

A THEORY OF USABLE INFORMATION UNDER COMPUTATIONAL CONSTRAINTS

Yilun Xu

CFCS, Peking University
xuyilun@pku.edu.cn

Shengjia Zhao

Stanford University
sjzhao@stanford.edu

Jiaming Song

Stanford University
tsong@cs.stanford.edu

Russell Stewart

russell.sb.nebel@gmail.com

Stefano Ermon

Stanford University
ermon@cs.stanford.edu

ABSTRACT

We propose a new framework for reasoning about information in complex systems. Our foundation is based on a variational extension of Shannon’s information theory that takes into account the modeling power and computational constraints of the observer. The resulting *predictive \mathcal{V} -information* encompasses mutual information and other notions of informativeness such as the coefficient of determination. Unlike Shannon’s mutual information and in violation of the data processing inequality, \mathcal{V} -information can be created through computation. This is consistent with deep neural networks extracting hierarchies of progressively more informative features in representation learning. Additionally, we show that by incorporating computational constraints, \mathcal{V} -information can be reliably estimated from data even in high dimensions with PAC-style guarantees. Empirically, we demonstrate predictive \mathcal{V} -information is more effective than mutual information for structure learning and fair representation learning.

1 INTRODUCTION

Extracting actionable *information* from noisy, possibly redundant, and high-dimensional data sources is a key computational and statistical challenge at the core of AI and machine learning. Information theory, which lies at the foundation of AI and machine learning, provides a conceptual framework to characterize information in a mathematically rigorous sense (Shannon & Weaver, 1948; Cover & Thomas, 1991). However, important computational aspects are not considered in information theory. To illustrate this, consider a dataset of encrypted messages intercepted from an opponent. According to information theory, these encrypted messages have high mutual information with the opponent’s plans. Indeed, with infinite computation, the messages can be decrypted and the plans revealed. Modern cryptography originated from this observation by Shannon that perfect secrecy is (essentially) impossible if the adversary is computationally unbounded (Shannon & Weaver, 1948). This motivated cryptographers to consider restricted classes of adversaries that have access to limited computational resources (Pass & Shelat, 2010). More generally, it is known that information theoretic quantities can be expressed in terms of betting games (Cover & Thomas, 1991). For example, the (conditional) entropy of a random variable X is directly related to how predictable X is in a certain betting game, where an agent is rewarded for correct guesses. Yet, the standard definition unrealistically assumes agents are computationally unbounded, i.e., they can employ arbitrarily complex prediction schemes.

Leveraging modern ideas from variational inference and learning (Ranganath et al., 2013; Kingma & Welling, 2013; LeCun et al., 2015), we propose an alternative formulation based on realistic computational constraints that is in many ways closer to our intuitive notion of information, which we term *predictive \mathcal{V} -information*. Without constraints, predictive \mathcal{V} -information specializes to classic mutual information. Under natural restrictions, \mathcal{V} -information specializes to other well-known notions of predictiveness, such as the coefficient of determination (R^2). A consequence of this new formulation is that computation can “create usable information” (e.g., by decrypting the intercepted messages), invalidating the famous data processing inequality. This generalizes the idea that clever

feature extraction enables prediction with extremely simple (e.g., linear) classifiers, a key notion in modern representation and deep learning (LeCun et al., 2015).

As an additional benefit, we show that predictive information can be estimated with statistical guarantees using the Probably Approximately Correct framework (Valiant, 1984). This is in sharp contrast with Shannon information, which is well known to be difficult to estimate for high dimensional or continuous random variables (Battiti, 1994). Theoretically we show that the statistical guarantees of estimating information translate to statistical guarantees for a variant of the Chow-Liu algorithm for structure learning. In practice, when the observer employs deep neural networks as a prediction scheme, \mathcal{V} -information outperforms methods that approximate Shannon information in various applications, including Chow-Liu tree construction in high dimension and gene regulatory network inference.

2 DEFINITIONS AND NOTATIONS

To formally define the predictive \mathcal{V} -information, we begin with a formal model of a computationally bounded agent trying to predict the outcome of a real-valued random variable. The agent is either provided another real-valued random variable as side information, or provided no side information. We use X and Y to denote the sample spaces \mathcal{X} and \mathcal{Y} respectively (while assuming they are separable), and use $\mathcal{P}(X)$ to denote the set of all probability measures over the Borel algebra \mathcal{A} on $(\mathcal{X}, \mathcal{A})$ (similarly defined for Y).

Definition 1 (Predictive Family)¹ Let $\mathcal{V} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f \in \mathcal{V}\}$. We say that \mathcal{V} is a predictive family if it satisfies

$$\forall f \in \mathcal{V}; \forall P \in \mathcal{P}(\mathcal{X}) \text{ range}(f); \forall Q \in \mathcal{P}(\mathcal{Y}); \exists s: \exists x \in \mathcal{X}; f(x) = P; f(Q) = P \quad (1)$$

A predictive family is a set of predictive models the agent is allowed to use, e.g., due to computational or statistical constraints. We refer to the additional condition in Eq.(1) as *optional ignorance*. Intuitively, it means that the agent can, in the context of the prediction game we define next, ignore the side information if she chooses to.

Definition 2 (Predictive conditional \mathcal{V} -entropy) Let X, Y be two random variables taking values in \mathcal{X}, \mathcal{Y} , and \mathcal{V} be a predictive family. Then the predictive conditional entropy is defined as

$$H_{\mathcal{V}}(Y|X) = \inf_{f \in \mathcal{V}} E_{x,y} [-\log f(x)(y)]$$

$$H_{\mathcal{V}}(Y|?) = \inf_{f \in \mathcal{V}} E_y [-\log f(?) (y)]$$

We additionally call $H_{\mathcal{V}}(Y|?)$ the \mathcal{V} -entropy, and also denote it as $H_{\mathcal{V}}(Y)$.

In our notation f is a function $\mathcal{X} \rightarrow \mathcal{Y}$, so $f(x) \in \mathcal{P}(\mathcal{Y})$ is a probability measure on \mathcal{Y} chosen based on the received side information x (we use $f(x)$ instead of the more conventional $f(x|y)$); and $f(x)(y) \in \mathbb{R}$ is the value of the density evaluated at $y \in \mathcal{Y}$. Intuitively, \mathcal{V} (conditional) entropy is the smallest expected negative log-likelihood that can be achieved predictively on observation (side information) X (or no side information $?$), using models from \mathcal{V} . Eq.(1) means that whenever the agent can use to predict Y 's outcomes, it has the option to ignore the input, and it does not matter whether X is observed or not.

Definition 2 generalizes several known definitions of uncertainty. For example, as shown in proposition 2, if the \mathcal{V} is the largest possible predictive family that includes all possible models, i.e. $\mathcal{V} = \mathcal{P}(\mathcal{Y})$, then Definition 2 reduces to Shannon entropy: $H_{\mathcal{V}}(Y|X) = H(Y|X)$ and $H_{\mathcal{V}}(Y|?) = H(Y) = H(Y)$. By choosing more restrictive families, we recover several other notions of uncertainty such as trace of covariance, as will be shown in Proposition 1.

Shannon mutual information is a measure of changes in entropy when conditioning on new variables:

$$I(X; Y) = H(Y) - H(Y|X) = H(Y) - H(Y|X) \quad (2)$$

Here, we will use predictive \mathcal{V} -entropy to define an analogous quantity, $I_{\mathcal{V}}(X; Y)$, to represent the change in predictability of an output variable Y when given side information X .

¹Regularity Conditions: To minimize technical overhead we restrict our discussion only to distributions with probability density functions (PDF) or probability mass functions (PMF) with respect to the underlying measure. Also $\mathcal{P}(\mathcal{X})$

Definition 3 (Predictive V-information). Let X, Y be two random variables taking values in \mathcal{X}, \mathcal{Y} , and V be a predictive family. The predictive information from X to Y is defined as

$$I_V(X \rightarrow Y) = H_V(Y|?) - H_V(Y|X) \quad (3)$$

2.1 IMPORTANT SPECIAL CASES

Several important notions of uncertainty and predictiveness are special cases of our definition. Note that when we are defining V-entropy of a random variable in sample space $\mathcal{X} \subseteq \mathbb{R}^d$ (without side information), out of convenience we can assume the empty $X = ?$ (this does not violate our requirement that $\mathcal{X} \subseteq \mathbb{R}^d$)

Proposition 1. For V-entropy and V-information, we have

1. Let $\mathcal{X} = \mathcal{Y}$ be as in Def. 1. Then $H_V(Y)$ is the Shannon entropy, $H_V(Y|X)$ is the Shannon conditional entropy, and $I_V(X \rightarrow Y)$ is the Shannon mutual information.
2. Let $\mathcal{Y} = \mathbb{R}^d$ and $V = \{f : f(y) = \int_{\mathcal{X}} p_j(y) dy\}$, where P is the distribution with density $p_j(y) = \frac{1}{Z} e^{-k_j y}$ where $Z = \int_{\mathcal{X}} e^{-k_j y} dy$, then the V-entropy of a random variable Y equals its mean absolute deviation, up to an additive constant.
3. Let $\mathcal{Y} = \mathbb{R}^d$ and $V = \{f : f(y) = \sum_{j=1}^d y_j\}$, $\mathcal{X} = \mathbb{R}^d$, $\mathcal{X} = \mathbb{R}^d$, $\mathcal{X} = \mathbb{R}^d$, then the V-entropy of a random variable Y equals the trace of its covariance matrix $\text{tr}(\text{Cov}(Y))$, up to an additive constant.
4. Let $V = \{f : f(y) = \sum_{i=1}^d Q_i(y)\}$, where Q_i is a distribution in a minimal exponential family with sufficient statistics $s : \mathcal{Y} \rightarrow \mathbb{R}^d$ and set of natural parameters. For a random variable Y with expected sufficient statistics $\bar{y} = E[s(Y)]$, the V-entropy of Y is the maximum Shannon entropy over all random variables with identical expected sufficient statistics, i.e. $E[s(\tilde{Y})] = \bar{y}$.
5. Let $\mathcal{Y} = \mathbb{R}^d$, \mathcal{X} be any vector space, and $V = \{f : f(x) = \sum_{i=1}^d \lambda_i x_i\}$, $\mathcal{X} = \mathbb{R}^d$, $\mathcal{X} = \mathbb{R}^d$, $\mathcal{X} = \mathbb{R}^d$, where \mathcal{L} is the set of linear functions $f : \mathcal{X} \rightarrow \mathbb{R}^d$, then V-information $I_V(X \rightarrow Y)$ equals the (unnormalized) maximum coefficient of determination $R^2 = \text{tr}(\text{Cov}(Y))$ for linear regression.

The trace of covariance represents a natural notion of uncertainty – for example, a random variable with zero variance (when $\mathcal{X} = \mathbb{R}^d$, $\text{tr}(\text{Cov}(Y)) = \text{Var}(Y)$) is trivial to predict. Proposition 1.3 shows that the trace of covariance corresponds to a notion of surprise (in the Shannon sense) for an agent restricted to make predictions using certain Gaussian models. More broadly, a similar analogy can be drawn for other exponential families of distributions. In the same spirit, the coefficient of determination, also known as the fraction of variance explained, represents a natural notion of informativeness for computationally bounded agents. Also note that in the case of Proposition 1.4, the V-entropy is invariant if the expected sufficient statistics remain the same.

3 PROPERTIES OF V-INFORMATION

3.1 ELEMENTARY PROPERTIES

We first show several elementary properties of V-entropy and V-information. In particular, V-information preserves many properties of Shannon information that are desirable in a machine learning context. For example, mutual information (and information) should be non-negative as conditioning on additional side information should not reduce an agent's ability to predict.

Proposition 2. Let Y and X be any random variables on \mathcal{Y} and \mathcal{X} , and V and U be any predictive families, then we have

1. Monotonicity. If $V \subseteq U$, then $H_V(Y) \geq H_U(Y)$, $H_V(Y|X) \geq H_U(Y|X)$.
2. Non-Negativity $I_V(X \rightarrow Y) \geq 0$.
3. Independence. If X is independent of Y , $I_V(X \rightarrow Y) = I_V(Y \rightarrow X) = 0$.

The optional ignorance requirement in Eq.(1) is a technical condition needed for these properties to hold. Intuitively, it guarantees that conditioning on side information does not restrict the class of densities the agent can use to predict. This property is satisfied by many existing machine learning models, often by setting some weights to zero so that an input is effectively ignored.

3.2 ON THE PRODUCTION OF INFORMATION THROUGH PREPROCESSING

The Data Processing Inequality guarantees that computing on data cannot increase its mutual information with other random variables. Formally, letting $f: X \rightarrow X'$ be any function, $I(f(X); Y)$ cannot have higher mutual information with Y than $X: I(X; Y) = I(X; Y)$. But is this property desirable? In analyzing optimal communication, yes - it demonstrates a fundamental limit to the number of bits that can be transmitted through a communication channel. However, we argue that in machine learning settings this property is less appropriate.

Consider an RSA encryption scheme where the public key is known. Given plain text and its corresponding encrypted text, if we have in finite computation, we can perfectly compute one from the other. Therefore, the plain text and the encrypted text should have identical Shannon mutual information with respect to any label we want to predict. However, to any human (or machine learning algorithm), it is certainly easier to predict the label from the plain text than the encrypted text. In other words, decryption increases a human's ability to predict the label: processing increases the "usable information". More formally, denoting \mathcal{F} as the decryption algorithm and \mathcal{G} as a class of natural language processing functions, we have $I_{\mathcal{F}}(X; Y) > I_{\mathcal{G}}(X; Y) = 0$.

As another example, consider the mutual information between an image's pixels and its label. Due to data processing inequality, we cannot expect to use a function to map raw pixels to "features" that have higher mutual information with the label. However, the fundamental principle of representation learning is precisely the ability to learn predictive features — functions of the raw inputs that enable predictions with higher accuracy. Because of this key difference between information and Shannon information, machine learning practices such as representation learning can be justified in the information theoretic context.

3.3 ON THE ASYMMETRY OF PREDICTIVE V-INFORMATION

V-information also captures the intuition that sometimes, it is easy to predict Y from X but not vice versa. In fact, modern cryptography is founded on the assumption that certain functions $h: X \rightarrow Y$ are one-way, meaning that there exists a polynomial algorithm to compute $h(x)$ but no polynomial algorithm to compute $h^{-1}(y)$. This means that if \mathcal{V} contains all polynomial-time computable functions, then $I_{\mathcal{V}}(X; h(X)) = I_{\mathcal{V}}(h(X); X)$.

This property is also reasonable in the machine learning context. For example, several important methods for causal discovery (Peters et al., 2017) rely on this asymmetry: if X causes Y , then usually it is easier to predict Y from X than vice versa; another commonly used assumption is that X can be accurately modeled by a Gaussian distribution, while Y cannot (Pearl, 2000).

4 PAC GUARANTEES FOR V-INFORMATION ESTIMATION

For many practical applications of mutual information (e.g., structure learning), we do not know the joint distribution of $X; Y$, so cannot directly compute the mutual information. Instead we only have samples $(x_i; y_i)_{i=1}^N$ of $X; Y$ and need to estimate mutual information from data.

Shannon information is notoriously difficult to estimate for high dimensional random variables. Although non-parametric estimators of mutual information exist (Kraskov et al., 2004; Darbellay & Vajda, 1999; Gao et al., 2017), these estimators do not scale to high dimensions. Several variational estimators for Shannon information have been recently proposed (van den Oord et al., 2018; Nguyen et al., 2010; Belghazi et al., 2018), but have two shortcomings: due to their variational assumptions, their bias/variance tradeoffs are poorly understood and they are still not efficient enough for high dimensional problems. For example, the CPC estimator suffers from large bias, since its estimates saturate at $\log N$ where N is the batch size (van den Oord et al., 2018; Poole et al., 2019); the NWJ estimator suffers from large variance that grows at least exponentially in the ground-truth mutual information (Song & Ermon, 2019). Please see Appendix B for more details and proofs.

On the other hand, V -information is explicit about the assumptions (as a feature instead of a bug). V -information is also easy to estimate with guarantees if we can bound the complexity of its Rademacher or covering number complexity. As we will show, bounds on the complexity of directly translate to PAC (Valiant, 1984) bounds for information estimation. In practice, we can efficiently optimize over V , e.g., via gradient descent. In this paper we will present the Rademacher complexity version; other complexity measures (such as covering number) can be derived similarly.

Definition 4 (Empirical V -information) Let X, Y be two random variables taking values in \mathcal{X} and \mathcal{Y} and $D = \{(x_i, y_i)\}_{i=1}^n$. \mathcal{X}, \mathcal{Y} denotes the set of samples drawn from the joint distribution \mathbb{P} and V is a predictive family. The empirical V -information (under D) is the following V -information under the empirical distribution defined via:

$$\hat{I}_V(X \llcorner Y; D) = \inf_{f \in \mathcal{V}} \frac{1}{n} \sum_{(x_i, y_i) \in D} \log \frac{1}{f(x_i|y_i)} - \inf_{f \in \mathcal{V}} \frac{1}{n} \sum_{(x_i, y_i) \in D} \log \frac{1}{f(x_i, y_i)} \quad (4)$$

Then we have the following PAC bound over the empirical V -information:

Theorem 1. Assume $f \in \mathcal{V}; x \in \mathcal{X}; y \in \mathcal{Y}; \log f(x|y) \in [-B, B]$. Then for any $\epsilon \in (0, 0.5)$, with probability at least $1 - \epsilon$, we have:

$$I_V(X \llcorner Y) - \hat{I}_V(X \llcorner Y; D) \leq 4R_{n, D}(\mathcal{G}_V) + 2B \frac{\log \frac{1}{\epsilon}}{n} \quad (5)$$

where we define the function family $\mathcal{G}_V = \{f \circ g; g(x, y) = \log f(x|y); f \in \mathcal{V}\}$, and $R_N(\mathcal{G})$ denotes the Rademacher complexity of \mathcal{G} with sample size n .

Typically, the Rademacher complexity term satisfies $R_{n, D}(\mathcal{G}_V) = O(n^{-\frac{1}{2}})$ (Bartlett & Mendelson, 2001; Gao & Zhou, 2016). It's worth noticing that a complex function family (i.e., with large Rademacher complexity) could lead to overfitting. On the other hand, an overly-simple family may not be expressive enough to capture the relationship between X and Y . As an example of the theorem, we provide a concrete estimation bound where \mathcal{V} is chosen to be linear functions mapping to the mean of a Gaussian distribution. This was shown in Proposition 1 to lead to the coefficient of determination.

Corollary 1.1. Assume $X = \{x \in \mathbb{R}^{d_x}; \|x\|_2 \leq k_x\}$ and $Y = \{y \in \mathbb{R}^{d_y}; \|y\|_2 \leq k_y\}$. If

$$\mathcal{V} = \{f : f(x) = N(Wx + b; I); f(y) = N(c; I); W \in \mathbb{R}^{d_y \times d_x}; b, c \in \mathbb{R}^{d_y}; \|W\|_2 \leq k; \|b\|_2 \leq 1\}$$

Denote $M = (k_x + k_y)^2 + \log 2$, then for any $\epsilon \in (0, 0.5)$, with probability at least $1 - \epsilon$:

$$I_V(X \llcorner Y) - \hat{I}_V(X \llcorner Y; D) \leq \frac{M}{4n} + 4 \frac{\log \frac{1}{\epsilon}}{n}$$

Similar results can be obtained using other classes of machine learning models with known (Rademacher) complexity.

5 STRUCTURE LEARNING WITH V -INFORMATION

Among many possible applications of V -information, we show how to use it to perform structure learning with provable guarantees. The goal of structure learning is to learn a directed graphical model (Bayesian network) or undirected graphical model (Markov network) that best captures the (conditional) independence structure of an underlying data generating process. Structure learning is difficult in general, but if we restrict ourselves to certain set of graphs there are efficient algorithms. In particular, the Chow-Liu algorithm (Chow & Liu, 1968) can efficiently learn tree graphs (i.e. the set of trees). Chow & Liu (1968) show that the problem can be reduced to:

$$g^* = \arg \max_{g \in \mathcal{G}_{\text{tree}}} \sum_{(X_i, X_j) \in \text{edge}(g)} I(X_i; X_j) \quad (6)$$

where $I(X_i; X_j)$ is the Shannon mutual information between variables X_i and X_j . In other words, it suffices to construct the maximal weighted spanning tree where the weight between two vertices is

their Shannon mutual information. [Chow & Wagner \(1973\)](#) show that the Chow-Liu algorithm is consistent, i.e, it recovers the true solution as the dataset size goes to infinity. However, the finite sample behavior of the Chow-Liu algorithm for high dimensional problems is much less studied, due to the difficulty of estimating mutual information. In fact, we show in our experiments that the empirical performance is often poor, even with state-of-the-art estimators. Additionally, methods based on mutual information cannot take advantage of intrinsically asymmetric relationships, which are common for example in gene regulatory networks ([Meyer et al., 2007](#)).

To address these issues, we propose a new structure learning algorithm based on \mathcal{V} -information instead of Shannon information. The idea is that we can associate to each edge in G (i.e., each pair of variables) a suitable predictive family (cf. Def 1). The main challenge is that we cannot simply replace mutual information with \mathcal{V} -information in Eq. 6 because \mathcal{V} -information is asymmetric – we now have to optimize over directed trees:

$$g = \arg \max_{g \in \mathcal{G}_{d \text{ tree}}} \sum_{i=2}^n I_{V_{t(g)(i)}}(X_{t(g)(i)} \mid X_i) \quad (7)$$

where $\mathcal{G}_{d \text{ tree}}$ is the set of directed trees, $at(g) : N \rightarrow N$ is the function mapping each non-root node of directed tree g to its parent, and $V_{i,j}$ is the predictive family for random variables X_i and X_j . After estimating \mathcal{V} -information on each edge, we use the Chu-Liu algorithm ([Chu & Liu, 1965](#)) to construct the maximal directed spanning tree. This allows us to solve (7) exactly, even though there is a combinatorially large number of trees to consider. Pseudocode is summarized in Algorithm 1 in Appendix. Denote $C(g) = \sum_{i=2}^m I_{V_{t(g)(i)}}(X_{t(g)(i)} \mid X_i)$, we show in the following theorem that unlike the original Chow-Liu algorithm, our algorithm has guarantees in the finite samples regime, even in continuous settings:

Theorem 2. Let $\{X_i\}_{i=1}^m$ be the set of m random variables $D_{i,j}$ (resp. D_j) be the set of samples drawn from $P(X_i \mid X_j)$ (resp. $P(X_j)$). Denote the optimal directed tree with maximum expected edge weights $\sum C(g)$ as g^* and the optimal directed tree constructed on the data as \hat{g} . Then with the assumption in theorem 1, for any $\epsilon \in (0, \frac{1}{2m(m-1)})$, with probability at least $1 - 2m(m-1)\epsilon$, we have:

$$C(\hat{g}) - C(g^*) \leq 2(m-1) \max_{i,j} \left(2R_{D_{i,j}}(G_{V_{i,j}}) + 2R_{D_j}(G_{V_j}) + B \frac{r}{2 \log \frac{1}{(jD_j)^{\frac{1}{2}} + jD_{i,j}^{\frac{1}{2}}}} \right) \quad (8)$$

Theorem 2 shows that the total edge weights of the maximal directed spanning tree constructed by algorithm 1 would be close to the optimal total edge weights if the Rademacher term is small. Although larger $C(g)$ does not necessarily lead to better Chow-Liu trees, empirically we find that the optimal tree in the sense of equation (7) is consistent with the optimal tree in equation (6) under commonly used \mathcal{V} .

6 EXPERIMENTAL RESULTS

6.1 STRUCTURE LEARNING WITH CONTINUOUS HIGH-DIMENSIONAL DATA

We generate synthetic data using various ground-truth tree structures between 7 and 20 variables, where each variable is 10-dimensional. We use Gaussians, Exponentials, and Uniforms as ground truth edge-conditionals. We use \mathcal{V} -information (Gaussian) and \mathcal{W} -information (Logistic) to denote Algorithm 1 with two different \mathcal{V} -families. Please refer to Appendix D.1 for more details. We compare with the original Chow-Liu algorithm equipped with state-of-the-art mutual information estimators CPC ([van den Oord et al., 2018](#)), WJ ([Nguyen et al., 2010](#)) and MINE ([Belghazi et al., 2018](#)), with the same neural network architecture as \mathcal{V} -families for fair comparison. All the experiments are repeated for 10 times. As a performance metric, we use the wrong-edges-ratio (the ratio of edges that are different from ground truth) as a function of the amount of training data.

We show two illustrative experiments in figure 1a; please refer to Appendix D.1 for all simulations. We can see that although the two \mathcal{V} -families used are misspecified with respect to the true underlying (conditional) distributions, the estimated Chow-Liu trees are much more accurate across all data regimes, with CPC (blue) being the best alternative. Surprisingly, \mathcal{V} -information (Gaussian) works consistently well in all cases and only requires about 100 samples to recover the ground-truth Chow-Liu tree in simulation-A.

(a) Chow-Liu tree Construction

(b) Gene network inference

(c) V-information of frames

Figure 1: (a) The expected wrong-edges-ratio of algorithm 1 with different mutual information estimators-based algorithms from sample size 5×10^3 . (b) AUC curve for gene regulatory network inference. (c) The predictive information versus frame distance.

6.2 GENE REGULATORY NETWORK INFERENCE

Mutual information between pairs of gene expressions is often used to construct gene regulatory networks. We evaluate V -information on the in-silico dataset from the DREAM5 challenge (Marbach et al., 2012) and use the setup of Gao et al. (2017), where 20 genes with 660 datapoints are utilized to evaluate all methods. We compare with state-of-the-art non-parametric Shannon mutual information estimators in this low dimensional setting: KDE, the traditional kernel density estimator; KSG estimator (Kraskov et al., 2004); mixed KSG estimator (Gao et al., 2017) and partitioning, an adaptive partitioning estimator (Darbellay & Vajda, 1999) implemented by Song (2014). For fair comparison with these low dimensional estimators, we select $f : f[x] = N(g(x); \frac{1}{2}); x \in X; f[y] = N(\cdot; \frac{1}{2})$ range(g), where g is a 3-rd order polynomial.

The task is to predict whether a directed edge between genes exists in the ground-truth gene network. We use the estimated mutual information and V -information for gene pairs as the test statistic to obtain the AUC for various methods. As shown in Figure 1b, our method outperforms all other methods in network inference under different fractions of data used for estimation. The natural information measure in this task is asymmetry since the goal is to find the pairs of genes (A_i, B_i) in which A_i regulates B_i , thus V -information is more suitable for such case than mutual information.

6.3 RECOVERING THE ORDER OF VIDEO FRAMES

Let X_1, \dots, X_{20} be random variables each representing a frame in videos from the Moving-MNIST dataset, which contains 10,000 sequences each of length 20 showing two digits moving with stochastic dynamics. Can Algorithm 1 be used to recover the natural (causal) order of the frames? Intuitively, predictability should be inversely related with frame distance, thus enabling structure learning. Using a conditional PixelCNN++ (Salimans et al., 2017) as predictive family we shown in Figure 1c that predictive V-information does indeed decrease with frame distance, despite some fluctuations when the frame distances are large. Using Algorithm 1 to construct a Chow-Liu tree, we find that tree perfectly recovers the relative order of the frames

We also generate a Deterministic-Moving-MNIST dataset, where digits move according to deterministic dynamics. From the perspective of Shannon mutual information, every pair of frames has the same mutual information. Hence, standard Chow-Liu tree learning algorithm would fail to discover the natural ordering of the frames (causal structure). In contrast, once we constrain the observer to PixelCNN++ models, algorithm 1 with predictive V-information can still recover the order of different frames when the frame distances are relatively small (less than 9). Compared to the stochastic dynamics case, V-information is more irregular with increasing frame distance, since the PixelCNN++ tends to overfit.

6.4 INFORMATION THEORETIC APPROACHES TO FAIRNESS

The goal of fair representation learning is to map input X to a feature space Z such that the mutual information between Z and some sensitive attribute U (such as race or gender) is minimized. The motivation is that using Z (instead of X) as input we can no longer use the sensitive attributes U to make decisions, thus ensuring some notion of fairness. Existing methods obtain fair representations by optimizing against an “adversarial” discriminator so that the discriminator cannot predict U from Z (Edwards & Storkey, 2015; Louizos et al., 2015; Madras et al., 2018; Song et al., 2018). Under some assumptions and V , we show in Appendix D.2 that these works actually use V-information minimization as part of their objective, where V depends on the functional form of the discriminator.

However, it is clear from the V-information perspective that features trained with V-information minimization might not generalize to V-information and vice versa. To illustrate this, we use a function family V_j as the attacker to extract information from features trained with V_i (with $Z \perp U$) minimization, where all the V s are neural nets. On three datasets commonly used in the fairness literature (Adult, German, Heritage), previous methods work well at preventing information “leak” against the class of adversary they’ve been trained on, but fail when we consider different ones. As shown in Figure 3b in Appendix, the diagonal elements in the matrix are usually the smallest in rows, indicating that the attacker function family extracts more information on features trained with V_j (for $j \neq i$) V-information minimization. This challenges the generalizability of fair representations in previous works. Please refer to Appendix D.2 for details.

7 RELATED WORK

Alternative definitions of Information Several alternative definitions of mutual information are available in the literature. Renyi entropy and Renyi mutual information (Lenzi et al., 2000) extend Shannon information by replacing KL divergence with divergences. However, they have the same difficulty when applied to high dimensional problems as Shannon information.

The line of work most related to ours is the entropy and H-mutual information (DeGroot et al., 1962; Günwald et al., 2004), which associate a definition of entropy to every prediction loss. However, there are two key differences. First, literature on entropy only consider a few special types of prediction functions that serve unique theoretical purposes; for example, (Duchi et al., 2018) considers the set of all functions on a feature space to prove surrogate risk consistency, and (Gret et al., 2004) only considers the entropy to prove the duality between maximum entropy and worst-case loss minimization. In contrast, our definition takes a completely different perspective — emphasizing bounded computation and intuitive properties of “usable” information. Furthermore, entropy still suffers from difficulty of estimation in high dimension because the definitions do not restrict to functions with small complexity (e.g. Rademacher complexity).

Mutual information estimation The estimation of mutual information in the machine learning field is often on the continuous underlying distribution. For non-parametric mutual information estimators, many methods have exploited the principle to calculate the mutual information, such as the Kernel density estimator (Paninski & Yajima, 2008), k-Nearest-Neighbor estimator and the KSG estimator (Kraskov et al., 2004). However, these non-parametric estimators usually aren't scalable to high dimension. Recently, several works utilize the variational lower bounds of MI to design MI estimator based on deep neural network in order to estimate MI of high dimension continuous random variables (Nguyen et al., 2010; van den Oord et al., 2018; Belghazi et al., 2018).

8 CONCLUSION

We defined and investigated \mathcal{V} -information, a variational extension to classic mutual information that incorporates computational constraints. Unlike Shannon mutual information, \mathcal{V} -information attempts to capture usable information, and has very different properties, such as invalidating the data processing inequality. In addition, \mathcal{V} -information can be provably estimated, and can thus be more effective for structure learning and fair representation learning.

ACKNOWLEDGEMENTS

This research was supported by AFOSR (FA9550-19-1-0024), NSF (#1651565, #1522054, #1733686), ONR, and FLI.

REFERENCES

- Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *J. Mach. Learn. Res.* 3:463–482, 2001.
- Roberto Battiti. Using mutual information for selecting features in supervised neural net learning. *IEEE Transactions on neural networks*, 5(4):537–550, 1994.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, R. Devon Hjelm, and Aaron C. Courville. Mutual information neural estimation. *ICML*, 2018.
- C Chow and Cong Liu. Approximating discrete probability distributions with dependence trees. *IEEE transactions on Information Theory*, 14(3):462–467, 1968.
- C. K. Chow and Terry J. Wagner. Consistency of an estimate of tree-dependent probability distributions (corresp.). *IEEE Trans. Information Theory*, 19:369–371, 1973.
- Yau Chu and T. Liu. On the shortest arborescence of a directed graph. *Scientia Sinica* 14:1396–1400, 1965.
- Thomas M. Cover and Joy A. Thomas. *Elements of information theory*. 1991.
- Georges A. Darbellay and Igor Vajda. Estimation of the information by an adaptive partitioning of the observation space. *IEEE Trans. Information Theory*, 45:1315–1321, 1999.
- Morris H DeGroot et al. Uncertainty, information, and sequential experiments. *The Annals of Mathematical Statistics*, 33(2):404–419, 1962.
- John Duchi, Khashayar Khosravi, Feng Ruan, et al. Multiclass classification, information, divergence and surrogate risk. *The Annals of Statistics*, 46(6B):3246–3275, 2018.
- Harrison A Edwards and Amos J. Storkey. Censoring representations with an adversarial. [abs/1511.05897](https://arxiv.org/abs/1511.05897), 2015.
- Wei Gao and Zhi-Hua Zhou. Dropout rademacher complexity of deep neural networks. *Science China Information Sciences*, 59(7):072104, 2016.
- Weihaio Gao, Sreeram Kannan, Sewoong Oh, and Pramod Viswanath. Estimating mutual information for discrete-continuous mixtures. *NIPS*, 2017.

- Peter D Grünwald, A Philip Dawid, et al. Game theory, maximum entropy, minimum discrepancy and robust bayesian decision theory. *The Annals of Statistics*, 32(4):1367–1433, 2004.
- Edwin T Jaynes. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(9):939–952, 1982.
- Sham M. Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. *NIPS*, 2008.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- Alexander Kraskov, Harald Gsbauer, and Peter Grassberger. Estimating mutual information. *Phys. Rev. E* 69:066138, Jun 2004. doi: 10.1103/PhysRevE.69.066138. <https://link.aps.org/doi/10.1103/PhysRevE.69.066138>
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–444, 2015.
- Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.
- EK Lenzi, RS Mendes, and LR Da Silva. Statistical mechanics based on renyi entropy. *Physica A: Statistical Mechanics and its Applications*, 280(3-4):337–345, 2000.
- Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard S. Zemel. The variational fair autoencoder. *CoRR*, abs/1511.00830, 2015.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard S. Zemel. Learning adversarially fair and transferable representations. *ArXiv*, abs/1802.06309, 2018.
- Daniel Marbach, James C. Costello, Robert Küfner, N. Vega, Robert J. Prill, Diogo M Camacho, Kyle R. Allison, Manolis Kellis, James J. Collins, and Gustavo Stolovitzky. Wisdom of crowds for robust gene network inference. *Nature Methods*, 2012.
- Patrick E. Meyer, Kevin Kontos, Frédéric La tte, and Gianluca Bontempi. Information-theoretic inference of large transcriptional regulatory networks. *EURASIP J. Bioinformatics and Systems Biology*, 2007.
- XuanLong Nguyen, Martin J. Wainwright, and Michael I. Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56:5847–5861, 2010.
- Liam Paninski and Masanao Yajima. Undersmoothed kernel entropy estimation. *IEEE Transactions on Information Theory*, 54:4384–4388, 2008.
- Rafael Pass and Abhi Shelat. *A course in cryptography*. 2010.
- Judea Pearl. *Causality: Models, reasoning, and inference*. 2000.
- Jonas Peters, Dominik Janzing, and Bernhard Schöpf. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.
- Ben Poole, Sherjil Ozair, Aaron van den Oord, Alexander A Alemi, and George Tucker. On variational bounds of mutual information. *arXiv preprint arXiv:1905.06922*, 2019.
- Rajesh Ranganath, Sean Gerrish, and David M. Blei. Black box variational inference. *AISTATS*, 2013.
- Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. Pixelcnn++: Improving the pixel-cnn with discretized logistic mixture likelihood and other modifications. *ArXiv*, abs/1701.05517, 2017.
- Claude E. Shannon and Warren Weaver. *The mathematical theory of communication*. 1948.
- Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information estimators. *arXiv preprint arXiv:1910.06222*, 2019.

Jiaming Song, Pratyusha Kalluri, Aditya Grover, Shengjia Zhao, and Stefano Ermon. Learning controllable fair representations. *AI STAT* 2018.

Zoltán Szabó. Information theoretical estimators toolbox. *Mach. Learn. Res* 15:283–287, 2014.

Leslie G Valiant. A theory of the learnable. *Communications of the ACM* 27(11):1134–1142, 1984.

Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *ArXiv*, abs/1807.03748, 2018.

Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning* 1(1–2):1–305, 2008.

A PROOFS

A.1 PROOF OF PROPOSITION 1

Proposition 1. For V -entropy and V -information, we have

1. Let \mathcal{X} be as in Def. 1. The $H(Y)$ is the Shannon entropy, $H(Y|X)$ is the Shannon conditional entropy, and $I(Y; X)$ is the Shannon mutual information.
2. Let $Y = \mathbb{R}^d$ and $V = \{f : f \in \mathcal{G}, P \in \mathcal{P}(\mathbb{R}^d)\}$, where P is the distribution with density $y \mapsto \frac{1}{Z} e^{k \cdot y - k_2}$ where $Z = \int_{\mathbb{R}^d} e^{k \cdot y - k_2} dy$, then the V -entropy of a random variable Y equals its mean absolute deviation, up to an additive constant.
3. Let $Y = \mathbb{R}^d$ and $V = \{f : f \in \mathcal{G}, N(\cdot; \mu) \in \mathcal{N}(\mathbb{R}^d; \Sigma) : \mu \in \mathbb{R}^d, \Sigma \in \mathcal{S}_d^+\}$, then the V -entropy of a random variable Y equals the trace of its covariance matrix ($\text{tr}(\text{Cov}(Y))$), up to an additive constant.
4. Let $V = \{f : f \in \mathcal{G}, Q_{t; \theta} \in \mathcal{Q}_t(\mathbb{R}^d; \theta) : \theta \in \Theta\}$, where $Q_{t; \theta}$ is a distribution in a minimal exponential family with sufficient statistics $s : Y \rightarrow \mathbb{R}^d$ and set of natural parameters. For a random variable Y with expected sufficient statistics $\bar{s} = E[s(Y)]$, the V -entropy of Y is the maximum Shannon entropy over all random variables with identical expected sufficient statistics, i.e. $E[s(\hat{Y})] = \bar{s}$.
5. Let $Y = \mathbb{R}^d$, X be any vector space, and $V = \{f : x \mapsto N(\langle x, \phi \rangle; \sigma^2) : \phi \in X, \sigma^2 > 0\}$, where X is the set of linear functions $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, then V -information $I_V(X; Y)$ equals the (unnormalized) maximum coefficient of determination $R^2 = \text{tr}(\text{Cov}(Y))$ for linear regression.

Proof. (1)

Let $P_{Y|X}$ denote the density function of random variable Y conditioned on $X = x$ (we denote this random variable as $Y|x$).

$$\begin{aligned} H(Y|X) &= \int_{\mathcal{X}} E_{x; Y} \log \frac{1}{f[X](y)} = \int_{\mathcal{X}} E_x \int_{\mathcal{Y}} P_{Y|X}(y) \log \frac{P_{Y|X}(y)}{f[X](y)P_{Y|X}(y)} \\ &= \int_{\mathcal{X}} E_x \text{KL}(P_{Y|X} \| f[X]) + H(Y|X) \\ &= E_x \text{KL}(P_{Y|X} \| f[X]) = H(Y|X) \end{aligned} \quad (9)$$

where in sum is achieved for f where $f[x] = P_{Y|X}$ and H is the Shannon (conditional) entropy. The same proof technique can be used to show that $H(Y) = H(Y)$, with the in sum achieved by f where $f[\cdot] = P_Y$. Hence we have

$$I(Y; X) = H(Y) - H(Y|X) = H(Y) - H(Y|X) = I(Y; X) \quad (10)$$

(2)

$$\begin{aligned} H_V(Y) &= \int_{\mathcal{V}} E_{y; Y} [\log f_V(y)] = \int_{\mathbb{R}^d} E_{y; Y} \log \frac{1}{Z} e^{k \cdot y - k_2} \\ &= \int_{\mathbb{R}^d} E_{y; Y} [k \cdot y - k_2] + \log Z \\ &= \text{MAD}(Y) + \log Z \end{aligned} \quad (11)$$

where MAD denotes mean absolute deviation $\int_{\mathbb{R}^d} \|y - E[Y]\| dy$.

(3)

$$\begin{aligned}
 H_V(Y) &= \inf_{f \in \mathcal{F}_V} E_{y \sim Y} [\log f(y)] \\
 &= \inf_{2^{\mathbb{R}^d}} E_{y \sim Y} \log \frac{1}{(2\pi)^{\frac{d}{2}} |j|^{-\frac{1}{2}}} e^{-\frac{1}{2}(y - \mu)^T j^{-1} (y - \mu)} \\
 &= \inf_{2^{\mathbb{R}^d}} E_{y \sim Y} [-(y - \mu)^T j^{-1} (y - \mu)] + \frac{d}{2} \log |j| \\
 &= \inf_{2^{\mathbb{R}^d}} E_{y \sim Y} [\text{tr}((y - \mu)(y - \mu)^T) j^{-1}] + \frac{d}{2} \log |j| \quad (\text{Cyclic property of trace}) \\
 &= \text{tr}(\text{Cov}(Y)) + \frac{d}{2} \log |j| \quad (\text{Linearity of trace})
 \end{aligned}$$

(4) The density function of an exponential family distribution with sufficient statistics $t(y)$ is $\exp(-t(y) \cdot \eta - A(\eta))$ where $A(\eta)$ is the partition function.

$$\begin{aligned}
 H_V(Y) &= \inf_{f \in \mathcal{F}_V} E_{y \sim Y} [\log f(y)] = \inf_{\eta} E_{y \sim Y} [\log \exp(-t(y) \cdot \eta - A(\eta))] \\
 &= \sup_{\eta} (E_{y \sim Y} [t(y) \cdot \eta] - A(\eta)) \\
 &= A^*(E_{y \sim Y} [t(y)]) \quad (12)
 \end{aligned}$$

where A^* is the Fenchel dual of the log-partition function $A(\eta)$. Under mild conditions (Wainwright et al., 2008)

$$A^*(\eta) = H(P_\eta)$$

where P_η is the maximum entropy distribution out of all distributions satisfying $E_{y \sim P_\eta} [t(y)] = \eta$ (Jaynes, 1982), and $H(\cdot)$ is the Shannon entropy.

(5) Assume random variables $X \in \mathbb{R}^d, Y \in \mathbb{R}^d, V = \text{ff} : x \mapsto N(x; \mu, \Sigma); x \in \mathbb{R}^d; ? \mapsto N(\cdot; \mu, \Sigma) \in \mathbb{R}^d$; $\Sigma = \frac{1}{2} I_d$; $\mu = 0$. Then the V -information from X to Y is

$$\begin{aligned}
 I_V(X \rightarrow Y) &= H_V(Y) - H_V(Y|X) \\
 &= \inf_{2^{\mathbb{R}^d}} E_{y \sim Y} \log \frac{1}{(2\pi)^{\frac{d}{2}} |j|^{-\frac{1}{2}}} e^{-\frac{1}{2} y^T j^{-1} y} - \inf_{2^{\mathbb{R}^d}} E_{x,y \sim X,Y} \log \frac{1}{(2\pi)^{\frac{d}{2}} |j|^{-\frac{1}{2}}} e^{-\frac{1}{2} (y - \mu)^T j^{-1} (y - \mu)} \\
 &= \inf_{2^{\mathbb{R}^d}} E_{x,y \sim X,Y} \left[-\frac{1}{2} y^T j^{-1} y + \frac{1}{2} (y - \mu)^T j^{-1} (y - \mu) \right] \\
 &= \text{tr}(\text{Cov}(Y)) - \frac{\text{tr}(\text{Cov}(Y))}{\text{tr}(\text{Cov}(Y))} \text{tr}(\text{Cov}(Y)) \\
 &= \text{tr}(\text{Cov}(Y)) - \text{tr}(\text{Cov}(Y)) = 0
 \end{aligned} \quad (13)$$

□

A.2 PROOF OF PROPOSITION 2

Proposition 2. Let Y and X be any random variables on \mathcal{Y} and \mathcal{X} , and V and U be any predictive families, then we have

1. Monotonicity: If $V \subseteq U$, then $H_V(Y) \geq H_U(Y)$, $H_V(Y|X) \geq H_U(Y|X)$.
2. Non-Negativity: $I_V(X \rightarrow Y) \geq 0$.
3. Independence: If X is independent of Y , $I_V(X \rightarrow Y) = I_V(Y \rightarrow X) = 0$.

Proof. (1)

$$H_V(Y) = \inf_{f \in \mathcal{F}_V} E_{y \sim Y} \log \frac{1}{f(y)} = \inf_{f \in \mathcal{F}_U} E_{y \sim Y} \log \frac{1}{f(y)} = H_U(Y) \quad (14)$$

$$H_V(Y|X) = \inf_{f \in \mathcal{F}_V} E_{x,y \sim X,Y} \log \frac{1}{f(x,y)} = \inf_{f \in \mathcal{F}_U} E_{x,y \sim X,Y} \log \frac{1}{f(x,y)} = H_U(Y|X) \quad (15)$$

The inequalities (14) and (15) are because we are taking the in mum over a larger set.

(2)

Denote $V_?$ as the subset of V that satisfy $[x] = f[?]$, $8x \in X$.

$$\begin{aligned} H_V(Y) &= \inf_{f \in V} E_{x,y} [-\log f[?](y)] \\ &= \inf_{f \in V_?} E_{x,y} [-\log f[?](y)] && \text{(By Optional Ignorance)} \\ &= \inf_{f \in V_?} E_{x,y} [-\log f[x](y)] \\ &= \inf_{f \in V} E_{x,y} [-\log f[x](y)] = H_V(Y | X) \end{aligned}$$

Therefore

$$I_V(Y \perp X) = H_V(Y) - H_V(Y|X) \geq 0$$

(3)

Denote $V_?$ as the subset of V that satisfy $[x] = f[?]$, $8x \in X$.

$$\begin{aligned} H_V(Y | X) &= \inf_{f \in V} E_{x,y} [-\log f[x](y)] \\ &= \inf_{f \in V} E_x \times E_y [-\log f[x](y)] && \text{(Independence)} \\ &= E_x \times \inf_{f \in V} E_y [-\log f[x](y)] && \text{(Jensen)} \\ &= E_x \times \inf_{f \in V_?} E_y [-\log f[x](y)] && \text{(Optional Ignorance)} \\ &= \inf_{f \in V_?} E_y [-\log f[?](y)] && \text{(No dependence on } x) \\ &= \inf_{f \in V} E_y [-\log f[?](y)] = H_V(Y) \end{aligned}$$

Therefore $I_V(Y \perp X) = H_V(Y) - H_V(Y|X) \geq 0$. Combined with the Proposition 2.2 that $I_V(X \perp Y)$ must be non-negative, $I_V(X \perp Y)$ must be 0.

□

A.3 PROOF OF THEOREM 1

Theorem 1. Assume $f \in V; x \in X; y \in Y; \log f[x](y) \in [B; B]$. Then for any $\epsilon \in (0; 0.5)$, with probability at least $1 - \epsilon$, we have:

$$I_V(X \perp Y) - \hat{I}_V(X \perp Y; D) \leq 4R_{|D|}(G_V) + 2B \frac{\sqrt{s}}{|D|} \tag{5}$$

where we define the function family $G_V = \{g(x; y) = \log f[x](y); f \in V\}$, and $R_N(G)$ denotes the Rademacher complexity of G with sample size N .

Before proving theorem 1, we introduce two lemmas. Proofs for these Lemmas follow the same strategy as theorem 8 in [Bartlett & Mendelson \(2001\)](#):

Lemma 3. Let $X; Y$ be two random variables taking values in X and Y and D denotes the set of samples drawn from the joint distribution over $X \times Y$. Assume $f \in V; x \in X; y \in Y; \log f[x](y) \in [B; B]$.

Take $\hat{f} = \arg \min_{f \in V} \frac{1}{|D|} \sum_{x_i, y_i \in D} \log f[x_i](y_i)$, then $\epsilon \in (0; 1)$, with probability at least $1 - \epsilon$, we have:

$$H_V(Y|X) - \frac{1}{|D|} \sum_{x_i, y_i \in D} \log \hat{f}[x_i](y_i) \leq 2R_{|D|}(G_V) + 2B \frac{\sqrt{s}}{|D|} \tag{16}$$

Proof. We apply McDiarmid's inequality to the function defined for any sample D by

$$(\mathcal{D}) = \sup_{f \in \mathcal{F}} E_{x,y} [\log f(x)(y)] - \frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) \quad (17)$$

Let D and D^0 be two samples differing by exactly one point, then since the difference of suprema does not exceed the supremum of the differences and $\mathcal{F} = \{x \in \mathcal{X}; y \in \mathcal{Y}; \log f(x)(y) \in [B, B]\}$, we have:

$$\begin{aligned} (\mathcal{D}) - (\mathcal{D}^0) &= \sup_{f \in \mathcal{F}} E_{x,y} [\log f(x)(y)] - \frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) - \left(\sup_{f \in \mathcal{F}} E_{x,y} [\log f(x)(y)] - \frac{1}{jD^0{}^j} \sum_{x_i, y_i \in \mathcal{D}^0} \log f(x_i)(y_i) \right) \\ &= \sup_{f \in \mathcal{F}} \left(\frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) - \frac{1}{jD^0{}^j} \sum_{x_i, y_i \in \mathcal{D}^0} \log f(x_i)(y_i) \right) \\ &\leq \frac{2B}{jD^j} \end{aligned} \quad (18)$$

then by McDiarmid's inequality, for any $\epsilon \in (0, 1)$, with probability at least $1 - \epsilon$, the following holds:

$$(\mathcal{D}) \leq E_D[(\mathcal{D})] + B \frac{\sqrt{\epsilon}}{2 \log \frac{1}{\epsilon}} \quad (18)$$

Then we bound the $E_D[(\mathcal{D})]$ term:

$$E_D[(\mathcal{D})] = E_D \left[\sup_{f \in \mathcal{F}} E_{x,y} [\log f(x)(y)] - \frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) \right] \quad (19)$$

$$= E_D \left[\sup_{f \in \mathcal{F}} E_{D^0} \left[\frac{1}{jD^0{}^j} \sum_{x_i^0, y_i^0 \in \mathcal{D}^0} \log f(x_i^0)(y_i^0) - \frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) \right] \right] \quad (20)$$

$$E_D \left[\sup_{f \in \mathcal{F}} E_{D^0} \left[\frac{1}{jD^0{}^j} \sum_{x_i^0, y_i^0 \in \mathcal{D}^0} \log f(x_i^0)(y_i^0) - \frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) \right] \right] \quad (21)$$

$$E_{D; D^0} \left[\sup_{f \in \mathcal{F}} \frac{1}{jD^0{}^j} \sum_{x_i^0, y_i^0 \in \mathcal{D}^0} \log f(x_i^0)(y_i^0) - \frac{1}{jD^j} \sum_{x_i, y_i \in \mathcal{D}} \log f(x_i)(y_i) \right] \quad (22)$$

$$= E_{D; D^0} \left[\sup_{f \in \mathcal{F}} \frac{1}{jD^j} \sum_{i=1}^j (\log f(x_i^0)(y_i^0) - \log f(x_i)(y_i)) \right] \quad (23)$$

$$E_{D; D^0} \left[\sup_{f \in \mathcal{F}} \frac{1}{jD^j} \sum_{i=1}^j (\log f(x_i^0)(y_i^0) - \log f(x_i)(y_i)) \right] \quad (24)$$

$$E_{D; D^0} \left[\sup_{f \in \mathcal{F}} \frac{1}{jD^j} \sum_{i=1}^j \log f(x_i)(y_i) \right] + E_{D^0} \left[\sup_{f \in \mathcal{F}} \frac{1}{jD^0{}^j} \sum_{i=1}^j \log f(x_i^0)(y_i^0) \right] \quad (25)$$

$$= 2 E_{D; D^0} \left[\sup_{f \in \mathcal{F}} \frac{1}{jD^j} \sum_{i=1}^j \log f(x_i)(y_i) \right] \quad (26)$$

$$= 2 E_{D; \mathcal{G}} \sup_{g \in \mathcal{G}} \frac{1}{jDj} \sum_{i=1}^j \mathcal{X}^i g(x_i; y_i) = 2 R_{jDj}(\mathcal{G}_V) \quad (27)$$

where \mathcal{X} is a Rademacher variable that is uniform on $\{-1, +1\}$. Inequality (22) follows from the convexity of \sup , inequality (24) follows from the symmetrization argument for Rademacher random variables (Ledoux & Talagrand (2013), Section 6.1), inequality (21) follows from the convexity of \log . (27) follows from the definition of \mathcal{G} and Rademacher complexity.

Finally, combining inequality (18) and (27) yields for all V , with probability at least

$$E_{x;Y} [\log f[x](y)] - \frac{1}{jDj} \sum_{x_i; y_i \in D} \log f[x_i](y_i) \leq 2 R_{jDj}(\mathcal{G}_V) + B \frac{\sqrt{s}}{jDj} \frac{2 \log \frac{1}{\delta}}{jDj} \quad (28)$$

In particular, the inequality holds for $\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{jDj} \sum_{x_i; y_i \in D} \log f[x_i](y_i)$ and $\hat{f} = \arg \min_{f \in \mathcal{F}} E_{x;Y} [\log f[x](y)]$. Then we have:

$$E_{x;Y} [\log \hat{f}[x](y)] - \frac{1}{jDj} \sum_{x_i; y_i \in D} \log \hat{f}[x_i](y_i) \leq H_V(Y|X) + \frac{1}{jDj} \sum_{x_i; y_i \in D} \log \hat{f}[x_i](y_i)$$

$$E_{x;Y} [\log \hat{f}[x](y)] - \frac{1}{jDj} \sum_{x_i; y_i \in D} \log \hat{f}[x_i](y_i)$$

Hence the bound (16) holds. \square

Similar bounds can be derived for $H_V(Y)$ when we choose the domain \mathcal{X} to be $\mathcal{X} = \{f\}$:

Lemma 4. Let Y be random variable taking values in \mathcal{Y} and D denotes the set of samples drawn from the underlying distribution $\mathbb{P}(Y)$. Assume $f \in \mathcal{F}$, $y \in \mathcal{Y}$; $\log fy \in [-B, B]$. Take $\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{jDj} \sum_{x_i; y_i \in D} \log fy_i$, then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have:

$$H_V(Y) - \frac{1}{jDj} \sum_{y_i \in D} \log \hat{f}y_i \leq 2 R_{jDj}(\mathcal{G}_{V^?}) + B \frac{\sqrt{s}}{jDj} \frac{2 \log \frac{1}{\delta}}{jDj} \quad (29)$$

$$2 R_{jDj}(\mathcal{G}_V) + B \frac{\sqrt{s}}{jDj} \frac{2 \log \frac{1}{\delta}}{jDj} \quad (30)$$

where $\mathcal{G}_{V^?} = \{g(y) = \log fy; f \in \mathcal{F}\}$.

Proof. The first inequality (29) can be derived similarly as Lemma 3. Since a predictive family, hence there exists a function $h: \mathcal{V} \rightarrow \mathcal{V}$, such that $h(f) = f$ and $h(x) = f[x] = f[y]$.

$$R_{jDj}(\mathcal{G}_{V^?}) = E_{D; \mathcal{F}} \sup_{f \in \mathcal{F}} \frac{1}{jDj} \sum_{i=1}^j \mathcal{X}^i \log fy_i$$

$$= E_{D; \mathcal{F}} \sup_{f \in \mathcal{F}} \frac{1}{jDj} \sum_{i=1}^j \mathcal{X}^i \log h(f)[x_i](y_i)$$

$$= E_{D; \mathcal{F}} \sup_{f \in \mathcal{F}} \frac{1}{jDj} \sum_{i=1}^j \mathcal{X}^i \log f[x_i](y_i) \quad (31)$$

$$= R_{jDj}(\mathcal{G}_V)$$

The inequality (31) holds because $h(f) = f$. \square

Now we prove theorem 1:

Theorem 1. Assume $\epsilon \in (0, 1/2]$, $\delta \in (0, 1/2]$, $\log f(x)(y) \in [B, B]$, for any $\epsilon \in (0, 1/2]$, with probability at least $1 - \delta$, we have:

$$I_V(X \parallel Y) - \hat{I}_V(X \parallel Y; D) \leq 4R_{|D|}(G_V) + 2B \frac{\epsilon}{|D|}$$

Proof. Define $\hat{x} = \arg \min_{x_i, y_i \in D} \log f(x_i)(y_i)$ and $\hat{y} = \arg \min_{y_i \in D} \log f(\hat{x})(y_i)$. Using the triangular inequality we have:

$$\begin{aligned} I_V(X \parallel Y) - \hat{I}_V(X \parallel Y; D) &= (H_V(Y) - H_V(Y|X)) - (H_V(Y) - H_V(Y|\hat{x})) \\ &\leq H_V(Y) - H_V(Y|X) + H_V(Y|\hat{x}) - H_V(Y) \\ &= H_V(Y|\hat{x}) - H_V(Y|X) \\ &= H_V(Y|\hat{x}) - H_V(Y|\hat{x}) + H_V(Y|\hat{x}) - H_V(Y|X) \\ &= H_V(Y|\hat{x}) - H_V(Y|X) \end{aligned} \quad (32)$$

For simplicity let

$$D_{Y|X} = H_V(Y|X) - H_V(Y|\hat{x})$$

and

$$D_Y = H_V(Y) - H_V(Y|\hat{y})$$

With inequality (32), Lemma 3 and Lemma 4, we have:

$$\begin{aligned} \Pr \left[I_V(X \parallel Y) - \hat{I}_V(X \parallel Y; D) > 4R_{|D|}(G_V) + 2B \frac{\epsilon}{|D|} \right] &\leq \Pr \left[D_{Y|X} + D_Y > 4R_{|D|}(G_V) + 2B \frac{\epsilon}{|D|} \right] \\ &\leq \Pr \left[D_{Y|X} > 2R_{|D|}(G_V) + B \frac{\epsilon}{|D|} \right] + \Pr \left[D_Y > 2R_{|D|}(G_V) + B \frac{\epsilon}{|D|} \right] \\ &\leq \Pr \left[D_{Y|X} > 2R_{|D|}(G_V) + B \frac{\epsilon}{|D|} \right] + \Pr \left[D_Y > 2R_{|D|}(G_V) + B \frac{\epsilon}{|D|} \right] \end{aligned}$$

(Union bound)
(Lemma 3 and Lemma 4)

Hence we have:

$$\Pr \left[I_V(X \parallel Y) - \hat{I}_V(X \parallel Y; D) > 4R_{|D|}(G_V) + 2B \frac{\epsilon}{|D|} \right] \leq \delta$$

which completes the proof. \square

A.4 PROOF OF COROLLARY 1.1

Corollary 1.1. Assume $X = f x \in \mathbb{R}^{d_x}; k_x k_2 \leq k_x g$ and $Y = f y \in \mathbb{R}^{d_y}; k_y k_2 \leq k_y g$. If

$$V = f f : f [x] = N(Wx + b; I); f [y] = N(c; I); W \in \mathbb{R}^{d_y \times d_x}; b, c \in \mathbb{R}^{d_y}; k(W; b) k_2 \leq 1/g$$

Denote $M = (k_x + k_y)^2 + \log 2$, then $\mathbb{E} I_V(X \parallel Y) \leq (0; 0.5)$, with probability at least $1 - \delta$:

$$I_V(X \parallel Y) \leq \hat{I}_V(X \parallel Y; D) + \frac{M}{4jD_j} + \frac{r}{1 + 4} \frac{1}{2 \log 2}$$

The proof is an adaptation of the proof for theorem 3 in [Kakade et al. \(2008\)](#).

Proof. From theorem 1 we have:

$$I_V(X \parallel Y) \leq \hat{I}_V(X \parallel Y; D) + 4R_{jD_j}(G_V) + 2B \frac{s}{jD_j} \frac{2 \log 1}{jD_j}$$

In the following $k(W; b) k_2$ is the matrix 2-norm of $(W; b)$, then the Rademacher term can be bounded as follows:

$$\begin{aligned} R_{jD_j}(G_V) &= \frac{1}{jD_j} \mathbb{E} \sup_{W; b; k(W; b) k_2 \leq 1} \sum_{i=1}^{jD_j} \mathbb{R}^i \left(\log p \frac{1}{2} \left(\frac{1}{2} k y_i \right) W x_i + b k_2^2 \right) \\ &\leq \frac{1}{jD_j} \mathbb{E} \sup_{W; b; k(W; b) k_2 \leq 1} \sum_{i=1}^{jD_j} \mathbb{R}^i \left(\frac{1}{2} k y_i \right) W x_i + b k_2^2 + \frac{1}{jD_j} \mathbb{E} \sum_{i=1}^{jD_j} \mathbb{R}^i \log p \frac{1}{2} \end{aligned} \quad (33)$$

The second term in RHS can be bounded as follows:

$$\begin{aligned} \frac{1}{jD_j} \mathbb{E} \sum_{i=1}^{jD_j} \mathbb{R}^i \log p \frac{1}{2} &\leq \frac{1}{jD_j} \mathbb{E} \sum_{i=1}^{jD_j} \mathbb{R}^i \log p \frac{1}{2} \quad (\text{concavity of } \log \frac{1}{2}) \\ &= \frac{1}{jD_j} \sum_{i=1}^{jD_j} \mathbb{E} \mathbb{R}^i \log p \frac{1}{2} \quad (\text{Independence of } \mathbb{R}^i) \\ &= \frac{s}{jD_j} \frac{(\log p \frac{1}{2})^2}{jD_j} \end{aligned} \quad (34)$$

The first term in RHS can be bounded as follows:

$$\begin{aligned} \frac{1}{jD_j} \mathbb{E} \sup_{W; b; k(W; b) k_2 \leq 1} \sum_{i=1}^{jD_j} \mathbb{R}^i \left(\frac{1}{2} k y_i \right) W x_i + b k_2^2 &\leq \frac{1}{2jD_j} \mathbb{E} \sup_{W; b; k(W; b) k_2 \leq 1} \sum_{i=1}^{jD_j} \mathbb{R}^i \left(k y_i \right) W x_i + b k_2^2 \\ &\leq \frac{\max_i k y_i k_2^2}{2} \frac{1}{jD_j} + \max_i k x_i k_2 \frac{\max_i k y_i k_2^2}{jD_j} \\ &\quad + \frac{1}{2jD_j} \mathbb{E} \sup_{W; b; k(W; b) k_2 \leq 1} \sum_{i=1}^{jD_j} \mathbb{R}^i \left(k W x_i + b k_2^2 \right) \end{aligned} \quad (35)$$

$$\frac{\max_i k_y k_2^2}{2} \frac{1}{|D_j|} + \max_i k_x k_2 \frac{\max_i k_y k_2^2}{|D_j|} + \frac{\max_i k_x k_2}{2|D_j|} E_{D_j} \sup_{W; b; k} \sum_{i=1}^j (k W x_i + b k)^2 \quad (36)$$

$$\frac{\max_i k_y k_2^2}{2} \frac{1}{|D_j|} + \max_i k_x k_2 \frac{\max_i k_y k_2^2}{|D_j|} + \frac{\max_i k_x k_2}{2} \frac{\max_i k_x k_2^2}{|D_j|} \quad (37)$$

$$p \frac{M}{4|D_j|}$$

The inequalities (36) and (35) follow the same proof in (34).

Hence we have:

$$R_{|D_j|}(G_v) \leq p \frac{M}{4|D_j|} \quad (38)$$

In this example, we can bound the upper bound of function G_v by

$$B = \sup_{x \in X, y \in Y; k(W; b)_{k_2} \leq 1} \log p \frac{1}{2} \frac{1}{2} k_y W x + b k_2^2$$

$$\sup_{x \in X, y \in Y; k(W; b)_{k_2} \leq 1} \log p \frac{1}{2} + \frac{1}{2} k_y k_2^2 + k W x + b k_2^2 + 2 k_y k W x + b k_2^2$$

$$\log p \frac{1}{2} + \frac{1}{2} (k_x + k_y)^2 < M$$

Combining inequality (38) we arrive at the theorem. □

A.5 PROOF OF THEOREM 2

Theorem 2. Let $\{X_i\}_{i=1}^m$ be the set of m random variables $D_{i,j}$ (resp. D_j) be the set of samples drawn from $P(X_i; X_j)$ (resp. $P(X_j)$). Denote the optimal directed tree with maximum expected edge weights $\sum_{ij} C(g)$ as g^* and the optimal directed tree constructed on the data set \mathcal{D} . Then with the assumption in theorem 1, for any $\epsilon \in (0; \frac{1}{2m(m-1)})$, with probability at least $1 - \epsilon$, we have:

$$C(\hat{g}) - C(g^*) \leq 2(m-1) \max_{i,j} 2R_{D_{i,j}}(G_{v_{i,j}}) + 2R_{D_j}(G_{v_j}) + B \frac{r}{2 \log \frac{1}{|D_j|} + |D_{i,j}|} \quad (8)$$

Proof. Let $C_D(g)$ be the estimated sum of edge weights on data set \mathcal{D} of the tree g , i.e.,

$$C_D(g) = \sum_{i=2}^n \hat{I}_{V_{t(g)(i)}, i}(X_{t(g)(i)}; X_i; D)$$

where $t(g) : N \rightarrow N$ is the function mapping each non-root node of directed tree g to its parent. The same notation for tree \hat{g} . Let

$$= \max_{i,j} \hat{I}_V(X_i; X_j) - \hat{I}_V(X_i; X_j; D)$$

be the maximum absolute estimation error of single edge weight. By the definition we have $8g; jC(\hat{g}) - C_D(\hat{g}) \leq (m-1)\epsilon$, then:

$$C(\hat{g}) + (m-1)\epsilon \leq C_D(\hat{g}) \leq C_D(g) + C(g) - (m-1)\epsilon \quad (39)$$

From lemma 4 and lemma 3 we have:

$$\begin{aligned}
 & \Pr \left(\max_{i,j} 2R_{D_{ij}}(G_{i,j}) + 2R_{D_j}(G_j) + B \frac{r}{2 \log^{-1}(jD_j j^{\frac{1}{2}} + jD_{ij} j^{\frac{1}{2}})} \right) \\
 & \Pr \left(\max_{i,j} \left| I_{V_{ij}}(X_i \setminus X_j) - \hat{I}_{V_{ij}}(X_i \setminus X_j; D) \right| > 2R_{D_{ij}}(G_{i,j}) + 2R_{D_j}(G_j) + B \frac{r}{2 \log^{-1}(jD_j j^{\frac{1}{2}} + jD_{ij} j^{\frac{1}{2}})} \right) \\
 & \Pr \left(\max_{i,j} \left| \mathbb{E}_{H_{V_{ij}}(X_j)} \left[\frac{1}{jD_j j^{\frac{1}{2}}} \sum_{x_j \in \mathcal{X}_j} \log \hat{f}_{V_{ij}}(x_j) \right] - \mathbb{E}_{H_{V_{ij}}(X_j; X_i)} \left[\frac{1}{jD_{ij} j^{\frac{1}{2}}} \sum_{x_i, x_j \in \mathcal{X}_{i,j}} \log \hat{f}_{V_{ij}}(x_j) \right] \right| \right) \\
 & > 2R_{D_{ij}}(G_{i,j}) + 2R_{D_j}(G_j) + B \frac{r}{2 \log^{-1}(jD_j j^{\frac{1}{2}} + jD_{ij} j^{\frac{1}{2}})} \\
 & \Pr \left(\max_{i,j} \left| \mathbb{E}_{H_{V_{ij}}(X_j)} \left[\frac{1}{jD_j j^{\frac{1}{2}}} \sum_{x_j \in \mathcal{X}_j} \log \hat{f}_{V_{ij}}(x_j) \right] - \mathbb{E}_{H_{V_{ij}}(X_j; X_i)} \left[\frac{1}{jD_{ij} j^{\frac{1}{2}}} \sum_{x_i, x_j \in \mathcal{X}_{i,j}} \log \hat{f}_{V_{ij}}(x_j) \right] \right| \right) \\
 & > 2R_{D_j}(G_j) + B \frac{r}{2 \log^{-1} jD_j j^{\frac{1}{2}}} \\
 & \Pr \left(\max_{i,j} \left| \mathbb{E}_{H_{V_{ij}}(X_j; X_i)} \left[\frac{1}{jD_{ij} j^{\frac{1}{2}}} \sum_{x_i, x_j \in \mathcal{X}_{i,j}} \log \hat{f}_{V_{ij}}(x_j) \right] - \mathbb{E}_{H_{V_{ij}}(X_j)} \left[\frac{1}{jD_j j^{\frac{1}{2}}} \sum_{x_j \in \mathcal{X}_j} \log \hat{f}_{V_{ij}}(x_j) \right] \right| \right) \\
 & > 2R_{D_{ij}}(G_{i,j}) + B \frac{r}{2 \log^{-1} jD_{ij} j^{\frac{1}{2}}} A \\
 & m(m-1)2 \quad \text{(By lemma 3, 4 and union bound)}
 \end{aligned}$$

Hence

$$\Pr \left(\max_{i,j} 2R_{D_{ij}}(G_{i,j}) + 2R_{D_j}(G_j) + B \frac{r}{2 \log^{-1}(jD_j j^{\frac{1}{2}} + jD_{ij} j^{\frac{1}{2}})} > 1 - m(m-1)2 \right) \quad (40)$$

Then combining inequality (39) and (40) we arrive at the result. \square

B ANALYSIS OF APPROXIMATE ESTIMATORS FOR SHANNON INFORMATION

We consider two approximate estimators for Shannon information. The first is the CPC (or InfoNCE in [Poole et al. \(2019\)](#)) estimator (I_{CPC}) proposed by [van den Oord et al. \(2018\)](#):

$$I_{CPC} = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \log \frac{f(x_i; y_i)}{\frac{1}{N} \sum_{j=1}^N f(x_i; y_j)} \right] \quad (41)$$

where the expectation is over N independent samples from the joint distribution $p(x; y)$.

The second is the NWJ estimator (I_{NWJ}) proposed by [Nguyen et al. \(2010\)](#):

$$I_{NWJ} = \mathbb{E}_{x,y \sim p(x,y)} [f(x; y)] - \mathbb{E}_{x,y \sim p(x)p(y)} [f(x; y)] \quad (42)$$

In both cases, f is a parameterized function, and the objectives are to maximize these lower bounds parameterized by f to approximate mutual information. Ideally, with sufficiently flexible models and data, we would be able to recover the true mutual information. However, these ideal cases do not carry over to practical scenarios.

For I_{CPC} , [van den Oord et al. \(2018\)](#) show that I_{CPC} is no larger than $\log N$, where N is the batch size. This means that the I_{CPC} estimator will incur large bias when $I(X; Y) \ll \log N$. We provide a proof for completeness as follows.

Proposition 3. If $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$,

$$I_{CPC} \leq \log N \quad (43)$$

Proof. We have:

$$I_{\text{CPC}} := E \frac{1}{N} \sum_{i=1}^N \log \frac{f(x_i; y_i)}{\frac{1}{N} \sum_{j=1}^N f(x_i; y_j)} \quad (44)$$

$$E \frac{1}{N} \sum_{i=1}^N \log \frac{f(x_i; y_i)}{\frac{1}{N} \sum_{j=1}^N f(x_i; y_j)} = E \frac{1}{N} \sum_{i=1}^N \log N = \log N \quad (45)$$

which completes the proof. \square

For NWJ, we note that the NWJ involves a term denoted as $\frac{e^{f(x;y)}}{p(x)p(y)}$, which could be dominated by rare data-points that have high values. Intuitively, this would make it a poor mutual information estimator by optimizing. The NWJ estimator may suffer from high variance when the estimator is optimal (Song & Ermon, 2019), this is also empirically observed in Poole et al. (2019). We provide a proof for completeness as follows.

Proposition 4. Assume that \hat{f}_{NWJ} achieves the optimum value for NWJ. Then the variance of the empirical NWJ estimator satisfies $\text{Var} \hat{f}_{\text{NWJ}} = \frac{e^{I(X;Y)} - 1}{N}$, where

$$\hat{f}_{\text{NWJ}} = \frac{1}{N} \sum_{i=1}^N [f(x_i; y_i)] \frac{e^{\frac{1}{N} \sum_{i=1}^N f(x_i; y_i)}}{e^{f(x_i; y_i)}}$$

is the empirical NWJ estimator with N i.i.d. samples $(x_i; y_i)_{i=1}^N$ from $p(x; y)$ and N i.i.d. samples $f(x_i; y_i)_{i=1}^N$ from $p(x)p(y)$.

Proof. Let us denote $z_i = \frac{p(x_i; y_i)}{p(x_i)p(y_i)}$. Clearly $E_{p(x)p(y)} [z_i] = 1$. Then we have:

$$\begin{aligned} \text{Var}(z_i) &= E_{p(x)p(y)} z_i^2 - (E_{p(x)p(y)} [z_i])^2 \\ &= E_{p(x)p(y)} z_i^2 - 1 \\ &= E_{p(x)p(y)} \frac{p(x_i; y_i)^2}{p(x_i)p(y_i)} - 1 \\ &= E_{p(x; y)} \frac{p(x_i; y_i)}{p(x_i)p(y_i)} - 1 \end{aligned} \quad (46)$$

$$= E_{p(x; y)} \log \frac{p(x_i; y_i)}{p(x_i)p(y_i)} = I(X; Y) - 1 \quad (47)$$

where we use Jensen's inequality for \log at the last step.

From Nguyen et al. (2010), we have:

$$f(x; y) = 1 + \log \frac{p(x; y)}{p(x)p(y)} \quad (48)$$

for all $x; y$. Since $(x_i; y_i)_{i=1}^N$ (resp. $f(x_i; y_i)_{i=1}^N$) are N datapoints independently sampled from the distribution $p(x; y)$ (resp. $p(x)p(y)$), we have

$$\begin{aligned} \text{Var} \hat{f}_{\text{NWJ}} &= \text{Var} \frac{1}{N} \sum_{i=1}^N [f(x_i; y_i)] \frac{e^{\frac{1}{N} \sum_{i=1}^N f(x_i; y_i)}}{e^{f(x_i; y_i)}} \\ &= \text{Var} \frac{e^{\frac{1}{N} \sum_{i=1}^N f(x_i; y_i)}}{N} \\ &= \text{Var} \frac{1}{N} \sum_{i=1}^N z_i \frac{e^{I(X; Y)} - 1}{N} \end{aligned} \quad (49)$$

which completes the proof. \square

Algorithm 1 Construct Chow-Liu Trees with \mathcal{V} -Information

Require: $D = \{X_i\}_{i=1}^m$, with each X_i being a set of datapoints sampled from the underlying distribution of random variable X_i . The set of function families $\mathcal{F}_{i,j} = \{g_{i,j}^m\}_{i,j \in [m]}$ between all the nodes.

- 1: for $i = 1; \dots; m$ do
- 2: for $j = 1; \dots; m$ do
- 3: if $i \neq j$ then
- 4: Calculate the edge weight $w_{i,j} = \hat{I}_{\mathcal{V}_{i,j}}(X_i \parallel X_j; \mathcal{F}_{i,j})$.
- 5: end if
- 6: end for
- 7: end for
- 8: Construct the fully connected graph $G = (V; E)$, with node set $V = \{X_1; \dots; X_m\}$ and edge set $E = \{e_{i,j}\}_{i,j \in [m]}$.
- 9: Construct the maximal directed spanning tree on G by Chow-Liu algorithm, where mutual information is replaced by \mathcal{V} -information.
- 10: return g

C THE NEW ALGORITHM FOR CHU-LIU TREE CONSTRUCTION

See Algorithm 1. $\hat{I}_{\mathcal{V}_{i,j}}(X_i \parallel X_j; \mathcal{F}_{i,j})$ denotes the empirical \mathcal{V} -information.

D DETAILED EXPERIMENTS SETUP**D.1 CHU-LIU TREE CONSTRUCTION**

Figure 2 shows the Chu-Liu tree construction of Simulation-6. The Simulation-A and Simulation-B in the main body correspond to Simulation-1 and Simulation-4.

Simulation-1 Simulation-3 :

The ground-truth Chu-Liu tree is a star tree (i.e. all random variables are conditionally independent given X_1). We conduct all experiments for 10 times, each time with random simulated orthogonal matrices $\{W_i\}_{i=2}^{20}$. Simulation-1: $X_1 \sim U(0; 10)$ and $X_i \mid X_1 \sim N(W_i X_1; 6I); (2 \leq i \leq 20)$; Simulation-2: $X_1 \sim U(0; 10)$ and $X_i \mid X_1 \sim W_i E(X_1 + \cdot); (2 \leq i \leq 20), i \sim E(0; 1)$; Simulation-3 is a mixed version $X_1 \sim U(0; 10); X_i \mid X_1 \sim \frac{1}{2}N(W_i X_1; 6I) + \frac{1}{2}W_i E(X_1 + \cdot); (2 \leq i \leq 20)$.

Simulation-4 Simulation-6 :

The ground-truth Chu-Liu tree is a tree of depth two. We conduct all experiments for 10 times, each time with random simulated orthogonal matrices $\{W_i\}_{i=2}^7$. Simulation-4: $X_1 \sim U(0; 10), X_i \mid X_1 \sim N(W_i X_1; 2I) (i = 2; 3), X_i \mid X_2 \sim N(W_i X_2; 2I) (i = 4; 5), X_i \mid X_3 \sim N(W_i X_3; 2I) (i = 6; 7)$; Simulation-5: $X_1 \sim U(0; 10), X_i \mid X_1 \sim E(X_1 + \cdot) (i = 2; 3), X_i \mid X_2 \sim W_i E(X_2 + \cdot) (i = 4; 5), X_i \mid X_3 \sim W_i E(X_3 + \cdot) (i = 6; 7), i \sim E(0; 1)$; Simulation-6 is a mixed version $X_1 \sim U(0; 10), X_i \mid X_1 \sim W_i E(X_1 + \cdot) (i = 2; 3), X_i \mid X_2 \sim N(W_i X_2; 2I) (i = 4; 5), X_i \mid X_3 \sim N(W_i X_3; 2I) (i = 6; 7), i \sim E(0; 1)$.

D.2 FAIRNESS

We can adapt the \mathcal{V} -information perspective to fairness. Denote the random variable that represents sensitive data and the representation U and Z respectively. Assume U is discrete and \mathcal{V} belongs to predictive family 1. Then we have $I_{\mathcal{V}}(U) = H(U)$ as long as \mathcal{V} has softmax on the top and belongs to predictive family. In this case, minimizing $I_{\mathcal{V}}(Z \parallel U)$ equals to minimize $H_{\mathcal{V}}(Y \mid X)$. Let the joint distribution of Z and U be parameterized by θ . Hence the final objective is:

$$\min_{\theta} I_{\mathcal{V}}(u; z) = \min_{\theta} \sup_{f \in \mathcal{F}_{\mathcal{V}}} \mathbb{E}_{z;u \sim q(z;u)} [\log P_f(z|u)]$$

In [Edwards & Storkey \(2015\)](#); [Madras et al. \(2018\)](#); [Louizos et al. \(2015\)](#); [Song et al. \(2018\)](#), functions in \mathcal{V} are parameterized by a discriminator.

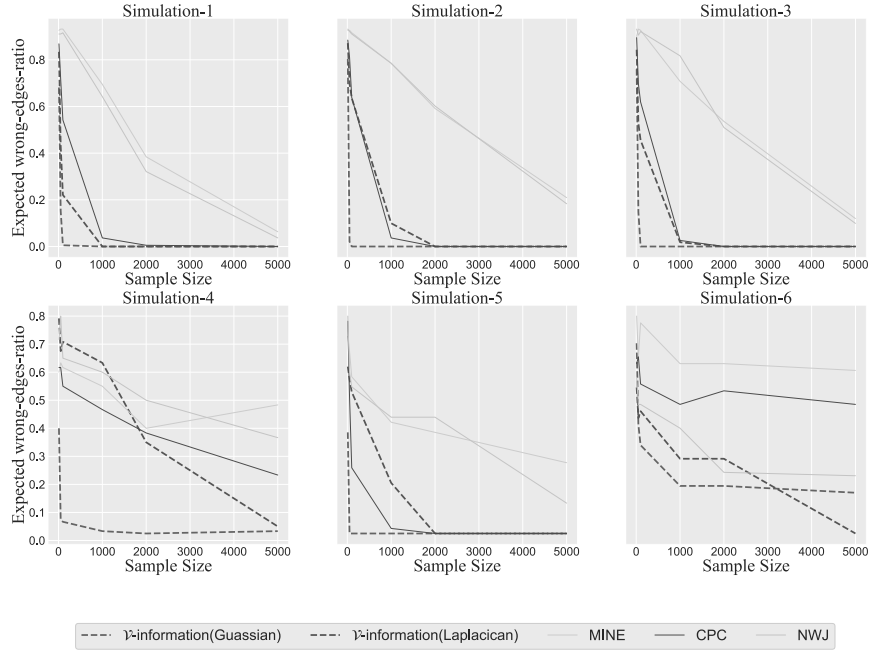


Figure 2: Chu-Liu Tree Construction: The expected wrong-edges-ratio of algorithm 1 with different \mathcal{V} and other mutual information estimators-based algorithms from sample size 10 to 5×10^3 .

For the (F_i, F_j) elements described in the main body, please refer to figure 3b. The three datasets are: the UCI Adult dataset² which has gender as the sensitive attribute; the UCI German credit dataset³ which has age as the sensitive attribute and the Heritage Health dataset⁴ which has the 18 configurations of ages and gender as the sensitive attribute.

The models in the figure are:

$$\mathcal{V}_A = \{f : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{U}) | f[z](u) = \prod_{(z_i; u_i) \in \mathcal{D}} \frac{e^{kz_i} z k_2^2 = h}{\prod_{(z_i; u_i) \in \mathcal{D}} e^{kz_i} z k_2^2 = h} * I(u_i = u), h \in \mathbb{R}\}, \text{ where } \mathcal{D} \text{ is the training set.}$$

$\mathcal{V}_B = \{f : f[z] = \text{softmax}(g(z))\}$, where g is a two-layer MLP with Relu as the activation function.

$\mathcal{V}_C = \{f : f[z] = \text{softmax}(g(z))\}$, where g is a three-layer MLP with LeakyRelu as the activation function.

We further visualize a special case of the $(\mathcal{V}_A, \mathcal{V}_B)$ pair in figure 3a, where the $\mathcal{V}_i = \{f : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{U}) | f[z](u) = \prod_{(z_i; u_i) \in \mathcal{D}} \frac{e^{kz_i} z k_2^2 = h}{\prod_{(z_i; u_i) \in \mathcal{D}} e^{kz_i} z k_2^2 = h} * I(u_i = u), h \in \mathbb{R}\}$ explicitly makes the features of different sensitivity attributes more evenly spread, and functions in \mathcal{V}_B is a simple two layers MLP with softmax at the top. The learned features by \mathcal{V}_A -information minimization appear more evenly spread as expected, however, the attacker functions in \mathcal{V}_B can still achieve a high AUC of 0.857.

The (i, j) elements of tables in Figure 3b stand for using function family \mathcal{V}_i to attack features trained with \mathcal{V}_j -information minimization. The diagonal elements in the matrix are usually the smallest in rows, indicating that the attacker function family \mathcal{V}_i extracts more information on featured trained with \mathcal{V}_j ($j \neq i$)-information minimization.

²<https://archive.ics.uci.edu/ml/datasets/adult>

³<https://archive.ics.uci.edu/ml/datasets>

⁴<https://www.kaggle.com/c/hhp>

