

Towards Robust Scale-Invariant Mutual Information Estimators

Anonymous authors

Paper under double-blind review

Abstract

Mutual information (MI) is hard to estimate for high dimensional data, and various estimators have been proposed over the years to tackle this problem. Here, we note that there exists another challenging problem, namely that many estimators of MI, which we denote as $I(X; T)$, are sensitive to scale, i.e., $I(X; \alpha T) \neq I(X; T)$ where $\alpha \in \mathbb{R}$. Although some normalization methods have been hinted at in previous works, there is no in-depth study of the problem. In this work, we study new normalization strategies for MI estimators to be scale-invariant, particularly for the Kraskov–Stögbauer–Grassberger (KSG) and the neural network-based MI (MINE) estimators. We provide theoretical and empirical results and show that the original un-normalized estimators are not scale-invariant and highlight the consequences of an estimator’s scale-dependence. We propose new global normalization strategies that are tuned to the corresponding estimator and scale invariant. We compare our global normalization strategies to existing local normalization strategies and provide intuitive and empirical arguments to support the use of global normalization. Extensive experiments across multiple distributions and settings are conducted, and we find that our proposed variants KSG-Global- L_∞ and MINE-Global-Corrected are most accurate within their respective approaches. Finally, we perform an information plane analysis of neural networks and observe clearer trends of fitting and compression using the normalized estimators compared to the original un-normalized estimators. Our work highlights the importance of scale awareness and global normalization in the MI estimation problem.

1 Introduction

Mutual information (MI), is a fundamental measure of dependency between two variables, which has become pivotal in various machine learning domains, including generalization (Xu & Raginsky, 2017; Bu et al., 2019; Russo & Zou, 2020), representation learning (Bachman et al., 2019; Tschannen et al., 2020) and fairness (Wang et al., 2023; Roh et al., 2020). Estimating MI for high-dimensional continuous variables (Xu et al., 2020) is particularly challenging, due to the hardness of accurately estimating the probability distribution in high dimensions (Goldfeld & Greenewald, 2021). For example, traditional estimators like Kraskov–Stögbauer–Grassberger (KSG) (Kraskov et al., 2004), rely on distance metrics, and for high dimensional data, the distances would have less variation due to the curse of dimensionality.

In this paper, we highlight a critical but underexplored factor that leads to inaccuracies in MI estimation: the scale of the variables (i.e., $|X|$). Specifically, when considering the mutual information $I(X; \alpha T)$, where $\alpha \in \mathbb{R}^+$ is a scaling factor, we demonstrate that when $\alpha \ll 1$ or $\alpha \gg 1$, the MI estimates can deviate significantly from the true value. This is problematic since by definition, $I(X; \alpha T) = I(X; T)$ for any two continuous random variables (RVs) X and T . Moreover, a stronger result states that $I(X; f(Y)) = I(X; T)$ for any continuous and invertible transformation f (Cover & Thomas, 2006). In this paper, we mainly focus on the specific impact of scale.

Most mutual information estimators, including the widely adopted KSG estimator (Kraskov et al., 2004) and its subsequent variants (Gao et al., 2017), lack scale invariance—a limitation that we rigorously demonstrate in this study. We provide a theoretical analysis explaining why this deficiency arises. We also show the

binning estimator (Paninski, 2003) for MI can be scale invariant when the number of bins used is fixed. However, the binning estimator itself is not well-suited for estimating high-dimensional continuous variables. Recently the mutual information neural estimator (MINE) (Belghazi et al., 2018) was proposed, which is a neural network-based estimator of MI that makes use of its Donsker-Varadhan (DV) representation. We demonstrate theoretically that ideally, MINE should be scale-invariant, but MINE fails in practice due to limitations introduced by stochastic gradient descent optimization.

Despite numerous surveys that have explored various methods of MI estimation (Walters-Williams & Li, 2009; McAllester & Stratos, 2020; Paninski, 2003), the critical importance of normalization (preprocessing) has been largely overlooked. A natural solution to ensure scale invariance is to pre-process the data using standard normalization, where each dimension is adjusted to have a variance of 1, and we refer as *local normalization*. This pre-processing step was hinted in (Kraskov et al., 2004) for the KSG estimator. Local normalization also has been commonly applied as a preprocessing step in many deep learning studies involving mutual information perspective (Hjelm et al., 2019; Xie et al., 2024). However, local normalization treats each dimension independently and normalizes them to have a variance of 1, which, as we demonstrate in Section 5.1.1, does not work well in the high-dimension setting especially in neural networks, across two separate experiments. This is because most high-dimensional feature representations in neural networks always contain some noisy dimensions, which are of low energy and contain irrelevant features. Thus, amplifying these low energy dimensions can lead to suboptimal MI estimates. We also note that the recent work by (Czyż et al., 2023), in addition to trying out local normalization approaches, also studied other preprocessing methods including the transformation of the margin distribution to uniform distribution (via converting to rank). We note that this conversion step also brings all individual dimensions to equal importance like local normalization, and thus would have the same pitfalls in this scenario. Note that we have included additional discussions on related work in Appendix B.

To address this issue, in our work, we propose a set of *global normalization* approaches. Unlike local normalization, global normalization preserves the relative energies between the different dimensions, and thus avoids scaling up low-energy noisy dimensions. Our proposed estimator modifications do not only include new normalization approaches, however, and often also have an additional maximization step, which helps bias our estimators better. It is well known that KSG and other MI estimators have a tendency to have negative bias Czyż et al. (2023), especially in high dimensions. Our normalization approaches for KSG incorporate this observation via an additional maximization step, which also follows intuitively from one of our theoretical observations in Proposition 3.

We now summarize our contributions:

- We propose novel scale-invariant extensions of KSG and MINE-based estimators that effectively address the one-sided scale-invariance issue and substantially improve estimator accuracy. To the best of our knowledge, our work is the first comprehensive analysis of the effect of scale and various normalization methodologies, some of which are introduced for the first time in this work.
- We demonstrate that the KSG-Global- L_∞ and MINE-Global-Corrected variants consistently produce the most accurate estimations within their respective approaches, across a broad range of experiments involving synthetic data, which are targeted towards the high-dimensional and low-data regime. These experiments include multiple types of transformations, noise injections, and changes in dimensionality.
- We explore the dynamics of MI between inputs X and hidden layers T during neural network training. Our results highlight that unnormalized estimators significantly confound the scale of T in their estimates, while our normalized approaches can often capture distinct phases of training, such as fitting and compression.

The rest of the paper is organized as follows. In Section 2, we first provide motivation for improving MI Estimators by including robustness to scale. In Section 3, we theoretically study scale-invariance behaviour of three estimators: binning, KSG and MINE. In Section 4, we introduce our proposed scale-invariant variants of KSG and MINE using new normalization strategies. In Section 5.1, we provide additional empirical motivation for each aspect of the proposed variants of KSG and MINE. In section 5.2, we conduct extensive experiments on synthetic datasets and rigorously test our proposed estimators against the non-normalized and standard locally-normalized estimators. In 5.3 we conduct experiments on three real datasets, which mainly includes monitoring the mutual information measures during neural network training. Comparisons

are made against both the original estimators and standard normalization approaches (local normalization). Finally, we summarize our findings and discuss their implications in Section 6. Note that a detailed background on MI Estimators studied in this work is presented in Appendix A.

2 Motivation

Estimating MI is fundamental to various domains, ranging from learning theory to practical applications such as medical analysis and wireless communication (Shwartz-Ziv & Tishby, 2017; Saxe et al., 2018). To motivate our proposed normalization strategy, this section outlines several desirable properties that effective MI estimators should possess. Let $S = \{(X_1, T_1), (X_2, T_2), \dots, (X_n, T_n)\}$ be the sampled data. With this, let $\hat{I}_{est}^n(X; T)$ represent an estimate of the MI between X and T using the estimator est , given N sampled points from the joint distribution $P(X, T)$. Ideally, we seek the estimator to have the following properties:

1. **Global Scale Invariance:** For any arbitrary $\alpha \in \mathbb{R}$ and $n \in \mathbb{Z}^+$, $\hat{I}_{est}^n(\alpha X; \alpha T) = \hat{I}_{est}^n(X; T)$
2. **One-Sided Scale Invariance** For any arbitrary $\alpha \in \mathbb{R}$ and $n \in \mathbb{Z}^+$, $\hat{I}_{est}^n(X; \alpha T) = \hat{I}_{est}^n(X; T)$

We emphasize the importance of these properties because true mutual information inherently satisfies them. By definition, $I(\alpha X; \alpha T) = I(X; T)$ and $I(\alpha X; T) = I(X; T)$ for a scalar α . In the case of neural networks, where X represents the input, Y represents the target, and T represents the features, estimation of $I(X; T)$ becomes important, as it was hypothesized that it can predict the generalization behavior of deep learning networks (Shwartz-Ziv & Tishby, 2017). Furthermore, (Shwartz-Ziv & Tishby, 2017) also predicts a two-phase behavior of $I(X; T)$ during training: (a) fitting, where $I(X; T)$ and $I(T; Y)$ increases, and (b) compression where $I(X; T)$ decreases. However, this is often not observed (Saxe et al., 2018). We hypothesize that it could be because of the scale-sensitivity of the estimators, as the scale of T changes significantly during training.

We note that the current estimators may not obey one-sided scale invariance. First, we study three estimators theoretically: KSG, MINE, and binning.

3 Testing One-sided Scale-Invariance of MI Estimators

In this section, we theoretically test whether the common MI estimators are global-scale invariant and one-sided scale invariant. In Section 5, we also present an experimental test of one-sided scale-invariance on MI estimators. Note that for all results that follow, we assume every random variable is bounded. That is, if X is bounded, we have that $|X| \leq B$ for some finite $B < \infty$. Also, for the following results, let $X \in \mathbb{R}^d$ and $T \in \mathbb{R}^m$. Note that although, some of the following results show that the estimators do not adhere to scale-invariant behaviour in the asymptotic regime, we eventually find that most estimators show significant disruption in response to scale changes of less than 10 times in magnitude. Also note that by \mathbb{R}^+ , we mean the set of all positive real numbers, excluding zero.

Binning: Let us denote the binning estimator described in (Paninski, 2003) by \hat{I}_{bin}^n . Then we have the following result.

Proposition 1. It holds that $\hat{I}_{bin}^n(\alpha X; \alpha T) = \hat{I}_{bin}^n(X; T)$ and $\hat{I}_{bin}^n(X; \alpha T) = \hat{I}_{bin}^n(X; T) \forall \alpha \in \mathbb{R}^+$.

Remark 1. We note that even though the binning estimator is scale-invariant, it is not a good estimator for MI, more so in the high dimension setting (Kraskov et al., 2004). This is because in high dimensions the data occupies the space very sparsely, and most bins will yield zero datapoints and thus a zero probability. Due to this, it is common practice to use fewer bins overall, which instead leads to less accurate estimates of MI as more information is lost.

KSG: Let us denote the KSG estimator proposed in (Kraskov et al., 2004) by \hat{I}_{KSG}^n . Then, we have the following results.

Proposition 2. It holds that $\hat{I}_{KSG}^n(\alpha X; \alpha T) = \hat{I}_{KSG}^n(X; T), \forall \alpha \in \mathbb{R}^+$.

This proof also leads to the following result KSG estimator may not follow one-sided invariance.

Proposition 3. It holds that $\lim_{\alpha, n \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = -\frac{1}{k}$ and $\lim_{\alpha \rightarrow 0^+, n \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = -\frac{1}{k}$, where k is the k -nearest neighbor parameter for the estimator. Thus, $\hat{I}_{KSG}^n(X; \alpha T)$ need not be equal to $\hat{I}_{KSG}^n(X; T)$.

MINE: We first define two variants of the MINE estimator as follows:

MINE-Opt: This estimator refers to the MINE estimator where instead of training the neural network on the loss function defined in equation 7 by stochastic gradient descent (SGD), we pick the best neural network configuration that directly maximizes equation 7. Thus, we pick the global optimum.

MINE-SGD: This estimator refers to the MINE estimator where optimization of the loss function defined in equation 7, is performed using conventional stochastic gradient descent. This is the standard approach proposed originally by (Belghazi et al., 2018).

We denote the MINE-based MI estimators by $\hat{I}_{MINE-opt}^n$ and $\hat{I}_{MINE-sgd}^n$. We then have the following results.

Proposition 4. It holds that $\hat{I}_{MINE-opt}^n(X; \alpha T) = \hat{I}_{MINE-opt}^n(X; T) \forall \alpha \in \mathbb{R}^+$, where $\mathbb{R}^+ = \{x \in \mathbb{R} | x > 0\}$.

Next, we outline a theoretical result regarding the limiting behaviour of the first layer weights for the MINE estimator’s neural network, when the scale of one of the variables approaches zero.

Proposition 5. Consider the MINE optimization problem with input data $S = \{(\alpha X_1, Y_1), \dots, (\alpha X_n, Y_n)\}$ where $X \in \mathbb{R}^{d_x}$, $Y \in \mathbb{R}^{d_y}$, $(X, Y) \sim P(X, Y)$ are bounded RVs and $\alpha \in \mathbb{R}^+$ is a scaling factor. We consider a neural network of depth $d_n + 1$ having h_1, h_2, \dots, h_{d_n} ReLU-activated hidden neurons in the respective layers. The network is trained via stochastic gradient descent on the MINE loss function in equation 8 for a finite number of epochs n_e . Let the *trained* weights between the j^{th} node of the $l + 1^{th}$ hidden layer and the i^{th} node of the l^{th} hidden layer be denoted by $w_{ji}^l \in \mathbb{R}^d$. We consider the case where the initialized weights are very close to zero but not exactly zero (to allow unsymmetrical learning). We assume that the network weights are bounded, such that every weight $|w_{ji}^l| \leq B$ for some $B \in \mathbb{R}$. Let $\eta(t)$ denote the learning rate used at epoch t . Then we have, $\forall i, j$,

$$\lim_{\alpha \rightarrow 0^+} |w_{ji}^1| = 0 \quad (1)$$

With this, we have the following result that explores scale invariance in neural network based MINE estimators.

Proposition 6. We consider the same setting as Proposition 5 for the MINE estimation problem. There, it holds that $\lim_{\alpha \rightarrow 0} \hat{I}_{MINE-sgd}^n(X; \alpha T) = 0$. Thus, $\hat{I}_{MINE-sgd}^n(X; \alpha T)$ need not be equal to $\hat{I}_{MINE-sgd}^n(X; T)$.

4 Methodology

4.1 Normalization Strategies

We consider a setting where we are given an RV $X = [x_1, x_2, \dots, x_d] \in \mathbb{R}^d$, where $x_i \in \mathbb{R}$ represents the i -th component of X . Suppose $X \sim P$, where P is a probability distribution, and let $S = \{X_1, X_2, \dots, X_n\}$ be an independent and identically distributed (i.i.d.) sample drawn from P , with $X_j \in \mathbb{R}^d$ for $j = 1, \dots, n$. With this, we outline three normalization strategies that form the basis of our studies in this work. We define them as follows.

Definition 1. (Local Normalization) The *locally normalized variable* $X_{\sigma|S} = [x'_1, \dots, x'_i, \dots, x'_d] \in \mathbb{R}^d$ is defined by normalizing each dimension i individually as $x'_i = \frac{x_i}{\sqrt{\mathbb{E}_S[(x_i - \mathbb{E}_S[x_i])^2]}}$, for $i = 1, \dots, d$, where

$\mathbb{E}_S[\cdot]$ denotes the empirical expectation over S .

Definition 2. (Global Normalization)

The *globally normalized variable* $X_{\Sigma|S} \in \mathbb{R}^d$ is defined as $X_{\Sigma|S} = \frac{X}{\sqrt{\mathbb{E}_S[\|X - \mathbb{E}_S[X]\|_2^2]}}$, where $\|\cdot\|_2$ denotes the L_2 -norm, and $\mathbb{E}_S[\cdot]$ is the empirical expectation over S .

Definition 3. (Global L_∞ Normalization) The globally L_∞ -normalized variable $X_{\Sigma_\infty|S} \in \mathbb{R}^d$ is $X_{\Sigma_\infty|S} = \frac{X}{\mathbb{E}_S[\|X - \mathbb{E}_S[X]\|_\infty]}$, where $\|\cdot\|_\infty$ denotes the L_∞ -norm and $\mathbb{E}_S[\cdot]$ is the empirical expectation over S .

Note that for any RV X , we denote by $X_{\sigma|S}$ and $X_{\Sigma|S}$ its locally and globally normalized versions respectively.

4.2 Studied Scale-Invariant Estimators

We are given the RVs $X \in \mathbb{R}^d$ and $T \in \mathbb{R}^m$, and sampled data $S = \{(X_1, T_1), (X_2, T_2), \dots, (X_n, T_n)\} \sim P_{XT}^n$. All following estimates are for the MI between X and T , given S . With this, we propose the following normalization approaches for KSG and MINE estimators. We outline our approaches for scale-invariant KSG and MINE extensions in Table 1.

Table 1: Proposed scale-Invariant KSG and MINE variants

KSG	MINE
KSG-Local: $\hat{I}_{KSG}^n(X_{\sigma S}; T_{\sigma S})$	MINE-Local: $\hat{I}_{MINE}^n(X_{\sigma S}; T_{\sigma S})$
KSG-Global: $\max_{c \in \{c_1, c_2, \dots, c_n\}} [\hat{I}_{KSG}^n(X_{\Sigma S}; cT_{\Sigma S})]$	MINE-Global: $\hat{I}_{MINE}^n(X_{\Sigma S}; T_{\Sigma S})$
KSG-Global-L_∞: $\max_{c \in \{c_1, c_2, \dots, c_n\}} [\hat{I}_{KSG}^n(X_{\Sigma_\infty S}; cT_{\Sigma_\infty S})]$	MINE-Global-Corrected: $\hat{I}_{MINE}^n(\sqrt{d_X}X_{\Sigma S}; \sqrt{d_T}T_{\Sigma S})$

Remark 2. In addition to the above approaches, we compare the standard baselines of KSG and MINE. Furthermore, we also include a recent variant of KSG in our comparisons, called BI-KSG (Gao et al., 2017), which has smaller bias levels for highly correlated data. We do not include binning-based measures in our experimental results, as we find that they fare poorly for almost all of our studied cases. Thus, we only study the KSG and MINE variants empirically in this work. Also, note that the range of multiplier scales c_1, c_2, \dots, c_n are tunable hyperparameters, and we fix $c_1 = 0.1, c_2 = 0.2, \dots, c_n = 2$ for all our experiments. Note that for the KSG-Global variants, this does slow down the MI estimation process, as it takes n times the computation time to estimate the measure when compared to KSG.

5 Experimental Studies

Our empirical studies can be categorized into roughly four broad sections:

- Empirical motivation for proposed normalization variants:** We provide in-depth empirical analyses for each normalization variant proposed in this work, and also the overall reasons for potentially choosing global normalization approaches over local.
- Scale dependence and Signal to Noise Ratio (SNR) analysis of estimators:** We perform some basic tests and analyses of all estimators. First, we study their overall responses to scale changes, and then we study their responses to changes in noise levels.
- Accuracy analysis of estimators:** We conduct an extensive accuracy-bias-correlation analysis of all estimators in two different settings where ground truth MI is known. In each setting, we generate synthetic data using a diverse set of transformations to simulate different distribution scenarios.
- Studying neural network training using estimators:** We study the MI dynamics of neural networks during training. Specifically, we analyze the MI between input and features and compare the trends resulting from various estimators.

For our experiments, we use two *base* distributions for generating the random variables X and T . We refer to them in various parts of the experiments. They are as follows:

- Correlated Gaussians:** Here, $X \in \mathbb{R}^d \sim \mathcal{N}(0, I_d)$ and $T \in \mathbb{R}^d \sim \mathcal{N}(0, I_d)$, and $\mathbb{E}[X_i T_i] = \rho$ for $1 \leq i \leq d$ and $\mathbb{E}[X_i T_j] = 0$ when $i \neq j$. I_d denotes the identity matrix of size $d \times d$. This is a standard setting used in many prior MI estimation works.
- Additive Gaussian Noise:** Here $X \in \mathbb{R}^d \sim \mathcal{N}(0, I_d)$ and $T = X + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2 I_d)$.

The details for our estimators are provided in **Appendix D**. We used the NPEET MI estimator toolbox for estimating KSG and KSG-based measures ¹. For MINE, we used the popular pytorch-based package ².

5.1 Additional Motivation for Normalization Variants

In this section, we provide both intuitive and empirical arguments for our proposed variants in the previous section. First, we provide intuitive and empirical reasons for when and why global normalization approaches could be preferred. Next, we provide a rationale for our proposed global normalization variants for KSG and MINE estimators.

5.1.1 Global Over Local

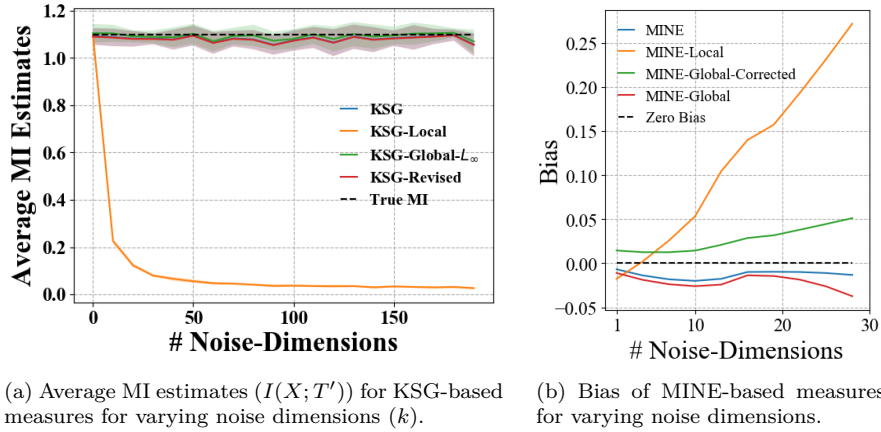


Figure 1: Comparative evaluation of MI estimators (KSG and MINE) under increasing noise dimensions. Please see Appendix C.1 for details.

Global normalization is advantageous for estimating mutual information (MI) in high-dimensional settings, such as neural network feature representations, where many features in T are sparse. In contrast, local normalization scales all variables equally, potentially overemphasizing irrelevant, low-energy dimensions, which degrades MI estimates.

KSG: Consider random variables $X, T \in \mathbb{R}^2$, where $T = X + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2 I_2)$. We augment X with k independent noise components $\epsilon = [\epsilon_1, \dots, \epsilon_k]$, $\epsilon_i \sim \mathcal{N}(0, \sigma'^2)$ with $\sigma' \ll \sigma$, forming $X' = [X, \epsilon]$. Although $I(X; T) = I(X'; T)$, as shown in Figure 1a, KSG and global normalization maintain stable MI estimates with increasing k , while local normalization underestimates MI by over-scaling the noise.

MINE: For MINE, we generate correlated Gaussian RVs $X, T \in \mathbb{R}^2$ (with a random correlation $\rho \in (0, 0.8)$) and similarly extend X with noise to form X' . Averaging over trials, Figure 1b reveals that local normalization causes the bias to increase significantly with k , whereas global normalization yields stable estimates. This occurs because, for MINE, the network tends to overfit the added noise when each variable is normalized equally, while global variants preserve the effective input dimensionality.

5.1.2 KSG-Global: Why the Maximization Step?

We summarize two main arguments for our maximization step in Table 1:

1. **Negative Bias in High Dimensions:** The KSG estimator shows increasing negative bias as data dimensionality grows. In one experiment, correlated Gaussians $X, T \in \mathbb{R}^d$ with fixed ground truth MI (≈ 0.8) were sampled (20 trials of 1000 points each) while d ranged from 1 to 9, yielding a clear negative bias in the estimators as seen in Figure 2(a). A second experiment with $d \in \{2, 4, 8, 16, 32, 64\}$ and randomly chosen

¹ <https://github.com/gregversteeg/NPEET> ² <https://github.com/gtegnr/mine-pytorch>

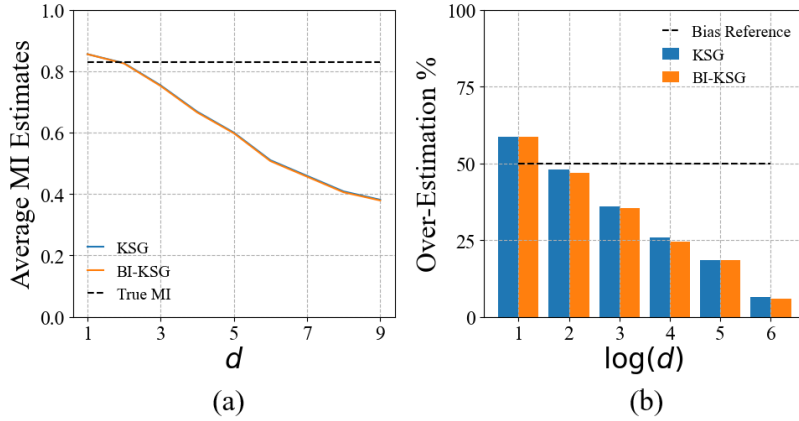


Figure 2: Dependency of KSG estimator bias on data dimension (using base-2 log). See Appendix C.1 for details.

correlation confirmed that the fraction of MI estimates below the ground truth rises with d (Figure 2(b)). Taking the maximum estimate over scales helps mitigate this bias.

2. Consequence of Proposition 3: Proposition 3 shows that for $I(X; \alpha T)$, the KSG estimator converges to negative values at extreme scales α . This observation motivates taking the maximum over a range of scales c for the globally normalized variables, i.e., $\hat{I}_{KSG}^n(X_{\Sigma|S}, cT_{\Sigma|S})$. Empirically, the estimator values exhibit an approximately Gaussian trend with respect to c (see Figure 4a), making the maximum both well-defined and meaningful.

Remark 3. Note that MINE implicitly maximizes over relative scales via the arbitrary first-layer weights. Since scaling these weights adjusts the effective input similarly (i.e., $(\alpha W)^T X = W^T(\alpha X)$), the network optimizes over affine transformations. Nonetheless, due to the tendency of gradient descent towards flatter minima Keskar et al. (2017), the full benefit of this invariance may not be achieved.

5.1.3 Motivation for Global Normalization Variants

We motivate our proposed global normalization variants: KSG-Global- L_∞ and MINE-Global-Corrected.

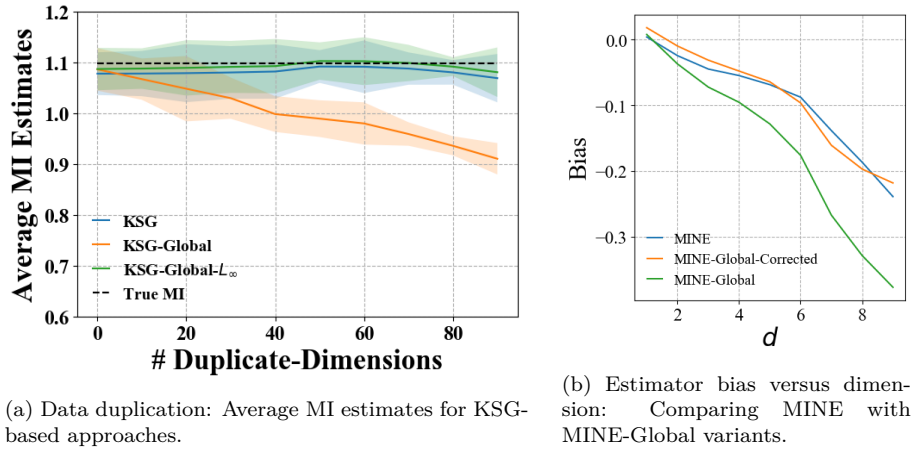


Figure 3: Analysis of normalization variants: (a) Impact of data duplication for KSG-based approaches, and (b) Estimator bias for MINE and MINE-Global variants across dimensions. Please see Appendix C.1 for details.

KSG: Global normalization aims to equate the scale of the nearest neighbor distances of X and T , to ensure both variables are equally regarded in the MI estimation. This also avoids biases that could lead to low or negative MI estimates (see Proposition 3 and Figure 4a). In contrast, local normalization scales

all dimensions equally, potentially amplifying noise. However, since KSG computes distances using the L_∞ -norm, when $d_X \gg d_T$, global normalization makes the individual dimensions of X significantly smaller than those of T , yielding smaller L_∞ distances. This motivates the proposed KSG-Global- L_∞ approach, which scales based on L_∞ norm distances instead of L_2 norm. To illustrate its impact, we duplicate X to form $X' = [X, X, \dots, X]$ (thus $I(X'; T) = I(X; T)$). As shown in Figure 3a, KSG-Global actually decreases its estimate with increasing duplicates, whereas both KSG and KSG-Global- L_∞ maintain consistent MI estimates. This confirms the necessity of employing the L_∞ -norm for scaling the variables using the global normalization approach.

MINE: Global normalization can lead to low per-dimension energy when $d_X \gg d_T$, since $\mathbb{E}[X_{\Sigma|S}(i)^2] = 1/d_X$ versus $\mathbb{E}[T_{\Sigma|S}(i)^2] = 1/d_T$. This imbalance may cause gradient descent to focus predominantly on T . To counter this, we rescale the normalized variables so that $\mathbb{E}[X'_{\Sigma|S}(i)^2] = \mathbb{E}[T'_{\Sigma|S}(j)^2] = 1$, i.e., $X'_{\Sigma|S} = \sqrt{d_X} X_{\Sigma|S}$ (and similarly for T). Experiments with d ranging from 1 to 9 (Figure 3b) reveal that while MINE-Global’s bias grows more negative with higher d , MINE-Global-Corrected has bias levels comparable to standard MINE. This demonstrates that rescaling effectively mitigates the effects of imbalances in terms of per-dimension energy of X and T .

5.1.4 Scale and SNR Analysis

Scale: Experiments were conducted to study the effect of scaling on MI estimators $I(\eta X; T)$, using correlated Gaussian variables X, T . Scaling factors η were sampled logarithmically between 10^{-2} and 10^3 for KSG, and 10^{-2} to 10 for MINE, as MINE estimates degrade sharply beyond $\eta = 10$. Results (Figure 4a, 4b) show KSG estimates converging to -0.33 (matching Proposition 3) as η increases, while MINE estimates drop to zero as $\eta \rightarrow 0$ (supporting Proposition 6). Global and local normalization variants remained robust to scaling, with significantly lower RMSE compared to vanilla estimators (Figure 4c, 4d). Notably, KSG exhibited tighter confidence intervals, whereas MINE estimates showed higher variance at extreme scales for the native estimator.

SNR: MI estimators were tested with Signal-to-Noise Ratios (SNR) ranging from 0 to 5, where $T = X + \epsilon$ ($\epsilon \sim \mathcal{N}(0, \sigma^2)$) and scaled to $T' = 0.1T$ (estimating $I(X; T')$). We expect the MI estimates to increase

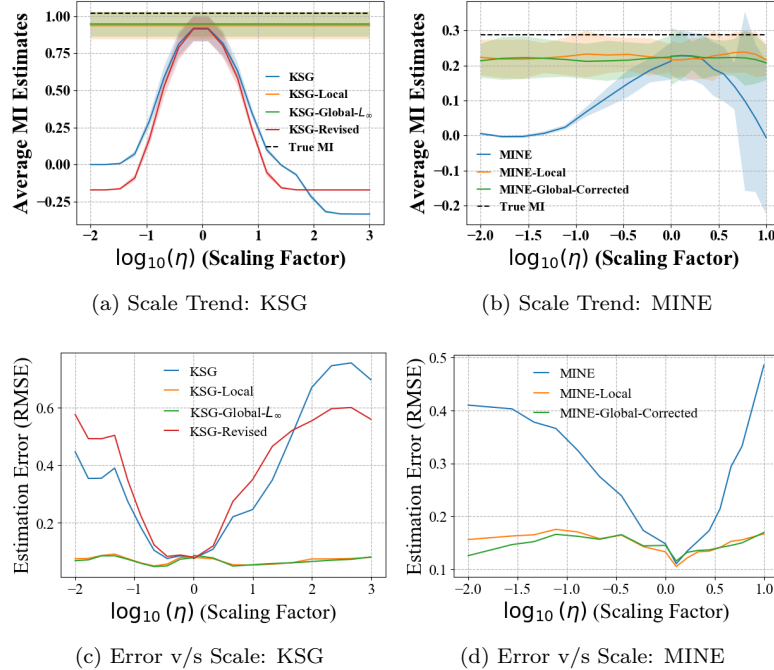


Figure 4: Analysis of MI Estimators in response to data scaling. Estimates are for $I(\eta X; T)$, where η is the scaling factor. Please see Appendix C.1 for details.

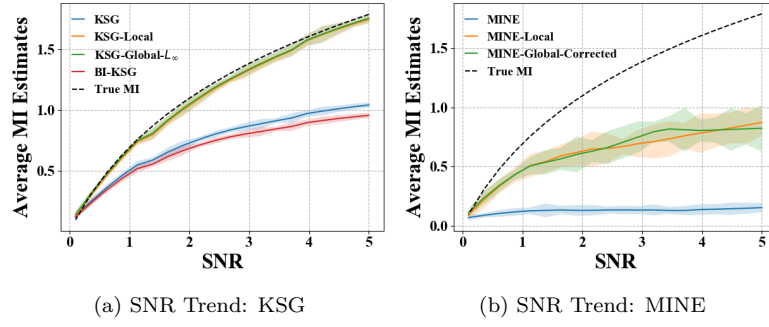


Figure 5: MI estimates across varying Signal-to-Noise Ratios (SNR). Estimates are for $I(X; \eta T)$, where T' is scaled by $\eta = 0.1$. See Appendix C.1 for details.

with SNR, following the ground truth. Results (Figure 5a, 5b) show global and local variants accurately following ground truth MI trends, increasing with SNR. However, vanilla KSG and MINE estimators failed to reflect the true MI trend, instead converging at higher SNR values due to their inherent scale dependence. This demonstrates the potential biases introduced by unnormalized estimators in scenarios involving variable scaling.

5.2 Comparing MI Estimators: Error Analysis

In this section, we undergo a comprehensive series of experiments, where we compute various error measures of all estimators on a diverse range of datasets.

5.2.1 Experiment Summary

Dataset creation: To create these datasets, we follow the two base distributions described in the beginning of this section. First, we generate X, T according to the two base distributions: Additive Gaussian and Correlated Gaussian as defined in Section 5. Then, we then make X undergo some (or none) of the following transformations, which are all MI preserving. For what follows, let $X \in \mathbb{R}^d$ and $T \in \mathbb{R}^d$.

1. **Randommat (rm):** $X' = \alpha W^T X$, where $\alpha \sim \text{Unif}(0, 1)$ and $W \in \mathbb{R}^{d \times d}$ where $W(i, j) \sim \text{Unif}(0, 1)$. $\text{Unif}(a, b)$ denotes a uniform distribution over $[a, b]$. If the randomly generated W is not invertible, we keep generating until we get an invertible W .
2. **Cube (cb):** $X' = X \circ X \circ X$, where \circ denotes element wise multiplication (Hadamard Product).
3. **Sigmoid (sg):** $X' = \sigma(X)$, where $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is such that $X'[i] = \frac{1}{1 + e^{-X[i]}}$, where $X[i]$ denotes the i^{th} dimension of X and similarly for X' .
4. **Duplicate-self (ds):** $X' = [X, X, \dots, X] \in \mathbb{R}^{Kd}$. We set $K = 20$ in our experiments.
5. **Duplicate-noise (dn):** $X' = [X, \epsilon] \in \mathbb{R}^{d+k}$, where $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_k]$ where $\epsilon_i \sim \mathcal{N}(0, \sigma'^2)$. We set $\sigma' = 0.2$ and $k = 20$.

Our objective is to evaluate the accuracy of the estimation of $I(X'; T)$.

Performance Measures: We study three different measures of performance in our experiments. For what follows, let $\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_k$ denote the estimated values of MI for any estimator across k trials, and let $\mu_1, \mu_2, \dots, \mu_k$ denote the ground truth values. With this, we summarize our performance measures as follows:

- **Normalized RMSE:** We first estimate the RMSE as $RMSE(\hat{\mu}, \mu) = \sqrt{\mathbb{E}_i[(\hat{\mu}_i - \mu_i)^2]}$. Then we estimate a baseline RMSE as $RMSE_Base(\mu) = \sqrt{\mathbb{E}_{i,j}[(\mu_i - \mu_j)^2]}$. With this, we can estimate the final measure as: $RMSE_Norm(\hat{\mu}, \mu) = \frac{RMSE(\hat{\mu}, \mu)}{RMSE_Base(\mu)}$.
- **Spearman Correlation:** The Spearman correlation measures the degree of monotonic relationship between $\hat{\mu}$ and μ (Zar, 2005). This is estimated as the Pearson's correlation coefficient between the rank values of $\hat{\mu}$ and μ .
- **Bias:** We estimate the bias as $\mathbb{E}_i[\mu_i - \hat{\mu}_i]$.

Table 2: Normalized RMSE of KSG-Based Estimators: Additive Gaussian Noise Base

Transformation					d	KSG-Based Measures				
rm	cb	sg	ds	dn		ksg	bi-ksg	ksg-loc	ksg-glo	ksg-glo- L_∞
	✓	✓	✓		2	0.351	0.404	0.147	0.208	0.110
		✓	✓		2	0.286	0.335	0.060	0.084	0.050
	✓	✓			2	0.458	0.533	0.145	0.113	0.113
		✓			4	1.275	1.424	0.312	0.300	0.301
✓				✓	4	0.862	0.932	0.445	0.594	0.396
			✓		4	0.332	0.342	0.304	0.520	0.297
				✓	4	0.332	0.342	1.327	0.298	0.297
✓	✓				4	1.131	1.233	0.977	0.959	0.931
		✓		✓	4	1.275	1.424	1.334	0.300	0.301
	✓	✓		✓	6	1.981	2.129	1.904	1.021	1.021
		✓			6	1.983	2.131	0.816	0.811	0.812
	✓		✓		6	1.377	1.408	1.343	1.605	1.290
✓			✓		6	1.643	1.730	1.275	1.660	1.153
		✓	✓	✓	6	1.983	2.131	1.905	0.814	0.816

Table 3: Normalized RMSE of MINE-Based Estimators: Additive Gaussian Noise Base

Transformation					d	MINE-Based Measures			
rm	cb	sg	ds	dn		mine	mine-loc	mine-glo	mine-glo-corr
	✓	✓	✓		2	0.470	0.292	0.337	0.278
		✓	✓		2	0.445	0.255	0.275	0.233
	✓	✓			2	1.036	0.560	0.658	0.565
		✓			4	1.302	0.720	0.968	0.684
✓				✓	4	0.930	0.438	0.803	0.381
			✓		4	0.276	0.369	0.642	0.375
				✓	4	0.622	0.423	0.895	0.269
✓	✓				4	1.574	1.099	1.219	1.162
		✓		✓	4	1.335	0.423	0.895	0.286
	✓	✓		✓	6	1.881	0.570	1.516	0.437
		✓			6	1.843	1.147	1.535	1.185
	✓		✓		6	1.213	0.991	1.444	1.004
✓			✓		6	1.014	1.088	1.441	1.058
		✓	✓	✓	6	1.831	0.517	1.499	0.545

Evaluation Process: We summarize the empirical process for the additive Gaussian noise base (full results, including the correlated Gaussian base, are in Appendix E) as reported in Tables 2 and 3. For each experiment, a specific MI-preserving transformation (indicated in the first column) is applied to X . Then, over 40 trials, we generate $N = 1000$ samples of $X, T \sim P(X, T)$, apply the transformation to obtain X' , and compute MI estimates $I(X'; T)$ across different data dimensions d . The performance measures: normalized RMSE, Spearman correlation, and bias are computed from these trials. (Red entries in the tables denote cases where the normalized RMSE exceeds one; we later note that even then, the estimates often retain significant Spearman correlation with the true MI.)

Remark 4. Since every transformation is MI preserving, combining them yields new distributions whose ground truth MI is unchanged. This flexible framework allows us to simulate high-dimensional data (up to 200 dimensions) with low intrinsic dimension (typically <10), reflecting the characteristics of neural network features. Our choice of transformations is motivated by the behaviors observed under different normalization strategies—for example, local normalization is adversely affected by added noise (duplicate-noise), while KSG-Global struggles with duplicate-self, and nonlinear transformations (e.g., sigmoid and cube) can alter nearest neighbor distances.

5.2.2 Takeaways

The main observations from the results are as follows:

- Overall, global and local normalization variants fare significantly better than the baseline measures.
- Our global normalization variants (MINE-global-corrected and KSG-Global- L_∞) overall fare better than other normalization strategies. In fact when the base distribution is additive Gaussian noise, we find that in most cases MINE-global-corrected and KSG-global- L_∞ outperform compared to the other normalization approaches.
- KSG-Global- L_∞ has very consistent performance, and across both settings, it seems to have the best performance in most cases. Even when the normalized RMSE estimates are insignificant (red entries), KSG-Global shows significant correlation with true MI in many of the cases (see Appendix E).
- As discussed in our motivation, we find that overall the global normalization variants (MINE-Global-Corrected and KSG-Global- L_∞) perform better than their vanilla global normalization counterparts. This is much more apparent in the case of MINE.

5.3 Application of MI Estimations in Deep Learning

Mutual information is a key measure for analyzing neural network behavior during training. We evaluate MI on IB (Shwartz-Ziv & Tishby, 2017), MNIST (Deng, 2012), and CIFAR-10 (Krizhevsky & Hinton, 2009) datasets. Network architectures, activations, and other details are in Appendix C. For each classification dataset $\{X, Y\}$, networks are trained and an intermediate layer’s output Z is extracted (third layer for IB and MNIST; Global Average Pooling for CIFAR-10). MI estimates $I(X; Z)$ and $I(Z; Y)$ are computed using KSG, KSG-Local, MINE, MINE-Local, and our proposed KSG-Global- L_∞ and MINE-Global-Corrected estimators.

We analyze the MI estimates from two perspectives:

- **Training Dynamics:** In Figure 6, we plot $I(X; Z)$ over epochs (averaged over 10 trials), along with the scale $|Z|$ (in blue). On MNIST and CIFAR-10, the original KSG estimates strongly follow the feature scale, indicating sensitivity to scaling. In contrast, KSG-Global- L_∞ does not mimic the scale curve; for IB and CIFAR-10, it first rises then decreases (with CIFAR-10 showing a drop after 3 epochs), consistent with the fitting-compression trend of (Shwartz-Ziv & Tishby, 2017).
- **Information Plane Visualization:** We also plot $I(X; Z)$ vs. $I(Z; Y)$ (details in Appendix G) to illustrate the trade-off between input representation and label relevance.

Figure 6 reveals distinct behaviors the proposed MI estimators across datasets:

IB and MNIST: On the IB dataset, both local and global variants of KSG and MINE estimators successfully exhibit the fitting (increase in $I(X; Z)$) and compression (decrease in $I(X; Z)$) phases during training, consistent with the information bottleneck theory (Shwartz-Ziv & Tishby, 2017). However, on the MNIST dataset, only the global variants, KSG-Global- L_∞ and MINE-Global-Corrected, clearly demonstrate these trends, while the vanilla variants primarily track the feature scale ($|Z|$) rather than capturing intrinsic mutual information dynamics. This suggests that the normalized variants, which have been tested in the previous section accuracy-wise, can yield different $I(X; Z)$ trends than their vanilla counterparts.

CIFAR-10: For CIFAR-10, an interesting divergence emerges. Vanilla KSG and MINE estimators show fitting and compression phases, but the fitting phase appears to strongly correlate with initial scale jumps of $|Z|$, indicating sensitivity to feature scale. In contrast, the normalized variants predominantly display compression-like trends throughout training, highlighting a different behavior. Notably, KSG-Global- L_∞ reveals a different dynamic, which only finds a compression phase from the start. This divergence emphasizes the ability of normalization variants to decouple scale effects from mutual information trends, providing deeper insights into network behavior.

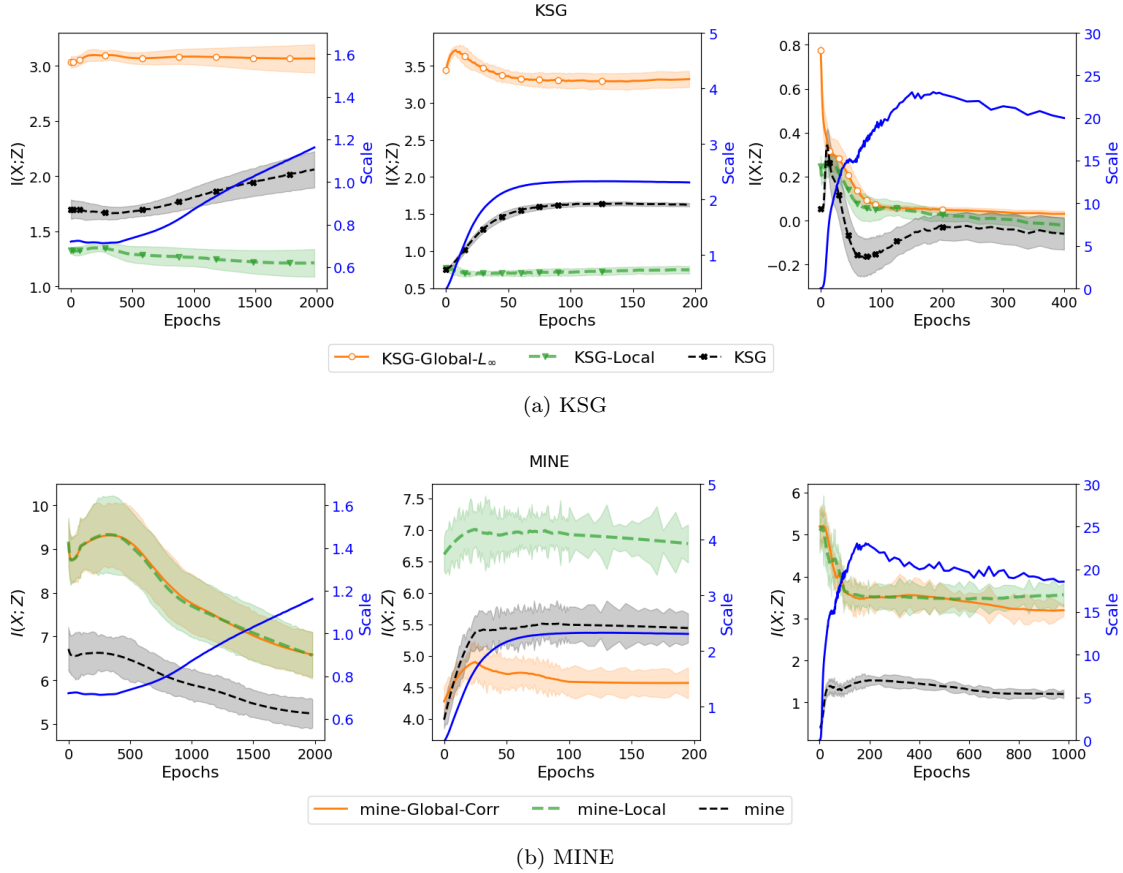


Figure 6: $I(X;Z)$ measures estimated after every epoch of training on IB, MNIST and CIFAR-10 datasets. Z represents the output of 3^{rd} layer for IB dataset and MNIST dataset, and 7^{th} layer for CIFAR-10 dataset. Details in Appendix D.3. The scale of the features ($|Z|$) is plotted in blue.

6 Conclusion

We presented a comprehensive study of scale invariance in MI estimators, and its impact on estimation accuracy, trends, and on MI-based analysis of neural network training. We outlined multiple normalization approaches to combat scale changes, centered around KSG and MINE, and discussed the pros and cons of each approach. Specifically targeting the high-dimensional and low-data regime, intuitive and empirical arguments were given for each normalization approach and the final choice of estimators. Overall we found that while both local normalization and global normalization have their own strengths, in most practical scenarios, global normalization variants fare better. Both normalization strategies lead to desirable behaviour in response to input scale changes. Extensive experiments across two broad settings were conducted to measure the overall performance of each estimator. In almost all cases, the local and global normalization approaches fare much better than their unnormalized counterparts, while global normalization variants have the best performance overall. Lastly, on three real datasets, we studied the information plane dynamics w.r.t the hidden layer feature representations during training, for the unnormalized and normalized estimator variants. More clear trends of fitting and compression were observed with global normalization approaches in two out of the three datasets, with KSG-Global variants showing clearer trends than MINE-Global variants. While KSG-Global normalization variants demonstrate superior performance, their computational cost is notably higher, as the MI estimation process requires n times the computation time compared to standard KSG variants. Our work highlights the importance of scale-awareness in the problem of MI estimation, and its potential impact on MI estimates.

References

- Milton Abramowitz. *Handbook of Mathematical Functions, With Formulas, Graphs, and Mathematical Tables*. Dover Publications, Inc., USA, 1974. ISBN 0-486-61272-4.
- Philip Bachman et al. Learning representations by maximizing mutual information across views. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Mohamed Ishmael Belghazi et al. Mutual information neural estimation. In Jennifer Dy and Andreas Krause (eds.), *International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 531–540. PMLR, July 2018.
- Yuheng Bu, Shaofeng Zou, and Venugopal V. Veeravalli. Tightening mutual information based bounds on generalization error. In *IEEE International Symposium on Information Theory*, pp. 587–591, 2019.
- Ivan Butakov, Alexander Tolmachev, Sofia Malanchuk, Anna Neopryatnaya, Alexey Frolov, and Kirill Andreev. Information bottleneck analysis of deep neural networks via lossy compression. In *International Conference on Learning Representations*, 2024.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, USA, 2006. ISBN 0-471-24195-4.
- Paweł Czyż et al. Beyond normal: On the evaluation of mutual information estimators. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Li Deng. The MNIST database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012. Publisher: IEEE.
- M. D. Donsker and S. R.S. Varadhan. Asymptotic evaluation of certain markov process expectations for large time. IV. *Communications on Pure and Applied Mathematics*, 36(2):183–212, March 1983. ISSN 0010-3640. doi: 10.1002/cpa.3160360204. Publisher: Wiley-Liss Inc.
- Shuyang Gao et al. Efficient Estimation of Mutual Information for Strongly Dependent Variables. In Guy Lebanon and S. V. N. Vishwanathan (eds.), *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pp. 277–286, San Diego, California, USA, May 2015. PMLR.
- Weihaio Gao, Sewoong Oh, and Pramod Viswanath. Demystifying fixed k-nearest neighbor information estimators. In *IEEE International Symposium on Information Theory*, pp. 1267–1271, 2017. doi: 10.1109/ISIT.2017.8006732.
- Ziv Goldfeld and Kristjan Greenewald. Sliced mutual information: A scalable measure of statistical dependence. *Advances in Neural Information Processing Systems*, 34:17567–17578, 2021.
- Gokul Gowri, Xiaokang Lun, Allon M Klein, and Peng Yin. Approximating mutual information of high-dimensional variables using learned representations. In *Neural Information Processing Systems (NeurIPS)*, 2024.
- R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2019.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017.
- Justin B. Kinney and Gurinder S. Atwal. Equitability, mutual information, and the maximal information coefficient. *Proceedings of the National Academy of Sciences*, 111(9):3354–3359, 2014.

- Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical Review E*, 69(6):066138, June 2004. Publisher: American Physical Society.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, Toronto, Ontario, 2009.
- David McAllester and Karl Stratos. Formal limitations on the measurement of mutual information. In Silvia Chiappa and Roberto Calandra (eds.), *International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 875–884, August 2020.
- Liam Paninski. Estimation of entropy and mutual information. *Neural computation*, 15(6):1191–1253, 2003. Publisher: MIT Press.
- Dávid Pál, Barnabás Póczos, and Csaba Szepesvári. Estimation of Rényi entropy and mutual information based on generalized nearest-neighbor graphs. In *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc., 2010.
- David N. Reshef et al. Equitability analysis of the maximal information coefficient, with comparisons. *arXiv preprint arXiv:1703.00810*, abs/1301.6314, 2013.
- Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Fr-train: A mutual information-based approach to fair and robust training. In *International Conference on Machine Learning*, 2020.
- Daniel Russo and James Zou. How much does your data exploration overfit? controlling bias via information usage. *IEEE Transactions on Information Theory*, 66(1):302–323, 2020.
- Andrew Michael Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan Daniel Tracey, and David Daniel Cox. On the information bottleneck theory of deep learning. In *International Conference on Learning Representations*, 2018.
- Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.
- Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information estimators. *arXiv preprint arXiv:1910.06222*, 2019.
- Michael Tschannen, Josip Djolonga, Paul K. Rubenstein, Sylvain Gelly, and Mario Lucic. On mutual information maximization for representation learning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- Janett Walters-Williams and Yan Li. Estimation of Mutual Information: A Survey. In Peng Wen, Yuefeng Li, Lech Polkowski, Yiyu Yao, Shusaku Tsumoto, and Guoyin Wang (eds.), *Rough Sets and Knowledge Technology*, pp. 389–396, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg. ISBN 978-3-642-02962-2.
- Rui Wang, Pengyu Cheng, and Ricardo Henao. Toward fairness in text generation via mutual information minimization based on importance sampling. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (eds.), *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 4473–4485. PMLR, 25–27 Apr 2023.
- Yaochen Xie, Ziqian Xie, Sheikh Muhammad Saiful Islam, Degui Zhi, and Shuiwang Ji. Genetic InfoMax: Exploring Mutual Information Maximization in High-Dimensional Imaging Genetics Studies. *Trans. Mach. Learn. Res.*, 2024, 2024.
- Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. A theory of usable information under computational constraints. In *International Conference on Learning Representations*, 2020.
- Jerrold H Zar. Spearman rank correlation. *Encyclopedia of biostatistics*, 7, 2005. Publisher: Wiley Online Library.

A Background on Mutual Information Estimators

A.1 Mutual Information

Mutual information of two variables is a statistical measure that quantifies the mutual dependence between two random variables. Specifically, it measures the amount of information obtained about one random variable through the observation of another. To understand mutual information, it is essential to first examine another foundational concept, Shannon entropy. Shannon entropy represents the intrinsic informational uncertainty associated with a probabilistic system. Given a continuous random variable X with a probability density function f from a set \mathcal{X} , the continuous entropy $h(X)$ is defined as:

$$h(X) := - \int_{\mathcal{X}} f(x) \log f(x) dx \quad (2)$$

Then, the mutual information between continuous random variables X and Y is given by:

$$I(X; Y) = h(X) + h(Y) - h(X, Y) \quad (3)$$

where $h(X, Y)$ represents the joint differential entropy of X and Y , defined as $h(X, Y) = - \int_{\mathcal{X}, \mathcal{Y}} f(x, y) \log f(x, y) dx dy$. Mutual information can be interpreted as the reduction in the uncertainty of X due to the knowledge of Y , or equivalently, as the amount of information that X and Y share.

In the case of jointly continuous random variables, the mutual information can be expressed in terms of Kullback–Leibler (KL-) divergence

$$I(X; Y) = D_{\text{KL}}(P(X, Y) \| P(X) \otimes P(Y)), \quad (4)$$

where $P(X) \otimes P(Y)$ is the dot product of two marginal distributions $P(X)$ and $P(Y)$, $P(X, Y)$ is their joint distribution. D_{KL} is defined as

$$D_{\text{KL}}(P \| Q) := \mathbb{E}_P \left[\log \frac{dP}{dQ} \right]. \quad (5)$$

In practice, estimating the true distribution of continuous random variables is challenging, especially for high-dimensional data. In the following section, we will discuss various non-parametric MI estimators, which estimate the distribution of random variables and subsequently compute estimated mutual information.

A.2 Mutual Information Estimators

The overall setting of the MI estimation problem is as follows. We are given two RVs $X, T \sim P(X, T)$ and sampled data $S = \{(X_1, T_1), \dots, (X_n, T_n)\}$. Our estimators are denoted in the form $\hat{I}_{est}^n(X; T)$, where *est* denotes the name of the estimator.

In this section, we present several widely-used nonparametric MI estimators that are studied in our work and have been extensively applied in other research.

Binning Estimator: Also called histogram based estimator in many research. This method represents the most elementary technique for estimating mutual information. To estimate MI, the continuous random variable is discretized into bins, counting the number of samples that fall into each bin, and computing the probability density (Paninski, 2003). The binning estimator for n samples can be expressed as: $\hat{I}_{bin}^n(X; T) = H_{bin}(X) + H_{bin}(T) - H_{bin}(X, T)$. where $H_{bin}(X)$ represents the binned entropy given a RV X , such that $H_{bin}(X) = - \sum_i P(X_i) \log P(X_i)$. Let $n(X_i)$ be the number of samples that fall in i -th bin of X , and N is the total number of data points. Then we have $P(X_i) \approx n(X_i)/N$ for binning method. Similarly, we represent binned joint entropy as $H_{bin}(X, T) = - \sum_{i,j} P(X_i, T_j) \log P(X_i, T_j)$, and $P(X_i, T_j) \approx n(X_i, T_j)/N$.

Kraskov–Stögbauer–Grassberger (KSG) Estimator: Another popular non-parametric approach to estimate MI in high dimensions is the KSG estimator in (Kraskov et al., 2004). Unlike the binning estimator, the KSG estimator uses the k -nearest neighbor (K -NN) statistic to estimate the probability function of

continuous random variables, which also uses the joint entropy decomposition method to estimate MI. The KSG estimator effectively uses the k -nearest neighbor distances to estimate the various entropies involved in the joint-entropy decomposition of $I(X; T)$. To define the KSG estimator, we will need some pre-requisites. First, let the k -NN distance $\rho_{k,i,p}$ be defined as the distance from (X_i, T_i) to k -th nearest neighbor in the joint space (X, T) as measured in l_p distance. We now denote $n_{x,i,p} = \sum_{j \neq i} \mathbb{I}\{\|X_j - X_i\|_p \leq \rho_{k,i,p}\}$ as the number of neighbors of the i -th sample X_i within a specified distance under the l_p norm, and similarly for T as $n_{t,i,p}$. Eventually, the KSG estimator yields the following estimate of $I(X; T)$:

$$\hat{I}_{KSG}^n(X; T) = \psi(k) + \psi(n) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n (\psi(n_{x,i,\infty}) + \psi(n_{t,i,\infty})), \quad (6)$$

where $\psi(x)$ is the digamma function (i.e., $\psi(x) = \Gamma(x)^{-1} d\Gamma(x)dx$).

In (Gao et al., 2017), authors proposed a bias-improved KSG (**BI-KSG**) that performs better than KSG when N is small and X and T are not independent. It is also important to note that many other variants of KSG and other estimators (Pál et al., 2010; Gao et al., 2015) use k -NN approach.

Mutual Information Neural Estimator (MINE): In our work, we utilize neural network based MI estimators, specifically Mutual Information Neural Estimation (Belghazi et al., 2018). This approach estimates mutual information by using a dual representation of the KL-divergence, known as Donsker-Varadhan (DV) representation (Donsker & Varadhan, 1983). Given RVs $X \sim P(X)$, $T \sim P(T)$, and $(X_i, T_i) \sim P(X, T)$, we express equation 4 in terms of DV representation as:

$$I(X; T) = \sup_{F: \mathcal{X} \times \mathcal{T} \rightarrow \mathbb{R}} \mathbb{E}_{X, T \sim P(X, T)}[F(X, T)] - \log \left(\mathbb{E}_{X, T \sim P(X) \times P(T)} \left[e^{F(X, T)} \right] \right), \quad (7)$$

where F can be any measurable function from $\mathcal{X} \times \mathcal{T} \rightarrow \mathbb{R}$ that satisfies the necessary integrability constraints of two expectations in equation 7 to be well-defined, and the supremum is taken over all functions F that contribute meaningfully to the optimization problem, ensuring the validity of equation 7.

To compute $I(X; T)$ in practice, assuming n independent and identically distributed samples (i.i.d.) are drawn from $(X_i, T_i) \sim P(X, T)$, and (X_i, \tilde{T}_i) is artificially constructed by choosing \tilde{T}_i as a randomly shuffled set of $(T_i)_{i=1}^n$. When n is large enough, inspired by the law of large numbers, the MINE estimator approximates the MI as:

$$\hat{I}_{\text{MINE}}(X; T) = \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n F_{\theta}(X_i, T_i) - \log \left(\frac{1}{n} \sum_{i=1}^n e^{F_{\theta}(X_i, \tilde{T}_i)} \right), \quad (8)$$

where $F_{\theta} : \mathcal{X} \times \mathcal{T} \rightarrow \mathbb{R}$ is parameterized by a deep neural network with parameters $\theta \in \Theta$, and Θ is the parameter space of the neural network. By training a neural network to optimize the above equation (i.e., finding the optimal F_{θ}), the final output will yield the MINE estimates of MI between X and T .

B Additional Discussions on Related Works

B.1 Equitability

There is a body of work on the equitability of mutual information estimators, of which the notion of self-consistent equitability is related to our work. It seems (Reshef et al., 2013) is the first that proposes the notion of self-consistent equitability, which is a generalized form of one-sided scale invariance discussed in our work, as we consider the specific case when the function is scale-related. In our work, however, we do not focus on the aspect of data transformation specifically, but mainly focus on the behaviour of estimators in response to scaling. In such works concerned with self-consistent equitability, there is always a very wide range of functions (Table 2 of (Reshef et al., 2013)) to get a broad overview of the behavior of MI in response to any type of transformation. However, it is notable that scaling has not been one of them, most likely due to the way the data was already preprocessed to be scale-invariant, using local normalization, which already makes them self-equitable to scaling.

Our work only indirectly tests self-consistent equitability in Tables 1 and 2, where instead of prioritizing the degree of self-consistent equitability, we prioritize the overall MI estimation accuracy. We also use a variety of transformations as functions on the data (in a cascaded manner) and see the resulting estimation errors for the various compared metrics. However, our choice of transformations (apart from cube and sigmoid) is very different when compared to Table 2 of Reshef et al. (2013)’s work, as they either encompass all dimensions (like invertible random matrix multiplication) or change the dimensionality of the input by adding noisy dimensions or duplicate dimensions. In contrast, most functions in self-consistent equitability literature work with one-dimensional transformations applied to each data dimension individually. Lastly, as we cascade a random subset of these operations in a random order, we get a richer set of MI-preserving transformations, which can potentially also be tested in an equitability setup using their choice of metrics (like in Figure 2 of (Kinney & Atwal, 2014)). We are considering this for future work.

B.2 Recent MI Estimation works

We outline three different aspects in which our work differs from (Czyż et al., 2023), which is an important recent study that also tests MI Estimators and their performance under various settings.

Choice of transformations: As the focus in (Czyż et al., 2023) is not particularly on the neural network use case where X is the input and Z represents a feature layer, their choice of transformations is motivated differently. Most of their transformations are dimension-wise, i.e., a transformation applied to each dimension. The only exception is their spiral diffeomorphism (Figure 5 of (Czyż et al., 2023)), which radially morphs the distribution in such a way that the MI is preserved. We note that as our focus is mainly on the natural use cases of MI in deep learning, we construct certain types of transformations relevant to this setting. Apart from the cubic transformation, which is dimension-wise, all our other transformations are motivated by the potential use case in deep learning. We outline each one as follows. The sigmoid transformation is motivated by potential uses of sigmoid in the network’s hidden layers. The random matrix multiplication (randmat) transformation is motivated via the features undergoing similar transformations through neural network layers, which are usually matrix multiplications followed by non-linearity. Note that the randmat also has an additional scaling term α (Section 5.2), which scales the resulting transformed vector as we wish to also focus on the robustness to scale. The duplicate-noise transformation, which adds dummy noise dimensions, is motivated by the fact that the number of hidden neurons can change through the layers and often many of these dimensions deeper within the network are usually very sparse and noisy. Similarly, the duplicate-self transformation is another approach to changing the dimensionality of the input while preserving the total information.

Preprocessing methods: Czyż et al. (2023)’s work indeed finds that the Gaussianization-based local normalization yields better results than uniform marginalization, when the base distribution is the multi-variate student distribution. However, they do note that the improvements are minor, and overall they end up choosing the standard variance-based local normalization over the other two. Our conjecture is that Gaussianization may work better when the distribution has long tails, as long-tail distributions are typically harder to estimate MI for, which was observed in (Czyż et al., 2023). This is potentially because the data

samples that are a part of the long tail may end up adding more noise to the final estimate, than in a typical case with Gaussian variables, and Gaussianizing the data preserves MI while avoiding long tails in the input, thereby leading to better performance. A potential direction of future work is thus incorporating similar considerations for our global normalization strategies to further enhance the accuracy of MI estimators.

Scale invariance: As the data was already preprocessed using local normalization approaches, scale invariance isn't a part of the analysis in (Czyż et al., 2023), similar to (Kinney & Atwal, 2014). In our case, all our normalization approaches are based on ensuring that the MI estimates are scale-invariant and while retaining other desired properties which may not hold for standard normalization approaches, such as robustness to noisy dimensions.

B.3 Neural network-based pre-processing

As one of our contributions is a set of new pre-processing strategies to ensure desirable scale-invariant behaviour of MI estimators, we discuss how this compares to other recent work that uses neural networks to learn pre-processing strategies for more accurate MI Estimation.

Gowri et al. (2024): In this work, the authors propose a two-step MI estimation procedure. First, they train compressed representations of the input which minimize an upper bound on the conditional entropies between the compressed representations and the input/output, after which they estimate the MI between the compressed representations using KSG. In the limiting case, these compressed representations would preserve the true MI, otherwise, normally they are upper-bounded by it. It is worth noting their compressed representations can be of any isomorphic form w.r.t the choice of compressor, and as they use neural networks, the scale of the compression can vary. So to alleviate the issue, they use the standard unit normalization approach (local). As such, our findings in this work, including the different variants proposed for global normalization, can be applied to their KSG estimation step to generate more robust MI estimates.

Butakov et al. (2024): It seemed to us that this work proposed a new estimator using normalizing flows, which in this case are function maps (which can be a neural network) which eventually transform the data distribution to a simpler one which has a closed form solution for MI estimation. In their work, they seem to focus on a Gaussian base distribution. In principle, their estimator could be scale invariant, as their final estimator is the analytical MI expression itself. However, we did not see any explicit theoretical study on this, as one will need to show that the correlation matrix at the end of applying multiple normalizing flows is invariant to one-sided scaling of one of the variables. Overall, our work points out that scale-invariance is an important consideration for any MI estimator to prevent scale confounding in the estimates.

C Appendix: Details for Experiments

C.1 Details for the Section 5: Experimental Studies

Key Parameters:

Figure 1: Average MI estimates for KSG for a varying number of noise dimensions

- **Setup:** Additive Gaussian ($X, T \in \mathbb{R}^2$ where $T = X + \epsilon$, with $\epsilon \sim \mathcal{N}(0, \sigma^2 I_2)$)
- **Number of Samples:** 1000
- **Number of Trials:** 10
- $\sigma = \frac{1}{\sqrt{2}}$
- $\sigma' = 0.04$

Figure 2: Bias of MINE-based measures for varying noise dimensions

- **Setup:** Correlated Gaussian $X, T \in \mathbb{R}^2$ with correlated coefficient ρ
- **Number of Samples:** 1000
- **Number of Trials:** 10
- **Correlation coefficient:** ρ : 0.2

Figure 3: Bias of the KSG estimator with the real data dimension

Figure 3(a):

- **Setup:** Correlated Gaussian $X, T \in \mathbb{R}^d$ with correlated coefficient ρ
- **Number of Samples:** 1000
- **Number of Trials:** 20
- **Dimensionality (d):** Evaluated over $d \in \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$

Figure 3(b):

- **Setup:** Correlated Gaussian $X, T \in \mathbb{R}^d$ with correlated coefficient ρ
- **Number of Samples:** 200
- **Number of Trials:** 20
- **Correlation coefficient:** $\rho \sim \text{Unif}[0.4, 0.9]$, where *Unif* denotes the uniform distribution
- **Dimensionality (d):** Evaluated over $d \in \{2, 4, 8, 16, 32, 64\}$, and $\log_2 d \in \{1, 2, 3, 4, 5, 6\}$.

Figure 4: Data duplication: Average MI estimates for KSG-based approaches.

- **Setup:** Additive Gaussian ($X, T \in \mathbb{R}^2$ where $T = X + \epsilon$, with $\epsilon \sim \mathcal{N}(0, \sigma^2 I_2)$), and $X' = [X, X, X, \dots X]$.
- **Number of Samples:** 1000
- **Number of Trials:** 10
- $\sigma = \frac{1}{\sqrt{2}}$
- **Dimensionality of Final Input:** $d_{X'} \in \{2, 22, 42, 62, 82, 102, 122, 142, 162, 182\}$

Figure 5: Estimator bias versus dimension: Comparing MINE with MINE-Global variants

- **Setup:** Correlated Gaussian $X, T \in \mathbb{R}^d$ with correlated coefficient ρ
- **Number of Samples:** 1000
- **Number of Trials:** 10
- **Correlation coefficient:** ρ : 0.2
- **Dimensionality (d):** Evaluated over $d_X, d_T \in \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$

Figure 6: Analysis of MI Estimators in response to data scaling. Estimates are for $I(\eta X; T)$, where η is the scaling factor

Figure 6(a) & 6(b):

- **Setup:** Correlated Gaussian $X, T \in \mathbb{R}^2$ with correlated coefficient ρ , estimates are for $I(\eta X; T)$, where η is the scaling factor.
- **Number of Samples:** 1000
- **Number of Trials:** 20
- **Correlation Coefficient:** ρ : 0.8 (KSG); 0.5 (MINE)
- **Scaling Factor (η):**
 - For KSG: $\eta \in [10^{-2}, 10^3]$, equispaced on a \log_{10} scale.
 - For MINE: $\eta \in [10^{-2}, 10]$, equispaced on a \log_{10} scale.

Figure 6(c) & 6(d):

- **Setup:** Correlated Gaussian $X, T \in \mathbb{R}^2$ with correlated coefficient ρ , RMSE of MI estimates $I(\eta X; T)$, where η is the scaling factor.
- **Number of Samples:** 1000
- **Number of Trials:** 20
- **Correlation Coefficient:** $\rho \sim \text{Unif}[0, 0.8]$
- **Scaling Factor (η):**
 - For KSG: $\eta \in [10^{-2}, 10^3]$, equispaced on a \log_{10} scale.
 - For MINE: $\eta \in [10^{-2}, 10]$, equispaced on a \log_{10} scale.
- **Error Metric:** Root Mean Squared Error (RMSE) computed between MI estimates and ground truth.

Figure 7: Average MI estimates for various estimators across different values of SNR.

- **Setup:** Additive Gaussian noise base, where $T = X + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2 I_2)$. Additionally, T is scaled to $T' = 0.1T$, and the mutual information is estimated as $I(X; T')$.
- **Number of Samples:** 1000
- **Number of Trials:** 10
- $\sigma = \frac{1}{\sqrt{\text{SNR}}}$
- **Dimensionality (d):** $d_X, d_T = 2$

D MI Estimators: Configurations

D.1 KSG

We used the NPEET MI estimator toolbox for estimating KSG and KSG-based measures ³. We set $k = 3$ for all experiments. For the global KSG variants, we fix $c_1 = 0.1, c_2 = 0.2, \dots, c_n = 2$ for all our experiments.

D.2 MINE

We used the popular pytorch-based package ⁴ for the MINE implementation.

Overall MINE implementation: For estimating $I(X; T)$, we used single-hidden layer ReLU-activated neural networks of the configuration: $(d_X + d_T) \rightarrow H_1 \rightarrow H_2 \rightarrow \dots \rightarrow H_k \rightarrow 1$, where $d_X + d_T$ is the dimensionality of the input and H_1, \dots, H_k is the number of hidden neurons for each hidden layer. The last layer is a linear layer. We used the Adam optimizer with a learning rate of 0.001. The hidden neuron configuration varies depending on our experiment. We set the number of epochs to 50 for all experiments. Given a training dataset $S = \{(X_1, T_1), \dots, (X_n, T_n)\}$, the network F_θ effectively minimizes the following loss:

$$-\frac{1}{n} \sum_{i=1}^n F_\theta(X_i, T_i) + \log \left(\frac{1}{n} \sum_{i=1}^n e^{F_\theta(X_i, \tilde{T}_i)} \right), \quad (9)$$

Note that we estimated MINE in the standard manner as per its original definition, which is, once the networks are optimized, we estimate the the above loss on the training dataset Song & Ermon (2019). It is noteworthy that to avoid overestimation, (Czyż et al., 2023) proposes an approach where the MINE estimate

³ <https://github.com/gregversteeg/NPEET> ⁴ <https://github.com/gtegner/mine-pytorch>

is computed on the test data. However, in our case, we find that MINE does not overestimate for most of the cases in Tables 3 and 5, as it almost always has negative bias, and thus the training/test split approach is not necessary in our case.

Figures 1-7 and Tables 3 and 5: For the small dataset cases, we found that using a smaller number of hidden neurons yielded significantly better and more stable results. Thus, for Figures 1-7 and Tables 3 and 5, we found that using a single hidden layer with $H_1 = 20$ yielded the most stable results on average, and thus we set $H_1 = 20$. We found that for this small sample size setting, increasing H_1 led to unstable estimates and large variance of estimators.

Figure 8-10: For the larger dataset cases, which is the case for our real datasets, we were able to increase H , and the details are as follows. For the IB dataset, we have one hidden layer with $H = 30$ neurons. For MNIST and CIFAR-10 datasets, we have two hidden layers with $H = 30$ neurons each. Each dense layer uses ReLU activation.

D.3 Network Architecture for Neural Network Analysis in Section 6

Table 4: Model Architecture for IB Dataset

Layer	Dimension	Activation Function
Input	28×28	-
Flatten	12	-
Dense	10	ReLU
Dense	7	ReLU
Dense	5	ReLU
Dense	4	ReLU
Dense	4	ReLU
Dense	2	SoftMax

Table 5: Model Architecture for MNIST Dataset

Layer	Dimension	Activation Function
Input	28×28	-
Flatten	784	-
Dense	1024	ReLU
Dense	20	ReLU
Dense	20	ReLU
Dense	20	ReLU
Dense	10	SoftMax

For the MNIST and IB datasets, we replicate the network architectures from Saxe et al. (2018)’s work, using the widely-adopted *ReLU* activation function for the hidden layers. Specifically, for the IB dataset, we utilize a neural network with 7 hidden layers of dimensions 12-10-7-5-4-3-2. For the MNIST dataset, the neural network consists of 6 fully connected layers with dimensions 784-1024-20-20-20-10.

For the CIFAR-10 dataset, we adopt a neural network with 4 convolutional layers, 3 fully connected layers. The tasks for the MNIST and CIFAR-10 datasets involve classifying image inputs into their respective classes, while the task for the IB dataset involves training a binary decision rule based on 12 randomly distributed points. The networks are trained using SGD and cross-entropy loss. We train 2000 epochs for the IB dataset, 200 epochs for the MNIST dataset, and 1000 epochs for the CIFAR-10 dataset.

In Table 5, Table 4 and Table 6, we present the network architecture and output dimensions for each layer of the neural networks used in our study. The layers with bold text are the layers for extracted Z .

Table 6: Model Architecture for CIFAR-10 Dataset

Layer	Dimension	Activation Function
Input	$32 \times 32 \times 3$	-
Conv2D	$32 \times 32 \times 16$	ReLU
Conv2D	$32 \times 32 \times 16$	ReLU
MaxPooling	$16 \times 16 \times 16$	-
Conv2D	$16 \times 16 \times 32$	ReLU
Conv2D	$16 \times 16 \times 32$	ReLU
Global AveragePooling	32	-
Dense	64	ReLU
Dense	10	SoftMax

For the IB dataset, we trained for 2000 epochs with an SGD optimizer and a learning rate of 5×10^{-3} . For the MNIST dataset, we trained for 200 epochs with an SGD optimizer and a learning rate of 5×10^{-4} . For the CIFAR-10 dataset, we trained for 1000 epochs with an SGD optimizer and a learning rate of 1×10^{-3} . The batch sizes were 256 for the IB dataset, 128 for the MNIST dataset, and 512 for the CIFAR-10 dataset.

E Full Results on Synthetic Data

We provide the full results of Tables 2 and 3 of the main paper, in Table 7, and include the full results for the correlated Gaussian distribution base in Table 8. In the following tables, in addition to the normalized root mean squared error, we also report the Spearman correlation and Bias of each estimator, as defined in Section 5.2.1. Table 7 addresses the full KSG and MINE results in the same setting as Tables 2 and 3, and similarly, Table 8 addresses the full KSG and MINE results for the correlated Gaussian distribution base.

These are the specific takeaways from the Spearman correlation and bias results from Tables 7 and 8:

1. Very interestingly, we find that although the prediction error in terms of normalized RMSE quickly becomes higher than random guessing (red entries) for data in higher dimensions, the Spearman correlation of MI estimates with the ground truth MI still stays high in many cases.
2. A particularly notable example of this is for KSG variants, where we find that in spite of normalized RMSE exceeding 1, the KSG variants still have Spearman correlation measures very close to 1, showing that they are very highly rank-correlated with the true MI. It is therefore clear that the large negative bias of KSG variants in high dimensions affects all cases similarly, and thus the dependency between the estimated MI and the true MI remains. This also highlights that using some carefully calibrated ways of estimating MI in high dimensional settings may yield very accurate predictions in terms of normalized RMSE as well.
3. In fact, we find that over all cases in Tables 7 and 8, only the KSG-global variant stays very consistent in terms of high Spearman correlation. Barring only two cases, we see that the Spearman correlation of KSG-global- L_∞ is always greater than 0.9.
4. We find that overall our normalized variants of KSG and MINE showcase bias closest to 0. Interestingly, we see that for all variants, the bias roughly becomes increasingly negative as the dimensionality of the input increases.

Table 7: Comparing performance measures of MI Estimators: Additive Gaussian Noise Base (Full Results)

Transformation				d	measure	KSG-Based Measures				MINE-Based Measures			
rm	cb	sg	ds	dn		ksg	bi-ksg	ksg-loc	ksg-glo	ksg-glo- L_∞	mine	mine-loc	mine-glo-corr
✓	✓	✓	✓	2	RMSE-norm	0.351	0.404	0.147	0.208	0.110	0.470	0.292	0.337
					spearman bias	0.985	0.985	0.983	0.982	0.983	0.938	0.920	0.889
✓	✓	✓	✓	2	RMSE-norm	-0.191	-0.226	-0.072	-0.105	-0.047	-0.246	-0.130	-0.161
					spearman bias	0.286	0.335	0.060	0.084	0.050	0.445	0.255	0.275
✓	✓	✓	✓	2	RMSE-norm	-0.160	-0.192	-0.013	-0.035	0.004	0.938	0.929	0.917
					spearman bias	0.458	0.533	0.145	0.113	0.113	-0.233	-0.108	-0.126
✓	✓	✓	✓	2	RMSE-norm	0.982	0.981	0.986	0.987	0.987	1.036	0.560	0.658
					spearman bias	-0.253	-0.302	-0.071	-0.050	-0.050	0.548	0.816	0.853
✓	✓	✓	✓	4	RMSE-norm	1.275	1.424	0.312	0.300	0.301	1.302	0.720	0.968
					spearman bias	0.988	0.991	0.993	0.995	0.995	0.805	0.930	0.901
✓	✓	✓	✓	4	RMSE-norm	-1.078	-1.235	-0.254	-0.240	-0.240	-1.105	-0.566	-0.781
					spearman bias	0.862	0.932	0.445	0.594	0.396	0.930	0.438	0.803
✓	✓	✓	✓	4	RMSE-norm	0.599	0.549	0.995	0.990	0.994	0.660	0.897	0.944
					spearman bias	-0.653	-0.714	-0.363	-0.490	-0.322	-0.738	-0.296	-0.640
✓	✓	✓	✓	4	RMSE-norm	0.332	0.342	0.304	0.520	0.297	0.276	0.369	0.642
					spearman bias	0.994	0.994	0.995	0.994	0.994	0.931	0.922	0.922
✓	✓	✓	✓	4	RMSE-norm	-0.282	-0.298	-0.247	-0.426	-0.239	-0.122	-0.223	-0.494
					spearman bias	0.332	0.342	1.327	0.298	0.297	0.622	0.423	0.895
✓	✓	✓	✓	4	RMSE-norm	0.994	0.994	0.954	0.995	0.994	0.876	0.865	0.917
					spearman bias	-0.282	-0.298	-1.107	-0.239	-0.239	-0.486	0.002	-0.719
✓	✓	✓	✓	4	RMSE-norm	1.131	1.233	0.977	0.959	0.931	1.574	1.099	1.219
					spearman bias	0.588	0.534	0.975	0.961	0.969	0.542	0.905	0.837
✓	✓	✓	✓	4	RMSE-norm	-0.942	-1.045	-0.816	-0.796	-0.773	-1.198	-0.914	-1.012
					spearman bias	1.275	1.424	1.334	0.300	0.301	1.335	0.423	0.895
✓	✓	✓	✓	4	RMSE-norm	0.988	0.991	0.916	0.995	0.995	0.807	0.870	0.914
					spearman bias	-1.078	-1.235	-1.113	-0.239	-0.240	-1.119	-0.004	-0.718
✓	✓	✓	✓	6	RMSE-norm	1.981	2.129	1.904	1.021	1.021	1.881	0.570	1.516
					spearman bias	0.969	0.966	0.905	0.989	0.989	0.852	0.786	0.881
✓	✓	✓	✓	6	RMSE-norm	-1.999	-2.174	-1.919	-1.021	-1.021	-1.898	-0.179	-1.502
					spearman bias	1.983	2.131	0.816	0.811	0.812	1.843	1.147	1.535
✓	✓	✓	✓	6	RMSE-norm	0.959	0.957	0.995	0.995	0.995	0.842	0.771	0.862
					spearman bias	-2.001	-2.175	-0.816	-0.809	-0.810	-1.861	-1.116	-1.527
✓	✓	✓	✓	6	RMSE-norm	1.377	1.408	1.343	1.605	1.290	1.213	0.991	1.444
					spearman bias	0.983	0.981	0.979	0.908	0.989	0.694	0.783	0.831
✓	✓	✓	✓	6	RMSE-norm	-1.377	-1.410	-1.347	-1.609	-1.291	-1.158	-0.933	-1.432
					spearman bias	1.643	1.730	1.275	1.660	1.153	1.014	1.088	1.441
✓	✓	✓	✓	6	RMSE-norm	0.338	0.314	0.937	0.871	0.976	0.775	0.798	0.841
					spearman bias	-1.631	-1.730	-1.275	-1.666	-1.150	-0.967	-1.042	-1.431
✓	✓	✓	✓	6	RMSE-norm	1.983	2.131	1.905	0.814	0.816	1.831	0.517	1.499
					spearman bias	0.956	0.955	0.926	0.995	0.995	0.813	0.806	0.811
✓	✓	✓	✓	6	RMSE-norm	-2.001	-2.175	-1.921	-0.811	-0.814	-1.847	-0.068	-1.471
					spearman bias								

Table 8: Comparing performance measures of MI Estimators: Correlated Gaussian Base

Transformation				N	d	measure	KSG-Based Measures					MINE-Based Measures				
rm	cb	sg	ds	dn			ksg	bi-ksg	ksg-loc	ksg-glo	ksg-glo- L_{∞}	mine	mine-loc	mine-glo	mine-glo-corr	
					200	2	RMSE-norm spearman bias	0.125 0.918 -0.021	0.125 0.922 -0.023	0.122 0.924 -0.020	0.113 0.936 0.014	0.115 0.936 0.015	0.573 0.711 -0.238	0.579 0.711 -0.247	0.601 0.596 -0.257	0.555 0.801 -0.238
			✓		200	2	RMSE-norm spearman bias	0.138 0.923 -0.049	0.140 0.923 -0.051	0.140 0.914 -0.048	0.160 0.898 -0.042	0.116 0.878 -0.005	0.359 0.886 -0.114	0.321 0.929 -0.089	0.475 0.931 -0.192	0.315 0.940 -0.108
	✓		✓		200	2	RMSE-norm spearman bias	0.305 0.794 -0.124	0.332 0.798 -0.146	0.281 0.911 -0.117	0.233 0.916 -0.125	0.233 0.867 -0.077	0.350 0.889 -0.282	0.416 0.897 -0.149	0.544 0.907 -0.227	0.405 0.944 -0.156
	✓	✓	✓		1000	2	RMSE-norm spearman bias	0.129 0.961 -0.046	0.151 0.958 -0.064	0.072 0.911 -0.019	0.094 0.960 -0.023	0.060 0.952 0.000	0.278 0.951 -0.106	0.135 0.942 -0.029	0.195 0.969 -0.057	0.155 0.967 -0.038
		✓	✓		1000	2	RMSE-norm spearman bias	0.090 0.947 -0.031	0.103 0.950 -0.041	0.046 0.954 -0.004	0.052 0.943 -0.003	0.049 0.961 0.017	0.274 0.957 -0.104	0.130 0.941 -0.025	0.169 0.973 -0.042	0.144 0.973 -0.032
	✓	✓			1000	2	RMSE-norm spearman bias	0.126 0.940 -0.042	0.147 0.935 -0.060	0.065 0.939 -0.020	0.060 0.949 0.000	0.060 0.950 0.000	0.489 0.879 -0.213	0.281 0.950 -0.105	0.358 0.946 -0.145	0.290 0.950 -0.112
		✓			200	5	RMSE-norm spearman bias	0.962 0.736 -0.664	1.088 0.408 -0.822	0.481 0.964 -0.322	0.470 0.966 -0.293	0.469 0.969 -0.294	1.014 0.715 -0.703	0.949 0.982 -0.641	0.982 0.852 -0.673	0.934 0.789 -0.629
✓					200	5	RMSE-norm spearman bias	0.738 0.713 -0.442	0.783 0.569 -0.498	0.688 0.893 -0.406	0.628 0.895 -0.349	0.625 0.905 -0.346	0.974 0.606 -0.672	0.956 0.930 -0.648	0.993 0.765 -0.681	0.952 0.833 -0.649
	✓				200	5	RMSE-norm spearman bias	0.778 0.933 -0.539	0.819 0.931 -0.589	0.726 0.932 -0.498	0.695 0.937 -0.462	0.693 0.939 -0.460	1.084 0.605 -0.784	0.980 0.773 -0.670	0.999 0.805 -0.688	0.967 0.711 -0.660
				✓	1000	5	RMSE-norm spearman bias	0.258 0.990 -0.149	0.259 0.989 -0.152	0.793 0.812 -0.486	0.253 0.985 -0.134	0.253 0.986 -0.134	0.514 0.974 -0.304	0.489 0.847 0.171	0.680 0.968 -0.441	0.285 0.954 0.076
✓	✓				1000	5	RMSE-norm spearman bias	0.898 0.582 -0.523	0.981 0.433 -0.615	0.792 0.943 -0.472	0.767 0.963 -0.456	0.744 0.965 -0.440	0.952 0.511 -0.680	0.767 0.923 -0.526	0.819 0.924 -0.567	0.766 0.902 -0.533
		✓			200	10	RMSE-norm spearman bias	1.381 0.131 -1.351	1.473 0.170 -1.522	1.008 0.956 -0.985	0.999 0.962 -0.957	1.000 0.963 -0.960	1.441 0.512 -1.412	1.357 0.790 -1.288	1.424 0.806 -1.388	1.352 0.856 -1.290
	✓	✓	✓		200	10	RMSE-norm spearman bias	1.424 -0.127 -1.212	1.514 -0.083 -1.383	1.384 0.751 -1.175	1.122 0.935 -0.910	1.121 0.941 -0.909	1.446 0.590 -1.418	1.063 0.005 -0.451	1.411 0.771 -1.362	1.096 0.786 -0.784
		✓			200	10	RMSE-norm spearman bias	1.356 0.831 -1.164	1.412 0.823 -1.263	1.229 0.872 -1.043	1.321 0.885 -1.111	1.213 0.900 -1.011	1.287 0.652 -1.182	1.163 0.933 -0.990	1.382 0.942 -1.325	1.162 0.916 -0.976
✓		✓			1000	10	RMSE-norm spearman bias	1.217 0.815 -1.033	1.260 0.641 -1.120	1.158 0.977 -0.976	1.353 0.947 -1.145	1.063 0.986 -0.890	0.837 0.956 -0.869	0.877 0.954 -0.844	1.077 0.982 -1.030	0.912 0.969 -0.861
		✓			1000	10	RMSE-norm spearman bias	1.424 -0.065 -1.212	1.516 -0.450 -1.386	1.386 0.845 -1.177	0.875 0.992 -0.721	0.875 0.992 -0.721	1.185 0.748 -1.145	0.808 0.875 0.444	1.090 0.951 -0.981	0.723 0.898 0.560
✓		✓			1000	2	RMSE-norm spearman bias	0.141 0.968 -0.040	0.159 0.943 -0.050	0.093 0.967 -0.031	0.078 0.964 -0.014	0.061 0.967 0.000	0.280 0.793 -0.105	0.157 0.953 -0.017	0.235 0.960 -0.085	0.150 0.950 -0.027

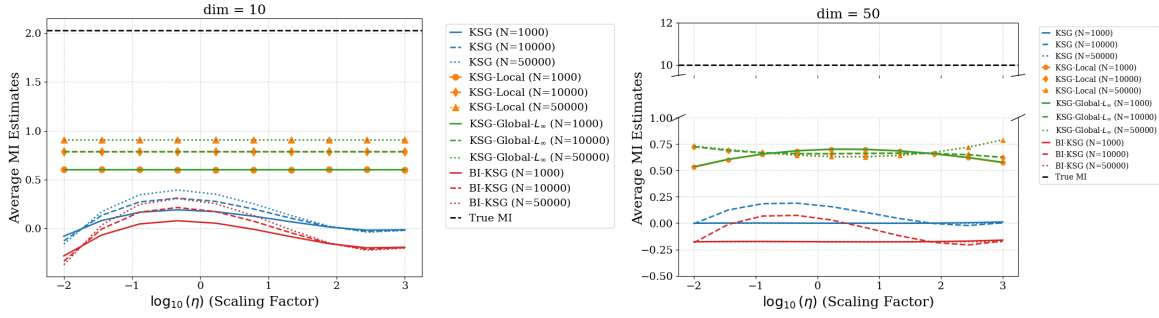


Figure 7: Analysis of the KSG Estimators in response to data scaling for varying number of data samples $N = \{1000, 10000, 50000\}$, and different data dimensionality $d = \{10, 50\}$. Estimates are for $I(\eta X; T)$, where η is the scaling factor.

F Additional Results

F.1 Impact of Noise

As an important measure of the dependence between two variables, mutual information is also widely used to analyze the behavior of neural networks during training. In this section, we present a comprehensive analysis of MI measures on network behaviors across datasets, including IB dataset (Shwartz-Ziv & Tishby, 2017), MNIST (Deng, 2012) and CIFAR-10 (Krizhevsky & Hinton, 2009). Network architecture architectures and activation functions and other details are provided in Appendix C.

F.2 Scale Invariance Testing

In the same setting as Section 5.1.4, we conduct more experiments to see if behaviour of the various estimators in response to data scaling remains unchanged for a different number of sampled datapoints $N = \{1000, 5000\}$ and $d = \{10, 50\}$. The results are shown in Figures 7 and 8.

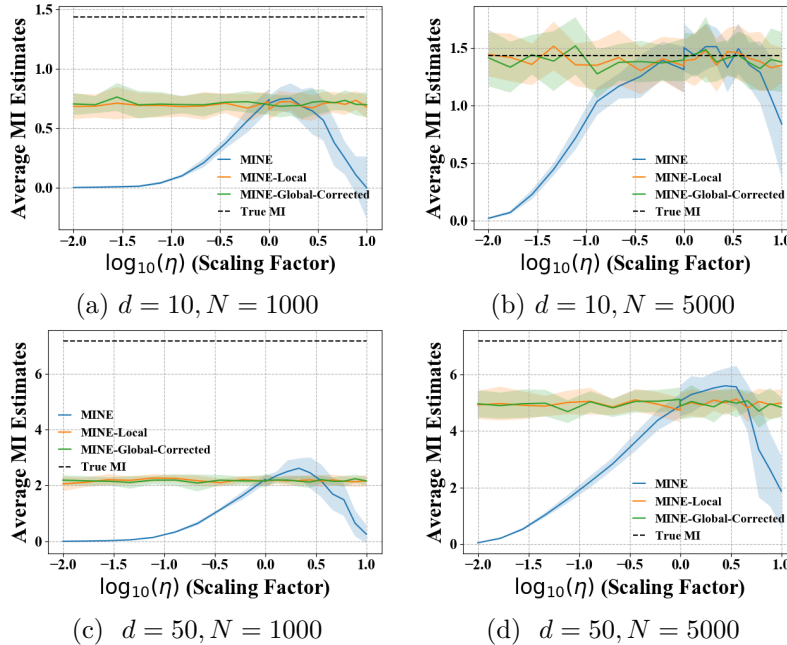


Figure 8: Analysis of the MINE Estimators in response to data scaling for varying number of data samples $N = \{1000, 5000\}$, and different data dimensionality $d = \{10, 50\}$. Estimates are for $I(\eta X; T)$, where η is the scaling factor.

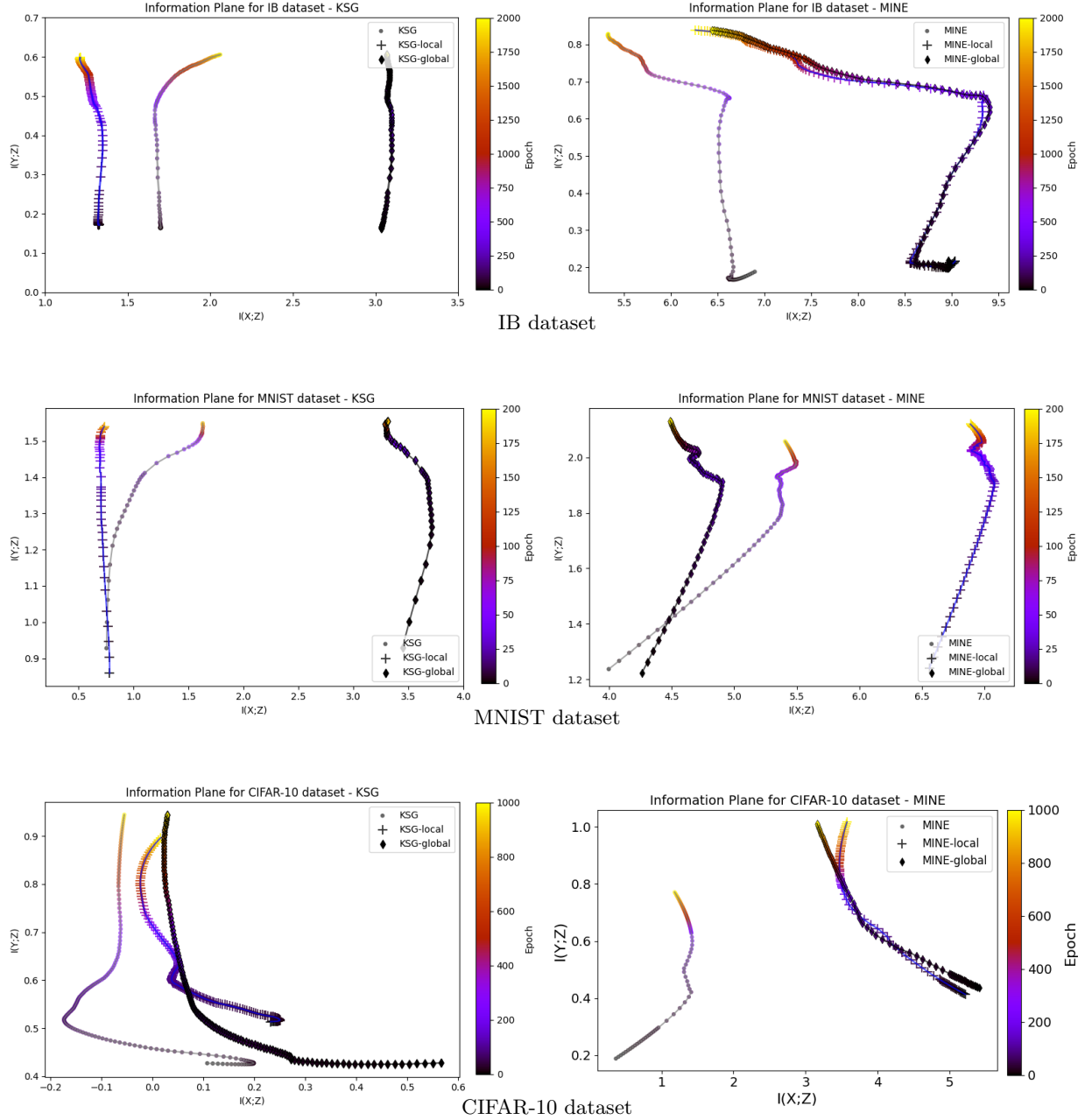


Figure 9: Information plane ($I(X; Z)$ against $I(Y; Z)$) for IB, MNIST and CIFAR-10 datasets. Z is the output of 3^{rd} layer for IB dataset and MNIST dataset, and 7^{th} layer for CIFAR-10 dataset. Details in Appendix D.3

KSG: Overall, for KSG (Figure 7), when we’re analyzing the behaviour of the native estimators, we find that they show similar response to scale, i.e., they converge to very low values near zero as the scale η goes to either extremity. In contrast, the scale-invariant estimators preserve the response across scales for $d = 10$. For higher dimensionality $d = 50$, we see some interesting behavioral changes for the scale-invariant KSG-Local and KSG-Global variants. We find that although they do not reach negligible values when scale reaches either extreme, the average MI estimates do not stay the same for all η , and there is some variation. Overall, the KSG-Global and Local variants behave similarly as before, and their average measures are relatively stable when η is no not small or too large. Only when η is either less than 0.1 or greater than 10, we see a slightly more pronounced change for the average MI estimates.

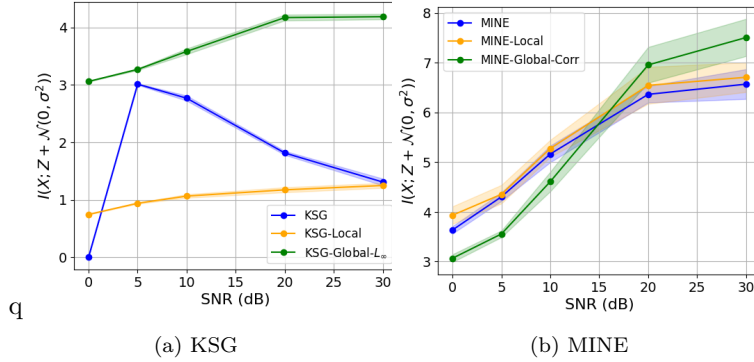


Figure 10: $I(X; Z + \mathcal{N}(0, \sigma^2))$ results with varying SNR. Shaded regions indicate 95% confidence intervals derived from 10 trials.

MINE: For the MINE results (Figure 8), we observe some interesting variations in the scale response depending on the dimensionality and the number of datapoints. We find that when we use a greater number of datapoints, then the MINE estimate converges to zero only for smaller scaling factors η . Our result in Proposition 5 essentially applies for the limiting case when $\eta \rightarrow 0^+$, and we can intuitively show that the value of η for which we see this limiting behaviour reduces with the greater sample number. The proof of Proposition 5 essentially finds that the weights associated with the scaled RV ηT is upper bounded in magnitude by a multiplicative factor of η and the number of updates and epochs. As a greater number of datapoints also implies a greater number of gradient descent updates, this also implies that we shall have potentially larger weights when the sample number increases. This implies that for larger datapoints N the MI estimate may stay non-zero for a larger range of η , and drop to near-zero values only near the extreme values of η , which is what we see in Figure 8. Lastly, we see that the local and global variants of MINE have relatively stable behaviour of their mean values across η for all d, N combinations. However, we do see that the variance of these estimators increases with more datapoints.

F.3 Impact of noise in neural network training

In this experiment, $I(X; Z + \mathcal{N}(0, \sigma^2))$ is measured at the last epoch of the training to evaluate how the MI changes when additive Gaussian noise is introduced. an additive noise layer $N \sim \mathcal{N}(0, \sigma^2)$ is introduced before the Z layer. In figure 10, the noise level is characterized using the signal-to-noise ratio (SNR) numbers, which quantifies the strength of a signal relative to the background noise. Specifically, an SNR of a dB implies that with unit signal power, the noise variance is $10^{-a/10}$. Consequently, as SNR increases, the noise level decreases, leading to a reduction in interference. Briefly, the mutual information $I(X; Z + \mathcal{N}(0, \sigma^2))$ is expected to increase as the noise level decreases. Since the SNR is expressed in decibels (dB), the noise reduction is more pronounced within the range of low SNRs (e.g. 0 to 15 dB), decreasing faster compared to the range of high SNRs (15 to 30 dB), and MI within the range of low SNR should also have a higher growing rate compared to high SNRs.

In figure 10, we present the change of $I(X; Z + \mathcal{N}(0, \sigma^2))$ during training, comparing the original KSG estimator, its local-normalized and global-normalized variants, and the MINE estimator and its variants. The displayed results represent the averages from 10 trials. For the original KSG estimator, we observe that as the SNR increases, the mutual information does not change as initially anticipated. Instead, the results for the original KSG estimator initially increase and then decline as SNR continues to rise. Notably, both estimators with global normalization exhibit the most consistent trend, reflecting the expected increase in dependence between X and Z as noise is reduced, and have a higher growing rate at low SNR regime.

G Information Plane Analysis

In figure 9, we present the information plane analysis for the IB, MNIST, and CIFAR-10 datasets, employing both KSG and MINE estimators to examine MI dynamics during neural network training. The displayed curves represent averages from 10 independent trials.

For the IB and MNIST datasets, the original KSG estimator shows that both $I(X; Z)$ and $I(Y; Z)$ generally increase as training progresses. However, the KSG-local and KSG-global- L_∞ estimators reveal a clearer information bottleneck pattern. Specifically, these estimators distinctly identify two phases: a fitting phase characterized by an initial increase in $I(X; Z)$ that eventually stabilizes, and a subsequent compression phase where $I(X; Z)$ decreases while $I(Y; Z)$ continues to increase. For CIFAR-10, the fitting phase is evident only during the first 10 epochs, after which a compression phase emerges. The original KSG estimator exhibits an irregular trend, initially increasing, then decreasing, and subsequently increasing again. We hypothesize that the initial compression phase is later overshadowed by scale-related effects. In contrast, the KSG-global- L_∞ estimator clearly and consistently demonstrates a monotonic compression phase. Among these methods, KSG-global- L_∞ consistently captures the most distinct fitting and compression phases across datasets, effectively highlighting the information bottleneck phenomenon.

Results from experiments with MINE estimators also yield insightful observations. For the IB dataset, the vanilla MINE estimator initially shows an unexpected reduction in both $I(X; Z)$ and $I(Y; Z)$, indicating a brief period where the model seems to discard label information, followed by clear fitting and compression phases. On the MNIST dataset, vanilla MINE captures only the fitting phase, with $I(X; Z)$ continuously increasing and no subsequent compression phase. In contrast, our proposed MINE-global-corrected variant clearly demonstrates both fitting and compression phases, outperforming the local-normalized MINE variant. For CIFAR-10, vanilla MINE only detects a fitting phase, while local-normalized and global-normalized MINE variants successfully reveal clear compression phases.

H Proofs of theoretical results

Proposition 1. It holds that $\hat{I}_{bin}^n(\alpha X; \alpha T) = \hat{I}_{bin}^n(X; T)$ and $\hat{I}_{bin}^n(X; \alpha T) = \hat{I}_{bin}^n(X; T) \forall \alpha \in \mathbb{R}^+$.

Proof. We note that the number of bins chosen for each dimension is fixed, and the locations of the bins are determined by the minimum and maximum values of the data in each dimension, i.e., they determine the edges of the bins. Let $X_{min} \in \mathbb{R}^d$ then denote the vector of minimum values across all dimensions, and vice-versa for X_{max} . When X scales to αX , as $\alpha > 0$, we have that the vector of minimum values for αX is simply αX_{min} and similarly for X_{max} , and the binning locations also get scaled by α . Thus, there is a bijection between the binning locations of X and αX . Since both the binning structure and the data points within each bin are scaled uniformly, the probability of data falling into any given bin remains unchanged. Therefore, the distribution of data across the bins is invariant under scaling, leading to the same binning estimator $\hat{I}_{bin}^n(X; T) = \hat{I}_{bin}^n(\alpha X; T) = \hat{I}_{bin}^n(\alpha X; \alpha T)$ for any scaling factor α . \square

Proposition 2. It holds that $\hat{I}_{KSG}^n(\alpha X; \alpha T) = \hat{I}_{KSG}^n(X; T)$, $\forall \alpha \in \mathbb{R}^+$.

Proof. We note the expression for the KSG estimator (equation 3 from (Kraskov et al., 2004)) as follows:

$$\hat{I}_{KSG}^n(\alpha X; \alpha T) = \psi(k) + \psi(n) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n (\psi(n_{\alpha x, i, \infty}) + \psi(n_{\alpha t, i, \infty})) \quad (10)$$

Here, ψ denotes the digamma function (Abramowitz, 1974), and $n_{\alpha x, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha X_i - \alpha X_j\|_\infty \leq \rho_{k, i, \infty}\}$, where $\rho_{k, i, \infty}$ is the k-NN distance of the joint sample i , $\{\alpha X, \alpha T\}$ (this distance is computed in $d + m$ dimensions). Furthermore, $\|\alpha X_i - \alpha X_j\|_\infty$ represents the X -dimensions only distance (i.e. in d dimensional space). Let $\rho'_{k, i, \infty}$ be the k-NN distance of the joint sample i for the unscaled variables $\{X, T\}$. It is trivial to see that $\rho_{k, i, \infty} = \alpha \rho'_{k, i, \infty}$. Thus, $n_{\alpha x, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha X_i - \alpha X_j\|_\infty \leq \alpha \rho'_{k, i, \infty}\} = \sum_{j \neq i} \mathbb{I}\{\|X_i - X_j\|_\infty \leq \rho'_{k, i, \infty}\} = n_{x, i, \infty}$, and similarly $n_{\alpha t, i, \infty} = n_{t, i, \infty}$. This shows that $\hat{I}_{KSG}^n(\alpha X; \alpha T) = \hat{I}_{KSG}^n(X; T)$. \square

Proposition 3. It holds that $\lim_{\alpha, n \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = -\frac{1}{k}$ and $\lim_{\alpha \rightarrow 0^+, n \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = -\frac{1}{k}$, where k is the k -nearest neighbor parameter for the estimator. Thus, $\hat{I}_{KSG}^n(X; \alpha T)$ need not be equal to $\hat{I}_{KSG}^n(X; T)$.

Proof. Following from the proof of Proposition 2, we note that as $\alpha \rightarrow 0$, we first show $n_{\alpha T, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha T_i - \alpha T_j\|_\infty \leq \rho_{k, i, \infty}\} \rightarrow n$. First, distances in the joint space $(X, \alpha T)$ can be expressed as $\|(X_i, \alpha T_i) - (X_j, \alpha T_j)\|_\infty$. Thus, as $\alpha \rightarrow 0$, $\rho_{k, i, \infty}$ becomes the k -nearest neighbor distance in the X -space only. Next, because X and T are bounded, and when $\alpha \rightarrow 0$, the distance $\|\alpha T_i - \alpha T_j\|_\infty \rightarrow 0$. Thus, $n_{\alpha T, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha T_i - \alpha T_j\|_\infty \leq \rho_{k, i, \infty}\} \rightarrow n$, as all points are essentially at a zero T -only distance between each other in T -space. Similarly, $n_{x, i, \infty} = k$ in this case, as the nearest neighbor distance in the joint $(X, \alpha T)$ space becomes X -only distance, and $n_{x, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|X_i - X_j\|_\infty \leq \rho_{k, i, \infty}\} = k$.

Recalling the expression for KSG, we have:

$$\hat{I}_{KSG}^n(X; \alpha T) = \psi(k) + \psi(n) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n (\psi(n_{x, i, \infty}) + \psi(n_{\alpha T, i, \infty})) \quad (11)$$

Thus we then have: $\lim_{\alpha, n \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = \psi(k) + \psi(n) - \frac{1}{k} + \frac{1}{n} \sum_{i=1}^n (\psi(k) + \psi(n)) = -\frac{1}{k}$.

Lastly, as KSG is global scale-invariant (Proposition 2), we have that $\lim_{\alpha \rightarrow 0, n \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = \lim_{\alpha \rightarrow 0, n \rightarrow \infty} \hat{I}_{KSG}^n(\frac{1}{\alpha} X; T) = \lim_{\alpha, n \rightarrow \infty} \hat{I}_{KSG}^n(\alpha X; T) = -\frac{1}{k}$. The final result follows from the fact that $\hat{I}_{KSG}^n(X; T) = \hat{I}_{KSG}^n(T; X)$. \square

Proposition 4. It holds that $\hat{I}_{MINE-opt}^n(X; \alpha T) = \hat{I}_{MINE-opt}^n(X; T) \forall \alpha \in \mathbb{R}^+$.

Proof. To demonstrate that $\hat{I}_{MINE-opt}^n(X; \alpha T) = \hat{I}_{MINE-opt}^n(X; T)$, we begin by considering any neural network function f that yields a specific value for the expression $\mathbb{E}_{X, Y \sim P(X, T)} [f(X, T)] - \mathbb{E}_{X, T \sim P(X) \times P(T)} [e^{f(X, T)}]$, there exists a corresponding neural network function f' such that $\mathbb{E}_{X, \alpha T \sim P(X, \alpha T)} [f'(X, \alpha T)] - \mathbb{E}_{X, \alpha T \sim P(X) \times P(\alpha T)} [e^{f'(X, \alpha T)}]$ has the same value of the expression involving f and vice-versa. To construct f' , let W_T be the weights of the first layer of the network f that are attached to T , and similarly W'_T for f' . Define a new network function f' with the same architecture as f except that $W'_T = W_T / \alpha$. By construction, the function f' satisfies $f'(X, \alpha T) = f(X, T)$, which implies that for every function f that optimizes the expression in equation 7 there is a corresponding function f' for the variables X and αT . This also shows that the optimization for $I(X; T)$ and $I(X; \alpha T)$ as expressed in equation 7 is equivalent. As a result, the mutual information estimator $\hat{I}_{MINE-opt}^n$, which corresponds to the supremum of the value of this expression over all possible neural network functions, is invariant under scaling of T . Therefore, we conclude that $\hat{I}_{MINE-opt}^n(X; \alpha T) = \hat{I}_{MINE-opt}^n(X; T)$. \square

Proposition 5. Consider the MINE optimization problem with input data $S = \{(\alpha X_1, Y_1), \dots, (\alpha X_n, Y_n)\}$ where $X \in \mathbb{R}^{d_x}$, $Y \in \mathbb{R}^{d_y}$, $(X, Y) \sim P(X, Y)$ are bounded RVs and $\alpha \in \mathbb{R}^+$ is a scaling factor. We consider a neural network of depth $d_n + 1$ having h_1, h_2, \dots, h_{d_n} ReLU-activated hidden neurons in the respective layers. The network is trained via gradient descent on the MINE loss function in equation 8 for a finite number of epochs n_e . Let the *trained* weights between the j^{th} node of the $l + 1^{th}$ hidden layer and the i^{th} node of the l^{th} hidden layer be denoted by $w_{ji}^l \in \mathbb{R}^d$. We consider the case where the initialized weights are very close to zero but not exactly zero (to allow unsymmetrical learning). We assume that the network weights are bounded, such that every weight $|w_{ji}^l| \leq B$ for some $B \in \mathbb{R}$. Let $\eta(t)$ denote the learning rate used at epoch t . Then we have, $\forall i, j$,

$$\lim_{\alpha \rightarrow 0^+} |w_{ji}^1| = 0 \quad (12)$$

Proof. Let X be represented in terms of its individual dimension RVs as $X = [x_1, x_2, \dots, x_{d_x}]$. The weight update rule for w_{ji}^1 , at epoch t is then

$$\Delta w_{ji}^1 = -\eta(t) \alpha x_i \delta_j^1(t). \quad (13)$$

Here $\delta_j^1(t)$ is the backprop error signal at the j^{th} node of the second layer. Let $a(\cdot)$ denote the ReLU activation function, and a' its derivative. We then note the following chain rule for the error signal:

$$\delta_j^{l-1}(t) = \sum_{i=1}^{h_l} \delta_j^l(t) w_{ji}^{l-1} a'(z_j^{l-1}(t)), \quad (14)$$

where $z_j^{l-1}(t)$ denotes the output of the j^{th} node of the $l-1^{th}$ layer itself. As the weights are bounded, and the ReLU derivative is bounded by 1, we can write:

$$\delta_j^{l-1}(t) \leq \sum_{i=1}^{h_l} \delta_j^l(t) B, \quad (15)$$

Let us assume the whole network function can be denoted as $f_W(X, Y) : \mathbb{R}^{d_x+d_y} \rightarrow \mathbb{R}$. Note that the network outputs a single real number, and the last layer does not have any activation function (which is ReLU for the other layers). In the context of MINE's optimization problem, the network minimizes the following loss function:

$$\hat{I}_{\text{MINE}}(X; Y) = -\frac{1}{n} \sum_{i=1}^n f_W(\alpha X_i, Y_i) + \log \left(\frac{1}{n} \sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)} \right), \quad (16)$$

The error signal at the last layer, $\delta^{d_n+1}(t)$ is the derivative of the loss w.r.t the network output. However, as the MINE optimization effectively has two distributions $P(X, Y)$ and $P(X)P(Y)$ of input, we consider the error signal of these distributions separately. The loss function for a datapoint $X_j, Y_j \sim P(X, Y)$ is $-\frac{1}{n} f_W(\alpha X_j, Y_j)$, which yields an error signal $\delta^{d_n+1}(t) = -1/n$. The loss function for a datapoint $X_j, \tilde{Y}_j \sim P(X)P(Y)$ is $\log \left(\frac{1}{n} \sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)} \right)$, which yields an error signal

$$\delta^{d_n+1}(t) = \frac{d \left(\log \left(\frac{1}{n} \sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)} \right) \right)}{df_W(\alpha X_j, \tilde{Y}_j)} = \frac{e^{f_W(\alpha X_j, \tilde{Y}_j)}}{\sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)}} \leq 1 \quad (17)$$

Thus, across both cases, we have that the error signal $|\delta^{d_n+1}(t)| \leq 1$.

This observation coupled with applying the result in equation 15 yields:

$$|\delta_j^{l-1}(t)| \leq B^{d_n} \prod_{i=2}^{d_n} h_i \quad (18)$$

Next, as X and Y are bounded, we can assume that every dimension of $|x_i| \leq K$ for some $K \in \mathbb{R}^+$. Let us set $C = KB^{d_n} \prod_{i=2}^{d_n} h_i$. With this, we have that $|\Delta w_{ji}^1| \leq \eta(t) \alpha C$. Let the number of mini-batch updates per epoch be b_u . With this, the total change in $|\Delta w_{ji}^1|$ can be bounded as:

$$|\Delta_{\text{total}} w_{ji}^1| \leq \alpha \left(\sum_{t=1}^{n_e} \eta(t) b_u C \right) \quad (19)$$

The trained weight $|w_{ji}^1| \leq \epsilon + \alpha (\sum_{t=1}^{n_e} \eta(t) b_u C)$ where $\epsilon \rightarrow 0$ is the initialized value of the corresponding weight. Thus, when $\alpha \rightarrow 0^+$, we have that $\epsilon + \alpha (\sum_{t=1}^{n_e} \eta(t) b_u C) \rightarrow 0$, yielding $\lim_{\alpha \rightarrow 0^+} |w_{ji}^1| = 0$. \square

Proposition 6. We consider the same setting as Proposition 5 for the MINE estimation problem. There, it holds that $\lim_{\alpha \rightarrow 0} \hat{I}_{\text{MINE-sgd}}^n(X; \alpha T) = 0$. Thus, $\hat{I}_{\text{MINE-sgd}}^n(X; \alpha T)$ need not be equal to $\hat{I}_{\text{MINE-sgd}}^n(X; T)$.

Proof. Let f^* represent the neural network function which optimizes the expression for $\hat{I}_{MINE-sgd}^n(X; \alpha T)$, which is $\mathbb{E}_{X, \alpha T \sim P(X, \alpha T)} [f'(X, \alpha T)] - \mathbb{E}_{X, \alpha T \sim P(X) \times P(\alpha T)} [e^{f'(X, \alpha T)}]$, via SGD. Let W_T be the weights of the first layer of the network f^* which are attached to αT . From Proposition 5, we have that $\lim_{\alpha \rightarrow 0^+} |w_{ji}^1| = 0$. Thus, as $\alpha \rightarrow 0$, the weights $W_T \rightarrow 0$ as well. This indicates that the contribution of T to the function f^* , as $\alpha \rightarrow 0$, will be negligible. Thus effectively, $\lim_{\alpha \rightarrow 0} \hat{I}_{MINE-sgd}^n(X; \alpha T) = \hat{I}_{MINE-sgd}^n(X; 0) = 0$. This holds mainly because when $W_T \rightarrow 0$, the network essentially ignores the changes in αT , and thus αT is essentially treated as a constant variable. This proves the result. \square