

Neural Estimators for Conditional Mutual Information Using Nearest Neighbors Sampling

Sina Molavipour, *Student Member, IEEE*, Germán Bassi, *Member, IEEE*, and Mikael Skoglund, *Fellow, IEEE*

Abstract—The estimation of mutual information (MI) or conditional mutual information (CMI) from a set of samples is a long-standing problem. A recent line of work in this area has leveraged the approximation power of artificial neural networks and has shown improvements over conventional methods. One important challenge in this new approach is the need to obtain, given the original dataset, a different set where the samples are distributed according to a specific product density function. This is particularly challenging when estimating CMI.

In this paper, we introduce a new technique, based on k nearest neighbors (k -NN), to perform the resampling and derive high-confidence concentration bounds for the sample average. Then the technique is employed to train a neural network classifier and the CMI is estimated accordingly. We propose three estimators using this technique and prove their consistency, make a comparison between them and similar approaches in the literature, and experimentally show improvements in estimating the CMI in terms of accuracy and variance of the estimators.

Index Terms—conditional mutual information, neural networks, nearest neighbors.

I. INTRODUCTION

CONDITIONAL mutual information is recognized as an important statistical metric since, for example, characterizes the capacity of communication channels such as channels with random state and the relay channel [1]; however, its relevance goes beyond communication scenarios. Directed information [2], which is a notion for quantifying causal impact in stochastic processes, is computed as a possible infinite sum of CMIs [3]. Additionally, CMI has been adopted in machine learning [4], [5] as a way to extract shared information in data, while in the information bottleneck method, it can be used as a regularizer [6], [7].

The estimation of information-theoretic quantities has been an important subject in statistical inference for many years. In general, conventional methods are categorized as parametric and non-parametric estimators. In [8] several of these methods for estimating entropy, mutual information, and relative entropy are reviewed. One well-known non-parametric method to estimate MI of continuous random variables is the KSG estimator [9], [10]; this estimator is based on the k nearest neighbors method (k -NN) and shows a favorable performance for data with small dimensions. This method has subsequently been extended to estimate CMI in [11]–[13]. However, observations show that as the dimension of the data increases, the

estimation accuracy deteriorates, and addressing this issue has remained a challenge.

Recent studies leverage the power of artificial neural networks to improve the estimation of information-theoretic quantities. In the recent work [14], the authors propose the use of neural networks to estimate MI and, according to their numerical experiments, promising improvements with respect to the conventional KSG method can be seen for high-dimensional data. The key idea in [14] is to estimate a lower bound for the MI—known as a variational bound—instead of directly estimating the MI; the network is trained to maximize this lower bound which results in a tight approximation of the MI. This approach has been followed by a series of other works such as [15]–[20]. In particular, in [15], the limits of estimation using variational bounds are investigated and the authors provide high confidence bounds for these constraints in terms of the number of samples. Similar arguments can be found in [18] where the authors address the bias–variance trade-off in the neural estimators for MI. A thorough comparison for MI estimators is done in [16] and different methods based on variational bounds are compared in terms of bias and variance. In [19], further insights and applications are provided for a neural estimation of MI.

Before proceeding, consider the definition of CMI for continuous random variables:

$$\begin{aligned} I(X; Y|Z) &:= \iiint p(x, y, z) \log \frac{p(x, y, z)}{p(x|z)p(y, z)} dx dy dz \\ &= \mathbb{E}_{p(y, z)} \left[D(p(x|Y, Z) || p(x|Z)) \right]. \end{aligned} \quad (1)$$

A lower bound on the CMI can thus be obtained employing the Donsker–Varadhan (DV) variational characterization of the divergence [21]:

$$\begin{aligned} I(X; Y|Z) &\geq \mathbb{E}_{p(x, y, z)} [f(x, y, z)] \\ &\quad - \log \mathbb{E}_{p(x|z)p(y, z)} [\exp f(x, y, z)], \end{aligned} \quad (2)$$

where $f(\cdot)$ is any function such that the two expectations exist and are finite. The lower bound (2) may be relaxed, as suggested by Nguyen, Wainwright, and Jordan (NWJ) in [22], resulting in the following lower bound:

$$\begin{aligned} I(X; Y|Z) &\geq \mathbb{E}_{p(x, y, z)} [f(x, y, z)] \\ &\quad - e^{-1} \mathbb{E}_{p(x|z)p(y, z)} [\exp f(x, y, z)]. \end{aligned} \quad (3)$$

These bounds are tight with the appropriate choice of $f(\cdot)$, and equality holds in (2) by choosing $f(\cdot)$ as

$$f_{DV}^*(x, y, z) := C + \log \frac{p(x, y, z)}{p(x|z)p(y, z)}, \quad \forall C \in \mathbb{R}, \quad (4)$$

The authors are with the school of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden 100 44. (e-mails: {sinmo, germanb, skoglund}@kth.se)

This work was supported in part by the Knut and Alice Wallenberg Foundation and the Swedish Foundation for Strategic Research.

while

$$f_{NWJ}^*(x, y, z) := 1 + \log \frac{p(x, y, z)}{p(x|z)p(y, z)} \quad (5)$$

yields equality in both bounds (2) and (3). If the joint probability density function $p(x, y, z)$ were known, it would be possible to compute the optimal functions (4) and (5), and respectively the bounds (2) and (3). Most importantly, we could derive the CMI directly:

$$I(X; Y|Z) = \mathbb{E}_{p(x, y, z)} [f_{LDR}^*(x, y, z)], \quad (6)$$

where

$$f_{LDR}^*(x, y, z) := \log \frac{p(x, y, z)}{p(x|z)p(y, z)} \quad (7)$$

is the logarithm of the density ratio (LDR). However, we only have access to a set of samples distributed according to $p(x, y, z)$. Using these samples, we will approximate the functions (4), (5), and (7), which will allow us to estimate the CMI according to (2), (3), or (6).

We note that, for any fixed function $f(\cdot)$, the NWJ bound (3) is looser than the DV bound (2) except for the case of (5); however, the former bound has the advantage of having a linear form, which may be useful when the bound is estimated empirically. As noted in [20], the average of several estimates of the DV bound is neither a lower bound nor an upper bound of the CMI due to the concavity of the $\log(\cdot)$ function and Jensen's inequality. This becomes of paramount importance if the estimation must not exceed the true value of the CMI. For instance, when estimating the capacity of a communication channel determined by a CMI, the estimated value must be below the true value of the CMI to ensure a reliable communication. It is worth noting that, although estimating with insufficient number of samples may also cause such violation, this should not be confused with the issue caused by the non-linearity of the terms. Nonetheless, if there is no constraint in the estimated value of the CMI being below the true value, we may safely use any of the aforementioned three estimators. In fact, we show in our experiments that, in some cases, estimations based on (6) are more accurate while being above the true value of CMI.

As previously mentioned, the authors of [14] introduced the idea of using artificial neural networks to estimate MI; in particular, they calculate the DV bound, where $f(\cdot)$ is substituted with a neural network and the right-hand side (RHS) of (2) is maximized with the gradient descent method. A new approach to estimate both the MI and the CMI is taken in [17], where a neural network classifier is first trained to distinguish whether samples are generated according to the joint or product density function. Then the authors show that the output of this classifier can be used to approximate the optimal functions in (4) and (5). However, instead of estimating the lower bounds on the CMI directly, they express the CMI as a difference of two MI terms, i.e.,

$$I(X; Y|Z) = I(X; Y, Z) - I(X; Z), \quad (8)$$

and estimate the DV (or NWJ) lower bound for each term separately.

In this paper, we adopt the classifier technique of [17] and introduce a new method to apply it directly to the estimation of CMI. Estimating CMI is more complicated than estimating MI since the technique relies on having samples that are distributed according to the product density $p(x|z)p(y, z)$ apart from the original samples distributed according to $p(x, y, z)$. The approach of [17], which estimates the two terms on the RHS of (8), only requires samples distributed according to $p(x)p(y, z)$ and $p(x)p(z)$, which are simple to obtain given the original samples. Here, we address this issue in Section II and show that the k -NN method can be employed to obtain the desired samples from the original data. In fact, this technique can be applied to any resampling problem where we want to enforce a more restrictive factorization for the density function of the new samples. In Section III, we establish concentration bounds for the empirical average with respect to data sampled according to our k -NN method, which is one of the main contributions of this paper. Next, the consistency of our proposed estimators is investigated by the approximation and generalization power of our setup. Experiments and simulation results are presented in Section IV. Finally, we conclude the paper in Section V where we discuss possible future direction.

II. PRELIMINARIES AND CHALLENGES

Consider a dataset of n triples $(X, Y, Z) \in \mathcal{X}^3$ where X , Y and Z are mappings $\Omega \rightarrow \mathcal{X} \subset \mathbb{R}^d$ with finite Lebesgue measure $\lambda(\mathcal{X})$. For simplicity, we assume the mappings have the same range, while the extension is straightforward when variables range over different sets. Each triple is generated i.i.d. according to $p(x, y, z)$. The classifier technique, which is used at the core of our estimators, relies on a binary neural classifier that distinguishes whether an input sample (x, y, z) is more likely to be generated from the joint density $p(x, y, z)$ or the product density $p(x|z)p(y, z)$.

As neither of these density functions is known, estimating the CMI based on (2), (3), or (6) encounters the following two challenges:

- 1) The optimal functions f_{DV}^* , f_{NWJ}^* , and f_{LDR}^* cannot be computed due to the unknown densities, and thus must be approximated.
- 2) Even if the previous point is solved, it is not possible to derive (2), (3), or (6) analytically, and thus the expectations must also be approximated using the samples.

In the following, we address these issues. We will see that the output of the binary neural classifier, with a proper loss function, can be used to solve the first challenge. However, this leads to a new problem; in order to train the neural classifier, we need samples distributed according to both the joint and the product density functions. The solution to this new issue, which also addresses the second challenge, is to generate sample batches according to $p(x, y, z)$ and $p(x|z)p(y, z)$, where providing the latter is not straightforward and is the main focus of this paper.

Notation: Throughout the paper, capital letters (e.g., X) mostly denote random variables, while their lower-case counterparts (e.g., x) denote instances of said random variables. We use the notation x^n to denote the sequence of x_1, \dots, x_n .

However, we may also use n in the superscript to emphasize the dependence on a quantity with n ; this will be clear in the context. Additionally, for an arbitrary set \mathcal{I} , $x_{\{1,\dots,n\}\setminus\mathcal{I}}$ indicates the sequence of x_i 's, where i iterates on $1, \dots, n$ excluding the elements in \mathcal{I} .

A. Resampling

In this section, we explain how to generate the said batches of samples from the dataset $\{(X_i, Y_i, Z_i)\}_{i=1}^n$. Define \mathcal{I}_b to be a set of b numbers picked uniformly at random (without replacement) from the set $\{1, \dots, n\}$. Let $\mathcal{B}_{\text{joint}}^b$ denote the joint batch, which consists of b samples distributed i.i.d. according to $p(x, y, z)$ and it is defined as:

$$\mathcal{B}_{\text{joint}}^b := \{(X_i, Y_i, Z_i) \mid i \in \mathcal{I}_b\}. \quad (9)$$

On the other hand, let $\mathcal{B}_{\text{prod}}^{b'}$ be the product batch such that it contains b' samples distributed according to $p(x|z)p(y, z)$. To construct this batch, we exploit the notion of k nearest neighbors (k -NN).

Definition 1. Assume the dataset $\{(x_i, y_i, z_i)\}_{i=1}^n$ of size n is given. Let \mathcal{I}_m be a set of m indices chosen uniformly at random without replacement from $\{1, \dots, n\}$, and $\mathcal{I}_m^c := \{1, \dots, n\} \setminus \mathcal{I}_m$. For any $\zeta \in \mathcal{X}$, define $\mathcal{A}^{m,k,n}(\zeta, z^n)$ as the set of indices of the k nearest neighbors of ζ (by Euclidean distance) among z_i , for $i \in \mathcal{I}_m^c$. In other words, let $\pi_\zeta : \{1, \dots, |\mathcal{I}_m^c|\} \rightarrow \mathcal{I}_m^c$ be the bijection such that

$$|z_{\pi_\zeta(1)} - \zeta| \leq \dots \leq |z_{\pi_\zeta(|\mathcal{I}_m^c|)} - \zeta|,$$

then $\mathcal{A}^{m,k,n}(\zeta, z^n) := \{\pi_\zeta(1), \dots, \pi_\zeta(k)\}$. Hereafter we use $\mathcal{A}^m(\zeta)$ instead as the remaining parameters can be understood from the context. In particular, we note that $\mathcal{A}^0(\zeta)$ implies that the neighbors are chosen from all points z^n since $\mathcal{I}_0 = \emptyset$ and $\mathcal{I}_0^c = \{1, \dots, n\}$.

According to the previous definition, the product batch with $b' = mk$ samples is defined as

$$\mathcal{B}_{\text{prod}}^{b'} := \{(X_{j(i)}, Y_i, Z_i) \mid i \in \mathcal{I}_m, j(i) \in \mathcal{A}^m(Z_i)\}. \quad (10)$$

We refer to this sampling technique as *isolated k-NN* in the sequel.

Similar to the k nearest neighbors method, the complexity of *isolated k-NN* relies on the particular implementation. With a brute force approach, the time complexity of computing distances from each sample in \mathcal{I}_m to \mathcal{I}_m^c is $\mathcal{O}(dn)$ when $m \ll n$, and a trivial search among these distances to find the k nearest neighbors costs an additional $\mathcal{O}(kn)$. Thus, the total time and storage complexities for all m samples in the isolated set \mathcal{I}_m are $\mathcal{O}((d+k)mn)$ and $\mathcal{O}(dn)$, respectively. However, if the dimension of the data is small ($d \ll \log n$) and the number of queries (in our case m) is large, we can use alternative methods such as k -d trees where pre-processing enables us to find neighbors more efficiently. Using k -d trees data structure with a nearly balanced tree, the time complexity becomes $\mathcal{O}((dn + mk) \log n)$ and the storage, $\mathcal{O}(dn)$.

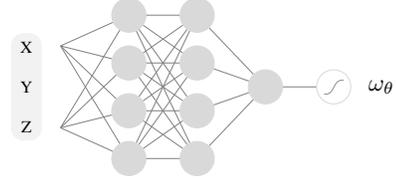


Fig. 1. The proposed neural network classifier that is parameterized with θ and consists of two hidden layers and concatenated with a sigmoid function.

B. Approximating f_{DV}^* , f_{NWJ}^* , and f_{LDR}^*

To estimate the optimal functions, it suffices to obtain the likelihood ratio $\frac{p(x,y,z)}{p(x|z)p(y,z)}$. As suggested in [17], [20], we use a feed-forward neural network to classify inputs from the joint and product batches. In this network, the input is a triple (x, y, z) and the last layer is concatenated with a sigmoid function. Let the network be parameterized with θ , then the output of the neural network is denoted as $\omega_\theta(x, y, z)$, see Fig. 1. As it will be clear later, in order to avoid unbounded values in the ratio of densities, the output of the sigmoid function is clipped¹ to the interval $[\tau, 1 - \tau]$ for $0 < \tau < \frac{1}{2}$.

The binary cross-entropy loss is chosen as the objective function to optimize the network. Define $q(x, y, z)$ (or simply q) as the batch type associated to an input, where $q = 1$ and $q = 0$ represent the joint and product batch, respectively. Then we have the following definition for the cross-entropy loss.

Definition 2. Given a function $\omega : \mathcal{X}^3 \rightarrow [0, 1]$, the expected cross-entropy loss is defined as:

$$L(\omega) := -\mathbb{E}_{p(q)p(x,y,z|q)} [Q \log \omega(X, Y, Z) + (1 - Q) \log(1 - \omega(X, Y, Z))]. \quad (11)$$

The pointwise minimizer of $L(\omega)$ can be used to compute the desired likelihood ratio and, accordingly, the optimal functions, as we see next. Note that unlike ω_θ , the function ω is not restricted by any parameterization.

Lemma 1. Let ω^* be the minimizer of the expected cross-entropy loss $L(\omega)$ and let $p(q = 1) = p_1$, then

$$\Gamma^*(x, y, z) := \frac{1 - p_1}{p_1} \frac{\omega^*(x, y, z)}{1 - \omega^*(x, y, z)} = \frac{p(x, y, z)}{p(x|z)p(y, z)}. \quad (12)$$

Using Lemma 1, the optimal functions (4), (5), and (7) can be evaluated as below,

$$\begin{aligned} f_{DV}^*(x, y, z) &:= C + \log \Gamma^*(x, y, z) \\ f_{NWJ}^*(x, y, z) &:= 1 + \log \Gamma^*(x, y, z) \\ f_{LDR}^*(x, y, z) &:= \log \Gamma^*(x, y, z). \end{aligned} \quad (13)$$

However, there are restrictions to obtain (13). First, the optimization to achieve ω^* is performed on $L(\omega_\theta)$ over the parameterized networks ω_θ , as searching over all functions is infeasible. Second, since the densities are not available, the expectations in $L(\omega_\theta)$ are approximated with sample averages.

¹It has been observed that such clipping also controls the bias–variance trade-off of the estimator [18]. In our case, choosing τ closer to zero decreases the bias and allows for the estimation of large values of CMI while it also increases the variance of the estimation.

Definition 3. Consider a neural-based classifier to be trained with sample batches $\mathcal{B}_{\text{joint}}^b$ and $\mathcal{B}_{\text{prod}}^{b'}$ such that

$$p_1 = \frac{b}{b + b'},$$

then the empirical cross-entropy loss is defined as:

$$L_{\text{emp}}(\omega_\theta) := p_1 L_b^1(\omega_\theta) + (1 - p_1) L_{b'}^2(\omega_\theta), \quad (14)$$

where

$$\begin{aligned} L_b^1(\omega_\theta) &:= -\frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \omega_\theta(x, y, z) \\ L_{b'}^2(\omega_\theta) &:= -\frac{1}{b'} \sum_{(x,y,z) \in \mathcal{B}_{\text{prod}}^{b'}} \log(1 - \omega_\theta(x, y, z)). \end{aligned} \quad (15)$$

Let $\hat{\theta}$ be the minimizer of $L_{\text{emp}}(\omega_\theta)$, according to the previous definition, and define

$$\hat{\Gamma}(x, y, z) := \frac{1 - p_1}{p_1} \frac{\omega_{\hat{\theta}}(x, y, z)}{1 - \omega_{\hat{\theta}}(x, y, z)}. \quad (16)$$

With a sufficiently large number of samples, n , and a proper tuning of the hyper-parameters of the network, $\hat{\Gamma}$ is close to Γ^* with high probability and the variational bounds for CMI can be estimated as:

$$\begin{aligned} \hat{I}_{DV}^{n, \hat{\theta}} &:= \frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \hat{\Gamma}(x, y, z) \\ &\quad - \log \frac{1}{b'} \sum_{(x,y,z) \in \mathcal{B}_{\text{prod}}^{b'}} \hat{\Gamma}(x, y, z), \\ \hat{I}_{NWJ}^{n, \hat{\theta}} &:= 1 + \frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \hat{\Gamma}(x, y, z) \\ &\quad - \frac{1}{b'} \sum_{(x,y,z) \in \mathcal{B}_{\text{prod}}^{b'}} \hat{\Gamma}(x, y, z). \end{aligned} \quad (17)$$

Similarly, the estimation based on LDR can be obtained as:

$$\hat{I}_{LDR}^{n, \hat{\theta}} := \frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \hat{\Gamma}(x, y, z). \quad (18)$$

There are two important caveats in computing these estimators. First, in training the classifier, it is desired to have $p_1 = \frac{1}{2}$ to avoid overfitting towards one of the classes. However, the prior can become biased due to the different resampling of the joint and product batches. Second, to implement the cross-validation, the final estimation is averaged over T trials where the train and test batches are re-sampled each time. This has been advocated in [14], [17] to control the variance of the estimation. Note that with the averaging over multiple trials, the final DV estimator is no longer a lower bound for the CMI [20]. The steps of our proposed method are stated in Algorithm 1. The functions `jointBatch` and `isolated_kNN` in the algorithm execute (9) and (10), respectively.

III. MAIN RESULTS

In this section, we discuss the consistency of the estimators $\hat{I}_{DV}^{n, \hat{\theta}}$, $\hat{I}_{NWJ}^{n, \hat{\theta}}$, and $\hat{I}_{LDR}^{n, \hat{\theta}}$, which in general relies on two things:

- The empirical sums in (15), (17), and (18) are concentrated around their expected values. For instance, this

Algorithm 1: Estimation of $I(X; Y|Z)$

Input: $Data = \{(x_i, y_i, z_i)\}_{i=1}^n, T, b, b', k$

- 1 Split $Data$ into $Train_set$ and $Test_set$
- 2 **for** $t=1, \dots, T$ **do**
- 3 $\mathcal{B}_{\text{joint,train}}^b \leftarrow \text{jointBatch}(Train_set, b)$
- 4 $\mathcal{B}_{\text{prod,train}}^{b'} \leftarrow \text{isolated_kNN}(Train_set, b', k)$
- 5 $\omega_{\hat{\theta}} \leftarrow \text{Train the classifier with } \mathcal{B}_{\text{joint,train}}^b, \mathcal{B}_{\text{prod,train}}^{b'}$
- 6 $\mathcal{B}_{\text{joint,test}}^b \leftarrow \text{jointBatch}(Test_set, b)$
- 7 $\mathcal{B}_{\text{prod,test}}^{b'} \leftarrow \text{isolated_kNN}(Test_set, b', k)$
- 8 Compute $\hat{I}_{DV}^{n, \hat{\theta}, t}$, $\hat{I}_{NWJ}^{n, \hat{\theta}, t}$, and $\hat{I}_{LDR}^{n, \hat{\theta}, t}$ using $\omega_{\hat{\theta}}$, $\mathcal{B}_{\text{joint,test}}^b$, and $\mathcal{B}_{\text{prod,test}}^{b'}$ as in (17) and (18)
- 9 **end**
- 10 $\hat{I}_{\text{est}}^{n, \hat{\theta}} \leftarrow \frac{1}{T} \sum_{t=1}^T \hat{I}_{\text{est}}^{n, \hat{\theta}, t}$ for $\text{est} = \text{'DV'}, \text{'NWJ'}, \text{'LDR'}$
- 11 **return** $\hat{I}_{DV}^{n, \hat{\theta}}, \hat{I}_{NWJ}^{n, \hat{\theta}}, \hat{I}_{LDR}^{n, \hat{\theta}}$

implies that for any θ , $L_{\text{emp}}(\omega_\theta)$ falls in the neighborhood of $L(\omega_\theta)$ with high probability, if certain conditions hold.

- The hyper parameters of the feed-forward neural network can be found such that with a perfect optimizer over parameters θ , one can desirably approximate ω^* , and accordingly f_{DV}^* , f_{NWJ}^* , and f_{LDR}^* .

In the following, we first show that the empirical average (using *isolated* k -NN re-sampled data) of any function g with bounded codomain converges to the expectation with respect to $p(x|z)p(y, z)$. Then, for our neural estimators of CMI, we exploit this result and the universal functional approximation theorem [23] to show that our estimators are consistent.

A. Concentration results

To obtain a high confidence concentration bound, let us make the following assumption on the value of k .

Assumption 1. Consider $k(n) = \Theta(n^{\frac{1}{2} + \epsilon_0})$ for some $\epsilon_0 > 0$ and $n - k(n) \geq m(n) \geq k(n)$. Then we select $b'(n) = m(n)k(n)$ samples to create the product batch² in the isolated k -NN method, using $k(n)$ neighbors and isolation set of size $m(n)$, as described in Definition 1. On the other hand, for the joint batch, assume $b(n) = \Theta(n)$. Hereafter we continue to use the notation b, b', m , and k , except where the dependency with n is important.

Remark 1. Assumption 1 is required to prove convergence. If we choose $m(n) = k(n) = n^{\frac{1}{2} + \epsilon_0}$, the size of the product batch becomes larger than n , i.e., $b'(n) = n^{1+2\epsilon_0}$. This is not an issue because the isolated k -NN technique enables us to construct batches of size larger than n , since the technique re-samples and mixes the original data, thus creating new samples. Nevertheless, we will show in the experimental results that even a smaller choice of k can yield a good estimation performance (e.g., see Fig. 4, where $n = 8e4$ while $k = 2$).

²It is worth noting that b' is then upper bounded by $m(n - m)$, as by choosing m indices for \mathcal{I}_m , we are left with at most $n - m$ samples from which to choose the neighbors, i.e., $k \leq n - m$.

Additionally, to balance the size of the joint and product batch, we can adjust the number of samples used to create the product batch. So in the example above, if we only use $\tilde{n} = n^{1/(1+2\epsilon_0)}$ samples and choose $m(\tilde{n}) = k(\tilde{n}) = \tilde{n}^{\frac{1}{2}+\epsilon_0}$, then $b'(n) = n$.

Now to address the concentration of the empirical average over samples taken with the isolated k -NN technique, we introduce the following theorem.

Theorem 1. Let $g(x, y, z) : \mathcal{X}^3 \rightarrow \mathbb{R}$ be any function such that $g^{\min} \leq g(x, y, z) \leq g^{\max}$, and $M := \max\{|g^{\min}|, |g^{\max}|\}$ is finite. Consider

$$\hat{g}(x^n, y^n, z^n) := \frac{1}{m} \sum_{i \in \mathcal{I}_m} \frac{1}{k} \sum_{j \in \mathcal{A}^m(z_i)} g(x_j, y_i, z_i), \quad (19)$$

with k and the set \mathcal{A}^m as defined in Assumption 1 and Definition 1, respectively. Then, for any $\epsilon > 0$ there exists an integer n_0 such that for $n > n_0$ and $m \leq n$,

$$\mathbb{P}\left(\left|\hat{g}(x^n, y^n, z^n) - \mathbb{E}_{p(x|z)p(y,z)}[g(X, Y, Z)]\right| \geq 3\epsilon\right) \leq \delta_1^n(\epsilon, c, M), \quad (20)$$

where δ_1^n is defined in Table I, $c := g^{\max} - g^{\min}$, and γ_d is the minimal number of cones centered at the origin, of angle $\pi/6$, that cover \mathbb{R}^d .

Proof. See Appendix A. \square

Remark 2. Note that $\lim_{n \rightarrow \infty} k(n)^2/n = \infty$, according to Assumption 1. Additionally, $n - m(n) = \Theta(n)$, which concludes that $\lim_{n \rightarrow \infty} \delta_1^n(\epsilon, c, M) = 0$.

Theorem 1, in conjunction with Hoeffding's inequality, leads to the concentration bound on $L_{\text{emp}}(\omega_\theta)$ found in the following proposition. This result is crucial in order to later show that $\omega_{\hat{\theta}}$ is close to ω^* , where we recall that $\hat{\theta}$ is the minimizer of the empirical loss.

Proposition 1. Let Assumption 1 hold. Then, for any $\mu > 0$ and any θ there exists n_0 such that for $n > n_0$,

$$\mathbb{P}\left(\left|L_{\text{emp}}(\omega_\theta) - L(\omega_\theta)\right| \geq \mu\right) \leq \delta_3^n(\mu), \quad (21)$$

where δ_3^n is defined in Table I.

Proof. See Appendix B. \square

B. Consistency of the estimators

To study the consistency of $\hat{I}_{DV}^{n, \hat{\theta}}$, $\hat{I}_{NWJ}^{n, \hat{\theta}}$, and $\hat{I}_{LDR}^{n, \hat{\theta}}$ in estimating $I(X; Y|Z)$, we make some further assumptions.

Assumption 2. There exist $0 < \alpha < \beta < \infty$ such that for any finite input $x, y, z \in \mathcal{X}^3$, the values of $p(x, y, z)$ and $p(x|z)p(y, z)$ are both constrained to the interval $[\alpha, \beta]$.

Remark 3. The constraint stated in Assumption 2 helps us show the consistency of our proposed estimators. However, most of our experiments are with Gaussian data, which does not fulfill the requirements of this assumption. Nonetheless, given the good simulation results, we believe that another

less restrictive assumption on the densities could replace Assumption 2.

Assumption 3. The classifier is parameterized with $\theta \in \Theta$ where $\Theta \subset \mathbb{R}^h$ and h is the number of parameters in the neural network. Also $\|\theta\|_2 \leq K$ for a constant K and the output of the classifier is B -Lipschitz with respect to θ . The hyper-parameters $h, K < \infty$ depend on the approximation power of the neural network to classify the dataset, and they can be determined according to the complexity of the dataset, the structure of the neural network, and others.

Remark 4. By restricting the output of the classifier to be Lipschitz continuous with respect to θ , the activation functions of the neural network must be differentiable (e.g., softplus). Nonetheless, we use rectified linear units (ReLU) in our experiments, similar to [14], [16], [17], and obtain a desirable estimation performance. Note that the softplus function, defined as $f(x) = \frac{1}{t} \ln(1 + e^{tx})$, is equivalent to ReLU, asymptotically as $t \rightarrow \infty$. However, the use of the ReLU function is encouraged over the softplus [24].

While Assumption 1 guarantees that $\delta_1^n(\epsilon, c, M), \delta_2^n(\epsilon), \dots, \delta_7^n(\epsilon)$ tend to zero asymptotically as $n \rightarrow \infty$, in order to obtain a concentration bound, the sample size n needs to be larger than a certain threshold. This value is determined by the true density $p(x, y, z)$ and the hyper-parameters of our setup, and it is stated in the following assumption.³

Assumption 4. For given $\epsilon^* > 0$ and $\delta^* > 0$, we assume that n is large enough such that the following conditions hold:

$$\delta_4^n\left(\frac{\epsilon}{8}\right) + \delta_i^n(\epsilon^*) \leq \delta^*, \quad \text{for } i \in \{5, 6, 7\},$$

where ϵ and $\delta_i^n(\cdot)$ are defined in Table I. Note that by the continuity of the functions $\delta_i^n(\cdot)$ and their asymptotic behavior, finding such an n is feasible.

Now we are able to express the consistency of our estimators in terms of concentration bounds in the following theorem.

Theorem 2. Let Assumptions 1, 2, 3, and 4 hold and $0 < \tau < \min\{\frac{1}{2}, p_1\}$. Then, given $\epsilon^*, \delta^* > 0$, there exists an integer n^* such that for all $n > n^*$,

$$\mathbb{P}\left(\left|\hat{I}_{\text{est}}^{n, \hat{\theta}} - I(X; Y|Z)\right| \geq \epsilon^*\right) \leq \delta^*, \quad (22)$$

where 'est' can be replaced with 'DV', 'NWJ', or 'LDR'.

Proof. See Appendix C. \square

Remark 5. Note that the proper choice of the hyper-parameters of the network and the value of n^* crucially relies on the true underlying density, and thus the bounds are not universal in that sense. This has been emphasized in [25, Remark 7] where the authors discuss required precautions for the neural estimator in [14].

³The sample complexity for the neural estimator of MI has been discussed in [14, Theorem 3] and recently revisited in the fourth online version of [15]. Similar results exist for the classifier estimator for the CMI in [17, Lemma 5].

TABLE I
TABLE OF PARAMETERS.

$$\begin{aligned}
\delta_1^n(\epsilon, c, M) &:= 2 \exp\left(\frac{-2\epsilon^2 k^2}{nc^2}\right) + 2 \exp\left(\frac{-2\epsilon^2 k^2}{(n-m)c^2}\right) + \exp\left(\frac{-(n-m)\epsilon^2}{8M^2\gamma_d^2}\right) \\
\delta_2^n(\epsilon) &:= \delta_1^n(\epsilon, \log \frac{1-\tau}{\tau}, -\log \tau) \\
\delta_3^n(\epsilon) &:= \delta_2^n\left(\frac{\epsilon}{3-2p_1}\right) + 2 \exp\left(\frac{-2b(n)\epsilon^2}{((3-2p_1)\log \frac{1-\tau}{\tau})^2}\right) \\
\delta_4^n(\epsilon) &:= \left(\frac{4BK\sqrt{h}}{\tau\epsilon}\right)^h \delta_3^n(\epsilon) \\
\delta_5^n(\epsilon) &:= \delta_1^n\left(\frac{(1-p_1)\epsilon\tau}{2\tau+6p_1-8p_1\tau}, \frac{1-p_1}{p_1} \frac{1-2\tau}{\tau(1-\tau)}, \frac{1-p_1}{p_1} \frac{1-\tau}{\tau}\right) \\
&\quad + 2 \exp\left(\frac{-b(n)(1-p_1)^2\epsilon^2\tau^2}{2((2\tau+6p_1-8p_1\tau)\log \frac{1-\tau}{\tau})^2}\right) \\
\delta_6^n(\epsilon) &:= \delta_1^n\left(\frac{\epsilon}{8}, \frac{1-p_1}{p_1} \frac{1-2\tau}{\tau(1-\tau)}, \frac{1-p_1}{p_1} \frac{1-\tau}{\tau}\right) + 2 \exp\left(\frac{-b(n)\epsilon^2}{128(\log \frac{1-\tau}{\tau})^2}\right) \\
\delta_7^n(\epsilon) &:= 2 \exp\left(\frac{-b(n)\epsilon^2}{8(\log \frac{1-\tau}{\tau})^2}\right) \\
\eta &:= \frac{\tau^3(1-\tau)\epsilon^*}{2(2\tau^2-2\tau+1)\beta} \\
\epsilon &:= \left(\frac{\eta}{1-\tau}\right)^2 \frac{\alpha}{2\lambda(\mathcal{X})}
\end{aligned}$$

IV. EXPERIMENTS

In this section, we compare our technique⁴ with the state-of-the-art approach proposed in [17], which we refer to as *MI-Diff* (also known as the CCMi method) since the method is based on computing the CMI with the difference of two MI terms, $I(X; Y, Z)$ and $I(X; Z)$, as in (8). Each MI term is then estimated by utilizing a neural network classifier with a similar structure as our method and training with the proper joint and product batches. In contrast to the *isolated k-NN* method, the construction of batches in *MI-Diff* is straightforward. The joint batches are created similar to (9), while the product batches for $I(X; Y, Z)$ and $I(X; Z)$ are constructed by taking b random indices separately from x^n and (y^n, z^n) . In particular for $I(X; Z)$, the product batch in the *MI-Diff* method is:

$$\mathcal{B}_{\text{prod}}^b := \{(X_i, Z_j) \mid i \in \mathcal{I}_b^{(1)}, j \in \mathcal{I}_b^{(2)}\}, \quad (23)$$

where $\mathcal{I}_b^{(1)}$ and $\mathcal{I}_b^{(2)}$ are random indices selected separately in $\{1, \dots, n\}$. Similarly for $I(X; Y, Z)$ the product batch is:

$$\mathcal{B}_{\text{prod}}^b := \{(X_i, Y_j, Z_j) \mid i \in \mathcal{I}_b^{(1)}, j \in \mathcal{I}_b^{(2)}\}. \quad (24)$$

Initially, we verify the approximation power and consistency of the estimators in two scenarios where the CMI is either non-zero (part A) or zero (part B). Additionally, we investigate the performance of our method as the dimension of data grows (part C). The generative model that we use is defined as:

$$\begin{aligned}
X &\sim \mathcal{N}(0, \sigma_x^2 \Sigma_d), \\
Y &\sim \mathcal{N}(X, \sigma_y^2 \Sigma_d), \\
Z &\sim \mathcal{N}(Y, \sigma_z^2 \Sigma_d).
\end{aligned} \quad (25)$$

In order to meet Assumption 2, we could use a truncated normal distribution by bounding the ℓ_2 norm of the random variables. However, slight deviations from this assumption do not significantly change the statistics of the generated dataset since the likelihood of observing a very large or low value is negligible.

⁴Code: https://github.com/smolavipour/CMI_Neural_Estimator

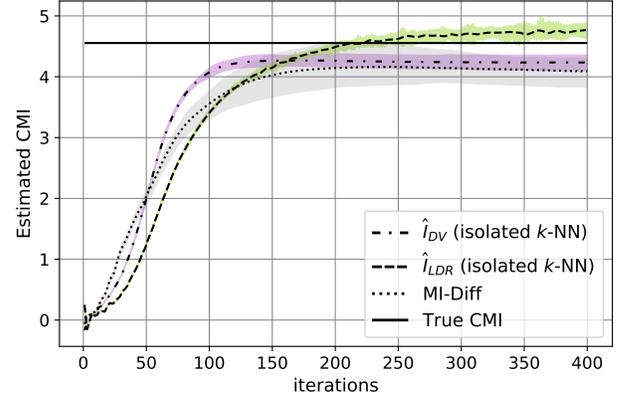


Fig. 2. Evolution of the estimation of $I(X; Y|Z)$ over training iterations, using our estimators $\hat{I}_{DV}^{n, \hat{\theta}}$ and $\hat{I}_{LDR}^{n, \hat{\theta}}$ compared with the *MI-Diff* method based on the DV bound. The input dataset contains $n = 8e4$ samples with $d = 3$. The shadows are based on the maximum and minimum values obtained at each iteration.

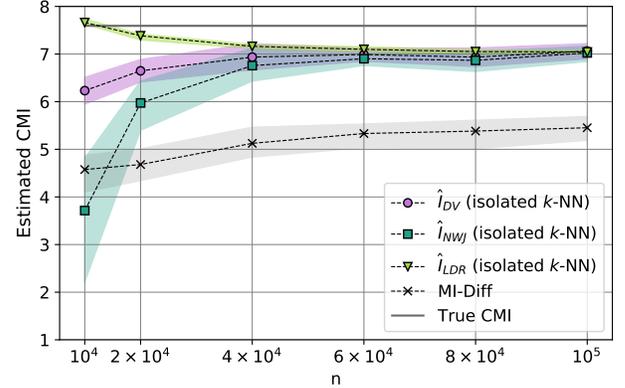


Fig. 3. Comparison between our proposed estimators using *isolated k-NN*, with $k = 2$, and the *MI-Diff* method to estimate $I(X; Y|Z)$, with $d = 5$, $E = 300$, and different values of n .

A. Estimating $I(X; Y|Z)$

Consider $\sigma_x = 10$, $\sigma_y = 1$, and $\sigma_z = 5$ in our model (25), and let $\Sigma_d = I_d$, i.e., the identity matrix of dimension d . According to the model, we compute $I(X; Y|Z)$ as follows:

$$\begin{aligned}
I(X; Y|Z) &= I(X; Y) - I(X; Z) \\
&= \frac{d}{2} \log \left(1 + \frac{\sigma_x^2}{\sigma_y^2} \right) - \frac{d}{2} \log \left(1 + \frac{\sigma_x^2}{\sigma_y^2 + \sigma_z^2} \right).
\end{aligned}$$

The evolution of the neural classifier's performance using a validation dataset during training is shown in Fig. 2, for $d = 3$ and $k = 2$. We made a comparison between the estimators $\hat{I}_{DV}^{n, \hat{\theta}}$ and $\hat{I}_{LDR}^{n, \hat{\theta}}$ according to Algorithm 1, with the *MI-Diff* method. Both the *MI-Diff* and $\hat{I}_{DV}^{n, \hat{\theta}}$ estimators converge after $E = 200$ epochs, while $\hat{I}_{LDR}^{n, \hat{\theta}}$ requires more iterations ($E \geq 300$) to converge. Comparing the range of the estimations suggests that the DV and LDR estimators have lower variance compared to the *MI-Diff*.

Next, we show significant improvements of our estimators compared with the *MI-Diff* method when the dimension increases. In Fig. 3, a comparison of the estimators for CMI

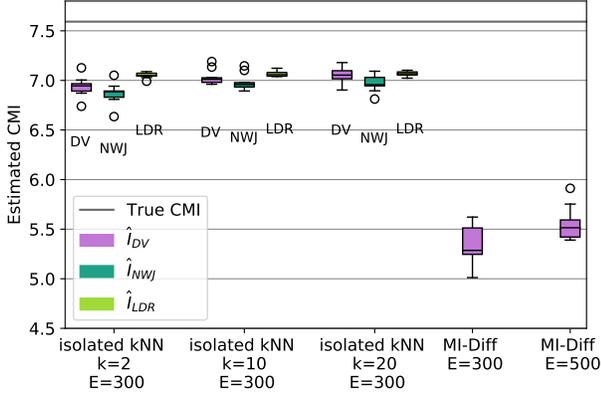


Fig. 4. Comparison between the *isolated k-NN* and the *MI-Diff* methods to estimate $I(X; Y|Z)$, for an input dataset with $n = 8e4$ samples and $d = 5$.

is depicted for $d = 5$ in terms of sample size n . Our LDR estimator performs better than both DV and NWJ estimators in terms of bias and variance. Note that the LDR estimator is averaging the density ratio over samples in the joint batch as

$$\frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \frac{p(x, y, z)}{p(x|z)p(y, z)}.$$

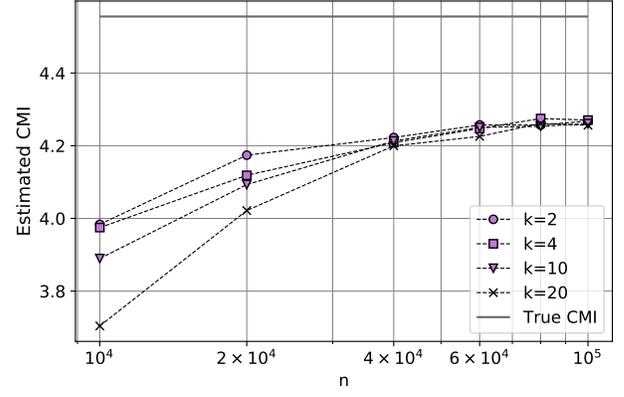
For a small number of samples, typically the samples with high probability density $p(x, y, z)$ appear while by having more samples, the chance of observing odd samples with low probability increases, which compensates the total average. This effect results in the LDR estimation to have a decreasing behavior by increasing n . On the other hand, the DV estimator is of the form

$$\frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \frac{p(x, y, z)}{p(x|z)p(y, z)} - \log \frac{1}{b'} \sum_{(x,y,z) \in \mathcal{B}_{\text{prod}}^{b'}} \frac{p(x, y, z)}{p(x|z)p(y, z)},$$

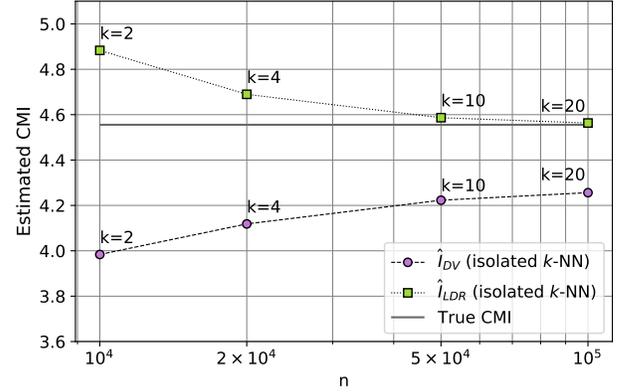
where the second sum is dominated by the unlikely events. So even by observing one odd event, the second term becomes very large. As n increases, more typical samples are collected in the product batch and this effect disappears.

In Fig. 4, $I(X; Y|Z)$ is estimated for $n = 8e4$ and dimension $d = 5$, with different choices of k . Then the results are compared with the *MI-Diff* method for estimating the DV bound. It can be observed that sampling batches using *isolated k-NN* improves the accuracy of the estimation. To ensure that the training in the *MI-Diff* method has been done with enough epochs, we repeated the experiment with more epochs as well. Despite leveraging additional learning iterations, both accuracy and variance of our estimators are more desirable.

As suggested in the *isolated k-NN* method, increasing k can improve the estimation if it is properly scaled with n and $\lim_{n \rightarrow \infty} k(n)/n = 0$. To investigate this, we compare the estimated CMI with $\hat{I}_{DV}^{n, \hat{\theta}}$ for $d = 3$ and different choices of k and n in Fig. 5a. In general, increasing the number of samples for a fixed k results in a more accurate estimation, as shown in Fig. 5a. However, when the number of samples is fixed, choosing a larger k worsens the estimation. The reason is that with n being fixed and $b = mk = \frac{n}{2}$, m becomes smaller and there are less samples of (y, z) to estimate the expectation $\mathbb{E}_{p(y,z)}[\cdot]$ with sample average. Nonetheless, this



(a) DV estimation for different values of k and n .



(b) DV and LDR estimations with a fixed k/n ratio.

Fig. 5. Estimated $I(X; Y|Z)$ using the *isolated k-NN* technique, with $d = 3$.

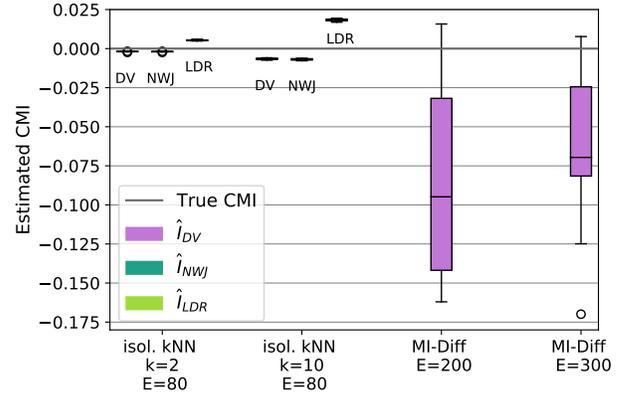


Fig. 6. Comparison between the *isolated k-NN* and the *MI-Diff* methods to estimate $I(X; Z|Y) = 0$, for $n = 8e4$ and $d = 3$.

behavior can be resolved if $k = k(n)$ increases with n , and as a result m can remain sufficiently large to obtain a desired accuracy. This is illustrated in Fig. 5b where with a fixed ratio of $k(n)/n$, the estimation improves by increasing k .

B. Estimating $I(X; Z|Y)$ (zero CMI)

A desirable estimator must be able to estimate both high and low values of CMI. In this scenario, we test the ability of our estimator for zero CMI. Due to the Markov chain $X \rightarrow Y \rightarrow Z$ in the model (25), $I(X; Z|Y) = 0$. Consider the same

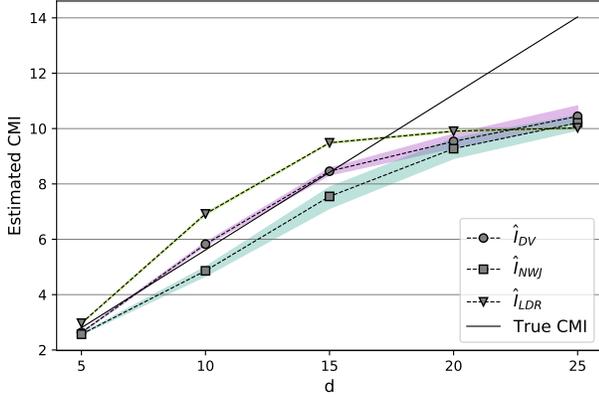
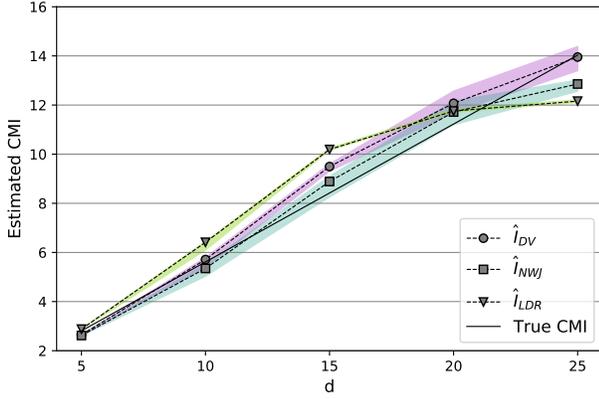
(a) Estimation with $n = 2e5$ and $\tau = 3e-5$ (b) Estimation with $n = 1e6$ and $\tau = 1e-6$

Fig. 7. Performance of different estimators using the *isolated k-NN* method to estimate $I(X; Y|Z)$ for different data dimensions d . The shadows are based on the maximum and minimum values obtained at each iteration.

choices for $\sigma_x, \sigma_y, \sigma_z$, and Σ_d as in previous part while $d = 3$. In Fig. 6 the box-plots are created by repeating Algorithm 1 and *MI-Diff* for 10 Monte Carlo trials and the sample size is $n = 8e4$. The results of *isolated k-NN* are shown for $k = 2$ and 10. The performance of the *MI-Diff* method is shown for different choices of number of epochs, while we exclude the case of $E = 80$ due to poor estimation (median < -0.3).

It can be observed that our technique has the advantage of lower bias and variance compared with *MI-Diff*. We note that, as explained earlier, the degrading performance by increasing k is due to n being fixed. Furthermore, although the CMI is non-negative, we see that the estimation can become negative if the density ratio is not estimated properly.

C. Effect of the data dimension

In order to test the performance of the estimators as the input dimension grows, we consider the generative model (25) with $\sigma_x = 10, \sigma_y = 3, \sigma_z = 5$, and $\Sigma_d = I_d$. In this experiment, we first estimate the CMI for several values of d using $n = 2e5$ number of samples and choosing $\tau = 3e-5$. As depicted in Fig. 7a, the estimation seems to saturate for dimension $d = 25$. This mainly occurs as the output of the neural network is clipped between $[\tau, 1 - \tau]$, so the density ratio and the estimated CMI are bounded accordingly.

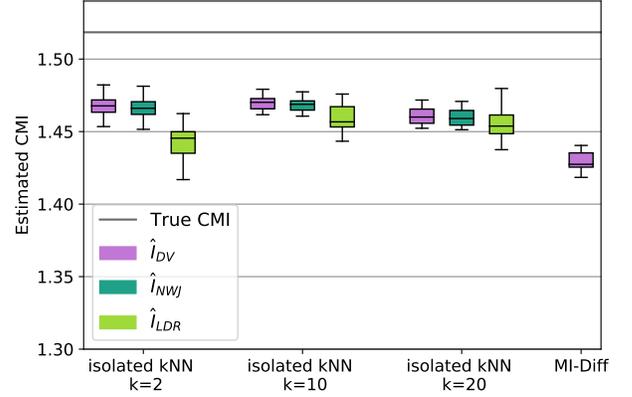


Fig. 8. Comparison between the *isolated k-NN* and the *MI-Diff* methods to estimate $I(f(X); Y|Z)$, where $f(x) = \tanh(0.05x)$, for an input dataset with $n = 8e4$ samples and $d = 1$.

Note that, although a smaller τ allows the estimator to reach higher values, it also makes the output more unstable (as remarked in [18]). To explain this behavior, consider the LDR estimation, $\frac{1}{b} \sum_{(x,y,z) \in \mathcal{B}_{\text{joint}}^b} \log \hat{\Gamma}(x, y, z)$. From (16) and assuming $p_1 = 0.5$, the estimated density ratio $\hat{\Gamma}(x, y, z)$ is bounded by $(1 - \tau)/\tau$. So if τ is very small, in some rare events, $\hat{\Gamma}(x, y, z)$ can become very large and affect the average. Nevertheless, by increasing n , more typical samples are included and we prevent the odd events from dominating the estimated value.

In Fig. 7b, we repeat the experiment with $n = 1e6$ and $\tau = 1e-6$. It can be seen that the estimations perform more accurately, in particular for higher dimensions. Additionally, it can be noted that the estimated CMI passes over the true value (even based on the NWJ bound) in some cases which is a consequence of choosing a small τ .

D. Non-linear model

To strengthen our justification on the proposed CMI neural estimator, we consider a non-linear scenario. First note that for any injective function $f(\cdot)$, $I(f(X); Y|Z) = I(X; Y|Z)$. This property allows us to test the performance of the CMI estimators when a non-linear function is applied on the data, while computing the true CMI remains tractable.

We thus estimate $I(f(X); Y|Z)$ with $n = 8e4$ samples for $f(x) = \tanh(0.05x)$ and the model defined in (25), assuming $\sigma_x, \sigma_y, \sigma_z$ equal to part A, while we choose the dimension $d = 1$ and $\Sigma_d = 1$; the coefficient inside $\tanh(\cdot)$ was chosen to avoid saturation of the output. The estimation results are plotted in Fig. 8, where it is clear that the non-linear function does not hinder the estimation performance. To compare our estimators with the *MI-Diff* method, we perform the Mann-Whitney U test. Since the box plots for the DV and NWJ estimator are superior to the result of the *MI-Diff* method and do not overlap, we only test the LDR estimation results. We obtained that our LDR estimator performs better than the *MI-Diff* method for all choices of k with p -values less than 0.01.

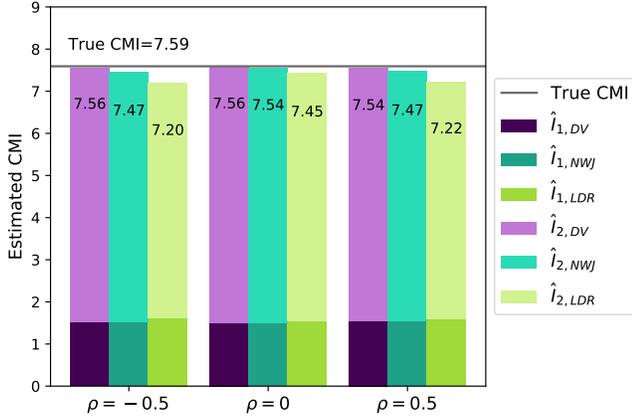


Fig. 9. Testing additivity property and DPI for $d = 5$, $d_1 = 1$, and $d_2 = 4$.

E. Additivity and data processing inequality

In this part, we test two properties of CMI in our estimations. We consider a slightly more complex model where we assume a component-wise dependency in (25) by choosing:

$$\Sigma_d = \begin{bmatrix} 1 & \rho & 0 & \dots & 0 \\ \rho & 1 & \rho & 0 & \dots & 0 \\ 0 & \rho & 1 & \rho & 0 & \dots & 0 \\ \vdots & & \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & \rho & 1 & & \end{bmatrix}.$$

Let Y , with dimension d , be decomposed into two parts Y_1 and Y_2 with dimensions d_1 and d_2 , respectively. Then the conditional mutual information can be split as below:

$$I(X; Y|Z) = I(X; Y_1|Z) + I(X; Y_2|Y_1, Z). \quad (26)$$

Denote the estimations of $I(X; Y_1|Z)$ and $I(X; Y_2|Y_1, Z)$ by $\hat{I}_{1,est}$ and $\hat{I}_{2,est}$, respectively. In this experiment, we test if:

- $\hat{I}_{1,est} + \hat{I}_{2,est}$ can estimate $I(X; Y|Z)$ (Additivity)
- $\hat{I}_{1,est}$ is smaller than $I(X; Y|Z)$ (data processing inequality (DPI))

Fig. 9 depicts the estimations $\hat{I}_{1,est}$ and $\hat{I}_{2,est}$ where *est* can be replaced with 'DV', 'NWJ', and 'LDR'. We use $n = 2e5$ samples with $d = 5$, $d_1 = 1$, and $d_2 = 4$, and the results show the averaged estimated values over all trials. It can be observed that both the additivity property and DPI hold for different choices of ρ .

F. Real-world scenario

One application of CMI is in the definition of directed information (DI), which quantifies the extent of causal effects in a network of processes. The directed information rate from X to Y causally conditioned on Z is defined as:

$$I(X \rightarrow Y||Z) := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n I(X^i; Y_i|Z^i, Y^{i-1}),$$

which if certain Markov properties hold, is simplified as

$$I(X \rightarrow Y||Z) = I(X^l; Y_l|Z^l, Y^{l-1}), \quad (27)$$

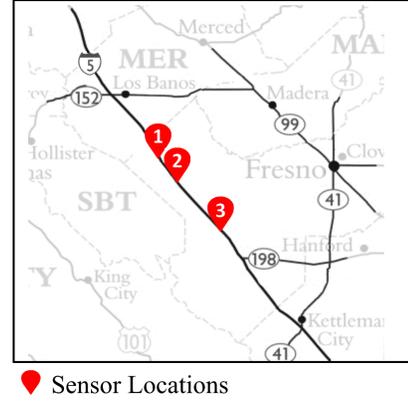


Fig. 10. Vehicular traffic sensors mounted on interstate 5 south, California.

where l is the Markov order in the data model (see [3] for more details). In other words, directed information captures the dependency of Y_l to the history of X^l , when the history of Y^{l-1} and the rest of the network (Z^l) is given. The causal effects are conventionally visualized via a directed graph where the weights of the links are the DI values; this is called a directed information graph (DIG). Assuming a network with three processes X , Y , and Z , the weight of the link $X \rightarrow Y$ in the DIG represents $I(X \rightarrow Y||Z)$.

In [26], the authors exploited the neural estimator in [17] to estimate the transfer entropy, which is a similar notion to DI. This motivates applying our neural estimators for DI, although the data is not i.i.d. but rather Markov. In this experiment, we apply our estimators on a real-world dataset which is collected from vehicular traffic sensors mounted on California highways⁵. Estimating DIG for vehicular traffic has been studied in [27] by quantizing the measurements' values and exploiting classical methods for discrete-value processes.

In this experiment, traffic flow data represents the number of cars passing the sensor for a time period of 5 minutes. We choose three consecutive sensors on the interstate 5 south, in Fresno (see Fig. 10), and the data is aggregated from January 2017 to June 2019, every day, for 24 hours. From the approximate distance between sensors and the average speed of the cars, we choose $l = 5$ in (27) to estimate the DI rate and construct the DIG graph. Since the sensors are in a row, the vehicles which are observed in one sensor will appear in the next sensor with a delay. However, due to the inbound and outbound traffic, it is possible to miss part of the flow.

The estimated DI between sensors is stated below in terms of weights of the adjacency matrix of the DIG:

$$G = \begin{bmatrix} 0 & 0.22 & 0.06 \\ 0.1 & 0 & 0.1 \\ 0.01 & 0.06 & 0 \end{bmatrix}.$$

The results indicate that the most significant causal effects are between sensor 1 and sensor 2, and from sensor 2 to sensor 3. This is in line with the physical placement of the sensors. Note that in (27) the dependency of X_l on Y_l is included in the DI.

⁵Accessed from California Department of Transportation (<http://pems.dot.ca.gov/>).

So if a car passes two sensors in the same time interval, it will be considered in both traffic flow signals. The backward link from sensor 2 to sensor 1 could be the result of this event since they are closer in distance (see Fig. 10).

V. DISCUSSION AND FINAL REMARKS

In this paper, we studied the use of a neural network classifier to estimate the conditional mutual information among some random variables, based on a dataset composed of i.i.d. samples of said random variables. Inspired by the k -NN method, we introduced a new technique for creating sample batches which re-samples the existing dataset; this re-sampling models a particular conditional independence in the distribution of the new samples. The classifier is then trained to distinguish between the original distribution of samples and the new one. This technique enabled us to estimate the CMI directly rather than estimating it as the difference of two MI terms. Our simulations showed that estimating with the proposed *isolated k-NN* method improved the accuracy of estimation for high and low values of CMI in several scenarios. However, extending the results to more complicated models or other real-world data requires further investigation in tuning the neural network and choosing its activation functions. This can be considered as a future direction of this work.

While the estimation based on the NWJ bound is always less than the DV estimation, we cannot make a general claim about the LDR estimation. For instance, in Fig. 5b the LDR estimation is higher than the DV estimation, while this does not hold in the non-linear experiment (Fig. 8) or the results in Fig. 9. In addition, we emphasize that a higher estimation does not necessarily reduce the bias as the estimations may overshoot the true value; for example in Fig. 6 or Fig. 7b.

Neural networks have been proposed in communication systems as part of the encoder/decoder blocks [28], [29]. However, learning the end-to-end communication system requires knowing the channel model, which might not be available in practice. While there exist approaches based on generative adversarial networks (GAN) [30], in [31], the authors optimize channel encoders by estimating the MI and advocate the use of neural estimators. This approach can be followed with our proposed estimators for channels with capacities characterized by CMI. However, as noted in [20], one should be careful to use an appropriate estimator for CMI; if the estimated value is used to determine the transmission rate and it is above the true value of CMI, the system will experience a catastrophic failure.

In some applications, a key step is to perform a threshold test on CMI rather than estimating its exact value. For instance in [3], the causal links of a network of random processes can be detected by checking if the CMI is above a certain threshold. In order to achieve a high accuracy in such tests (i.e., small type-I/II errors), it is not necessarily required to estimate the CMI accurately. Therefore, the performance of the tests can be investigated as a future direction of this work.

APPENDIX A PROOF OF THEOREM 1

A. Preliminaries

First let us review the lemmas that we require in this proof. Since each of the terms in $\hat{g}(x^n, y^n, z^n)$ is bounded, McDiarmid's inequality [32] is exploited to obtain concentration bounds.

Lemma 2 (McDiarmid's inequality). *Let V_1, \dots, V_n be independent random variables $V_i \in \mathcal{V}$ and assume $\phi : \mathcal{V}^n \rightarrow \mathbb{R}$ such that for all $i \in \{1, \dots, n\}$:*

$$\sup_{\substack{v_1, \dots, v_n \\ v'_i}} |\phi(v_1, \dots, v_n) - \phi(v_1, \dots, v'_i, \dots, v_n)| \leq c_i.$$

Then the following bound holds:

$$\mathbb{P}\left(|\phi(V^n) - \mathbb{E}[\phi(V^n)]| \geq \epsilon\right) \leq 2 \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}\right).$$

The inner sum in $\hat{g}(\cdot)$, defined in (19), resembles a k -NN regression which we leverage in our proof.

Lemma 3 ([33, Theorem 1]). *Let $(U_1, Z_1), \dots, (U_n, Z_n)$ be generated i.i.d. according to $p(u, z)$ and let $\zeta \in \mathcal{X}$. Using the same notation as in Definition 1, let $\mathcal{A}^m(\zeta)$ be the set of indices of the elements of z^n which are the k -NN of ζ , and define*

$$\psi_n^m(\zeta) := \frac{1}{k} \sum_{j \in \mathcal{A}^m(\zeta)} U_j,$$

and $\bar{\psi}(\zeta) := \mathbb{E}_{p(u|z)}[U|Z = \zeta]$. Further assume that $|U| \leq M$ and $\lim_{n \rightarrow \infty} k(n) = \infty$ and $\lim_{n \rightarrow \infty} k(n)/n = 0$. Then, if the neighbors are chosen in Z_1^n (i.e., by taking indices in $\mathcal{A}^0(\cdot)$), for any $\epsilon > 0$ there exists an integer n_0 such that for $n > n_0$:

$$\mathbb{P}\left(\int p(z) |\psi_n^0(z) - \bar{\psi}(z)| dz > \epsilon\right) \leq \exp\left(\frac{-n\epsilon^2}{8M^2\gamma_d^2}\right), \quad (28)$$

where γ_d is the minimal number of cones centered at the origin, of angle $\pi/6$, that cover \mathbb{R}^d .

Remark 6. *In Lemma 3, according to the definition of $\mathcal{A}^m(\zeta)$, we assume that $m = 0$. Nevertheless, if $0 < m < n$ and $\lim_{n \rightarrow \infty} k(n)/(n - m) = 0$, since the pairs are i.i.d., similar to (28), we have that*

$$\mathbb{P}\left(\int p(z) |\psi_n^m(z) - \bar{\psi}(z)| dz > \epsilon\right) \leq \exp\left(\frac{-(n - m)\epsilon^2}{8M^2\gamma_d^2}\right).$$

B. Proof of Theorem 1

To begin the proof, we make the following definitions:

$$g_n^k(y, z) := \frac{1}{k} \sum_{j \in \mathcal{A}^m(z)} g(x_j, y, z),$$

$$\bar{g}(y, z) := \mathbb{E}_{p(x|z)}[g(X, y, z)].$$

Note that $g_n^k(y, z)$ is in fact a function of x^n, z^n and y, z ; We use the simplified notation as the dependence on the data can be understood from the context. Moreover, since the pairs (y_i, z_i) are i.i.d. from the dataset, we may assume

$\mathcal{I}_m = \{1, \dots, m\}$ without loss of generality. We thus rewrite the estimator (19) as:

$$\hat{g}(x^n, y^n, z^n) = \frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i). \quad (29)$$

Now, using the triangle inequality, we have that

$$\begin{aligned} & \left| \hat{g}(x^n, y^n, z^n) - \mathbb{E}_{p(x|z)p(y,z)}[g(X, Y, Z)] \right| \\ &= \left| \frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i) - \int p(y, z) \bar{g}(y, z) dy dz \right| \\ &\leq \left| \frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i) - \int p(y, z) g_n^k(y, z) dy dz \right| \\ &\quad + \left| \int p(y, z) (g_n^k(y, z) - \bar{g}(y, z)) dy dz \right| \\ &\leq \left| \frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i) - \int p(y, z) g_n^k(y, z) dy dz \right| \\ &\quad + \left| \int p(z) \left| \int p(y|z) (g_n^k(y, z) - \bar{g}(y, z)) dy \right| dz \right| \\ &\leq \left| \frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i) - \mathbb{E}[g_n^k(Y, Z)] \right| \\ &\quad + \left| \int p(y, z) g_n^k(y, z) dy dz - \mathbb{E}[g_n^k(Y, Z)] \right| \\ &\quad + \left| \int p(z) \left| \int p(y|z) (g_n^k(y, z) - \bar{g}(y, z)) dy \right| dz \right|. \quad (30) \end{aligned}$$

To elaborate on the first two terms on the RHS of (30), we note that $\frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i)$ is a function of the random variables X_{m+1}^n, Y^m, Z^n , and thus random itself, while the randomness of $\int p(y, z) g_n^k(y, z) dy dz$ in the second term stems from X_{m+1}^n, Z_{m+1}^n . On the other hand, $\mathbb{E}[g_n^k(Y, Z)]$ is a deterministic term and the expectation is with respect to the density function $p(y, z)p(x_{m+1}^n, z_{m+1}^n)$.

In the following, we show the convergence of the first two terms in (30) according to Lemma 2. Next we show that the last term converges to zero according to Lemma 3 and Remark 6. Note that Assumption 1 guarantees the required assumption on k in Lemma 3 and Remark 6.

1) *First term in (30):* Define $w_i := (x_i, y_i, z_i)$ and let

$$\phi(w^n) = \frac{1}{m} \sum_{i=1}^m g_n^k(y_i, z_i),$$

which is a function of the random triples $\{(X_i, Y_i, Z_i)\}_{i=1}^n$. For any $i \in \{1, \dots, n\}$ we have that

$$\sup_{w^n, w'_i} \left| \phi(w_1, \dots, w_n) - \phi(w_1, \dots, w'_i, \dots, w_n) \right| \leq \frac{c}{\min\{m, k\}}, \quad (31)$$

where $c := g^{max} - g^{min}$. To see this, first consider a triple (x_i, y_i, z_i) is altered to (x'_i, y'_i, z'_i) for $i \leq m$. Then the largest difference that can happen is c/m . In case $i > m$, the extreme case is that z_i is the neighbor of all z_1, \dots, z_m , so in total the difference becomes c/k . By Assumption 1, $m > k$ and thus the RHS of (31) becomes c/k .

Since (31) holds, Lemma 2 implies the following bound:

$$\mathbb{P} \left(\left| \phi(W^n) - \mathbb{E}_{p(w^n)}[\phi(W^n)] \right| > \epsilon \right) \leq 2 \exp \left(\frac{-2\epsilon^2 k^2}{nc^2} \right).$$

The expectation inside the left hand side (LHS) of this equation may be rewritten as follows:

$$\begin{aligned} \mathbb{E}_{p(w^n)}[\phi(W^n)] &= \frac{1}{m} \mathbb{E}_{p(w^n)} \left[\sum_{i=1}^m g_n^k(Y_i, Z_i) \right] \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{p(y_i, z_i)p(x_{m+1}^n, z_{m+1}^n)} [g_n^k(Y_i, Z_i)] \\ &= \mathbb{E}_{p(y, z)p(x_{m+1}^n, z_{m+1}^n)} [g_n^k(Y, Z)], \quad (32) \end{aligned}$$

where the last equality holds since the pairs (Y_i, Z_i) are generated i.i.d. As a result,

$$\mathbb{P} \left(\left| \frac{1}{m} \sum_{i=1}^m g_n^k(Y_i, Z_i) - \mathbb{E}[g_n^k(Y, Z)] \right| > \epsilon \right) \leq 2 \exp \left(\frac{-2\epsilon^2 k^2}{nc^2} \right). \quad (33)$$

2) *Second term in (30):* Similarly, let

$$\phi'(w_{m+1}^n) = \int p(y, z) g_n^k(y, z) dy dz; \quad (34)$$

then, for any $i \in \{m+1, \dots, n\}$, we have that

$$\sup_{w_{m+1}^n, w'_i} \left| \phi'(w_{m+1}^n) - \phi'(w_{m+1}, \dots, w'_i, \dots, w_n) \right| \leq \frac{c}{k}.$$

Hence, Lemma 2 yields the following bound:

$$\mathbb{P} \left(\left| \phi'(W_{m+1}^n) - \mathbb{E}[\phi'(W_{m+1}^n)] \right| > \epsilon \right) \leq 2 \exp \left(\frac{-2\epsilon^2 k^2}{(n-m)c^2} \right). \quad (35)$$

The deviation of the second term in (30) can thus be bounded as below:

$$\mathbb{P} \left(\left| \int p(y, z) g_n^k(y, z) dy dz - \mathbb{E}[g_n^k(Y, Z)] \right| > \epsilon \right)$$

$$\begin{aligned}
&= \mathbb{P} \left(\left| \int p(y, z) g_n^k(y, z) dy dz \right. \right. \\
&\quad \left. \left. - \mathbb{E}_{p(y, z) p(x_{m+1}^n, z_{m+1}^n)} [g_n^k(Y, Z)] \right| > \epsilon \right) \\
&= \mathbb{P} \left(\left| \phi'(W_{m+1}^n) - \mathbb{E}[\phi'(W_{m+1}^n)] \right| > \epsilon \right) \\
&\leq 2 \exp \left(\frac{-2\epsilon^2 k^2}{(n-m)c^2} \right), \tag{36}
\end{aligned}$$

where the last step is due to (35).

3) *Third term in (30)*: Note that, for any z and any $j \in \mathcal{A}^m(z)$, and given the assumption of the theorem,

$$\left| \int p(y|z) g(x_j, y, z) dy \right| \leq M.$$

So if Assumption 1 holds, we know from Lemma 3 that, for any $\epsilon > 0$, there exists an integer n_0 such that for $n > n_0$

$$\begin{aligned}
&\mathbb{P} \left(\int p(z) \left| \int p(y|z) (g_n^k(y, z) - \bar{g}(y, z)) dy \right| dz > \epsilon \right) \\
&\leq \exp \left(\frac{-(n-m)\epsilon^2}{8M^2\gamma_d^2} \right). \tag{37}
\end{aligned}$$

Therefore, combining (30), (33), (36), and (37), we have that

$$\begin{aligned}
&\mathbb{P} \left(\left| \hat{g}(x^n, y^n, z^n) - \mathbb{E}_{p(x|z)p(y,z)} [g(X, Y, Z)] \right| \geq 3\epsilon \right) \\
&\leq \delta_1^n(\epsilon, c, M).
\end{aligned}$$

This concludes the proof of Theorem 1. \square

APPENDIX B PROOF OF PROPOSITION 1

The proof combines a concentration bound for $L_b^1(\omega_\theta)$, which follows from Hoeffding's inequality, and a concentration bound for $L_{b'}^2(\omega_\theta)$, which follows from Theorem 1.

By construction, $\mathcal{B}_{\text{joint}}^b$ consists of i.i.d. samples distributed according to $p(x, y, z)$. Moreover, the summands in (15) are bounded, i.e.,

$$\log \tau \leq \log \omega_\theta(x, y, z) \leq \log(1 - \tau),$$

given that the output of the classifier is clipped. Therefore, Hoeffding's inequality may be directly applied to find a concentration bound on $L_b^1(\omega_\theta)$, as seen in the lemma below.

Lemma 4. *For all θ and given $\epsilon > 0$, the following inequality holds:*

$$\begin{aligned}
&\mathbb{P} \left(\left| L_b^1(\omega_\theta) + \mathbb{E}_{p(x,y,z)} [\log \omega_\theta(X, Y, Z)] \right| \geq \epsilon \right) \\
&\leq 2 \exp \left(-\frac{2b\epsilon^2}{(\log \frac{1-\tau}{\tau})^2} \right).
\end{aligned}$$

Next, we show the convergence of $L_{b'}^2(\omega_\theta)$. According to Theorem 1, (15), and the definition of $\mathcal{B}_{\text{prod}}^b$, if we consider the function $g(x, y, z) = -\log(1 - \omega_\theta(x, y, z))$, we have that $\hat{g}(x^n, y^n, z^n) = L_{b'}^2(\omega_\theta)$. In this case, $g^{\max} = -\log(\tau)$,

$g^{\min} = -\log(1 - \tau)$, and $M = -\log(\tau)$, which implies that $c = \log \frac{1-\tau}{\tau}$. Then, the following Corollary is deduced.

Corollary 1. *Let Assumption 1 hold, then for any θ there exists n_0 such that for $n > n_0$:*

$$\mathbb{P} \left(\left| L_{b'}^2(\omega_\theta) + \mathbb{E}_{p(x|z)p(y,z)} [\log(1 - \omega_\theta(X, Y, Z))] \right| > 3\epsilon \right) \leq \delta_2^n(\epsilon),$$

where $\delta_2^n(\epsilon)$ is defined in Table I.

Now by the triangle inequality, we have that

$$\begin{aligned}
&|L_{\text{emp}}(\omega_\theta) - L(\omega_\theta)| \\
&\leq (1 - p_1) \left| L_b^2(\omega_\theta) + \mathbb{E}_{p(x|z)p(y,z)} [\log(1 - \omega_\theta(X, Y, Z))] \right| \\
&\quad + p_1 \left| L_b^1(\omega_\theta) + \mathbb{E}_{p(x,y,z)} [\log \omega_\theta(X, Y, Z)] \right|, \tag{38}
\end{aligned}$$

and the proof of Proposition 1 is complete by combining Lemma 4, Corollary 1, and choosing $\epsilon = \frac{\mu}{3-2p_1}$. \square

APPENDIX C PROOF OF THEOREM 2

The consistency of our proposed estimators is tied to the approximation power of the neural network, which is addressed in the following lemma.

Lemma 5. *For a given $\epsilon > 0$, $\exists \tilde{\theta} \in \Theta$ such that $\Theta \subset \mathbb{R}^h$ is compact and*

$$|L(\omega_{\tilde{\theta}}) - L^*| \leq \frac{\epsilon}{2}.$$

Proof. The proof can be shown similar to [17, Lemma 4]. \square

Remark 7. *The universal functional approximation introduced in [23] allows choosing parameters in a compact set of \mathbb{R}^h , while ϵ determines the number of neurons, and accordingly h , such that the desired approximation is achieved. Consider the network is approximating the function ω^* ; then, given $\epsilon > 0$, there exist a set Θ and a parameter $\tilde{\theta} \in \Theta$ such that $\omega_{\tilde{\theta}}$ is at an ϵ distance of ω^* .*

Adopting an optimizer such as *Adam*, we can minimize $L_{\text{emp}}(\omega_\theta)$ to find $\hat{\theta}$, and it is desired that $L(\omega_{\hat{\theta}})$ is close to L^* , which suggests we can use the neural network to approximate ω^* . As shown in Proposition 1, given a particular $\theta \in \Theta$, one can choose n such that $L_{\text{emp}}(\omega_\theta)$ falls in a the neighborhood of $L(\omega_\theta)$. Nevertheless, we need a more restrictive condition if we want to guarantee such convergence for all θ simultaneously. This is addressed in the lemma below.

Lemma 6. *Let Assumptions 1 and 3 hold, then for any $\mu > 0$, there exists n_1 such that, for $n > n_1$, we have that*

$$\mathbb{P} \left(\sup_{\theta \in \Theta} |L_{\text{emp}}(\omega_\theta) - L(\omega_\theta)| > 2\mu \right) \leq \delta_4^n(\mu), \tag{39}$$

where δ_4^n is defined in Table I.

Proof. The proof follows similar steps as the one for [17, Lemma 5] and we provide here only some details. Since $\Theta \subset \mathbb{R}^h$ and $\|\theta\|_2 \leq K, \forall \theta \in \Theta$, Θ can be covered with $N(\Theta, r)$

number of balls of radius r —the covering number with respect to ℓ_2 . The covering number is finite and bounded [34]:

$$N(\Theta, r) \leq \left(\frac{2K\sqrt{h}}{r} \right)^h. \quad (40)$$

Let $\{\theta_1, \dots, \theta_{N(\Theta, r)}\}$ denote the centers of the covering balls and $\{\Theta_1, \dots, \Theta_{N(\Theta, r)}\}$, the corresponding balls. We may then use the union bound on the LHS of (39) and take the supremum inside each ball Θ_j . By the triangle inequality, for any $\theta \in \Theta$, $\forall j$, and with probability at least $1 - \delta_4^n(\mu)$:

$$\begin{aligned} & |L_{\text{emp}}(\omega_\theta) - L(\omega_\theta)| \\ & \leq |L_{\text{emp}}(\omega_\theta) - L_{\text{emp}}(\omega_{\theta_j})| + |L_{\text{emp}}(\omega_{\theta_j}) - L(\omega_{\theta_j})| \\ & \quad + |L(\omega_{\theta_j}) - L(\omega_\theta)| \\ & \leq \frac{Br}{\tau} + \mu + \frac{Br}{\tau}, \end{aligned} \quad (41)$$

where the last step follows from Lipschitz continuity of $\log(\cdot)$ and ω_θ according to Assumption 3, and using Proposition 1. Finally, choosing $r = \frac{7\mu}{2B}$ concludes the proof. \square

The following proposition shows the convergence of $L(\omega_{\hat{\theta}})$.

Proposition 2. *Let Assumptions 1 and 3 hold. Given $\epsilon > 0$, there exists an integer n_1 such that, for $n > n_1$,*

$$\mathbb{P}(L(\omega_{\hat{\theta}}) - L^* \geq \epsilon) \leq \delta_4^n \left(\frac{\epsilon}{8} \right). \quad (42)$$

Proof. From Lemma 6, with probability at least $1 - \delta_4^n(\mu)$, we have that

$$|L_{\text{emp}}(\omega_{\hat{\theta}}) - L(\omega_{\hat{\theta}})| \leq 2\mu \quad \text{and} \quad |L_{\text{emp}}(\omega_{\hat{\theta}}) - L(\omega_{\hat{\theta}})| \leq 2\mu. \quad (43)$$

Since $\hat{\theta}$ minimizes $L_{\text{emp}}(\omega_\theta)$, by choosing $\mu = \frac{\epsilon}{8}$, we obtain:

$$\begin{aligned} L(\omega_{\hat{\theta}}) & \leq L_{\text{emp}}(\omega_{\hat{\theta}}) + \frac{\epsilon}{4} \leq L_{\text{emp}}(\omega_{\hat{\theta}}) + \frac{\epsilon}{4} \\ & \stackrel{(a)}{\leq} L(\omega_{\hat{\theta}}) + \frac{\epsilon}{2} \stackrel{(b)}{\leq} L^* + \epsilon, \end{aligned} \quad (44)$$

where the steps (a) and (b) are due to (43) and Lemma 5, respectively. \square

Proposition 2 implies that $L(\omega_{\hat{\theta}})$ is close to L^* and, due to the strong convexity of the cross-entropy loss, it can be shown that $\omega_{\hat{\theta}}$ is close to ω^* (in ℓ_1 norm). We continue the proof of the Theorem by defining the following terms:

$$\begin{aligned} I_{DV}^{\hat{\theta}} & := \mathbb{E}_{p(x,y,z)} [\log \hat{\Gamma}(X, Y, Z)] \\ & \quad - \log \mathbb{E}_{p(x|z)p(y,z)} [\hat{\Gamma}(X, Y, Z)], \\ I_{NWJ}^{\hat{\theta}} & := 1 + \mathbb{E}_{p(x,y,z)} [\log \hat{\Gamma}(X, Y, Z)] \\ & \quad - \mathbb{E}_{p(x|z)p(y,z)} [\hat{\Gamma}(X, Y, Z)], \\ I_{LDR}^{\hat{\theta}} & := \mathbb{E}_{p(x,y,z)} [\log \hat{\Gamma}(X, Y, Z)], \end{aligned} \quad (45)$$

where $\hat{\Gamma}(\cdot)$ is defined in (16). Then from the triangle inequality, we have that

$$\left| \hat{I}_{\text{est}}^{n, \hat{\theta}} - I(X; Y|Z) \right| \leq \left| \hat{I}_{\text{est}}^{n, \hat{\theta}} - I_{\text{est}}^{\hat{\theta}} \right| + \left| I_{\text{est}}^{\hat{\theta}} - I(X; Y|Z) \right|, \quad (46)$$

where ‘est’ can be replaced with ‘DV’, ‘NWJ’, or ‘LDR’. In the following, we show high-confidence convergence of the first and second terms on the RHS of the inequalities (46) due to Lemma 7 and Lemma 8, respectively.

Lemma 7. *Let Assumption 1 hold. For any $\epsilon^* > 0$, there exists n_2 such that for $n > n_2$ the following bounds hold:*

$$\mathbb{P} \left(\left| \hat{I}_{DV}^{n, \hat{\theta}} - I_{DV}^{\hat{\theta}} \right| \geq \frac{\epsilon^*}{2} \right) \leq \delta_5^n(\epsilon^*), \quad (47)$$

$$\mathbb{P} \left(\left| \hat{I}_{NWJ}^{n, \hat{\theta}} - I_{NWJ}^{\hat{\theta}} \right| \geq \frac{\epsilon^*}{2} \right) \leq \delta_6^n(\epsilon^*), \quad (48)$$

$$\mathbb{P} \left(\left| \hat{I}_{LDR}^{n, \hat{\theta}} - I_{LDR}^{\hat{\theta}} \right| \geq \frac{\epsilon^*}{2} \right) \leq \delta_7^n(\epsilon^*), \quad (49)$$

where δ_5^n , δ_6^n , and δ_7^n are defined in Table I.

Proof. See Appendix D. \square

Lemma 8. *Let Assumptions 1 and 3 hold and let $\tau < p_1$. Given $\epsilon^* > 0$, there exists an integer n_1 such that for $n > n_1$*

$$\mathbb{P} \left(\left| I_{\text{est}}^{\hat{\theta}} - I(X; Y|Z) \right| \geq \frac{\epsilon^*}{2} \right) \leq \delta_4^n \left(\frac{\epsilon^*}{8} \right), \quad (50)$$

where ‘est’ can be replaced with ‘DV’, ‘NWJ’, or ‘LDR’, and ϵ and δ_4^n are defined in Table I.

Proof. See Appendix E. \square

Therefore, if Assumption 4 also holds, we may combine Lemma 7 and Lemma 8 to yield a high-confidence bound for (46), which concludes the proof of the Theorem. \square

APPENDIX D PROOF OF LEMMA 7

Using the definitions (17) and (45), and the triangle inequality, we have that

$$\begin{aligned} \left| \hat{I}_{DV}^{n, \hat{\theta}} - I_{DV}^{\hat{\theta}} \right| & \leq \Delta_1 + \left| \log \frac{1}{b'} \sum_{\mathcal{B}_{\text{prod}}^{b'}} \hat{\Gamma}(x, y, z) \right. \\ & \quad \left. - \log \mathbb{E}_{p(x|z)p(y,z)} [\hat{\Gamma}(X, Y, Z)] \right| \\ & \leq \Delta_1 + \frac{p_1}{1-p_1} \frac{1-\tau}{\tau} \Delta_2, \end{aligned} \quad (51)$$

where

$$\begin{aligned} \Delta_1 & := \left| \frac{1}{b} \sum_{\mathcal{B}_{\text{joint}}^b} \log \hat{\Gamma}(x, y, z) - \mathbb{E}_{p(x,y,z)} [\log \hat{\Gamma}(X, Y, Z)] \right|, \\ \Delta_2 & := \left| \frac{1}{b'} \sum_{\mathcal{B}_{\text{prod}}^{b'}} \hat{\Gamma}(x, y, z) - \mathbb{E}_{p(x|z)p(y,z)} [\hat{\Gamma}(X, Y, Z)] \right|, \end{aligned}$$

and the last step in (51) follows since $\log(\cdot)$ is Lipschitz continuous given that $\hat{\Gamma}(\cdot)$ is bounded as below by definition,

$$\frac{1-p_1}{p_1} \frac{\tau}{1-\tau} \leq \hat{\Gamma}(x, y, z) \leq \frac{1-p_1}{p_1} \frac{1-\tau}{\tau}.$$

Similarly, for the NWJ estimator, we have that

$$\left| \hat{I}_{NWJ}^{n, \hat{\theta}} - I_{NWJ}^{\hat{\theta}} \right| \leq \Delta_1 + \Delta_2. \quad (52)$$

Finally, for the last estimator,

$$\left| \hat{I}_{LDR}^{n, \hat{\theta}} - I_{LDR}^{\hat{\theta}} \right| = \Delta_1. \quad (53)$$

We then use Hoeffding's inequality to bound the first terms on the RHS of (51), (52), and (53), which results in:

$$\mathbb{P}(\Delta_1 \geq \mu) \leq 2 \exp\left(-\frac{b\mu^2}{2(\log \frac{1-\tau}{\tau})^2}\right).$$

On the other hand, to show the concentration of the second terms on the RHS of (51) and (52), we leverage Theorem 1 with $c = \frac{1-p_1}{p_1} \frac{1-2\tau}{\tau(1-\tau)}$ and $M = \frac{1-p_1}{p_1} \frac{1-\tau}{\tau}$. Therefore, there exists an integer n_2 such that for all $n > n_2$,

$$\mathbb{P}(\Delta_2 \geq 3\mu) \leq \delta_1^n(\mu, c, M).$$

From (51), choosing $\mu = \frac{(1-p_1)\epsilon^*\tau}{2\tau+6p_1-8p_1\tau}$ yields (47), while for (52), we can choose $\mu = \frac{\epsilon^*}{8}$ to obtain (48). Finally for (53), we choose $\mu = \frac{\epsilon^*}{2}$ to obtain (49), and the proof of Lemma 7 is completed. \square

APPENDIX E PROOF OF LEMMA 8

To express the similarity between $\omega_{\hat{\theta}}$ and ω^* , we review a lemma from [17], which is based on the strong convexity of the cross-entropy loss and Assumption 2.

Lemma 9. ([17, Lemma 6]) *Let Assumption 2 hold. Given $\epsilon > 0$, if $L(\omega_{\theta}) \leq L^* + \epsilon$ for some $\theta \in \Theta$, then*

$$\int |\omega^*(x, y, z) - \omega_{\theta}(x, y, z)| dx dy dz \leq \eta,$$

where $\eta := (1-\tau)\sqrt{2\lambda(\mathcal{X})\epsilon/\alpha}$, and α is defined in Assumption 2 as the lower bound for the values of the joint and product density functions.

From Proposition 2, we know that, for any $\epsilon > 0$, with probability at least $1 - \delta_4^n(\epsilon/8)$

$$L(\omega_{\hat{\theta}}) \leq L^* + \epsilon.$$

Therefore, jointly with Assumption 2, the requirements of Lemma 9 are fulfilled. Let us further define $\Delta_{\omega}(x, y, z) := |\omega^*(x, y, z) - \omega_{\hat{\theta}}(x, y, z)|$, which leads to

$$\begin{aligned} \bar{\Delta}_{\omega} &:= \mathbb{E}_{p(x, y, z)} [\Delta_{\omega}(X, Y, Z)] \\ &= \int p(x, y, z) \Delta_{\omega}(x, y, z) dx dy dz \leq \eta\beta, \end{aligned} \quad (54)$$

and similarly

$$\bar{\Delta}'_{\omega} := \mathbb{E}_{p(x|z)p(y, z)} [\Delta_{\omega}(X, Y, Z)] \leq \eta\beta. \quad (55)$$

Next note that, from the continuity of $\Gamma^*(\cdot)$ and $\hat{\Gamma}(\cdot)$, defined in (12) and (16), respectively, we have that:

$$\begin{aligned} \bar{\Delta}_{\Gamma} &:= \mathbb{E}_{p(x, y, z)} \left| \log \Gamma^*(X, Y, Z) - \log \hat{\Gamma}(X, Y, Z) \right| \\ &\leq \frac{1}{\tau(1-\tau)} \bar{\Delta}_{\omega}, \\ \bar{\Delta}'_{\Gamma} &:= \mathbb{E}_{p(x|z)p(y, z)} \left| \Gamma^*(X, Y, Z) - \hat{\Gamma}(X, Y, Z) \right| \\ &\leq \frac{1-p_1}{p_1} \frac{1}{\tau^2} \bar{\Delta}'_{\omega}. \end{aligned} \quad (56)$$

So from the triangle inequality we have:

$$\begin{aligned} &\left| I_{DV}^{\hat{\theta}} - I(X; Y|Z) \right| \\ &\leq \left| \mathbb{E}_{p(x, y, z)} [\log \Gamma^*(X, Y, Z) - \log \hat{\Gamma}(X, Y, Z)] \right| \\ &\quad + \left| \log \mathbb{E}_{p(x|z)p(y, z)} [\Gamma^*(X, Y, Z)] \right. \\ &\quad \left. - \log \mathbb{E}_{p(x|z)p(y, z)} [\hat{\Gamma}(X, Y, Z)] \right| \\ &\stackrel{(a)}{\leq} \left| \mathbb{E}_{p(x, y, z)} [\log \Gamma^*(X, Y, Z) - \log \hat{\Gamma}(X, Y, Z)] \right| \\ &\quad + \frac{p_1(1-\tau)}{(1-p_1)\tau} \left| \mathbb{E}_{p(x|z)p(y, z)} [\Gamma^*(X, Y, Z) - \hat{\Gamma}(X, Y, Z)] \right| \\ &\leq \bar{\Delta}_{\Gamma} + \frac{p_1(1-\tau)}{(1-p_1)\tau} \bar{\Delta}'_{\Gamma} \stackrel{(b)}{\leq} \frac{1}{\tau(1-\tau)} \bar{\Delta}_{\omega} + \frac{1-\tau}{\tau^3} \bar{\Delta}'_{\omega} \\ &\leq \frac{\beta\eta(2\tau^2 - 2\tau + 1)}{\tau^3(1-\tau)}, \end{aligned} \quad (57)$$

where (a) and (b) are due to Lipschitz continuity of $\log(\cdot)$ and (56), respectively. Similarly, for the NWJ estimator:

$$\begin{aligned} &\left| I_{NWJ}^{\hat{\theta}} - I(X; Y|Z) \right| \\ &\leq \left| \mathbb{E}_{p(x, y, z)} [\log \Gamma^*(X, Y, Z) - \log \hat{\Gamma}(X, Y, Z)] \right| \\ &\quad + \left| \mathbb{E}_{p(x|z)p(y, z)} [\Gamma^*(X, Y, Z) - \hat{\Gamma}(X, Y, Z)] \right| \\ &\leq \bar{\Delta}_{\Gamma} + \bar{\Delta}'_{\Gamma} \leq \frac{1}{\tau(1-\tau)} \bar{\Delta}_{\omega} + \frac{1-p_1}{p_1\tau^2} \bar{\Delta}'_{\omega} \\ &\leq \frac{\beta\eta(1 + 2p_1\tau - p_1 - \tau)}{p_1\tau^2(1-\tau)}. \end{aligned} \quad (58)$$

Finally, for the LDR estimator, we have:

$$\begin{aligned} &\left| I_{LDR}^{\hat{\theta}} - I(X; Y|Z) \right| \\ &= \left| \mathbb{E}_{p(x, y, z)} [\log \Gamma^*(X, Y, Z) - \log \hat{\Gamma}(X, Y, Z)] \right| \\ &\leq \bar{\Delta}_{\Gamma} \leq \frac{1}{\tau(1-\tau)} \bar{\Delta}_{\omega} \leq \frac{\beta\eta}{\tau(1-\tau)}. \end{aligned} \quad (59)$$

Note that for $\tau < p_1$ we have the following:

$$\frac{1}{\tau(1-\tau)} \leq \frac{1 + 2p_1\tau - p_1 - \tau}{p_1\tau^2(1-\tau)} \leq \frac{2\tau^2 - 2\tau + 1}{\tau^3(1-\tau)}.$$

So by choosing $\eta = \frac{\tau^3(1-\tau)}{(2\tau^2-2\tau+1)\beta} \frac{\epsilon^*}{2}$, ϵ can be determined from η as defined in Lemma 9. This, together with the bounds (57), (58), and (59), concludes the proof of Lemma 8. \square

REFERENCES

- [1] A. El Gamal and Y.-H. Kim, *Network information theory*. Cambridge University Press, 2011.
- [2] J. Massey, "Causality, Feedback and Directed Information," in *Proc. Int. Symp. Inf. Theory Applic. (ISITA), Honolulu, HI, USA*, Nov. 1990, pp. 303-305.
- [3] S. Molavipour, G. Bassi, and M. Skoglund, "Testing for directed information graphs," in *55th Annual Allerton Conf. on Comm., Control, Comput. (Allerton)*, Oct. 2017, pp. 212-219.
- [4] F. Fleuret, "Fast binary feature selection with conditional mutual information," *Journal of Machine Learning Research*, vol. 5, pp. 1531-1555, 2004.
- [5] D. Loeckx, P. Slagmolen, F. Maes, D. Vandermeulen, and P. Suetens, "Nonrigid image registration using conditional mutual information," *IEEE Trans. Med. Imag.*, vol. 29, no. 1, pp. 19-29, 2009.

- [6] S. Mukherjee, "Machine learning using the variational predictive information bottleneck with a validation set," *arXiv preprint arXiv:1911.02210*, 2019.
- [7] B. Rodríguez-Gálvez, R. Thobaben, and M. Skoglund, "A variational approach to privacy and fairness," *arXiv preprint arXiv:2006.06332*, 2020.
- [8] Q. Wang, S. R. Kulkarni, and S. Verdú, "Universal estimation of information measures for analog sources," *Foundations and Trends® in Communications and Information Theory*, vol. 5, no. 3, pp. 265–353, 2009.
- [9] A. Kraskov, H. Stögbauer, and P. Grassberger, "Estimating mutual information," *Physical Review E*, vol. 69, no. 6, p. 066138, Jun. 2004.
- [10] W. Gao, S. Oh, and P. Viswanath, "Demystifying fixed k -nearest neighbor information estimators," *IEEE Trans. Inf. Theory*, vol. 64, no. 8, pp. 5629–5661, Aug. 2018.
- [11] J. Runge, "Conditional independence testing based on a nearest-neighbor estimator of conditional mutual information," in *21st Int. Conf. Art. Intell. Stats. (AISTATS)*, Apr. 2018, pp. 938–947.
- [12] M. Vejmelka and M. Paluš, "Inferring the directionality of coupling with conditional mutual information," *Physical Review E*, vol. 77, no. 2, p. 026214, Feb. 2008.
- [13] S. Frenzel and B. Pompe, "Partial mutual information for coupling analysis of multivariate time series," *Phys. Rev. Lett.*, vol. 99, no. 20, p. 204101, Nov. 2007.
- [14] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm, "MINE: Mutual information neural estimation," in *35th Int. Conf. Mach. Learn. (ICML)*, Jul. 2018, pp. 531–540.
- [15] D. McAllester and K. Statos, "Formal limitations on the measurement of mutual information," *arXiv:1811.04251*, Nov. 2018.
- [16] B. Poole, S. Ozair, A. Van Den Oord, A. Alemi, and G. Tucker, "On variational bounds of mutual information," ser. Proc. of Machine Learning Research, vol. 97. PMLR, Jun 2019, pp. 5171–5180.
- [17] S. Mukherjee, H. Asnani, and S. Kannan, "CCMI: Classifier based conditional mutual information estimation," in *Uncertainty in Artificial Intelligence*, Jul. 2019.
- [18] J. Song and S. Ermon, "Understanding the limitations of variational mutual information estimators," *arXiv preprint arXiv:1910.06222*, 2019.
- [19] Z. Qin, D. Kim, and T. Gedeon, "Rethinking softmax with cross-entropy: Neural network classifier as mutual information estimator," *arXiv preprint arXiv:1911.10688*, 2019.
- [20] S. Molavipour, G. Bassi, and M. Skoglund, "Conditional mutual information neural estimator," in *2020 IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, May 2020, pp. 5025–5029.
- [21] M. D. Donsker and S. S. Varadhan, "Asymptotic evaluation of certain Markov process expectations for large time. I," *Communications on Pure and Applied Mathematics*, vol. 28, no. 1, pp. 1–47, 1975.
- [22] X. Nguyen, M. J. Wainwright, and M. I. Jordan, "Estimating divergence functionals and the likelihood ratio by convex risk minimization," *IEEE Trans. Inf. Theory*, vol. 56, no. 11, pp. 5847–5861, Nov. 2010.
- [23] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural Networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [24] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [25] G. Pichler, P. Piantanida, and G. Koliander, "On the estimation of information measures of continuous distributions," *arXiv preprint arXiv:2002.02851*, 2020.
- [26] J. Zhang, O. Simeone, Z. Cvetkovic, E. Abela, and M. Richardson, "ITENE: Intrinsic transfer entropy neural estimator," *arXiv preprint arXiv:1912.07277*, 2019.
- [27] S. Molavipour, G. Bassi, M. Čičić, M. Skoglund, and K. H. Johansson, "Causality graph of vehicular traffic flow," *arXiv preprint arXiv:2011.11323*, 2020.
- [28] T. O'Shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Trans. on Cogn. Commun. Netw.*, vol. 3, no. 4, pp. 563–575, Dec. 2017.
- [29] S. Dörner, S. Cammerer, J. Hoydis, and S. ten Brink, "Deep learning based communication over the air," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 1, pp. 132–143, Feb. 2018.
- [30] T. J. O'Shea, T. Roy, N. West, and B. C. Hilburn, "Physical layer communications system design over-the-air using adversarial networks," in *2018 26th European Signal Process. Conf. (EUSIPCO)*, Sep. 2018, pp. 529–532.
- [31] R. Fritschek, R. F. Schaefer, and G. Wunder, "Deep learning for channel coding via neural mutual information estimation," in *2019 IEEE 20th Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, Jul. 2019.
- [32] C. McDiarmid, "On the method of bounded differences," *Surveys in Combinatorics*, vol. 141, no. 1, pp. 148–188, 1989.
- [33] L. Devroye, L. Györfi, A. Krzyżak, and G. Lugosi, "On the strong universal consistency of nearest neighbor regression function estimates," *Ann. Statist.*, vol. 22, no. 3, pp. 1371–1385, 1994.
- [34] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.