\gamma(P_{Y_j|s,x^{j-1}} \| P_{Y_j|x^{j-1}}), \end{aligned} \quad (21)$$

where the inequality follows from the convexity of \mathbb{E}_γ -divergence in each of its arguments (see, e.g., [52]). Notice that for any given $x^{j-1} \in \mathcal{X}^{j-1}$, we can write (with an abuse of notation)

$$\begin{aligned} \mathbb{E}_\gamma(P_{Y_j|s,x^{j-1}} \| P_{Y_j|x^{j-1}}) &= \int_B [P(dx_j|s, x^{j-1})\mathcal{N}(x_j, \lambda) - e^\varepsilon P(dx_j|x^{j-1})\mathcal{N}(x_j, \lambda)]_+ \\ &\quad + \int_{B^c} [P(dx_j|s, x^{j-1}) - e^\varepsilon P(dx_j|x^{j-1})]_+ \\ &= \int_B [P(dx_j|s, x^{j-1})\mathcal{N}(x_j, \lambda) - e^\varepsilon P(dx_j|x^{j-1})\mathcal{N}(x_j, \lambda)]_+ \end{aligned}$$

where we use B and B^c to write $B_j^\varepsilon(x^{j-1})$ and its complement. This demonstrates that the mass points corresponding to the event B^c do not contribute in the \mathbb{E}_γ -divergence.

Letting $P = P_{X_j|s,x^{j-1}}$ and $Q = P_{X_j|x^{j-1}}$ for a given x^{j-1} , it follows from above that

$$\sup_{x^{j-1}} \mathbb{E}_\gamma(P_{Y_j|s,x^{j-1}} \| P_{Y_j|x^{j-1}}) \leq \mathbb{E}_\gamma(P * \mathcal{N}(0, \lambda) \| Q * \mathcal{N}(0, \lambda)) = \mathbb{E} [\mathbb{E}_\gamma(\mathcal{N}(A, \lambda) \| \mathcal{N}(B, \lambda))], \quad (22)$$

where $*$ denotes the convolution operator and $A \sim P$ and $B \sim Q$ and the expectation is taken over any arbitrary coupling of P and Q (e.g., their product). It can be shown that

$$\mathbb{E}_\gamma(\mathcal{N}(\mu_1, \lambda^2 \mathbf{I}_r) \| \mathcal{N}(\mu_2, \lambda^2 \mathbf{I}_r)) = \mathbf{Q} \left(\frac{\log \gamma}{\beta} - \frac{1}{2} \beta \right) - \gamma \mathbf{Q} \left(\frac{\log \gamma}{\beta} + \frac{1}{2} \beta \right), \quad (23)$$

where $\mathbf{Q}(v) = \Pr(\mathcal{N}(0, 1) \geq v) = \int_v^\infty \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$ and $\beta = \frac{\|\mu_1 - \mu_2\|}{\lambda}$. Notice that the \mathbb{E}_γ -divergence between two Gaussian distributions depends on their means only through their differences.

$$\theta_\gamma(a, \lambda) \triangleq \mathbb{E}_\gamma(\mathcal{N}(\mu, \lambda^2 \mathbf{I}_r) \| \mathcal{N}(0, \lambda^2 \mathbf{I}_r)),$$

where $\|\mu\| = a$. According to (22), we can now write

$$\sup_{x^{j-1}} \mathbb{E}_\gamma(P_{Y_j|s,x^{j-1}} \| P_{Y_j|x^{j-1}}) \leq \sup_{a \in C} \theta_\gamma(\|a\|, \lambda) = \theta_\gamma(K, \lambda),$$

where the equality is due to the fact that $a \mapsto \theta_\gamma(a, \lambda)$ is increasing for a fixed λ . This, together with (21), implies

$$\mathbb{E}_\gamma(P_{Y_j|s,y^{j-1}} \| P_{Y_j|y^{j-1}}) \leq \theta_\gamma(K, \lambda),$$

and hence (14) is satisfied if $\theta_{e^\varepsilon}(K, \lambda) \leq \frac{\delta}{m}$.

A.4 Estimating Information Density using f -Divergences

Other f -divergence measures could also be used to estimate the information density by leveraging their dual representation [27]. Given a convex function f with $f(1) = 0$, the f -divergence $D_f(P \| Q) = \mathbb{E}_Q f\left(\frac{P}{Q}\right)$ can be expressed as

$$D_f(P \| Q) = \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(X)] - \mathbb{E}_Q[f^*(g(X))], \quad (24)$$

where $f^*(t) \triangleq \sup_{x \in \mathbb{R}} \{xt - f(t)\}$ is the Fenchel convex conjugate of f . It can be shown that the optimizer is the subdifferential $\partial f(\frac{P}{Q})$ which, in turn, is a non-decreasing function of $\frac{P}{Q}$. Thus, $D_f(P||Q)$ is also a candidate loss function in density ratio estimation problems.

A.5 Newey-McFadden Lemma

Lemma 1 ([55, Theorem 2.1]). *Given the extremum estimator $\hat{a} = \operatorname{argmax}_{a \in \mathcal{A}} \Lambda_n(a)$, if (i) \mathcal{A} is compact; (ii) there exists a limiting function $\Lambda(a)$ such that $\Lambda_n(a)$ converges to $\Lambda(a)$ in probability over \mathcal{A} ; (iii) $\Lambda(a)$ is continuous and has unique maximum at $a = a^*$, then \hat{a} is a consistent estimator of a^* .*

A.6 Proof of Theorem 4

Let the objective function of the extremum estimator be

$$\Lambda_n(g) \triangleq \mathbb{E}_{P_{S_n, X_n}} [g(S, X)] - \log \mathbb{E}_{P_{S_n} P_{X_n}} [e^{g(S, X)}]. \quad (25)$$

We prove this theorem by checking the properties of $\Lambda_n(g)$ according to Lemma 1. First, since Θ is compact and the mappings g_θ are continuous, the images $\mathcal{G}(\Theta)$ is also compact. Second, by triangular inequality, for $g \in \mathcal{G}(\Theta)$, we have

$$\begin{aligned} & |\Lambda_n(g) - (\mathbb{E}_{P_{S, X}} [g(S, X)] - \log \mathbb{E}_{P_S P_X} [e^{g(S, X)}])| \\ & \leq \sup_{g \in \mathcal{G}(\Theta)} |\mathbb{E}_{P_{S, X}} [g(S, X)] - \mathbb{E}_{P_{S_n, X_n}} [g(S, X)]| \\ & \quad + \sup_{g \in \mathcal{G}(\Theta)} |\log \mathbb{E}_{P_S P_X} [g(S, X)] - \log \mathbb{E}_{P_{S_n} P_{X_n}} [g(S, X)]|. \end{aligned} \quad (26)$$

Since the function g is uniformly bounded by M , i.e. $|g| \leq M$ for all θ, s and x , and logarithm is Lipschitz continuous with constant e^M in the interval $[e^{-M}, e^M]$, we have

$$|\log \mathbb{E}_{P_S P_X} [g(S, X)] - \log \mathbb{E}_{P_{S_n} P_{X_n}} [g(S, X)]| \leq e^M |\mathbb{E}_{P_S P_X} [g(S, X)] - \mathbb{E}_{P_{S_n} P_{X_n}} [g(S, X)]|. \quad (27)$$

Moreover, since \mathcal{G} is compact and g is continuous, the functions g and e^g satisfy the uniform law of large numbers [60]. Thus, Given $\eta > 0$, there exists an integer N such that for all $n \geq N$ and with probability one,

$$\sup_{g \in \mathcal{G}(\Theta)} |\mathbb{E}_{P_{S, X}} [g(S, X)] - \mathbb{E}_{P_{S_n, X_n}} [g(S, X)]| \leq \frac{\eta}{2}, \quad (28)$$

and

$$\sup_{g \in \mathcal{G}(\Theta)} |\log \mathbb{E}_{P_S P_X} [g(S, X)] - \log \mathbb{E}_{P_{S_n} P_{X_n}} [g(S, X)]| \leq \frac{\eta}{2} e^{-M}. \quad (29)$$

Summarizing (26)-(29), we have with probability one

$$|\Lambda_n(g) - (\mathbb{E}_{P_{S, X}} [g(S, X)] - \log \mathbb{E}_{P_S P_X} [e^{g(S, X)}])| \leq \eta. \quad (30)$$

In other words, there exists a limiting function $\Lambda(g) = \mathbb{E}_{P_{S, X}} [g(S, X)] - \log \mathbb{E}_{P_S P_X} [e^{g(S, X)}]$ such that $\Lambda_n(g)$ converges to $\Lambda(g)$ in probability.

Third, since $\Lambda(g) = \mathbb{E}_{P_{S, X}} [g(S, X)] - \log \mathbb{E}_{P_S P_X} [e^{g(S, X)}]$ consists of linear combinations (expectations) and continuous mappings (logarithm and exponential) of the continuous function g , $\Lambda(g)$ is continuous. Moreover, $\Lambda(g)$ has a unique optimizer g^* . Therefore, by Lemma 1, we know that with probability one,

$$|\hat{g}_n(s, x) - \hat{g}_\theta(s, x)| \leq \eta, \quad \forall s \in \mathcal{S}, x \in \mathcal{X}, \quad (31)$$

giving the consistency of the information density estimator.

A.7 Proof of Theorem 5

By Hoeffding's inequality [61], for all functions g bounded by M , i.e. $|g| \leq M$, we have

$$\Pr\{|\mathbb{E}_{P_{S_n, X_n}}[g(S, X)] - \mathbb{E}_{P_{S, X}}[g(S, X)]| > \frac{\eta}{4}\} \leq 2 \exp\left(-\frac{2n^2(\frac{\eta}{2})^2}{(2M)^2 n}\right) = 2 \exp\left(-\frac{n\eta^2}{32M^2}\right). \quad (32)$$

Moreover, since g_θ is parameterized by θ , we utilize the union bound [62, Lemma 2.2] to extend (32) for the parameters θ . For this purpose, recall that $\Theta \subset \mathbb{R}^d$ is compact and bounded by C , by the exterior covering number of bounded subspace [62, pp. 337], we know the r -covering number $N(r, \Theta)$ of Θ is upper bounded by

$$N(r, \Theta) \leq \left(\frac{2C\sqrt{d}}{r}\right)^d. \quad (33)$$

By (32) and (33), we have

$$\Pr\{\exists \theta_l \in \Theta \text{ s.t. } \sup_{g_\theta} |\mathbb{E}_{P_{S_n, X_n}}[g_\theta(S, X)] - \mathbb{E}_{P_{S, X}}[g_\theta(S, X)]| > \frac{\eta}{4}\} \leq 2N(r, \Theta) \exp\left(-\frac{n\eta^2}{32M^2}\right). \quad (34)$$

where θ_l is in the r -cover of Θ . Since $\mathcal{G}(\Theta)$ is compact, we can replace the supremum by maximum. To make $2N(r, \Theta) \exp\left(-\frac{n\eta^2}{32M^2}\right) < \delta$, we have

$$n > \frac{32M^2(\log N(r, \Theta) + \log \frac{2}{\delta})}{\eta^2}. \quad (35)$$

Now, let $r = \frac{\eta}{8L}$, and recall that g_θ is L -Lipschitz continuous with respect to θ , then for any $\theta \in \Theta$, we have with probability one

$$|g_\theta - g_{\theta_l}| \leq L|\theta - \theta_l| \leq Lr = L \times \frac{\eta}{8L} = \frac{\eta}{8}. \quad (36)$$

By triangular inequality, for any $\theta \in \Theta$, whenever $n > \frac{32M^2(d \log \frac{16LC\sqrt{d}}{\eta} + \log \frac{2}{\delta})}{\eta^2}$, we have with probability at least $1 - \delta$,

$$\begin{aligned} \max_{g_\theta} |\mathbb{E}_{P_{S_n, X_n}}[g_\theta(S, X)] - \mathbb{E}_{P_{S, X}}[g_\theta(S, X)]| &\leq \max_{g_\theta} |\mathbb{E}_{P_{S_n, X_n}}[g_\theta(S, X)] - \mathbb{E}_{P_{S_n, X_n}}[g_{\theta_l}(S, X)]| \\ &\quad + \max_{g_\theta} |\mathbb{E}_{P_{S_n, X_n}}[g_{\theta_l}(S, X)] - \mathbb{E}_{P_{S, X}}[g_{\theta_l}(S, X)]| \\ &\quad + \max_{g_\theta} |\mathbb{E}_{P_{S, X}}[g_\theta(S, X)] - \mathbb{E}_{P_{S, X}}[g_{\theta_l}(S, X)]| \\ &\leq \frac{\eta}{8} + \frac{\eta}{4} + \frac{\eta}{8} = \frac{\eta}{2} \end{aligned} \quad (37)$$

Therefore, we have

$$\Pr\{\max_{g_\theta} |\mathbb{E}_{P_{S_n, X_n}}[g_\theta(S, X)] - \mathbb{E}_{P_{S, X}}[g_\theta(S, X)]| \leq \frac{\eta}{2}\} \geq 1 - \delta. \quad (38)$$

Similarly, starting from

$$\begin{aligned} &\Pr\{\exists \theta_l \in \Theta \text{ s.t. } |\log \mathbb{E}_{P_{S_n} P_{X_n}}[e^{g_{\theta_l}(S, X)}] - \log \mathbb{E}_{P_S P_X}[e^{g_{\theta_l}(S, X)}]|\geq \frac{\eta}{4}\} \\ &\leq 2N(r, \Theta) \exp\left(-\frac{n\eta^2}{32M^2}\right), \end{aligned} \quad (39)$$

we also conclude that for any $\theta \in \Theta$, whenever $n > \frac{32M^2(d \log \frac{16LC\sqrt{d}}{\eta} + \log \frac{2}{\delta})}{\eta^2}$, we have with probability at least $1 - \delta$,

$$\Pr\{\max_{g_\theta} |\log \mathbb{E}_{P_{S_n, X_n}} [E^{g_\theta(S, X)}] - \log \mathbb{E}_{P_{S, X}} [e^{g_\theta(S, X)}]| \leq \frac{\eta}{2}\} \geq 1 - \delta. \quad (40)$$

Summarizing (38) and (40), whenever $n > \frac{32M^2(d \log \frac{16LC\sqrt{d}}{\eta} + \log \frac{2}{\delta})}{\eta^2}$, for any $\theta \in \Theta$, we have

$$\begin{aligned} & \Pr\{|\max_n \Lambda_n(\hat{g}_n(s, x)) - \max \Lambda(g(s, x))| \leq \eta\} \\ & \geq \Pr\{\max_{g_\theta} |\mathbb{E}_{P_{S_n, X_n}} [g_\theta(S, X)] - \mathbb{E}_{P_{S, X}} [g_\theta(S, X)]| + \max_{g_\theta} |\log \mathbb{E}_{P_{S_n, X_n}} [E^{g_\theta(S, X)}] \\ & \quad - \log \mathbb{E}_{P_{S, X}} [e^{g_\theta(S, X)}]| \leq \eta\} \\ & \geq 1 - \delta. \end{aligned} \quad (41)$$

The trimmed information density estimator, in this sense, gives a trimmed (clipped) information density, i.e. $|\hat{g}_n(s, x) - g^*(s, x)| \leq \eta$ if $g^*(s, x) \leq M$ and $|\hat{g}_n(s, x) - g^*(s, x)| \geq \eta$ otherwise. By the concentration of the information density [63], we also know the probability that the information density is clipped is upper bounded, i.e.

$$\Pr\{|g^*(s, x)| \geq M\} \leq e^{-M}. \quad (42)$$

Therefore, whenever $n > \frac{32M^2(d \log \frac{16LC\sqrt{d}}{\eta} + \log \frac{2}{\delta})}{\eta^2}$, for all $s \in \mathcal{S}$ and $x \in \mathcal{X}$, we have

$$\Pr\{|\hat{g}_n(s, x) - g^*(s, x)| \leq \eta\} \geq 1 - \delta \geq 1 - e^{-M}, \quad (43)$$

by choosing $\delta \leq e^{-M}$, and the desired result follows.

B Experimental Details

In this section, we provide detailed experimental setups including architecture of the function g in TIDE, training details for the experiments shown in the main text.

B.1 GENKI-4K Smiling Dataset

The GENKI-4K smiling dataset [31] contains 2400 colorful images for training and 600 for test, where each image, viewed as X , is a 64×64 pixels face that is smiling ($S = 1$) or not ($S = 0$).

Since the inputs of the encoder TIDE are images, we use adopt a convolutional neural net with three convolutional layers, two fully-connected layers, and a readout layer. The convolutional layers have kernels with dimension (5, 5, 3, 64), (5, 5, 64, 64), and (3, 3, 64, 128) respectively. After flattening the output of the third convolutional layer, we feed the output to two fully-connected layers with 384 and 192 neurons respectively. We train for 100 epochs using `AdagradOptimizer` with learning rate 0.0001 and batch size 256, and achieve $I(S, X) = 0.594 < H(S) = 1$ bits.

The adversary we used here is also a convolutional neural net with identical structure as the TIDE with the difference that the objective is the cross-entropy loss for classification, and is trained for 150 epochs using `AdagradOptimizer` with learning rate 0.005 and batch size 256.

Table 3: WMAE of the information density estimation on Gaussian synthetic data ($M = 5$).

		Empirical Plug-In Estimator				Kernel Density Estimator				TIDE			
$d \backslash \rho$		0.0	0.1	0.2	0.5	0.0	0.1	0.2	0.5	0.0	0.1	0.2	0.5
1		0.466	0.509	1.092	1.821	0.252	0.434	0.973	1.395	0.005	0.007	0.011	0.057
10		5.305	7.613	9.704	18.245	2.869	4.076	6.698	11.496	1.010	1.216	1.884	2.503

B.2 Celebrity Attributes (CelebA) Dataset

The CelebA dataset [32] contains 202599 colorful images, where each image is a 218×178 pixels face of a celebrity with 40 distinct binary labels, including **smiling**, **gender**, **Arched Eyebrows**, etc. We select 100000 face images as X and the private attribute S as smiling or not.

Since the inputs of the encoder TIDE are images, we use adopt a convolutional neural net with five convolutional layers, two fully-connected layers, and a readout layer. The convolutional layers have kernels with dimension $(5, 5, 3, 64)$, $(5, 5, 64, 64)$, $(3, 3, 128, 128)$, $(3, 3, 128, 128)$, and $(3, 3, 64, 128)$ respectively. After flattening the output of the third convolutional layer, we feed the output to two fully-connected layers with 384 and 192 neurons respectively. We train for 100 epochs using **AdagradOptimizer** with learning rate 0.005 and batch size 64, and achieve $I(S, X) = 0.967 \approx H(S) = 1$ bits.

The adversaries we used for emotion and gender detection here are also convolutional neural nets with identical structure as the TIDE with the difference that the objective is the cross-entropy loss for classification. We train the adversaries for 300 epochs using **AdagradOptimizer** with learning rate 0.001 and batch size 2000.

B.3 Politically-Biased Tweets

We collect 75946 tweets from more than 20 online publishers (e.g. CNN, Bloomberg, New York Times) using the Twitter API, and determine its private attribute S as the political bias of being right-wing ($S = 0$) and left-wing ($S = 1$) according to [33]. We clean up the tweets to only keep meaningful terms (i.e. pieces of words), and use bag-of-words representation [58] to tokenize all the pieces of words for each tweet according to term frequency, ending up with 24657 words (x_j). We order the x_j by the order it appears in the training texts of the Tweets.

The TIDE is a simple feed-forward neural network consists of three hidden layers with ReLU activation with 100 neurons for each hidden layer, and a readout layer with 32 neurons. We train for 50 epochs using **AdagradOptimizer** with learning rate 0.005 and batch size 128, and achieve $I(S; X) = 0.645$ bits.

C Additional Experiments on Synthetic Data

We apply the TIDE in Section 3 on Gaussian synthetic data to estimate the trimmed information density with limited number of samples and $M = 5$. We consider two d -dimensional multivariate standard Normal random variables S and X , with pairwise correlation $\text{corr}(S_i, X_j) = \rho \mathbf{1}_{\{i=j\}}$, $\rho \in (-1, 1)$, $1 \leq i, j \leq d$. Since the KL divergence is invariant to continuous bijective transformations of the considered variables, it is sufficient to consider S and X with standard Normal marginals. We generate 3000 samples with 70% – 30% train-test split accordingly. The TIDE is a simple feed-forward neural network consists of three hidden layers with ReLU activation with 100, 50, 50 neurons for each hidden layer, and a readout layer with 50 neurons. We jointly train over the entire training set for 3000 epochs using **AdagradOptimizer** with learning rate 0.005.

We compare the plug-in estimator using empirical distributions (with 30 bins for quantization), the Gaussian kernel density estimator [64], and the TIDE using 3k samples, and report the Weighted Mean Absolute Error (WMAE) of the information density in Table 3, where the weights are chosen as the ground true joint distributions and each number in the table is averaged over 10 repeated experiments. Note that since the Normal random variable is continuous, quantized empirical distribution gives loose estimate. The kernel density estimator performs better than the plug-in estimator but worse than the TIDE due to limited number of samples.

References

- [1] M. Bertran, N. Martinez, A. Papadaki, Q. Qiu, M. Rodrigues, G. Reeves, and G. Sapiro, “Adversarially learned representations for information obfuscation and inference,” in *Proc. of International Conference on Machine Learning (ICML)*, 2019.
- [2] X. Chen, T. Navidi, S. Ermon, and R. Rajagopal, “Distributed generation of privacy preserving data with user customization,” *arXiv preprint arXiv:1904.09415*, 2019.
- [3] C. Huang, P. Kairouz, X. Chen, L. Sankar, and R. Rajagopal, “Generative adversarial privacy,” *arXiv preprint arXiv:1807.05306*, 2018.
- [4] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, “Learning fair representations,” in *International Conference on Machine Learning*, 2013, pp. 325–333.
- [5] F. du Pin Calmon and N. Fawaz, “Privacy against statistical inference,” in *Proc. of IEEE Allerton Conference on Communication, Control, and Computing (Allerton)*, 2012, pp. 1401–1408.
- [6] S. Asoodeh, M. Diaz, F. Alajaji, and T. Linder, “Estimation efficiency under privacy constraints,” *IEEE Transactions on Information Theory*, vol. 65, no. 3, pp. 1512–1534, 2018.
- [7] I. Issa, A. B. Wagner, and S. Kamath, “An operational approach to information leakage,” *arXiv preprint arXiv:1807.07878*, 2018.
- [8] H. Hsu, S. Asoodeh, S. Salamatian, and F. P. Calmon, “Generalizing bottleneck problems,” in *Proc. of IEEE International Symposium on Information Theory (ISIT)*, 2018.
- [9] M. Diaz, H. Wang, F. P. Calmon, and L. Sankar, “On the robustness of information-theoretic privacy measures and mechanisms,” *arXiv preprint arXiv:1811.06057*, 2018.
- [10] Y. O. Basciftci, Y. Wang, and P. Ishwar, “On privacy-utility tradeoffs for constrained data release mechanisms,” in *2016 Information Theory and Applications Workshop (ITA)*. IEEE, 2016, pp. 1–6.
- [11] C. Huang, P. Kairouz, X. Chen, L. Sankar, and R. Rajagopal, “Context-aware generative adversarial privacy,” *Entropy*, vol. 19, no. 12, p. 656, 2017.
- [12] M. S. Pinsker, *Information and information stability of random variables and processes*. San Francisco: Holden-Day, 1964.
- [13] T. S. Han and S. Verdú, “Approximation theory of output statistics,” *IEEE Transactions on Information Theory*, vol. 39, no. 3, pp. 752–772, 1993.

- [14] A. Evfimievski, J. Gehrke, and R. Srikant, “Limiting privacy breaches in privacy preserving data mining,” in *Proc. of ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, 2003.
- [15] M. Bun and T. Steinke, “Concentrated differential privacy: Simplifications, extensions, and lower bounds,” in *Theory of Cryptography*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2016, pp. 635–658.
- [16] C. Dwork and G. N. Rothblum, “Concentrated differential privacy,” *arXiv preprint arXiv:1603.01887*, 2016.
- [17] B. Balle and Y.-X. Wang, “Improving the Gaussian mechanism for differential privacy: Analytical calibration and optimal denoising,” in *in Proc. of the International Conference on Machine Learning*, 2018.
- [18] A. D. Sarwate and K. Chaudhuri, “Signal processing and machine learning with differential privacy: Algorithms and challenges for continuous data,” *IEEE signal processing magazine*, vol. 30, no. 5, pp. 86–94, 2013.
- [19] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate, “Differentially private empirical risk minimization,” *Journal of Machine Learning Research*, vol. 12, no. Mar, pp. 1069–1109, 2011.
- [20] I. Csiszár, “Information-type measures of difference of probability distributions and indirect observations,” *Studia Sci. Math. Hungar.*, vol. 2, pp. 299–318, 1967.
- [21] G. Valiant and P. Valiant, “Estimating the unseen: an $n/\log(n)$ -sample estimator for entropy and support size, shown optimal via new clts,” in *Proc. of ACM symposium on Theory of computing (STOC)*, 2011.
- [22] Y. Wu and P. Yang, “Minimax rates of entropy estimation on large alphabets via best polynomial approximation,” *IEEE Transactions on Information Theory*, vol. 62, no. 6, pp. 3702–3720, 2016.
- [23] W. Gao, S. Kannan, S. Oh, and P. Viswanath, “Estimating mutual information for discrete-continuous mixtures,” in *Proc. of Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5986–5997.
- [24] V. Vapnik, *The nature of statistical learning theory*. Springer science & business media, 2013.
- [25] I. Belghazi, S. Rajeswar, A. Baratin, R. D. Hjelm, and A. Courville, “Mine: mutual information neural estimation,” *arXiv preprint arXiv:1801.04062*, 2018.
- [26] S. Liu, A. Takeda, T. Suzuki, and K. Fukumizu, “Trimmed density ratio estimation,” in *Proc. of Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [27] X. Nguyen, M. J. Wainwright, and M. I. Jordan, “Estimating divergence functionals and the likelihood ratio by convex risk minimization,” *IEEE Transactions on Information Theory*, vol. 56, no. 11, pp. 5847–5861, 2010.
- [28] R. McPherson, R. Shokri, and V. Shmatikov, “Defeating image obfuscation with deep learning,” *arXiv preprint arXiv:1609.00408*, 2016.

- [29] S. J. Oh, M. Fritz, and B. Schiele, “Adversarial image perturbation for privacy protection a game theory perspective,” in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 1491–1500.
- [30] Z. Wu, Z. Wang, Z. Wang, and H. Jin, “Towards privacy-preserving visual recognition via adversarial training: A pilot study,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 606–624.
- [31] T. MPLab, “The MPLab GENKI Database, GENKI-4K Subset,” <http://mplab.ucsd.edu>, 2009.
- [32] Z. Liu, P. Luo, X. Wang, and X. Tang, “Deep learning face attributes in the wild,” in *Proc. of International Conference on Computer Vision (ICCV)*, 2015.
- [33] A. Rachez, “Predicting political bias with python,” <https://medium.com/linalgo/predict-political-bias-using-python-b8575eedef13>, 2017, accessed: 2019-03-21.
- [34] J. Liao, O. Kosut, L. Sankar, and F. P. Calmon, “A tunable measure for information leakage,” *arXiv preprint arXiv:1806.03332*, 2018.
- [35] L. Sankar, S. R. Rajagopalan, and H. V. Poor, “Utility-privacy tradeoffs in databases: An information-theoretic approach,” *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 6, pp. 838–852, 2013.
- [36] A. Makhdoumi, S. Salamatian, N. Fawaz, and M. Médard, “From the information bottleneck to the privacy funnel,” in *Proc. of IEEE Information Theory Workshop (ITW)*, 2014.
- [37] K. Nissim, S. Raskhodnikova, and A. Smith, “Smooth sensitivity and sampling in private data analysis,” in *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*. ACM, 2007, pp. 75–84.
- [38] P. Cuff and L. Yu, “Differential privacy as a mutual information constraint,” in *Proc. of the ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2016.
- [39] I. Mironov, “Rényi differential privacy,” in *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. IEEE, 2017, pp. 263–275.
- [40] D. Kifer and A. Machanavajjhala, “No free lunch in data privacy,” in *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’11. New York, NY, USA: ACM, 2011, pp. 193–204. [Online]. Available: <http://doi.acm.org/10.1145/1989323.1989345>
- [41] M. Yamada, T. Suzuki, T. Kanamori, H. Hachiya, and M. Sugiyama, “Relative density-ratio estimation for robust distribution comparison,” in *Proc. of Advances in neural information processing systems (NeurIPS)*, 2011.
- [42] A. Smola, L. Song, and C. H. Teo, “Relative novelty detection,” in *Proc. of International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2009.
- [43] M. Sugiyama, M. Krauledat, and K.-R. Mäñzler, “Covariate shift adaptation by importance weighted cross validation,” *Journal of Machine Learning Research (JMLR)*, vol. 8, no. May, pp. 985–1005, 2007.

- [44] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [45] M. Sugiyama, T. Suzuki, and T. Kanamori, *Density ratio estimation in machine learning*. Cambridge University Press, 2012.
- [46] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.
- [47] J. Liu, P. Cuff, and S. Verdú, “ e_γ -resolvability,” *IEEE Transactions on Information Theory*, vol. 63, no. 5, pp. 2629–2658, May 2017.
- [48] Y. Polyanskiy, H. V. Poor, and S. Verdú, “Channel coding rate in the finite blocklength regime,” *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2307–2359, 2010.
- [49] G. Barthe and F. Olmedo, “Beyond differential privacy: Composition theorems and relational logic for f-divergences between probabilistic programs,” in *Automata, Languages, and Programming*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 49–60.
- [50] B. Balle, G. Barthe, M. Gaboardi, and J. Geumlek, “Privacy amplification by mixing and diffusion mechanisms,” *ArXiv*, vol. abs/1905.12264, 2019.
- [51] Y. Polyanskiy and S. Verdú, “Arimoto channel coding converse and Rényi divergence,” in *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, Sep. 2010, pp. 1327–1333.
- [52] I. Sason and S. Verdú, “f-divergence inequalities,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 5973–6006, 2016.
- [53] M. D. Donsker and S. S. Varadhan, “Asymptotic evaluation of certain markov process expectations for large time. iv,” *Communications on Pure and Applied Mathematics*, vol. 36, no. 2, pp. 183–212, 1983.
- [54] T. Amemiya, “Asymptotic properties of extremum estimators,” *Advanced econometrics, Harvard university press*, 1985.
- [55] W. K. Newey and D. McFadden, “Large sample estimation and hypothesis testing,” *Handbook of econometrics*, vol. 4, pp. 2111–2245, 1994.
- [56] K. W. Church and P. Hanks, “Word association norms, mutual information, and lexicography,” *Computational linguistics*, vol. 16, no. 1, pp. 22–29, 1990.
- [57] M. Alzantot, Y. Sharma, A. Elgohary, B.-J. Ho, M. Srivastava, and K.-W. Chang, “Generating natural language adversarial examples,” *arXiv preprint arXiv:1804.07998*, 2018.
- [58] C. Manning, P. Raghavan, and H. Schütze, “Introduction to information retrieval,” *Natural Language Engineering*, vol. 16, no. 1, pp. 100–103, 2010.
- [59] I. Žliobaitė and B. Custers, “Using sensitive personal data may be necessary for avoiding discrimination in data-driven decision models,” *Artificial Intelligence and Law*, vol. 24, no. 2, pp. 183–201, 2016.
- [60] S. van de Geer, *Empirical Processes in M-estimation*. Cambridge university press, 2000, vol. 6.

- [61] W. Hoeffding, “Probability inequalities for sums of bounded random variables,” in *The Collected Works of Wassily Hoeffding*. Springer, 1994, pp. 409–426.
- [62] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [63] Y. Polyanskiy and Y. Wu, “Lecture notes on information theory,” *Lecture Notes for ECE563 (UIUC) and*, vol. 6, pp. 2012–2016, 2014.
- [64] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.