

# 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  $\alpha$ , and the number of bins remains fixed. Thus, each bin in  $X$  corresponds bijectively to a bin in  $\alpha X$  with identical counts. Therefore, the empirical distribution over bins, and hence  $\hat{I}_{\text{bin}}^n$ , 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  $\psi$  denote 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  $\alpha T$  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  $w_{ji}^0[X]$  denote the first layer weights that are attached to  $X$ . We then have,  $\forall (i, j)$ ,

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

*Proof.* Let  $Z_j^l(\alpha x, y; t)$  denote the output of the  $j^{th}$  node at layer  $l$  in response to  $(\alpha x, y)$  as the input at iteration  $t$ , where  $(x, y)$  is a sampled instance of  $P(X, Y)$ . Note that when  $l = 0$ ,  $Z_j^l$  will be the inputs to the network  $(\alpha x, y)$  itself. As we are given the dataset  $S$ , the weight update rule for  $w_{ji}^l$ , at iteration  $t$  is then  $w_{ji}^l(t + 1) = w_{ji}^l(t) + \Delta w_{ji}^l(t)$ , where

$$\Delta w_{ji}^l(t) = -\eta(t) \frac{1}{n} \sum_{k=1}^n Z_j^l(\alpha X_k, Y_k; t) \delta_i^{l+1}(\alpha X_k, Y_k; t), \quad (12)$$

Here,  $\delta_i^{l+1}(\alpha X_k, Y_k; t)$  denotes the backpropagation error signal at the  $i^{th}$  neuron of the  $(l + 1)^{th}$  layer at iteration  $t$ , computed in response to the input pair  $(\alpha X_k, Y_k)$ . For notational simplicity, we omit the input instance arguments in the following expressions, and write the error signal as  $\delta_i^{l+1}(t)$ . Let  $a(\cdot)$  be the ReLU activation with its derivative  $a'(\cdot) \in \{0, 1\}$ . Then we have,

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

where  $z_j^{l-1}(t)$  denotes the pre-activation output of the  $j^{th}$  node of the  $(l - 1)^{th}$  layer itself. Note that  $Z_j^l(t) = a(z_j^l(t))$  for  $l \geq 1$ . As  $|a'| \leq 1$ , this yields:

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

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  $|\delta^{L+1}(t)| \leq 1$ .

Next, we consider the case when  $|\alpha| \leq A$ , where  $A \in \mathbb{R}^+$  is a finite real. Within this setting, using the principle of induction, we now will show that the error signal  $|\delta_j^l(t)| \leq D_t \forall (l, j)$  and the weights  $|w_{ji}^l(t)| \leq B_t \forall (l, i, j)$  for some finite reals  $D_t, B_t$ .

First, let us assume weights at iteration  $t$  satisfy  $|w_{ji}^l(t)| \leq B_t \forall (l, i, j)$  for some  $B_t \in \mathbb{R}^+$ . Then, it first follows from equation 14 that

$$|\delta_j^{l-1}(t)| \leq \sum_{i=1}^{h_l} B_t |\delta_i^l(t)|, \quad (17)$$

and as  $|\delta^{L+1}(t)| \leq 1$ , this ultimately yields

$$|\delta_j^l(t)| \leq (B_t)^{L+1-l} \prod_{v=l+1}^L h_v. \quad (18)$$

If we set  $D_t = \max_l \{(B_t)^{L+1-l} \prod_{v=l+1}^L h_v\}$ , we have that  $|\delta_j^l(t)| \leq D_t \forall (l, j)$ . Note that here  $D_t$  only depends on  $B_t$  and the network architecture. Following the weight update rule at iteration  $t$  from equation 12, we have:

$$\begin{aligned} |w_{ji}^l(t+1)| &= \left| w_{ji}^l(t) - \frac{\eta(t)}{n} \sum_{k=1}^n Z_j^l(\alpha X_k, Y_k; t) \delta_i^{l+1}(t) \right|, \\ &\leq |w_{ji}^l(t)| + \frac{\eta(t)}{n} \sum_{k=1}^n |Z_j^l(\alpha X_k, Y_k; t)| |\delta_i^{l+1}(t)|, \\ &\leq B_t + \eta(t) D_t \left( \frac{1}{n} \sum_{k=1}^n |Z_j^l(\alpha X_k, Y_k; t)| \right). \end{aligned} \quad (19)$$

Next, as  $X$  and  $Y$  are bounded, we can assume that every dimension  $|x_i| \leq K$  and  $|y_i| \leq K$  for some  $K \in \mathbb{R}^+$ . Now, as  $|\alpha| \leq A$ , the inputs  $\alpha X$  and  $Y$  are bounded. Furthermore, as the weights  $|w_{ji}^l(t)| \leq B_t$ , the outputs of the layers  $|Z_j^l(\alpha X_k, Y_k; t)|$  will be bounded by a finite scalar which can be written as a fixed function of  $B_t, K$  and the hidden layer counts  $h_1, \dots, h_L$ . Thus,

$$|w_{ji}^l(t+1)| \leq B_t + \eta(t) D_t g(B_t, K, h_1, \dots, h_L). \quad (20)$$

We can set  $B_{t+1} = B_t + \eta(t) D_t g(B_t, K, h_1, \dots, h_L)$ , and the above then implies that  $|w_{ji}^l(t+1)| \leq B_{t+1} \forall (l, i, j)$ .

Lastly, 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  $t$ ,  $|w_{ji}^l(t)| \leq B_t$  for some finite  $B_t$ , when  $|\alpha| \leq A$ . Let us consider  $B = \max_{t \in \{0, 1, \dots, T\}} B_t$ . Then for the entire duration of training and for any  $|\alpha| \leq A$ , we have  $|w_{ji}^l(t)| \leq B$ .

Thus, plugging this into equation 14, we obtain,  $\forall (l, j, t)$ ,

$$|\delta_j^{l-1}(t)| \leq \sum_{i=1}^{h_l} B |\delta_i^l(t)|. \quad (21)$$

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

$$|\delta_i^1(t)| \leq B^L \prod_{v=2}^L h_v. \quad (22)$$

Let  $C = KB^L \prod_{v=2}^L h_v$ . With this, following equation 12, we have that  $|\Delta w_{ji}^0[X](t)| \leq \eta(t)\alpha C$ . This yields

$$|w_{ji}^0[X](T)| \leq |w_{ji}^0[X](0)| + \alpha \left( \sum_{t=1}^T |\eta(t)C| \right), \quad (23)$$

$$\leq \epsilon + \alpha C \left( \sum_{t=1}^T \eta(t) \right). \quad (24)$$

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

**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  $\alpha T$ . From Proposition 5, as  $\alpha \rightarrow 0^+$ , we have  $|W_T|_{ij} \leq \epsilon$ , thus yielding finite weights. Hence, the contribution of  $\alpha T$  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  $\hat{I}_{MINE-sgd}^n(X; T)$ .  $\square$

**Proposition 7.** 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  $\alpha T$ . From Proposition 5, as  $\alpha \rightarrow 0^+$ , we have  $|W_T|_{ij} \leq \epsilon$ , thus yielding finite weights. Hence, the contribution of  $\alpha T$  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.