

Towards Robust Scale-Invariant Mutual Information Estimators

Anonymous authors
Paper under double-blind review

Appendix K

Proofs of theoretical results

Proposition 1. Let $\hat{I}_{\text{bin}}^n(X; T)$ denote the mutual information estimated using a fixed number of bins per dimension, with bin edges defined by the minimum and maximum values of the data. Then, for all $\alpha \in \mathbb{R}^+$,

$$\hat{I}_{\text{bin}}^n(\alpha X; \alpha T) = \hat{I}_{\text{bin}}^n(X; T) \quad \& \quad \hat{I}_{\text{bin}}^n(X; \alpha T) = \hat{I}_{\text{bin}}^n(X; T).$$

Proof. For $X \in \mathbb{R}^d$, define $X_{\min}, X_{\max} \in \mathbb{R}^d$ as the componentwise extrema. Under scaling by $\alpha > 0$, we have:

$$\min(\alpha X) = \alpha X_{\min} \quad \& \quad \max(\alpha X) = \alpha X_{\max},$$

so the bin edges scale by α , and the number of bins remains fixed. Thus, each bin in X corresponds bijectively to a bin in αX with identical counts. Therefore, the empirical distribution over bins, and hence \hat{I}_{bin}^n , is invariant under such scaling. \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. Let ψ denotes the digamma function (Abramowitz, 1974), and

$$n_{\alpha x, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha X_i - \alpha X_j\|_{\infty} \leq \rho_{k, i, \infty}\}, \quad (1)$$

where $\rho_{k, i, \infty}$ is the sup-norm distance of $\{\alpha X_i, \alpha T_i\}$ to its k -nearest neighbor in the joint $\{\alpha X, \alpha T\}$ space. Here $\|\alpha X_i - \alpha X_j\|_{\infty}$ represents the X -dimensions only distance.

Then, the KSG estimate of mutual information (equation 3 from (Kraskov et al., 2004)) is

$$\hat{I}_{KSG}^n(\alpha X; \alpha T) = \psi(k) + \psi(n) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n (\psi(n_{\alpha x, i, \infty}) + \psi(n_{\alpha t, i, \infty})). \quad (2)$$

Let $\rho'_{k, i, \infty}$ denote the sup-norm distance of $\{X_i, T_i\}$ to its k -nearest neighbor in the joint $\{X, T\}$ space for the unscaled variables. It is trivial to see that $\rho_{k, i, \infty} = \alpha \rho'_{k, i, \infty}$.

We then have,

$$n_{\alpha x, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha X_i - \alpha X_j\|_{\infty} \leq \alpha \rho'_{k, i, \infty}\} = \sum_{j \neq i} \mathbb{I}\{\|X_i - X_j\|_{\infty} \leq \rho'_{k, i, \infty}\} = n_{x, i, \infty}. \quad (3)$$

One can similarly show that $n_{\alpha t, i, \infty} = n_{t, i, \infty}$.

Thus, $\psi(n_{\alpha x, i, \infty}) = \psi(n_{x, i, \infty})$ and $\psi(n_{\alpha t, i, \infty}) = \psi(n_{t, i, \infty})$, and it follows that $\hat{I}_{KSG}^n(\alpha X; \alpha T) = \hat{I}_{KSG}^n(X; T)$. \square

Proposition 3. Let $\{(X_i, T_i)\}_{i=1}^n$ be drawn i.i.d. from a bounded distribution on $\mathbb{R}^{d_X} \times \mathbb{R}^{d_T}$ that is absolutely continuous. Then, it holds almost surely that

$$\lim_{\alpha \rightarrow 0^+} \hat{I}_{KSG}^n(X; \alpha T) = \lim_{\alpha \rightarrow \infty} \hat{I}_{KSG}^n(X; \alpha T) = -\frac{1}{k}.$$

Thus, $\hat{I}_{KSG}^n(X; \alpha T)$ need not be equal to $\hat{I}_{KSG}^n(X; T)$.

Proof. Following the proof of Proposition 2, let $\rho_{k,i,\infty}$ denote the sup-norm distance from $(X_i, \alpha T_i)$ to its k -th nearest neighbor in the joint space $(X, \alpha T)$. Define

$$n_{\alpha T, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha T_i - \alpha T_j\|_\infty \leq \rho_{k,i,\infty}\}, \quad (4)$$

and similarly for $n_{x,i,\infty}$. The KSG estimator's estimate here will be:

$$\hat{I}_{KSG}^n(X; \alpha T) = \psi(k) + \psi(n) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n (\psi(n_{x,i,\infty}) + \psi(n_{\alpha T, i, \infty})). \quad (5)$$

First, we observe that

$$\lim_{\alpha \rightarrow 0^+} \|(X_i, \alpha T_i) - (X_j, \alpha T_j)\|_\infty = \lim_{\alpha \rightarrow 0^+} \max\{\|X_i - X_j\|_\infty, \alpha \|T_i - T_j\|_\infty\} = \|X_i - X_j\|_\infty. \quad (6)$$

Furthermore, as X and T are absolutely continuous, it holds almost surely that the k -nearest neighbor distance in X space will be greater than zero and finite.

In the limit $\alpha \rightarrow 0^+$, all pairwise distances $\|\alpha T_i - \alpha T_j\|_\infty \rightarrow 0$, and thus $n_{\alpha T, i, \infty} = \sum_{j \neq i} \mathbb{I}\{\|\alpha T_i - \alpha T_j\|_\infty \leq \rho_{k,i,\infty}\} \rightarrow n$.

Next, from equation 6 it follows that $\rho_{k,i,\infty}$ becomes the k -nearest neighbor distance in X space as $\alpha \rightarrow 0^+$, and thus we can write

$$\lim_{\alpha \rightarrow 0^+} n_{x,i,\infty} = \sum_{j \neq i} \mathbb{I}\{\|X_i - X_j\|_\infty \leq \rho_{k,i,\infty}\} = k. \quad (7)$$

This enables us to replace the terms in equation 5 yielding

$$\lim_{\alpha \rightarrow 0^+} \psi(k) + \psi(n) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n (\psi(k) + \psi(n)) = -\frac{1}{k}. \quad (8)$$

Lastly, by global scale invariance (Proposition 2), $\hat{I}_{KSG}^n(X; \alpha T) = \hat{I}_{KSG}^n(\frac{1}{\alpha} X; T)$, and thus,

$$\lim_{\alpha \rightarrow 0} \hat{I}_{KSG}^n(X; \alpha T) = \lim_{\alpha \rightarrow 0} \hat{I}_{KSG}^n(\frac{1}{\alpha} X; T) = \lim_{\beta \rightarrow \infty} \hat{I}_{KSG}^n(\beta X; T) = -\frac{1}{k}. \quad (9)$$

□

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. Let f be any neural network function used in the MINE objective:

$$\mathbb{E}_{(X,T) \sim P(X,T)}[f(X, T)] - \log \mathbb{E}_{(X,T) \sim P(X) \times P(T)}[e^{f(X,T)}]. \quad (10)$$

Define a corresponding function f' for the scaled variable αT by setting the first-layer weights on T in f' to $W'_T = W_T/\alpha$, and keeping all other parameters identical. Then $f'(X, \alpha T) = f(X, T)$, so both expectations in the MINE objective remain unchanged after the transformation.

As this holds for any f , the supremum over all such functions is preserved under scaling of T . Hence, $\hat{I}_{MINE-opt}^n(X; \alpha T) = \hat{I}_{MINE-opt}^n(X; T)$ for all $\alpha > 0$. \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 $L + 1$ having h_1, h_2, \dots, h_L ReLU-activated hidden neurons in the respective layers. The network is trained via gradient descent on the MINE loss function in Belghazi et al. (2018) for a fixed number of iterations T , with a learning rate schedule $\eta(t) < \infty$ for all $t \leq T$. Let the weights between the i^{th} node of the $l + 1^{th}$ hidden layer and the j^{th} node of the l^{th} hidden layer after t iterations be denoted by $w_{ji}^l(t)$. Assume that the initialized weights are bounded, i.e., $\forall (l, i, j)$, $|w_{ji}^l(0)| \leq \epsilon$ for some $\epsilon > 0$. Lastly, let $w_{ji}^1[X]$ denote the first layer weights that are attached to X . We then have, $\forall i, j$,

$$\lim_{\alpha \rightarrow 0^+} |w_{ji}^1[X](T)| \leq \epsilon. \quad (11)$$

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 $w_{ji}^1(X)$, at iteration t is then $w_{ji}^1[X](t + 1) = w_{ji}^1[X](t) + \Delta w_{ji}^1[X](t)$, where

$$\Delta w_{ji}^1[X](t) = -\eta(t) \alpha x_i \delta_j^1(t). \quad (12)$$

Here $\delta_j^1(t)$ is the backprop error signal at the j^{th} node of the second layer at iteration t . Let $a(\cdot)$ be the ReLU activation with its derivative $a'(\cdot) \in \{0, 1\}$. Then we have,

$$\delta_j^{l-1}(t) = \sum_{i=1}^{h_l} \delta_j^l(t) w_{ji}^{l-1}(t) a'(z_j^{l-1}(t)), \quad (13)$$

where $z_j^{l-1}(t)$ denotes the output of the j^{th} node of the $(l - 1)^{th}$ layer itself. We note that as $|a'| \leq 1$, this yields:

$$|\delta_j^{l-1}(t)| \leq \sum_{i=1}^{h_l} |\delta_j^l(t) w_{ji}^{l-1}(t)|. \quad (14)$$

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:

$$\hat{I}_{MINE}(X; Y) = -\frac{1}{n} \sum_{i=1}^n f_W(\alpha X_i, Y_i) + \log \left(\frac{1}{n} \sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)} \right). \quad (15)$$

The error signal at the last layer, $\delta^{L+1}(t)$ 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 the input $X_j, Y_j \sim P(X, Y)$ is $-\frac{1}{n} f_W(\alpha X_j, Y_j)$, which yields an error signal $\delta^{L+1}(t) = -1/n$. Whereas, the loss function for the input $X_j, \tilde{Y}_j \sim P(X)P(Y)$ 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

$$\delta^{L+1}(t) = \frac{d \left(\log \left(\frac{1}{n} \sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)} \right) \right)}{df_W(\alpha X_j, \tilde{Y}_j)} = \frac{e^{f_W(\alpha X_j, \tilde{Y}_j)}}{\sum_{i=1}^n e^{f_W(\alpha X_i, \tilde{Y}_i)}} \leq 1. \quad (16)$$

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

Now, we note that 12 implies that if the error signal $\delta_j^1(t)$ and the layer outputs are finite, then the network weights, after a finite number of iterations, will stay finite. As the input to the network is bounded, finite weights imply finite layer outputs. Thus, we can state that finite $\delta_j^1(t)$ ultimately implies finite weights. To that end, we now prove that the error signal $\delta_j^1(t)$ will be finite for all finite t using the principle of induction, as follows.

First, let us assume that the weights at iteration t $w_{ji}^{l-1}(t)$ are finite (i.e. $|w_{ji}^{l-1}(t)| \leq B_t$ for some finite B_t). Next, as $|\delta^{L+1}(t)| \leq 1$, using 17, we immediately see that $\delta_j^1(t)$ will be finite. Thus, the weights $w_{ji}^{l-1}(t+1) = w_{ji}^{l-1}(t) + \Delta w_{ji}^1(t)$ will be finite as well, i.e., $|w_{ji}^{l-1}(t)| \leq B_{t+1}$ for some finite B_{t+1} .

As the weights at iteration 0 are finite ($|w_{ji}^l(0)| \leq \epsilon$), via the principle of induction we have that for any finite t , $|w_{ji}^{l-1}(t)| \leq B_t$ for some finite B_t . Let us consider $B = \max_{t \in \{0,1,\dots,T\}} B_t$. Then for the entire duration of training, we can write $|w_{ji}^{l-1}(t)| \leq B$.

Lastly, we note that B will be a function of α as it is the only parameter that is subject to change in the problem setting. Thus we represent it as $B(\alpha)$. However, we note that $B(\alpha)$ will remain finite for all finite values of α as it ensures the inputs and the initial weights are finite and the prior proof by induction follows.

Thus, plugging this into 17, we obtain:

$$|\delta_j^{l-1}(t)| \leq \sum_{i=1}^{h_l} B(\alpha) |\delta_j^l(t)|. \quad (17)$$

Subsequently, following the fact that $|\delta^{L+1}(t)| \leq 1$, we get

$$|\delta_j^1(t)| \leq (B(\alpha))^L \prod_{i=2}^L h_i. \quad (18)$$

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 $C(\alpha) = KB(\alpha)^L \prod_{i=2}^L h_i$. With this, we have that $|\Delta w_{ji}^1[X](t)| \leq \eta(t)\alpha C(\alpha)$. This yields

$$|w_{ji}^1[X](T)| \leq |w_{ji}^1[X](0)| + \alpha \left(\sum_{t=1}^T |\eta(t)C(\alpha)| \right) \quad (19)$$

$$\leq \epsilon + \alpha \left(\sum_{t=1}^T |\eta(t)C(\alpha)| \right). \quad (20)$$

Thus, when $\alpha \rightarrow 0^+$, we have that $\epsilon + \alpha \left(\sum_{t=1}^T |\eta(t)C(\alpha)| \right) = \epsilon$ as $C(\alpha)$ and $\eta(t)$ are finite, yielding $\lim_{\alpha \rightarrow 0^+} |w_{ji}^1[X](T)| \leq \epsilon$. \square

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

Proof. Let f^* denote the neural network function learned via stochastic gradient descent (SGD) that maximizes the following MINE objective for $\hat{I}_{MINE-sgd}^n(X; \alpha T)$:

$$\mathbb{E}_{X, \alpha T \sim P(X, \alpha T)} [f^*(X, \alpha T)] - \log \mathbb{E}_{X, \alpha T \sim P(X) \times P(\alpha T)} \left[e^{f^*(X, \alpha T)} \right].$$

Let W_T be the first-layer weights in f^* connected to αT . From Proposition 5, as $\alpha \rightarrow 0^+$, we have $|W_T|_{ij} \leq \epsilon$, thus yielding finite weights. Hence, the contribution of αT to the output of a neuron of the first hidden layer of f^* becomes $\alpha T \epsilon d_T \rightarrow 0$ as $\alpha \rightarrow 0^+$, where d_T is the dimensionality of T . Therefore, $f^*(X, \alpha T) = f^*(X)$ becomes independent of T in the limit $\alpha \rightarrow 0$.

In this case, both expectations in the above objective reduce to empirical averages of the same function $f^*(X)$, and Jensen's inequality implies

$$\frac{1}{n} \sum f^*(X_i) - \log \left(\frac{1}{n} \sum e^{f^*(X_i)} \right) \leq 0,$$

with equality if and only if $f^*(X)$ is constant. Thus, the maximized objective converges to zero.

Hence, $\lim_{\alpha \rightarrow 0^+} \hat{I}_{MINE-sgd}^n(X; \alpha T) = 0$, showing that $\hat{I}_{MINE-sgd}^n(X; \alpha T)$ need not equal $\hat{I}_{MINE-sgd}^n(X; T)$. \square

References

- Milton Abramowitz. *Handbook of Mathematical Functions, With Formulas, Graphs, and Mathematical Tables*. Dover Publications, Inc., USA, 1974. ISBN 0-486-61272-4.
- 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.
- Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical Review E*, 69(6):066138, June 2004. Publisher: American Physical Society.