

---

# Estimating Total Correlation with Mutual Information Bounds

---

Pengyu Cheng, Weituo Hao, Lawrence Carin  
Department of Electrical and Computer Engineering  
Duke University  
pengyu.cheng@duke.edu

## Abstract

Total correlation (TC) is a fundamental concept in information theory to measure the statistical dependency of multiple random variables. Recently, TC has shown effectiveness as a regularizer in many machine learning tasks when minimizing/maximizing the correlation among random variables is required. However, to obtain precise TC values is challenging, especially when the closed-form distributions of variables are unknown. In this paper, we introduced several sample-based variational TC estimators. Specifically, we connect the TC with mutual information (MI) and constructed two calculation paths to decompose TC into MI terms. In our experiments, we estimated the true TC values with the proposed estimators in different simulation scenarios and analyzed the properties of the TC estimators.

## 1 Introduction

Statistical dependency measures the correlation of random variables or factors in models, which is often an important concern in various scientific domains including statistics [12, 15], robotics [16, 4], bioinformatics [18, 24], and machine learning [7, 1, 14]. In recent deep learning studies, statistical dependency has increasingly served as learning objectives or regularizers for neural network training, and has achieved improvement in terms of model robustness [25], generalizability [1], interpretability [7, 9], *etc.*

Among statistical dependency measurements, mutual information (MI) is commonly used in machine learning. Given two random variables  $\mathbf{x}, \mathbf{y}$ , the mutual information is defined as:

$$\mathcal{I}(\mathbf{x}; \mathbf{y}) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} \right]. \quad (1)$$

Recently, mutual information has shown significant improvement when applied as a training criterion on learning tasks, such as conditional generation [7], domain adaptation [11], representation learning [6], and fairness [23]. However, MI can only handle the statistical dependency between two variables. When considering optimization of correlation among multiple variables, MI requires computation of each variable pair, which leads to a quadratic increase in computation cost. To address this problem, total correlation (TC) has been proposed by extending MI to multi-variable cases:

$$\mathcal{TC}(\mathbf{X}) = \mathcal{TC}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \mathbb{E}_{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)} \left[ \log \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)}{p(\mathbf{x}_1)p(\mathbf{x}_2) \dots p(\mathbf{x}_n)} \right]. \quad (2)$$

TC has also proven effective to enhance machine learning models in many tasks, such as independent component analysis [3], and disentangled representation learning [5, 19, 17]. However, TC suffers from the same numerical problem as MI: the exact values of TC are difficult to calculate without the closed-form distribution  $p(\mathbf{x}_i)$  and with only samples accessible. Previous works on disentangled representation learning [5, 10] avoid the estimation problem by assuming that both the latent priors and the inference posteriors follow multi-variate Gaussian distributions. Poole, *et al.* [22] proposed an upper bound of TC by further introducing another variable  $\mathbf{y}$ . With a strong assumption that given

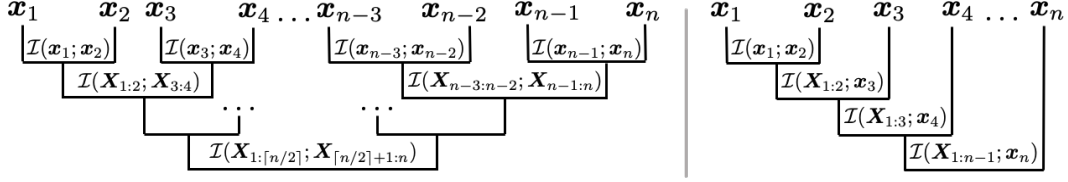


Figure 1: Two calculation paths of total correlation. **Left** (Tree-like calculation path): Divide the current variables into two subgroups with similar sizes. Calculate the MI between the subgroups and recursively calculate TC of both subgroups.  $\lceil n/2 \rceil$  is the smallest number larger than  $n/2$ . **Right** (Line-like calculation path): Calculate the MI between the current group of variables and the next variable, and then add the the next variable into current group.

$\mathbf{y}$ , all  $\mathbf{x}_i | \mathbf{y}$  are independent,  $p(\mathbf{X} | \mathbf{y}) = \prod_{i=1}^n p(\mathbf{x}_i | \mathbf{y})$ , Poole, *et al.* [22] concluded that  $\mathcal{TC}(\mathbf{X}) = \sum_{i=1}^n \mathcal{I}(\mathbf{x}_i; \mathbf{y}) - \mathcal{I}(\mathbf{X}; \mathbf{y})$ . All the aforementioned methods require additional assumptions to the distributions, which limits their application scenarios.

In this paper, we propose two TC estimation strategies based on mutual information variational bounds. More specifically, we decompose TC into the summation of MI terms along two different calculation paths: the tree-like path and the line-like path. Then the TC values are approximated by applying MI estimation to each decomposed term. In our experiments, we test the performance of the proposed TC estimators under multivariate Gaussian simulations.

## 2 Method

With the definition of total correlation (TC) and mutual information (MI) in (2) and (1), we find a connection between TC and MI summarized in Theorem 2.1.

**Theorem 2.1.** Suppose  $\mathcal{A} = \{i_1, i_2, \dots, i_m\} \subseteq \{1, 2, \dots, n\}$  is an index subset.  $\bar{\mathcal{A}} = \{j : j \notin \mathcal{A}\}$  is the complementary set of  $\mathcal{A}$ . Denote  $\mathbf{X}_{\mathcal{A}} = (\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_m})$  as the selected variables from  $\mathbf{X}$  with the indexes  $\mathcal{A}$ . Then we have  $\mathcal{TC}(\mathbf{X}) = \mathcal{TC}(\mathbf{X}_{\mathcal{A}}) + \mathcal{TC}(\mathbf{X}_{\bar{\mathcal{A}}}) + \mathcal{I}(\mathbf{X}_{\mathcal{A}}; \mathbf{X}_{\bar{\mathcal{A}}})$ .

**Corollary 2.1.1.** Given a variable group  $\mathbf{X}$  and another  $\mathbf{y}$ ,  $\mathcal{TC}(\mathbf{X} \cup \{\mathbf{y}\}) = \mathcal{TC}(\mathbf{X}) + \mathcal{I}(\mathbf{X}; \mathbf{y})$ .

**Corollary 2.1.2.** Given  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ , we have  $\mathcal{TC}(\mathbf{X}) = \sum_{i=1}^{n-1} \mathcal{I}(\mathbf{X}_{1:i}; \mathbf{x}_{i+1})$ .

The Theorem 2.1 provides insight that the TC of a group of variables  $\mathbf{X}$  can be decomposed into the TC of two subgroups  $\mathbf{X}_{\mathcal{A}}$  and  $\mathbf{X}_{\bar{\mathcal{A}}}$  and the MI between the two subgroups. Therefore, we can recursively represent the TC with MI terms. More specifically, we propose two schemes with different structures to calculate TC with different MI terms (as shown in Figure 1).

Let  $\mathbf{X}_{i:j} = (\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_j)$  denote a subset of variables with indexes from  $i$  to  $j$ . Based on Theorem 2.1, we propose two recursive TC calculation schemes: (1) **Line-like**:  $\mathcal{TC}(\mathbf{X}_{1:i+1}) = \mathcal{TC}(\mathbf{X}_{1:i}) + \mathcal{I}(\mathbf{X}_{1:i}; \mathbf{x}_{i+1})$ ; (2) **Tree-like**:  $\mathcal{TC}(\mathbf{X}_{i:j}) = \mathcal{TC}(\mathbf{X}_{i:\lfloor (i+j)/2 \rfloor}) + \mathcal{TC}(\mathbf{X}_{\lfloor (i+j)/2 \rfloor + 1:j}) + \mathcal{I}(\mathbf{X}_{i:\lfloor (i+j)/2 \rfloor}; \mathbf{X}_{\lfloor (i+j)/2 \rfloor + 1:j})$ , where  $\lfloor t \rfloor$  indicates the largest integer smaller than  $t$ . The line-like dynamic calculates the MI between a subgroup and a single variable, which leads to the representation of TC as the summation in Corollary 2.1.2. The tree-like dynamic divides the variables into balanced subgroups, so that the MI between two subgroups can be calculated with two variable parts in similar dimensions. Since the tree-like estimation is hard to summarize in an equation, we describe it in Algorithm 1. With the total correlation being decomposed into summation of MI terms, we can derive total correlation estimators based on the previous mutual information variational bounds.

---

### Algorithm 1 Tree-like TC estimation algorithm

---

**Prerequisite:** MI estimation method  $\hat{\mathcal{I}}$ , samples  $\{\mathbf{X}^{(i)}\}_{i=1}^M = \{(\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}, \dots, \mathbf{x}_n^{(i)})\}_{i=1}^M$   
**Function**  $\text{TC}_{\text{Tree-estimate}}(\mathbf{X}_{i:j})$ :  
**if**  $j - i \leq 0$  **then**  
    **return** 0  
**else**  
     $m = \lfloor (i + j) / 2 \rfloor$   
    **return**  $\text{TC}_{\text{Tree-estimate}}(\mathbf{X}_{i:m}) + \text{TC}_{\text{Tree-estimate}}(\mathbf{X}_{m+1:j}) + \hat{\mathcal{I}}(\mathbf{X}_{i:m}; \mathbf{X}_{m+1:j})$   
**end if**

---

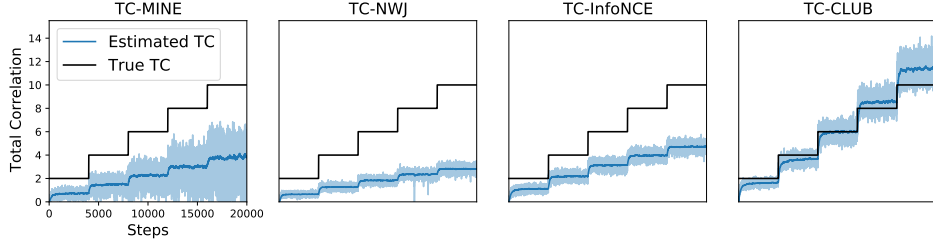


Figure 2: Simulation performance of TC **Line-like** estimators with different MI bounds.

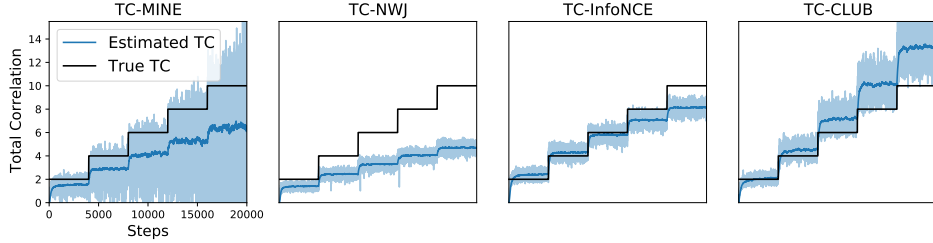


Figure 3: Simulation performance of TC **Tree-like** estimators with different MI bounds.

### 3 Experiments

We derive our TC estimators based on four MI bounds (MINE [2], NWJ [20], InfoNCE [21], and CLUB [8]) as TC-MINE, TC-NWJ, TC-InfoNCE, and TC-CLUB. The detailed description and implementation to the four MI estimators are shown in the Supplementary Material. Then we test the TC estimators with both tree-like and line-like strategies on simulations. The simulation data are drawn from four-dimensional Gaussian distributions  $(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , where  $\Sigma$  is a covariance matrix with all diagonal elements equal to 1. With this Gaussian assumption, the true TC value can be calculated as  $\mathcal{TC}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = -\frac{1}{2} \log \text{Det}(\Sigma)$ , where  $\text{Det}(\Sigma)$  is the determinant of  $\Sigma$ . Therefore, we can adjust the correlation coefficients in  $\Sigma$  to set the ground-truth TC values in the range  $\{2.0, 4.0, 6.0, 8.0, 10.0\}$ . At each TC true value, we sample data batches 4000 times, with batch size equal to 64, for the training of variational TC estimators.

In Figure 2, we report the performance of our TC estimators with different MI bounds at each training steps. In each figure, the true TC value is shown as a step function with black line. The estimation values are displayed among different steps with shadow blue curves. The dark blue curves shows the local averages of estimated TC, with a bandwidth equal to 200. Under both a tree-like and line-like path calculation, the TC-MINE, TC-NWJ and TC-InfoNCE remains a lower bound of the truth TC values, based on the fact that MINE, NWJ, and InfoNCE are lower bound of mutual information. CLUB is an MI upper bound, while the TC-CLUB also behaves as an upper bound of total correlation.

The upper bound method TC-CLUB achieves better performance with line-like calculation. This is because that CLUB requires a variational approximation  $q_\theta(\mathbf{v}|\mathbf{u})$  when estimating  $\mathcal{I}(\mathbf{v}; \mathbf{u})$ . When we use line-like calculation path,  $\mathbf{v} = \mathbf{x}_{i+1}$  is always a single variable, and  $\mathbf{u} = \mathbf{X}_{1:i}$  is the concatenation of  $(\mathbf{x}_1, \dots, \mathbf{x}_i)$ . The  $q_\theta(\mathbf{v}|\mathbf{u})$  as a neural network can have better performance with output  $\mathbf{v}$  in a fixed low dimension. In contrast, the lower bound methods show better estimation with tree-like calculation than line-like calculation. Because for all listed lower bound methods, the estimation of  $\mathcal{I}(\mathbf{v}; \mathbf{u})$  is based on  $\mathbf{v}$  and  $\mathbf{u}$  equally. With the tree-like strategy, each time the MI estimators are provided with samples in similar dimensions, which facilitates the learning of lower bound MI estimators. The bias and variance of the TC estimators are shown in the Supplementary Material.

### 4 Discussion

We have derived the line-like and tree-like calculation strategies to decompose the total correlation into summation of mutual information. By estimating mutual information terms with MI bounds, we introduced several TC estimators. The tree-like and line-like calculation strategies can bring advantages to TC estimation depending on different MI estimation processes. The proposed TC estimators can be further applied as learning criterion on many deep learning tasks, such as disentangled representation learning, ensemble learning, and model distillation.

## References

- [1] Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- [2] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Devon Hjelm, and Aaron Courville. Mutual information neural estimation. In *International Conference on Machine Learning*, pages 530–539, 2018.
- [3] Jean-François Cardoso. Dependence, correlation and gaussianity in independent component analysis. *Journal of Machine Learning Research*, 4(Dec):1177–1203, 2003.
- [4] Benjamin Charrow, Sikang Liu, Vijay Kumar, and Nathan Michael. Information-theoretic mapping using cauchy-schwarz quadratic mutual information. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4791–4798. IEEE, 2015.
- [5] Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, pages 2610–2620, 2018.
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- [7] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pages 2172–2180, 2016.
- [8] Pengyu Cheng, Weituo Hao, Shuyang Dai, Jiachang Liu, Zhe Gan, and Lawrence Carin. Club: A contrastive log-ratio upper bound of mutual information. *arXiv preprint arXiv:2006.12013*, 2020.
- [9] Pengyu Cheng, Martin Renqiang Min, Dinghan Shen, Christopher Malon, Yizhe Zhang, Yitong Li, and Lawrence Carin. Improving disentangled text representation learning with information-theoretic guidance. *arXiv preprint arXiv:2006.00693*, 2020.
- [10] Shuyang Gao, Rob Breckelmanns, Greg Ver Steeg, and Aram Galstyan. Auto-encoding total correlation explanation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1157–1166, 2019.
- [11] Behnam Gholami, Pritish Sahu, Ognjen Rudovic, Konstantinos Bousmalis, and Vladimir Pavlovic. Unsupervised multi-target domain adaptation: An information theoretic approach. *IEEE Transactions on Image Processing*, 2020.
- [12] Clive Granger and Jin-Lung Lin. Using the mutual information coefficient to identify lags in nonlinear models. *Journal of time series analysis*, 15(4):371–384, 1994.
- [13] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010.
- [14] 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. *arXiv preprint arXiv:1808.06670*, 2018.
- [15] Bo Jiang, Chao Ye, and Jun S Liu. Nonparametric k-sample tests via dynamic slicing. *Journal of the American Statistical Association*, 110(510):642–653, 2015.
- [16] Brian J Julian, Sertac Karaman, and Daniela Rus. On mutual information-based control of range sensing robots for mapping applications. *The International Journal of Robotics Research*, 33(10):1375–1392, 2014.
- [17] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning*, pages 2649–2658, 2018.

- [18] Alexander Lachmann, Federico M Giorgi, Gonzalo Lopez, and Andrea Califano. Aracne-ap: gene network reverse engineering through adaptive partitioning inference of mutual information. *Bioinformatics*, 32(14):2233–2235, 2016.
- [19] Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and Olivier Bachem. On the fairness of disentangled representations. In *Advances in Neural Information Processing Systems*, pages 14611–14624, 2019.
- [20] XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- [21] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [22] Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. On variational bounds of mutual information. In *International Conference on Machine Learning*, pages 5171–5180, 2019.
- [23] Jiaming Song, Pratyusha Kalluri, Aditya Grover, Shengjia Zhao, and Stefano Ermon. Learning controllable fair representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2164–2173, 2019.
- [24] Diego J Zea, Diego Anfossi, Morten Nielsen, and Cristina Marino-Buslje. Mitos. jl: mutual information tools for protein sequence analysis in the julia language. *Bioinformatics*, 33(4):564–565, 2016.
- [25] Sicheng Zhu, Xiao Zhang, and David Evans. Learning adversarially robust representations via worst-case mutual information maximization. *arXiv preprint arXiv:2002.11798*, 2020.

## A Proofs

*Proof of Theorem 2.1.* Note that  $\mathbf{X}_A := (\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_m})$  and  $\mathbf{X}_{\hat{A}} = \mathbf{X}/\mathbf{X}_A$ . Denote  $\mathbf{X}_{\hat{A}} = (\mathbf{x}_{j_1}, \mathbf{x}_{j_2}, \dots, \mathbf{x}_{j_l})$ . Then

$$\begin{aligned} \mathcal{TC}(\mathbf{X}) &= \mathbb{E}_{p(\mathbf{X})} \left[ \log \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)}{p(\mathbf{x}_1)p(\mathbf{x}_2) \dots p(\mathbf{x}_n)} \right] \\ &= \mathbb{E}_{p(\mathbf{X})} \left[ \log \left( \frac{p(\mathbf{X}_A)}{p(\mathbf{x}_{i_1})p(\mathbf{x}_{i_2}) \dots p(\mathbf{x}_{i_m})} \cdot \frac{p(\mathbf{X}_{\hat{A}})}{p(\mathbf{x}_{j_1})p(\mathbf{x}_{j_2}) \dots p(\mathbf{x}_{j_l})} \cdot \frac{p(\mathbf{X})}{p(\mathbf{X}_A)p(\mathbf{X}_{\hat{A}})} \right) \right] \\ &= \mathcal{TC}(\mathbf{X}_A) + \mathcal{TC}(\mathbf{X}_{\hat{A}}) + \mathcal{I}(\mathbf{X}_A; \mathbf{X}_{\hat{A}}) \end{aligned}$$

□

*Proof of Corollary 2.1.2.* We denote  $\mathbf{X}_{i:j} := (\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_{j-1}, \mathbf{x}_j)$ . Note that

$$\begin{aligned} \mathcal{TC}(\mathbf{X}_{1:n}) &= \mathbb{E}_{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)} \left[ \log \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)}{p(\mathbf{x}_1)p(\mathbf{x}_2) \dots p(\mathbf{x}_n)} \right] \\ &= \mathbb{E}_{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)} \left[ \log \left( \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}, \mathbf{x}_n)}{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1})p(\mathbf{x}_n)} \cdot \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1})}{p(\mathbf{x}_1)p(\mathbf{x}_2) \dots p(\mathbf{x}_{n-1})} \right) \right] \\ &= \mathcal{I}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}; \mathbf{x}_n) + \mathcal{TC}(\mathbf{X}_{1:n-1}) \\ &= \mathcal{I}(\mathbf{X}_{1:n-1}; \mathbf{x}_n) + \mathcal{TC}(\mathbf{X}_{1:n-1}) \end{aligned}$$

Similarly,

$$\mathcal{TC}(\mathbf{X}_{1:n}) = \mathcal{I}(\mathbf{X}_{1:n-1}; \mathbf{x}_n) + \mathcal{I}(\mathbf{X}_{1:n-2}; \mathbf{x}_{n-1}) + \mathcal{TC}(\mathbf{X}_{1:n-2}) = \sum_{i=1}^{n-1} \mathcal{I}(\mathbf{X}_{1:i}; \mathbf{x}_{i+1}) \quad (3)$$

□

## B Experiment Results

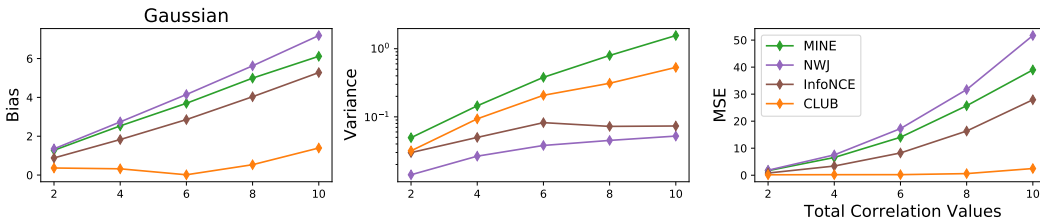


Figure 4: Bias, variance and MSE of **line-like** TC estimators

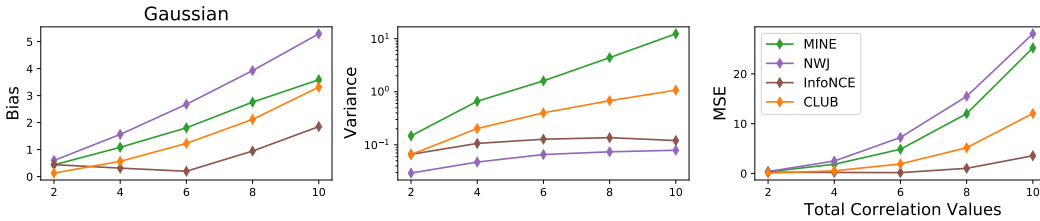


Figure 5: Bias, variance and MSE of **tree-like** TC estimators

## C MI Estimators

The Mutual Information Neural Estimator (MINE) [2] is defined as

$$\mathcal{I}_{\text{MINE}} := \mathbb{E}_{p(\mathbf{x}, \mathbf{y})}[f(\mathbf{x}, \mathbf{y})] - \log(\mathbb{E}_{p(\mathbf{x})p(\mathbf{y})}[e^{f(\mathbf{x}, \mathbf{y})}]), \quad (4)$$

where  $f(\cdot, \cdot)$  is a value function (or, a critic) approximated by a neural network.

The NWJ [20] lower bound is based on the  $f$ -divergence representation of MI:

$$\mathcal{I}_{\text{NWJ}} := \mathbb{E}_{p(\mathbf{x}, \mathbf{y})}[f(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{p(\mathbf{x})p(\mathbf{y})}[e^{f(\mathbf{x}, \mathbf{y})-1}]. \quad (5)$$

The InfoNCE [21] lower bound is based on Noise Contrastive Estimation (NCE) [13]:

$$\mathcal{I}_{\text{NCE}} := \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \log \frac{e^{f(\mathbf{x}_i, \mathbf{y}_i)}}{\frac{1}{N} \sum_{j=1}^N e^{f(\mathbf{x}_i, \mathbf{y}_j)}} \right], \quad (6)$$

where the expectation is over  $N$  samples  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$  drawn from the joint distribution  $p(\mathbf{x}, \mathbf{y})$ .

The MI contrastive log-ratio upper bound (CLUB) estimator [8] is based on a parameterized distribution  $q_\theta(\mathbf{y}|\mathbf{x})$ :

$$\mathcal{I}(\mathbf{x}; \mathbf{y}) \leq \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N [\log p(\mathbf{x}_i|\mathbf{y}_i) - \frac{1}{N} \sum_{j=1}^N \log p(\mathbf{x}_j|\mathbf{y}_i)] \right]. \quad (7)$$

All the MI lower bounds require learning of a value function  $f(\mathbf{x}, \mathbf{y})$ ; the CLUB upper bound requires learning of a network approximation  $q_\theta(\mathbf{y}|\mathbf{x})$ . To make fair comparison, we set the value function and the neural approximation with one hidden layer and the same hidden units. For the multivariate Gaussian setup, the number of hidden units is 20. On the top of hidden layer outputs, we add the ReLU activation function. The learning rate for all estimators is set to  $1 \times 10^{-4}$ .