

CLUB: A Contrastive Log-ratio Upper Bound of Mutual Information

Pengyu Cheng¹ Weituo Hao¹ Shuyang Dai¹ Jiachang Liu¹ Zhe Gan² Lawrence Carin¹

Abstract

Mutual information (MI) minimization has gained considerable interests in various machine learning tasks. However, estimating and minimizing MI in high-dimensional spaces remains a challenging problem, especially when only samples, rather than distribution forms, are accessible. Previous works mainly focus on MI lower bound approximation, which is not applicable to MI minimization problems. In this paper, we propose a novel Contrastive Log-ratio Upper Bound (CLUB) of mutual information. We provide a theoretical analysis of the properties of CLUB and its variational approximation. Based on this upper bound, we introduce a MI minimization training scheme and further accelerate it with a negative sampling strategy. Simulation studies on Gaussian distributions show the reliable estimation ability of CLUB. Real-world MI minimization experiments, including domain adaptation and information bottleneck, demonstrate the effectiveness of the proposed method. The code is at <https://github.com/Linear95/CLUB>.

1. Introduction

Mutual information (MI) is a fundamental measure of the dependence between two random variables. Mathematically, the definition of MI between variables x and y is

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

This important tool has been applied in a wide range of scientific fields, including statistics (Granger & Lin, 1994; Jiang et al., 2015), bioinformatics (Lachmann et al., 2016; Zea et al., 2016), robotics (Julian et al., 2014; Charrow et al., 2015), and machine learning (Chen et al., 2016; Alemi et al., 2016; Hjelm et al., 2018; Cheng et al., 2020).

¹Department of Electrical and Computer Engineering, Duke University, Durham, North Carolina, USA ²Microsoft, Redmond, Washington, USA. Correspondence to: Pengyu Cheng <pengyu.cheng@duke.edu>.

In machine learning, especially in deep learning frameworks, MI is typically utilized as a criterion or a regularizer in loss functions, to encourage or limit the dependence between variables. MI maximization has been studied extensively in various tasks, *e.g.*, representation learning (Hjelm et al., 2018; Hu et al., 2017), generative models (Chen et al., 2016), information distillation (Ahn et al., 2019), and reinforcement learning (Florensa et al., 2017). Recently, MI minimization has obtained increasing attention for its applications in disentangled representation learning (Chen et al., 2018), style transfer (Kazemi et al., 2018), domain adaptation (Gholami et al., 2018), fairness (Kamishima et al., 2011), and the information bottleneck (Alemi et al., 2016).

However, only in a few special cases can one calculate the exact value of mutual information, since the calculation requires closed forms of density functions and a tractable log-density ratio between the joint and marginal distributions. In most machine learning tasks, only samples from the joint distribution are accessible. Therefore, sample-based MI estimation methods have been proposed. To approximate MI, most previous works focused on lower-bound estimation (Chen et al., 2016; Belghazi et al., 2018; Oord et al., 2018), which is inconsistent to MI minimization tasks. In contrast, MI upper bound estimation lacks extensive exploration in the literature. Among the existing MI upper bounds, Alemi et al. (2016) fixes one of the marginal distribution ($p(y)$ in (1)) to a standard Gaussian, and obtains a variational upper bound in closed form. However, the Gaussian marginal distribution assumption is unduly strong, which makes the upper bound fail to estimate MI with low bias. Poole et al. (2019) points out a leave-one-out upper bound, which provides tighter MI estimation when sample size is large. However, it suffers from high numerical instability in practice when applied to MI minimization models.

To overcome the defects of previous MI estimators, we introduce a Contrastive Log-ratio Upper Bound (CLUB). Specifically, CLUB bridges mutual information estimation with contrastive learning (Oord et al., 2018), where MI is estimated by the difference of conditional probabilities between positive and negative sample pairs. Further, we develop a variational form of CLUB (vCLUB) into scenarios where the conditional distribution $p(y|x)$ is unknown, by approximating $p(y|x)$ with a neural network. We theoretically prove that, with good variational approximation,

vCLUB can either provide reliable MI estimation or remain a valid MI upper bound. Based on this new bound, we propose an MI minimization algorithm, and further accelerate it via a negative sampling strategy. The main contributions of this paper are summarized as follows.

- We introduce a Contrastive Log-ratio Upper Bound (CLUB) of mutual information, which is not only reliable as a mutual information estimator, but also trainable in gradient-descent frameworks.
- We extend CLUB with a variational network approximation, and provide theoretical analysis to the good properties of this variational bound.
- We develop a CLUB-based MI minimization algorithm, and accelerate it with a negative sampling strategy.
- We compare CLUB with previous MI estimators on both simulation studies and real-world applications, which demonstrate CLUB is not only better in the bias-variance estimation trade-off, but also more effective when applied to MI minimization.

2. Background

Although widely used in numerous applications, mutual information (MI) remains challenging to estimate accurately, when the closed-forms of distributions are unknown or intractable. Earlier MI estimation approaches include non-parametric binning (Darbellay & Vajda, 1999), kernel density estimation (Härdle et al., 2004), likelihood-ratio estimation (Suzuki et al., 2008), and K -nearest neighbor entropy estimation (Kraskov et al., 2004). These methods fail to provide reliable approximations when the data dimension increases (Belghazi et al., 2018). Also, the gradient of these estimators is difficult to calculate, which makes them inapplicable to back-propagation frameworks for MI optimization tasks.

To obtain differentiable and scalable MI estimation, recent approaches utilize deep neural networks to construct variational MI estimators. Most of these estimators focus on MI maximization problems, and provide MI lower bounds. Specifically, Barber & Agakov (2003) replaces the conditional distribution $p(\mathbf{y}|\mathbf{x})$ with an auxiliary distribution $q(\mathbf{y}|\mathbf{x})$, and obtains the Barber-Agakov (BA) bound:

$$I_{\text{BA}} := H(\mathbf{x}) + \mathbb{E}_{p(\mathbf{x},\mathbf{y})}[\log q(\mathbf{x}|\mathbf{y})] \leq I(\mathbf{x}; \mathbf{y}), \quad (2)$$

where $H(\mathbf{x})$ is the entropy of variable \mathbf{x} . Belghazi et al. (2018) introduces a Mutual Information Neural Estimator (MINE), which treats MI as the Kullback-Leibler (KL) divergence (Kullback, 1997) between the joint and marginal distributions, and converts it into the dual representation:

$$I_{\text{MINE}} := \mathbb{E}_{p(\mathbf{x},\mathbf{y})}[f(\mathbf{x}, \mathbf{y})] - \log(\mathbb{E}_{p(\mathbf{x})p(\mathbf{y})}[e^{f(\mathbf{x},\mathbf{y})}]), \quad (3)$$

where $f(\cdot, \cdot)$ is a score function (or, a critic) approximated by a neural network. Nguyen, Wainwright, and Jordan (NWJ) (Nguyen et al., 2010) derives another lower bound based on the MI f -divergence representation:

$$I_{\text{NWJ}} := \mathbb{E}_{p(\mathbf{x},\mathbf{y})}[f(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{p(\mathbf{x})p(\mathbf{y})}[e^{f(\mathbf{x},\mathbf{y})-1}]. \quad (4)$$

More recently, based on Noise Contrastive Estimation (NCE) (Gutmann & Hyvärinen, 2010), an MI lower bound, called InfoNCE, was introduced in Oord et al. (2018):

$$I_{\text{NCE}} := \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \log \frac{e^{f(\mathbf{x}_i, \mathbf{y}_i)}}{\frac{1}{N} \sum_{j=1}^N e^{f(\mathbf{x}_i, \mathbf{y}_j)}} \right], \quad (5)$$

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

Unlike the above MI lower bounds that have been studied extensively, MI upper bounds are still lacking extensive published exploration. Most existing MI upper bounds require the conditional distribution $p(\mathbf{y}|\mathbf{x})$ to be known. For example, Alemi et al. (2016) introduces a variational marginal approximation $r(\mathbf{y})$ to build a variational upper bound (VUB):

$$\begin{aligned} I(\mathbf{x}; \mathbf{y}) &= \mathbb{E}_{p(\mathbf{x},\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y})} \right] \\ &= \mathbb{E}_{p(\mathbf{x},\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \right] - \text{KL}(p(\mathbf{y}) \| r(\mathbf{y})) \\ &\leq \mathbb{E}_{p(\mathbf{x},\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \right] = \text{KL}(p(\mathbf{y}|\mathbf{x}) \| r(\mathbf{y})). \end{aligned} \quad (6)$$

The inequality is based on the fact that the KL-divergence is always non-negative. To be a good MI estimation, this upper bound requires a well-learned density approximation $r(\mathbf{y})$ to $p(\mathbf{y})$, so that the difference $\text{KL}(p(\mathbf{y}) \| r(\mathbf{y}))$ could be small. However, learning a good marginal approximation $r(\mathbf{y})$ without any additional information, recognized as the distribution density estimation problem (Magdon-Ismail & Atiya, 1999), is challenging, especially when variable \mathbf{y} is in a high-dimensional space. In practice, Alemi et al. (2016) fixes $r(\mathbf{y})$ as a standard normal distribution, $r(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{I})$, which results in a high-bias MI estimation. With N sample pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, Poole et al. (2019) replaces $r(\mathbf{y})$ with a Monte Carlo approximation $r_i(\mathbf{y}) = \frac{1}{N-1} \sum_{j \neq i} p(\mathbf{y}|\mathbf{x}_j) \approx p(\mathbf{y})$ and derives a leave-one-out upper bound (L1Out):

$$I_{\text{L1Out}} := \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \left[\log \frac{p(\mathbf{y}_i|\mathbf{x}_i)}{\frac{1}{N-1} \sum_{j \neq i} p(\mathbf{y}_i|\mathbf{x}_j)} \right] \right]. \quad (7)$$

This bound does not require any additional parameters, but highly depends on a sufficient sample size to achieve satisfying Monte Carlo approximation. In practice, L1Out suffers from numerical instability when applied to real-world MI minimization problems.

To compare our method with the aforementioned MI upper bounds in more general scenarios (*i.e.*, $p(\mathbf{y}|\mathbf{x})$ is unknown), we use a neural network $q_\theta(\mathbf{y}|\mathbf{x})$ to approximate $p(\mathbf{y}|\mathbf{x})$, and develop variational versions of VUB and L1Out as :

$$I_{\text{vVUB}} = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \frac{q_\theta(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \right], \quad (8)$$

$$I_{\text{vL1Out}} = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \left[\log \frac{q_\theta(\mathbf{y}_i|\mathbf{x}_i)}{\frac{1}{N-1} \sum_{j \neq i} q_\theta(\mathbf{y}_i|\mathbf{x}_j)} \right] \right]. \quad (9)$$

We discuss theoretical properties of these two variational bounds in the Supplementary Material. In a simulation study (Section 4.1), variational L1Out reaches better performance than previous lower bounds for MI estimation. However, the numerical instability problem remains for variational L1Out in real-world applications (Section 4.4). To the best of our knowledge, we provide the first variational version of VUB and L1Out upper bounds, and study their properties on both the theoretical analysis and the empirical performance.

3. Proposed Method

Suppose we have sample pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ drawn from an unknown or intractable distribution $p(\mathbf{x}, \mathbf{y})$. We aim to derive a upper bound estimator of the mutual information $I(\mathbf{x}; \mathbf{y})$ based on the given samples. In a range of machine learning tasks (*e.g.*, information bottleneck), one of the conditional distributions between variables \mathbf{x} and \mathbf{y} (as $p(\mathbf{x}|\mathbf{y})$ or $p(\mathbf{y}|\mathbf{x})$) can be known. To efficiently utilize this additional information, we first derive a mutual information (MI) upper bound with the assumption that one of the conditional distribution is provided (suppose $p(\mathbf{y}|\mathbf{x})$ is provided, without loss of generality). Then, we extend the bound into more general cases where no conditional distribution is known. Finally, we develop a MI minimization algorithm based on the derived bound.

3.1. CLUB with $p(\mathbf{y}|\mathbf{x})$ Known

With the conditional distribution $p(\mathbf{y}|\mathbf{x})$, our MI Contrastive Log-ratio Upper Bound (CLUB) is defined as:

$$I_{\text{CLUB}}(\mathbf{x}; \mathbf{y}) := \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})]. \quad (10)$$

To show that $I_{\text{CLUB}}(\mathbf{x}; \mathbf{y})$ is an upper bound of $I(\mathbf{x}; \mathbf{y})$, we calculate the gap Δ between them:

$$\begin{aligned} \Delta &:= I_{\text{CLUB}}(\mathbf{x}; \mathbf{y}) - I(\mathbf{x}; \mathbf{y}) \\ &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})] \\ &\quad - \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x}) - \log p(\mathbf{y})] \\ &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y})] - \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})] \\ &= \mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y}) - \mathbb{E}_{p(\mathbf{x})} [\log p(\mathbf{y}|\mathbf{x})]]. \end{aligned} \quad (11)$$

By the definition of the marginal distribution, we have $p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x} = \mathbb{E}_{p(\mathbf{x})}[p(\mathbf{y}|\mathbf{x})]$. Note that

$\log(\cdot)$ is a concave function, by Jensen's Inequality, we have $\log p(\mathbf{y}) = \log (\mathbb{E}_{p(\mathbf{x})}[p(\mathbf{y}|\mathbf{x})]) \geq \mathbb{E}_{p(\mathbf{x})}[\log p(\mathbf{y}|\mathbf{x})]$. Applying this inequality to equation (11), we conclude that the gap Δ is always non-negative. Therefore, $I_{\text{CLUB}}(\mathbf{x}; \mathbf{y})$ is an upper bound of $I(\mathbf{x}; \mathbf{y})$. The bound is tight when $p(\mathbf{y}|\mathbf{x})$ has the same value for any \mathbf{x} , which means variables \mathbf{x} and \mathbf{y} are independent. Consequently, we summarize the above discussion into the following Theorem 3.1.

Theorem 3.1. For two random variables \mathbf{x} and \mathbf{y} ,

$$I(\mathbf{x}; \mathbf{y}) \leq I_{\text{CLUB}}(\mathbf{x}; \mathbf{y}). \quad (12)$$

Equality is achieved if and only if \mathbf{x} and \mathbf{y} are independent.

With sample pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, $I_{\text{CLUB}}(\mathbf{x}; \mathbf{y})$ has an unbiased estimation as:

$$\begin{aligned} \hat{I}_{\text{CLUB}} &= \frac{1}{N} \sum_{i=1}^N \log p(\mathbf{y}_i|\mathbf{x}_i) - \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \log p(\mathbf{y}_j|\mathbf{x}_i) \\ &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N [\log p(\mathbf{y}_i|\mathbf{x}_i) - \log p(\mathbf{y}_j|\mathbf{x}_i)]. \end{aligned} \quad (13)$$

In the estimator \hat{I}_{CLUB} , $\log p(\mathbf{y}_i|\mathbf{x}_i)$ provides the conditional log-likelihood of positive sample pair $(\mathbf{x}_i, \mathbf{y}_i)$; $\{\log p(\mathbf{y}_j|\mathbf{x}_i)\}_{i \neq j}$ provide the conditional log-likelihood of negative sample pair $(\mathbf{x}_i, \mathbf{y}_j)$. The difference between $\log p(\mathbf{y}_i|\mathbf{x}_i)$ and $\log p(\mathbf{y}_j|\mathbf{x}_i)$ is the contrastive probability log-ratio between two conditional distributions. Therefore, we name this novel MI upper bound estimator as Contrastive Log-ratio Upper Bound (CLUB). Compared with previous MI neural estimators, CLUB has a simpler form as a linear combination of log-ratios between positive and negative sample pairs. The linear form of log-ratios improves the numerical stability for calculation of CLUB and its gradient, which we discuss in details in Section 3.3.

3.2. CLUB with Conditional Distributions Unknown

When the conditional distributions $p(\mathbf{y}|\mathbf{x})$ or $p(\mathbf{x}|\mathbf{y})$ is provided, the MI can be directly upper-bounded by equation (13) with samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$. Unfortunately, in a large number of machine learning tasks, the conditional relation between variables is unavailable.

To further extend the CLUB estimator into more general scenarios, we use a variational distribution $q_\theta(\mathbf{y}|\mathbf{x})$ with parameter θ to approximate $p(\mathbf{y}|\mathbf{x})$. Consequently, a variational CLUB term (vCLUB) is defined by:

$$I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y}) := \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q_\theta(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{p(\mathbf{y})} [\log q_\theta(\mathbf{y}|\mathbf{x})]. \quad (14)$$

Similar to the MI upper bound estimator \hat{I}_{CLUB} in (13), the

unbiased estimator for vCLUB with samples $\{\mathbf{x}_i, \mathbf{y}_i\}$ is:

$$\begin{aligned}\hat{I}_{\text{vCLUB}} &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left[\log q_{\theta}(\mathbf{y}_i | \mathbf{x}_i) - \log q_{\theta}(\mathbf{y}_j | \mathbf{x}_i) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\log q_{\theta}(\mathbf{y}_i | \mathbf{x}_i) - \frac{1}{N} \sum_{j=1}^N \log q_{\theta}(\mathbf{y}_j | \mathbf{x}_i) \right]. \quad (15)\end{aligned}$$

Using the variational approximation $q_{\theta}(\mathbf{y} | \mathbf{x})$, vCLUB no longer guarantees an upper bound of $I(\mathbf{x}; \mathbf{y})$. However, the vCLUB shares good properties with CLUB. We claim that with good variational approximation $q_{\theta}(\mathbf{y} | \mathbf{x})$, vCLUB can still hold a MI upper bound or become a reliable MI estimator. The following analyses support this claim.

Let $q_{\theta}(\mathbf{x}, \mathbf{y}) = q_{\theta}(\mathbf{y} | \mathbf{x})p(\mathbf{x})$ be the variational joint distribution induced by $q_{\theta}(\mathbf{y} | \mathbf{x})$. Generally, we have the following Theorem 3.2. Note that when \mathbf{x} and \mathbf{y} are independent, I_{vCLUB} has exactly the same value as $I(\mathbf{x}; \mathbf{y})$, without requiring any additional assumption on $q_{\theta}(\mathbf{y} | \mathbf{x})$. However, unlike in Theorem 3.1 as a sufficient and necessary condition, the ‘‘independence between \mathbf{x} and \mathbf{y} ’’ becomes sufficient but not necessary to conclude ‘‘ $I(\mathbf{x}; \mathbf{y}) = I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y})$ ’’, due to the variation approximation $q_{\theta}(\mathbf{y} | \mathbf{x})$.

Theorem 3.2. Denote $q_{\theta}(\mathbf{x}, \mathbf{y}) = q_{\theta}(\mathbf{y} | \mathbf{x})p(\mathbf{x})$. If

$$KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})) \leq KL(p(\mathbf{x})p(\mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})),$$

then $I(\mathbf{x}; \mathbf{y}) \leq I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y})$. The equality holds when \mathbf{x} and \mathbf{y} are independent.

Theorem 3.2 provides insight that vCLUB remains a MI upper bound if the variational joint distribution $q_{\theta}(\mathbf{x}, \mathbf{y})$ is ‘‘closer’’ to $p(\mathbf{x}, \mathbf{y})$ than to $p(\mathbf{x})p(\mathbf{y})$. Therefore, minimizing $KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y}))$ will facilitate the condition in Theorem 3.2 to be achieved. We show that $KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y}))$ can be minimized by maximizing the log-likelihood of $q_{\theta}(\mathbf{y} | \mathbf{x})$, because of the following equation:

$$\begin{aligned}& \min_{\theta} KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})) \\ &= \min_{\theta} \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log(p(\mathbf{y} | \mathbf{x})p(\mathbf{x})) - \log(q_{\theta}(\mathbf{y} | \mathbf{x})p(\mathbf{x}))] \\ &= \min_{\theta} \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y} | \mathbf{x})] - \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q_{\theta}(\mathbf{y} | \mathbf{x})]. \quad (16)\end{aligned}$$

Equation (16) equals $\min_{\theta} KL(p(\mathbf{y} | \mathbf{x}) || q_{\theta}(\mathbf{y} | \mathbf{x}))$, in which the first term has no relation with parameter θ . Therefore, $\min_{\theta} KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y}))$ is equivalent to the maximization of the second term, $\max_{\theta} \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q_{\theta}(\mathbf{y} | \mathbf{x})]$. With samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, we can maximize the log-likelihood function $\mathcal{L}(\theta) := \frac{1}{N} \sum_{i=1}^N \log q_{\theta}(\mathbf{y}_i | \mathbf{x}_i)$, which is the unbiased estimation of $\mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q_{\theta}(\mathbf{y} | \mathbf{x})]$.

In practice, the variational distribution $q_{\theta}(\mathbf{y} | \mathbf{x})$ is usually implemented with neural networks. By enlarging the network capacity (*i.e.*, adding layers and neurons)

and applying gradient-ascent to the log-likelihood $\mathcal{L}(\theta)$, we can obtain far more accurate approximation $q_{\theta}(\mathbf{y} | \mathbf{x})$ to $p(\mathbf{y} | \mathbf{x})$, thanks to the high expressiveness of neural networks (Hu et al., 2019; Oymak & Soltanolkotabi, 2019). Therefore, to further discuss the properties of vCLUB, we assume the neural network approximation q_{θ} achieves $KL(p(\mathbf{y} | \mathbf{x}) || q_{\theta}(\mathbf{y} | \mathbf{x})) \leq \varepsilon$ with a small number $\varepsilon > 0$. In the Supplementary Material, we quantitatively discuss the reasonableness of this assumption. Consider the KL-divergence between $p(\mathbf{x})p(\mathbf{y})$ and $q_{\theta}(\mathbf{x}, \mathbf{y})$. If $KL(p(\mathbf{x})p(\mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})) \geq KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y}))$, by Theorem 3.2, vCLUB is already a MI upper bound. Otherwise, if $KL(p(\mathbf{x})p(\mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})) < KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y}))$, we have the following corollary:

Corollary 3.3. Given $KL(p(\mathbf{y} | \mathbf{x}) || q_{\theta}(\mathbf{y} | \mathbf{x})) \leq \varepsilon$, if

$$KL(p(\mathbf{x}, \mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})) > KL(p(\mathbf{x})p(\mathbf{y}) || q_{\theta}(\mathbf{x}, \mathbf{y})),$$

then $|I(\mathbf{x}; \mathbf{y}) - I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y})| < \varepsilon$.

Combining Corollary 3.3 and Theorem 3.2, we conclude that with a good variational approximation $q_{\theta}(\mathbf{y} | \mathbf{x})$, vCLUB can either remain a MI upper bound, or become a MI estimator whose absolute error is bounded by the approximation performance $KL(p(\mathbf{y} | \mathbf{x}) || q_{\theta}(\mathbf{y} | \mathbf{x}))$.

3.3. CLUB in MI Minimization

One of the major applications of MI upper bounds is for mutual information minimization. In general, MI minimization aims to reduce the correlation between two variables \mathbf{x} and \mathbf{y} by selecting an optimal parameter σ of the joint variational distribution $p_{\sigma}(\mathbf{x}, \mathbf{y})$. Under some application scenarios, additional conditional information between \mathbf{x} and \mathbf{y} is known. For example, in the information bottleneck task, the joint distribution between input \mathbf{x} and bottleneck representation \mathbf{y} is $p_{\sigma}(\mathbf{x}, \mathbf{y}) = p_{\sigma}(\mathbf{y} | \mathbf{x})p(\mathbf{x})$. Then the MI upper bound I_{CLUB} can be calculated directly based on Eqn. (13).

Algorithm 1 MI Minimization with vCLUB

```

for each training iteration do
  Sample  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  from  $p_{\sigma}(\mathbf{x}, \mathbf{y})$ 
  Log-likelihood  $\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \log q_{\theta}(\mathbf{y}_i | \mathbf{x}_i)$ 
  Update  $q_{\theta}(\mathbf{y} | \mathbf{x})$  by maximizing  $\mathcal{L}(\theta)$ 
  for  $i = 1$  to  $N$  do
    if use sampling then
      Sample  $k'_i$  uniformly from  $\{1, 2, \dots, N\}$ 
       $U_i = \log q_{\theta}(\mathbf{y}_i | \mathbf{x}_i) - \log q_{\theta}(\mathbf{y}_{k'_i} | \mathbf{x}_i)$ 
    else
       $U_i = \log q_{\theta}(\mathbf{y}_i | \mathbf{x}_i) - \frac{1}{N} \sum_{j=1}^N \log q_{\theta}(\mathbf{y}_j | \mathbf{x}_i)$ 
    end if
  end for
  Update  $p_{\sigma}(\mathbf{x}, \mathbf{y})$  by minimize  $\hat{I}_{\text{vCLUB}} = \frac{1}{N} \sum_{i=1}^N U_i$ 
end for
    
```

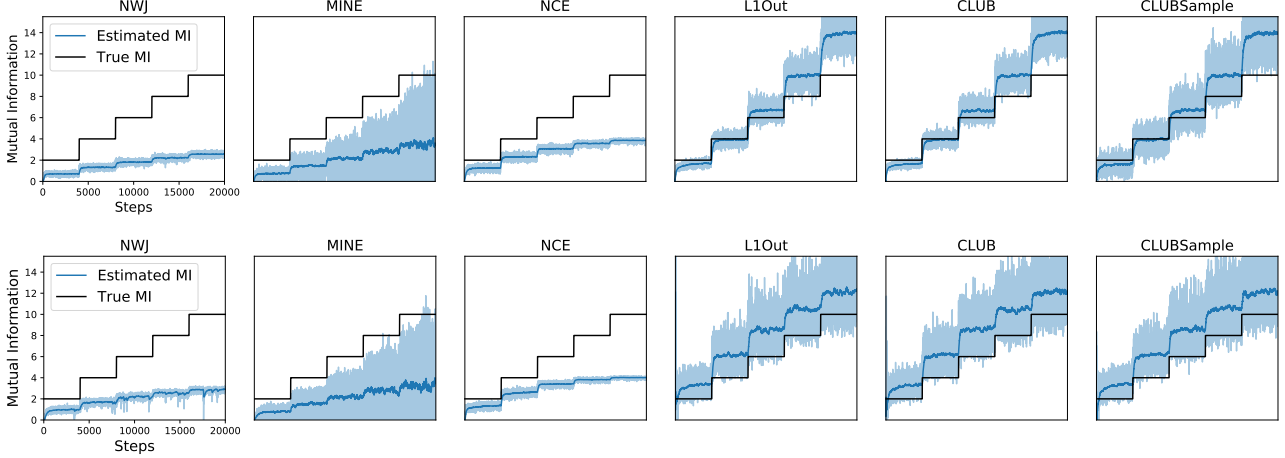


Figure 1. Simulation performance of MI estimators. In the **top** row, data are from joint Gaussian distributions with the MI true value stepping over time. In the **bottom** row, a cubic transformation is further applied to the Gaussian samples as \mathbf{y} . In each figure, the true MI values is a step function shown as the black line. The estimated values are displayed as shadow blue curves. The dark blue curves shows the local averages of estimated MI, with a bandwidth equal to 200.

For cases in which the conditional information between \mathbf{x} and \mathbf{y} remains unclear, we propose an MI minimization algorithm using the vCLUB estimator. At each training iteration, we first obtain a batch of samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}$ from $p_\sigma(\mathbf{x}, \mathbf{y})$. Then we update the variational approximation $q_\theta(\mathbf{y}|\mathbf{x})$ by maximizing the log-likelihood $\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \log q_\theta(\mathbf{y}_i|\mathbf{x}_i)$. After $q_\theta(\mathbf{y}|\mathbf{x})$ is updated, we calculate the vCLUB estimator as described in (15). Finally, the gradient of $\hat{\mathbf{I}}_{\text{vCLUB}}$ is calculated and back-propagated to parameters of $p_\sigma(\mathbf{x}, \mathbf{y})$. The reparameterization trick (Kingma & Welling, 2013) ensures the gradient back-propagates through the sampled embeddings $(\mathbf{x}_i, \mathbf{y}_i)$. Updating joint distribution $p_\sigma(\mathbf{x}, \mathbf{y})$ will lead to the change of conditional distribution $p_\sigma(\mathbf{y}|\mathbf{x})$. Therefore, we need to update the approximation network $q_\theta(\mathbf{y}|\mathbf{x})$ again. Consequently, $q_\theta(\mathbf{y}|\mathbf{x})$ and $p_\sigma(\mathbf{x}, \mathbf{y})$ are updated alternately during the training (as shown in Algorithm 1 without *sampling*).

In each training iteration, the vCLUB estimator requires calculation of all conditional distributions $\{p_\sigma(\mathbf{y}_j|\mathbf{x}_i)\}_{i,j=1}^N$, which leads to $\mathcal{O}(N^2)$ computational complexity. To further accelerate the calculate, for each positive sample pair $(\mathbf{x}_i, \mathbf{y}_i)$, instead of calculating the mean of the probabilities of all negative pairs as $\frac{1}{N} \sum_{i=1}^N \log q_\theta(\mathbf{y}_j|\mathbf{x}_i)$ in (15), we randomly sample a negative pair $(\mathbf{x}_i, \mathbf{y}_{k'_i})$ and use $\log q_\theta(\mathbf{y}_{k'_i}|\mathbf{x}_i)$ as an unbiased estimation, with k'_i uniformly selected from indices $\{1, 2, \dots, N\}$. Then we obtain the sampled vCLUB (vCLUB-S) MI estimator:

$$\hat{\mathbf{I}}_{\text{vCLUB-S}} = \frac{1}{N} \sum_{i=1}^N \left[\log q_\theta(\mathbf{y}_i|\mathbf{x}_i) - \log q_\theta(\mathbf{y}_{k'_i}|\mathbf{x}_i) \right],$$

with the property of unbiasedness that $\mathbb{E}[\hat{\mathbf{I}}_{\text{vCLUB-S}}] = \mathbb{E}[\hat{\mathbf{I}}_{\text{vCLUB}}] = \mathbf{I}_{\text{vCLUB}}(\mathbf{x}; \mathbf{y})$. By this sampling strategy, the computational complexity in each iteration can be reduced

to $\mathcal{O}(N)$ (as Algorithm 1 with *sampling*). A similar sampling strategy can also be applied to CLUB when $p(\mathbf{y}|\mathbf{x})$ is known. Besides the acceleration, the vCLUB-S estimator bridges the MI minimization with *negative sampling*, a commonly used training strategy for learning word embeddings (e.g., Word2Vec (Mikolov et al., 2013)) and node embeddings (e.g., Node2Vec (Grover & Leskovec, 2016)), in which a positive data pair $(\mathbf{x}_i, \mathbf{y}_i)$ includes two nodes with an edge connection or two words in the same sentence, and a negative pair $(\mathbf{x}_i, \mathbf{y}_{k'_i})$ is uniformly sampled from the whole graph or vocabulary. Although previous MI upper bounds also utilize the negative data pairs (such as L1Out in (7)), they cannot hold an unbiased estimation when accelerated with the sampling strategy, because of the non-linear log function applied after the linear probability summation. The unbiasedness of our sampled CLUB thanks to the form of linear log-ratio summation. In the experiments, we find the sampled vCLUB estimator not only provides comparable MI estimation performance, but also improves the model generalization abilities as a learning critic.

4. Experiments

In this section, we first show the performance of CLUB as a MI estimator on tractable toy (simulated) cases, with samples drawn from Gaussian and Cubic distributions. Then we evaluate the minimization ability of CLUB on two real-world applications: Information Bottleneck (IB) and Unsupervised Domain Adaptation (UDA). In the information bottleneck, the conditional distribution $p(\mathbf{y}|\mathbf{x})$ is known, so we compare performance of both CLUB and variational CLUB (vCLUB) estimators and their sampled versions. In the other experiments for which $p(\mathbf{y}|\mathbf{x})$ is unknown, all the tested upper bounds require variational approximation. Without ambiguity, in experiments except the Information Bottle-

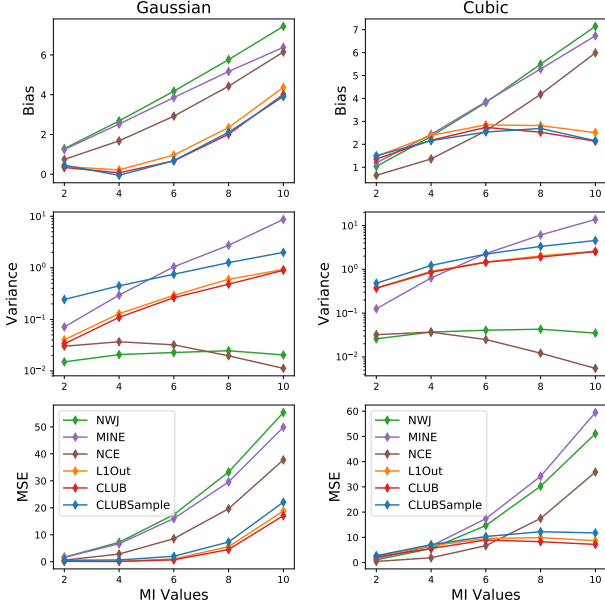


Figure 2. Estimation quality comparison of MI estimators. The **left** column shows the results of estimations under Gaussian distribution, while the **right** column is under Cubic setup. In each column, estimation metrics are reported as bias, variance, and mean-square-error (MSE). In each plot, the evaluation metric is reported with different true MI values varying from 2 to 10.

neck, we abbreviate all variational bounds (e.g., vCLUB) with their original names (e.g., CLUB) for simplicity.

4.1. MI Estimation Quality

Following the setup from Poole et al. (2019), we apply CLUB as an MI estimator in two toy tasks: (i) estimating MI with samples $\{(x_i, y_i)\}$ drawn jointly from a multivariate Gaussian distribution with correlation ρ ; (ii) estimating MI with samples $\{(x_i, (W y_i)^3)\}$, where (x_i, y_i) still comes from a Gaussian with correlation ρ , and W is a full-rank matrix. Since the transformation $y \rightarrow (W y)^3$ is smooth and bijective, the mutual information is invariant (Kraskov et al., 2004), $I(x; y) = I(x; (W y)^3)$. For both of the tasks, the dimension of samples x and y is set to $d = 20$. Under Gaussian distributions, the MI true value can be calculated as $I(x, y) = -\frac{d}{2} \log(1 - \rho^2)$, and therefore we set the MI true value in the range $\{2.0, 4.0, 6.0, 8.0, 10.0\}$ by varying the value of ρ . At each MI true value, we sample data batches 4000 times, with batch size equal to 64, for the training of variational MI estimators.

We compare our method with baselines including MINE (Belghazi et al., 2018), NWJ (Nguyen et al., 2010), InfoNCE (Oord et al., 2018), VUB (Alemi et al., 2016) and L1Out (Poole et al., 2019). Since the conditional distribution $p(y|x)$ is unknown in this simulation setup, all upper bounds (VUB, L1Out, CLUB) are calculated with an auxiliary approximation network $q_\theta(y|x)$. The

approximation network has the same structure for all upper bounds, parameterized in a Gaussian family, $q_\theta(y|x) = \mathcal{N}(y|\mu(x), \sigma^2(x) \cdot \mathbf{I})$ with mean $\mu(x)$ and variance $\sigma^2(x)$ inferred by neural networks. On the other hand, all the MI lower bounds (MINE, NWJ, InfoNCE) require learning of a value function $f(x, y)$. To make fair comparison, we set the value function and the neural approximation with one hidden layer and the same hidden units. For both the Gaussian and Cubic setups, the number of hidden units of our CLUB estimator is set to 15. On the top of hidden layer outputs, we add the ReLU activation function. The learning rate for all estimators is set to 5×10^{-3} .

We report in Figure 1 the estimated MI values in each training step. The estimation of VUB has incomparably large bias, so we provide its results in the Supplementary Material. Lower bound estimators, such as NWJ, MINE, and InfoNCE, provide estimated values mainly under the true MI values step function, while L1Out, CLUB and Sampled CLUB (CLUBSample) estimate values above the step function, which supports our theoretical analysis about CLUB with variational approximation. The numerical results of bias and variance in the estimation are reported in Figure 2. Among these methods, CLUB and CLUBSample have the lowest bias. The bias difference between CLUB and CLUBSample is insignificant, supporting our claim in Section 3.3 that CLUBSample is an unbiased stochastic approximation of CLUB. L1Out also provides small bias estimation which is slightly worse than CLUB. NWJ and InfoNCE have the lowest variance under both setups. CLUBSample has larger variance than CLUB and L1Out due to the use of the sampling strategy. When considering the bias-variance trade-off as the mean square estimation error (MSE, equals $\text{bias}^2 + \text{variance}$), CLUB outperforms other estimators, while L1Out and CLUBSample also provide competitive performance.

Although L1Out estimator reaches similar estimation performance as our CLUB on toy examples, we find L1Out fails to effectively reduce the MI when applied as a critic in real-world MI minimization tasks. The numerical results in Section 4.3 and Section 4.4 support our claim.

4.2. Time Efficiency of MI Estimators

Besides the estimation quality comparison, we further study the time efficiency of different MI estimators. We conduct the comparison under the same experimental setup as the Gaussian case in Section 4.1. Each MI estimator is tested with different batch size from 32 to 512. We count the total time cost of the whole estimation process and average it into each estimation step. In Figure 3, we report the average estimation time costs of different MI estimators. MINE and CLUBSample have the highest computational efficiency; both have $\mathcal{O}(N)$ computational complexity with respect to the sample size N , because of the negative sampling strat-

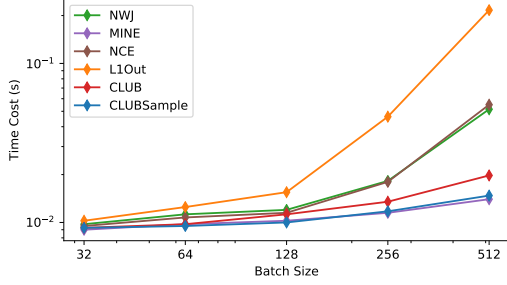


Figure 3. Estimator speed comparison with different batch size. Both the axes have a logarithm scale.

egy. Among other computational $\mathcal{O}(N^2)$ methods, CLUB has the highest estimation speed, thanks to its simple form as mean of log-ratios, which can be easily accelerated by matrix multiplication. Leave-one-out (L1out) has the highest time cost, because it requires “leaving out” the positive sample pair each time in the denominator of equation (7).

4.3. MI Minimization in Information Bottleneck

The Information Bottleneck (Tishby et al., 2000) (IB) is an information-theoretical method for latent representation learning. Given an input source $\mathbf{x} \in \mathcal{X}$ and a corresponding output target $\mathbf{y} \in \mathcal{Y}$, the information bottleneck aims to learn an encoder $p_\sigma(\mathbf{z}|\mathbf{x})$, such that the compressed latent code \mathbf{z} is highly relevant to the target \mathbf{y} , with irrelevant source information from \mathbf{x} being filtered. In other words, IB seeks to find the sufficient statistics of \mathbf{x} with respect to \mathbf{y} (Alemi et al., 2016), with minimum information used from \mathbf{x} . To address this task, an objective is introduced as

$$\min_{p_\sigma(\mathbf{z}|\mathbf{x})} -\mathcal{I}(\mathbf{y}; \mathbf{z}) + \beta \mathcal{I}(\mathbf{x}; \mathbf{z}) \quad (17)$$

where hyper-parameter $\beta > 0$. Following the same setup from Alemi et al. (2016), we apply the IB technique in the permutation-invariant MNIST classification. The input \mathbf{x} is a vector converted from a 28×28 image of a hand-written number, and the output \mathbf{y} is the class label of this number. The stochastic encoder $p_\sigma(\mathbf{z}|\mathbf{x})$ is implemented in a Gaussian variational family, $p_\sigma(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_\sigma(\mathbf{x}), \Sigma_\sigma(\mathbf{x}))$, where μ_σ and Σ_σ are two fully-connected neural networks.

For the first part of the IB objective (17), the MI between target \mathbf{y} and latent code \mathbf{z} is maximized. We use the same strategy as in the deep variational information bottleneck (DVB) (Alemi et al., 2016), where a variational classifier $q_\phi(\mathbf{y}|\mathbf{z})$ is introduced to implement a Barber-Agakov MI lower bound (Eqn. (2)) of $\mathcal{I}(\mathbf{y}; \mathbf{z})$. The second term in the IB objective requires the MI minimization between input \mathbf{x} and the latent representation \mathbf{z} . DVB (Alemi et al., 2016) utilizes the MI variation upper bound (VUB) (Eqn. (6)) for the minimization of $\mathcal{I}(\mathbf{x}; \mathbf{z})$. Since the closed form of $p_\sigma(\mathbf{z}|\mathbf{x})$ is already known as a Gaussian distribution parameterized by neural networks, we can directly apply our CLUB esti-

Method	Misclass. rate(%)
NWJ (Nguyen et al., 2010)	1.29
MINE (Belghazi et al., 2018)	1.17
InfoNCE (Oord et al., 2018)	1.24
DVB (VUB) (Alemi et al., 2016)	1.13
L1Out (Poole et al., 2019)	-
CLUB	1.12
CLUB (Sample)	1.10
vCLUB	1.10
vCLUB (Sample)	1.06

Table 1. Performance on the Permutation invariant MNIST classification. Different MI estimators are applied for the minimization of $\mathcal{I}(\mathbf{x}; \mathbf{z})$ in the Information Bottleneck. Misclassification rates of learned latent representation \mathbf{z} are reported. The top three methods are MI lower bounds, while the rest are MI upper bounds.

mator for minimizing $\mathcal{I}(\mathbf{x}; \mathbf{z})$. Alternatively, the variational CLUB can be also applied under this scenario. Besides CLUB and vCLUB, we compare previous methods such as MINE, NWJ, InfoNCE, and L1Out. The misclassification rates for different MI estimators are reported in Table 1.

MINE achieves the lowest misclassification error among lower bound estimators. Although providing good MI estimation in the Gaussian simulation study, L1Out suffers from numerical instability in MI optimization and fails during training. Both CLUB and vCLUB estimators outperform previous methods in bottleneck representation learning, with lower misclassification rates. Note that sampled versions of CLUB and vCLUB improve the accuracy compared with original CLUB and vCLUB, respectively, which verify the claim the negative sampling strategy improves model’s generalization ability. Besides, using variational approximation $q_\theta(\mathbf{y}|\mathbf{x})$ even attains higher accuracy than using ground truth $p_\sigma(\mathbf{y}|\mathbf{x})$ for CLUB. Although $p_\sigma(\mathbf{y}|\mathbf{x})$ provides more accurate MI estimation, the variational approximation $p_\sigma(\mathbf{y}|\mathbf{x})$ can add noise into the gradient of CLUB. Both the sampling and the variational approximation increase the randomness in the model, which helps to increase the model generalization ability (Hinton et al., 2012; Belghazi et al., 2018).

4.4. MI Minimization in Domain Adaptation

Another important application of MI minimization is disentangled representation learning (DRL) (Kim & Mnih, 2018; Chen et al., 2018; Locatello et al., 2019). Specifically, we aim to encode the data into several separate embedding parts, each with different semantic meanings. The semantically disentangled representations help improve the performance of deep learning models, especially in the fields of conditional generation (Ma et al., 2018), style transfer (John et al., 2019), and domain adaptation (Gholami et al., 2018). To learn (ideally) independent disentangled representations, one effective solution is to minimize the mutual information

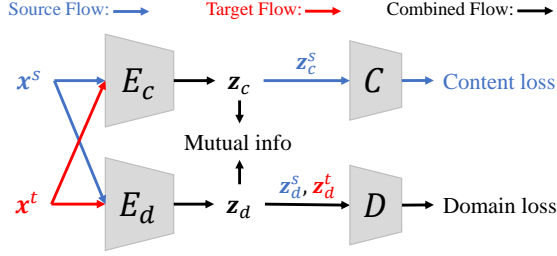


Figure 4. The information-theoretical framework for unsupervised domain adaptation. The input data \mathbf{x} (including \mathbf{x}^s and \mathbf{x}^t) are passed to a content encoder E_c and a domain encoder E_d , with output feature \mathbf{z}_c and \mathbf{z}_d , respectively. C is the content classifier, and D is the domain discriminator. The mutual information between \mathbf{z}_c and \mathbf{z}_d is minimized.

among different latent embedding parts.

We compare performance of MI estimators for learning disentangled representations in unsupervised domain adaptation (UDA) tasks. In UDA, we have images $\mathbf{x}^s \in \mathcal{X}^s$ from the source domain \mathcal{X}^s and $\mathbf{x}^t \in \mathcal{X}^t$ from the target domain \mathcal{X}^t . While each source image \mathbf{x}^s has a corresponding label y^s , no label information is available for observations in the target domain. The objective is to learn a model based on data $\{\mathbf{x}^s, y^s\}$ and $\{\mathbf{x}^t\}$, which not only performs well in source domain classification, but also provides satisfying predictions in the target domain.

To solve this problem, we use the information-theoretical framework inspired from [Gholami et al. \(2018\)](#). Specifically, two feature extractors are introduced: the domain encoder E_d and the content encoder E_c . The former encodes the domain information from an observation \mathbf{x} into a domain embedding $\mathbf{z}_d = E_d(\mathbf{x})$; the latter outputs a content embedding $\mathbf{z}_c = E_c(\mathbf{x})$ based on an input data point \mathbf{x} . As shown in Figure 4, the content embedding \mathbf{z}_c^s from the source domain is further used as an input to a content classifier $C(\cdot)$ to predict the corresponding class label, with a content loss defined as $\mathcal{L}_c = \mathbb{E}[-y^s \log C(\mathbf{z}_c^s)]$. The domain embedding \mathbf{z}_d (including \mathbf{z}_d^s and \mathbf{z}_d^t) is input to a domain discriminator $D(\cdot)$ to predict whether the observation comes from the source domain or target domain, with a domain loss defined as $\mathcal{L}_d = \mathbb{E}_{\mathbf{x} \in \mathcal{X}^s} [\log D(\mathbf{z}_d)] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}^t} [\log(1 - D(\mathbf{z}_d))]$. Since the content information and the domain information should be independent, we minimize the mutual information $I(\mathbf{z}_c, \mathbf{z}_d)$ between the content embedding \mathbf{z}_c and domain embedding \mathbf{z}_d . The final objective is (shown in Figure 4):

$$\min_{E_c, E_d, C, D} I(\mathbf{z}_c, \mathbf{z}_d) + \lambda_c \mathcal{L}_c + \lambda_d \mathcal{L}_d, \quad (18)$$

where $\lambda_c, \lambda_d > 0$ are hyper-parameters.

We apply different MI estimators to the framework (18), and evaluate the performance on several DA benchmark datasets, including MNIST, MNIST-M, USPS, SVHN, CIFAR-10, and STL. Detailed description to the datasets and model se-

Method	M→MM	M→U	U→M	SV→M	C→S	S→C
Source-Only	59.9	76.7	63.4	67.1	-	-
MI-based Disentangling Framework						
NWJ	83.3	98.3	91.1	86.5	78.2	71.0
MINE	88.4	98.1	94.8	83.4	77.9	70.5
InfoNCE	85.5	98.3	92.7	84.1	77.4	69.4
VUB	76.4	97.1	96.3	81.5	-	-
L1Out	76.2	96.3	93.9	-	77.8	69.2
CLUB	93.7	98.9	97.7	89.7	78.7	71.8
CLUB-S	94.6	98.9	98.1	90.6	79.1	72.3
Other Frameworks						
DANN	81.5	77.1	73.0	71.1	-	-
DSN	83.2	91.3	-	76.0	-	-
MCD	93.5	94.2	94.1	92.6	78.1	69.2

Table 2. Performance comparison on UDA. Datasets are MNIST (M), MNIST-M (MM), USPS (U), SVHN (SV), CIFAR-10 (C), and STL (S). Classification accuracy on target domain is reported. Among results in MI-based disentangling framework, the top three are MI lower bounds, while the rest are MI upper bounds. CLUB-S refers to Sampled CLUB.

tups is in the Supplementary Material. Besides the proposed information-theoretical UDA model, we also compare the performance with other UDA frameworks: DANN ([Ganin et al., 2016](#)), DSN ([Bousmalis et al., 2016](#)), and MCD ([Saito et al., 2018](#)). The numerical results are shown in Table 2. From the results, we find our MI-based disentangling shows competitive results with previous UDA methods. Among different MI estimators, the Sampled CLUB uniformly outperforms other competitive methods on four DA tasks. The stochastic sampling in CLUBSample improves the model generalization ability and preserves the model from overfitting. The other two MI upper bounds, VUB and L1Out, fail to train a satisfying UDA model, whose results are worse than the MI lower bound estimators. With L1Out, the training loss cannot even decrease on the most challenging SVHN→MNIST task, due to the numerical instability.

5. Conclusions

We have introduced a novel mutual information upper bound called Contrastive Log-ratio Upper Bound (CLUB). This novel MI estimator can be extended to a variational version for general scenarios when only samples of the joint distribution are obtainable. Based on the variational CLUB, we have proposed a new MI minimization algorithm, and further accelerated it with a negative sampling strategy. We have studied the good properties of CLUB both theoretically and empirically. Experimental results on simulation studies and real-world applications show the attractive performance of CLUB on both MI estimation and MI minimization tasks. This work provides an insight on the connection between mutual information and widespread machine learning training strategies, including contrastive learning and negative

sampling. We believe the proposed CLUB estimator will have vast applications for reducing the correlation of different model parts, especially in the domains of interpretable machine learning, controllable generation, and fairness.

Acknowledgements

Thanks to Dongruo Zhou from UCLA for helpful discussions on network expressiveness. The portion of this work performed at Duke University was supported in part by DARPA, DOE, NIH, NSF and ONR.

References

- Ahn, S., Hu, S. X., Damianou, A., Lawrence, N. D., and Dai, Z. Variational information distillation for knowledge transfer. In *CVPR*, 2019.
- Alemi, A. A., Fischer, I., Dillon, J. V., and Murphy, K. Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- Barber, D. and Agakov, F. V. The im algorithm: a variational approach to information maximization. In *NeurIPS*, 2003.
- Belghazi, M. I., Baratin, A., Rajeshwar, S., Ozair, S., Bengio, Y., Hjelm, D., and Courville, A. Mutual information neural estimation. In *ICML*, 2018.
- Bousmalis, K., Trigeorgis, G., Silberman, N., Krishnan, D., and Erhan, D. Domain separation networks. In *NeurIPS*, 2016.
- Charrow, B., Liu, S., Kumar, V., and Michael, N. Information-theoretic mapping using cauchy-schwarz quadratic mutual information. In *ICRA*, 2015.
- Chen, T. Q., Li, X., Grosse, R. B., and Duvenaud, D. K. Isolating sources of disentanglement in variational autoencoders. In *NeurIPS*, 2018.
- Chen, X., Duan, Y., Houthoofd, R., Schulman, J., Sutskever, I., and Abbeel, P. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NeurIPS*, 2016.
- Cheng, P., Min, M. R., Shen, D., Malon, C., Zhang, Y., Li, Y., and Carin, L. Improving disentangled text representation learning with information-theoretic guidance. *arXiv preprint arXiv:2006.00693*, 2020.
- Dai, S., Cheng, Y., Zhang, Y., Gan, Z., Liu, J., and Carin, L. Contrastively smoothed class alignment for unsupervised domain adaptation. *arXiv preprint arXiv:1909.05288*, 2019.
- Darbellay, G. A. and Vajda, I. Estimation of the information by an adaptive partitioning of the observation space. *IEEE Transactions on Information Theory*, 1999.
- Florensa, C., Duan, Y., and Abbeel, P. Stochastic neural networks for hierarchical reinforcement learning. *arXiv preprint arXiv:1704.03012*, 2017.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. Domain-adversarial training of neural networks. *JMLR*, 2016.
- Gholami, B., Sahu, P., Rudovic, O., Bousmalis, K., and Pavlovic, V. Unsupervised multi-target domain adaptation: An information theoretic approach. *arXiv preprint arXiv:1810.11547*, 2018.
- Granger, C. and Lin, J.-L. Using the mutual information coefficient to identify lags in nonlinear models. *Journal of time series analysis*, 1994.
- Grover, A. and Leskovec, J. node2vec: Scalable feature learning for networks. In *KDD*, 2016.
- Gutmann, M. and Hyvärinen, A. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, 2010.
- Härdle, W. K., Müller, M., Sperlich, S., and Werwatz, A. *Nonparametric and Semiparametric Models*. Springer Science & Business Media, 2004.
- Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- Hjelm, R. D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., and Bengio, Y. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018.
- Hu, W., Miyato, T., Tokui, S., Matsumoto, E., and Sugiyama, M. Learning discrete representations via information maximizing self-augmented training. In *ICML*, 2017.
- Hu, W., Li, Z., and Yu, D. Understanding generalization of deep neural networks trained with noisy labels. *arXiv preprint arXiv:1905.11368*, 2019.
- Jiang, B., Ye, C., and Liu, J. S. Nonparametric k-sample tests via dynamic slicing. *Journal of the American Statistical Association*, 2015.
- John, V., Mou, L., Bahuleyan, H., and Vechtomova, O. Disentangled representation learning for non-parallel text style transfer. In *ACL*, 2019.
- Julian, B. J., Karaman, S., and Rus, D. On mutual information-based control of range sensing robots for mapping applications. *The International Journal of Robotics Research*, 2014.

- Kamishima, T., Akaho, S., and Sakuma, J. Fairness-aware learning through regularization approach. In *IEEE 11th International Conference on Data Mining Workshops*, 2011.
- Kazemi, H., Soleymani, S., Taherkhani, F., Iranmanesh, S., and Nasrabadi, N. Unsupervised image-to-image translation using domain-specific variational information bound. In *NeurIPS*, 2018.
- Kim, H. and Mnih, A. Disentangling by factorising. In *ICML*, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Kraskov, A., Stögbauer, H., and Grassberger, P. Estimating mutual information. *Physical review E*, 2004.
- Kullback, S. *Information theory and statistics*. Courier Corporation, 1997.
- Lachmann, A., Giorgi, F. M., Lopez, G., and Califano, A. Aracne-ap: gene network reverse engineering through adaptive partitioning inference of mutual information. *Bioinformatics*, 2016.
- Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., and Bachem, O. Challenging common assumptions in the unsupervised learning of disentangled representations. In *ICML*, 2019.
- Ma, L., Sun, Q., Georgoulis, S., Van Gool, L., Schiele, B., and Fritz, M. Disentangled person image generation. In *CVPR*, 2018.
- Magdon-Ismael, M. and Atiya, A. F. Neural networks for density estimation. In *NeurIPS*, 1999.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- Nguyen, X., Wainwright, M. J., and Jordan, M. I. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 2010.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Oymak, S. and Soltanolkotabi, M. Overparameterized non-linear learning: Gradient descent takes the shortest path? In *ICML*, 2019.
- Poole, B., Ozair, S., Van Den Oord, A., Alemi, A., and Tucker, G. On variational bounds of mutual information. In *ICML*, 2019.
- Saito, K., Watanabe, K., Ushiku, Y., and Harada, T. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 2018.
- Suzuki, T., Sugiyama, M., Sese, J., and Kanamori, T. Approximating mutual information by maximum likelihood density ratio estimation. In *New challenges for feature selection in data mining and knowledge discovery*, 2008.
- Tishby, N., Pereira, F. C., and Bialek, W. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.
- Zea, D. J., Anfossi, D., Nielsen, M., and Marino-Buslje, C. Mitos. jl: mutual information tools for protein sequence analysis in the julia language. *Bioinformatics*, 2016.

A. Proofs of Theorems

Proof of Theorem 3.2. We calculate the gap between I_{vCLUB} and $I(\mathbf{x}; \mathbf{y})$:

$$\begin{aligned}
 \tilde{\Delta} &:= I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y}) - I(\mathbf{x}; \mathbf{y}) \\
 &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q_\theta(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{p(\mathbf{y})} [\log q_\theta(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x}) - \log p(\mathbf{y})] \\
 &= [\mathbb{E}_{p(\mathbf{y})} [\log p(\mathbf{y})] - \mathbb{E}_{p(\mathbf{x})p(\mathbf{y})} [\log q_\theta(\mathbf{y}|\mathbf{x})]] - [\mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log q_\theta(\mathbf{y}|\mathbf{x})]] \\
 &= \mathbb{E}_{p(\mathbf{x})p(\mathbf{y})} [\log \frac{p(\mathbf{y})}{q_\theta(\mathbf{y}|\mathbf{x})}] - \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log \frac{p(\mathbf{y}|\mathbf{x})}{q_\theta(\mathbf{y}|\mathbf{x})}] \\
 &= \mathbb{E}_{p(\mathbf{x})p(\mathbf{y})} [\log \frac{p(\mathbf{x})p(\mathbf{y})}{q_\theta(\mathbf{y}|\mathbf{x})p(\mathbf{x})}] - \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{q_\theta(\mathbf{y}|\mathbf{x})p(\mathbf{x})}] \\
 &= \text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) - \text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})).
 \end{aligned}$$

Therefore, $I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y})$ is an upper bound of $I(\mathbf{x}; \mathbf{y})$ if and only if $\text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) \geq \text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y}))$.

If \mathbf{x} and \mathbf{y} are independent, $p(\mathbf{x})p(\mathbf{y}) = p(\mathbf{x}, \mathbf{y})$. Then, $\text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) = \text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y}))$ and $\tilde{\Delta} = 0$. Therefore, $I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y}) = I(\mathbf{x}; \mathbf{y})$, the equality holds. \square

Proof of Corollary 3.3. If $\text{KL}(p(\mathbf{y}|\mathbf{x}) \| q_\theta(\mathbf{y}|\mathbf{x})) \leq \epsilon$, then

$$\text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log \frac{p(\mathbf{x}, \mathbf{y})}{q_\theta(\mathbf{x}, \mathbf{y})}] = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log \frac{p(\mathbf{y}|\mathbf{x})}{q_\theta(\mathbf{y}|\mathbf{x})}] = \text{KL}(p(\mathbf{y}|\mathbf{x}) \| q_\theta(\mathbf{y}|\mathbf{x})) \leq \epsilon.$$

By the condition $\text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) > \text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y}))$, we have $\text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) < \epsilon$.

Note that the KL-divergence is always non-negative. From the proof of Theorem 3.2,

$$\begin{aligned}
 |I_{\text{vCLUB}}(\mathbf{x}; \mathbf{y}) - I(\mathbf{x}; \mathbf{y})| &= |\text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) - \text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y}))| \\
 &< \max \{ \text{KL}(p(\mathbf{x})p(\mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})), \text{KL}(p(\mathbf{x}, \mathbf{y}) \| q_\theta(\mathbf{x}, \mathbf{y})) \} \leq \epsilon,
 \end{aligned}$$

which supports the claim. \square

B. Network Expressiveness in Variational Inference

In Section 3.2, when analyze the properties of the vCLUB estimator, we claim a reasonable assumption that with high expressiveness of the neural network $q_\theta(\mathbf{y}|\mathbf{x})$, we can achieve $\text{KL}(p(\mathbf{y}|\mathbf{x}) \| q_\theta(\mathbf{y}|\mathbf{x})) < \epsilon$. Here we provide a analysis under the scenario that the conditional distribution is a Gaussian distribution, $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}^*(\mathbf{x}), \mathbf{I})$. The variational approximation $q_\theta(\mathbf{y}|\mathbf{x})$ is parameterized by $q_\theta(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}), \mathbf{I})$.

Then training samples pair $(\mathbf{x}_i, \mathbf{y}_i)$ can be treated as $(\mathbf{x}_i, \boldsymbol{\mu}^*(\mathbf{x}_i) + \boldsymbol{\xi}_i)$, where $\boldsymbol{\xi}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then

$$\begin{aligned}
 \log p(\mathbf{y}|\mathbf{x}) &= \log \prod_{d=1}^D [\frac{1}{\sqrt{2\pi}} e^{(y^{(d)} - \mu^{*(d)}(\mathbf{x}))^2/2}] = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \|\mathbf{y} - \boldsymbol{\mu}^*(\mathbf{x})\|^2, \\
 \log q_\theta(\mathbf{y}|\mathbf{x}) &= \log \prod_{d=1}^D [\frac{1}{\sqrt{2\pi}} e^{(y^{(d)} - \mu_\theta^{(d)}(\mathbf{x}))^2/2}] = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \|\mathbf{y} - \boldsymbol{\mu}_\theta(\mathbf{x})\|^2.
 \end{aligned}$$

The log-ratio between $p(\mathbf{y}_i|\mathbf{x}_i)$ and $q_\theta(\mathbf{y}_i|\mathbf{x}_i)$ is

$$\log \frac{p(\mathbf{y}_i|\mathbf{x}_i)}{q_\theta(\mathbf{y}_i|\mathbf{x}_i)} = \log p(\mathbf{y}_i|\mathbf{x}_i) - \log q_\theta(\mathbf{y}_i|\mathbf{x}_i) = [\boldsymbol{\mu}^*(\mathbf{x}_i) - \boldsymbol{\mu}_\theta(\mathbf{x}_i)]^T [\mathbf{y}_i - \boldsymbol{\mu}_\theta(\mathbf{x}_i) + \boldsymbol{\xi}_i].$$

We further assume $\|\boldsymbol{\mu}^*(\mathbf{x}) - \boldsymbol{\mu}_\theta(\mathbf{x})\| < A$ is bounded. Then $|\log p(\mathbf{y}_i|\mathbf{x}_i) - \log q_\theta(\mathbf{y}_i|\mathbf{x}_i)| < A \|\mathbf{y}_i - \boldsymbol{\mu}_\theta(\mathbf{x}_i) + \boldsymbol{\xi}_i\|$.

Denote a loss function $l(\boldsymbol{\mu}_\theta(\mathbf{x}_i), \mathbf{y}_i) = \|\mathbf{y}_i - \boldsymbol{\mu}_\theta(\mathbf{x}_i) + \boldsymbol{\xi}_i\|$. With all reasonable assumptions in Hu et al. (2019), and applying the Theorem 5.1 in Hu et al. (2019), we know that when the number of samples $n \rightarrow \infty$, the expected error $\mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [l(\boldsymbol{\mu}_\theta(\mathbf{x}), \mathbf{y})] \rightarrow \infty$ with probability $1 - \delta$.

$$\text{KL}(p(\mathbf{y}|\mathbf{x}) \| q_\theta(\mathbf{y}|\mathbf{x})) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x}) - \log q_\theta(\mathbf{y}|\mathbf{x})] < A \cdot \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [l(\boldsymbol{\mu}_\theta(\mathbf{x}), \mathbf{y})].$$

Therefore, when given a small number $\varepsilon > 0$, having the sample size n large enough, we can guarantee that $\text{KL}(p(\mathbf{y}|\mathbf{x})\|q_\theta(\mathbf{y}|\mathbf{x}))$ is smaller than ε .

C. Properties of Variational Upper Bounds

In the Section 2, we introduce two variational MI upper bounds with neural network approximation $q_\theta(\mathbf{y}|\mathbf{x})$ to $p(\mathbf{y}|\mathbf{x})$:

$$\begin{aligned} \text{I}_{\text{vVUB}}(\mathbf{x}; \mathbf{y}) &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \frac{q_\theta(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \right], \\ \text{I}_{\text{vL1Out}}(\mathbf{x}; \mathbf{y}) &= \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \left[\log \frac{q_\theta(\mathbf{y}_i|\mathbf{x}_i)}{\frac{1}{N-1} \sum_{j \neq i} q_\theta(\mathbf{y}_i|\mathbf{x}_j)} \right] \right]. \end{aligned}$$

With the neural approximation $q_\theta(\mathbf{y}|\mathbf{x})$, I_{vVUB} and I_{vL1Out} no longer guarantee to be the MI upper bounds. However, both of the two estimators have good properties with a good approximation $q_\theta(\mathbf{y}|\mathbf{x})$.

Theorem C.1. *If $q_\theta(\mathbf{y}|\mathbf{x})$ satisfies $\text{KL}(p(\mathbf{y}|\mathbf{x})\|q_\theta(\mathbf{y}|\mathbf{x})) \leq \text{KL}(p(\mathbf{y})\|r(\mathbf{y}))$, then $\text{I}(\mathbf{x}; \mathbf{y}) \leq \text{I}_{\text{vVUB}}(\mathbf{x}; \mathbf{y})$.*

Proof of Theorem C.1. With the conditional $\text{KL}(p(\mathbf{y}|\mathbf{x})\|q_\theta(\mathbf{y}|\mathbf{x})) \leq \text{KL}(p(\mathbf{y})\|r(\mathbf{y}))$,

$$\begin{aligned} \text{I}(\mathbf{x}; \mathbf{y}) &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y})} \right] = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \left(\frac{p(\mathbf{y}|\mathbf{x})}{q_\theta(\mathbf{y}|\mathbf{x})} \cdot \frac{q_\theta(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \cdot \frac{r(\mathbf{y})}{p(\mathbf{y})} \right) \right] \\ &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \frac{q_\theta(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \right] + \text{KL}(p(\mathbf{y}|\mathbf{x})\|q_\theta(\mathbf{y}|\mathbf{x})) - \text{KL}(p(\mathbf{y})\|r(\mathbf{y})) \leq \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \frac{q_\theta(\mathbf{y}|\mathbf{x})}{r(\mathbf{y})} \right]. \end{aligned}$$

□

Theorem C.2. *Given $N - 1$ samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}$ from the marginal $p(\mathbf{x})$, If*

$$\text{KL}(p(\mathbf{y}|\mathbf{x})\|q_\theta(\mathbf{y}|\mathbf{x})) \leq \mathbb{E}_{\mathbf{x}_i \sim p(\mathbf{x})} \left[\text{KL} \left(p(\mathbf{y}) \left\| \frac{1}{N-1} \sum_{i=1}^{N-1} q_\theta(\mathbf{y}|\mathbf{x}_i) \right\| \right) \right],$$

then $\text{I}(\mathbf{x}; \mathbf{y}) \leq \text{I}_{\text{vL1Out}}(\mathbf{x}; \mathbf{y})$.

Proof. Assume we have N sample pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ drawn from $p(\mathbf{x}, \mathbf{y})$, then

$$\begin{aligned} \text{I}(\mathbf{x}; \mathbf{y}) &= \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i) \sim p(\mathbf{x}, \mathbf{y})} \left[\frac{1}{N} \sum_{i=1}^N \left[\log \frac{p(\mathbf{y}_i|\mathbf{x}_i)}{p(\mathbf{y}_i)} \right] \right] \\ &= \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i) \sim p(\mathbf{x}, \mathbf{y})} \left[\frac{1}{N} \sum_{i=1}^N \left[\log \left(\frac{p(\mathbf{y}_i|\mathbf{x}_i)}{q_\theta(\mathbf{y}_i|\mathbf{x}_i)} \cdot \frac{q_\theta(\mathbf{y}_i|\mathbf{x}_i)}{\frac{1}{N-1} \sum_{j \neq i} q_\theta(\mathbf{y}_i|\mathbf{x}_j)} \cdot \frac{\frac{1}{N-1} \sum_{j \neq i} q_\theta(\mathbf{y}_i|\mathbf{x}_j)}{p(\mathbf{y}_i)} \right) \right] \right] \\ &= \text{KL}(p(\mathbf{y}|\mathbf{x})\|q_\theta(\mathbf{y}|\mathbf{x})) + \text{I}_{\text{vVUB}}(\mathbf{x}; \mathbf{y}) - \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \text{KL} \left(p(\mathbf{y}) \left\| \frac{1}{N-1} \sum_{j \neq i} q_\theta(\mathbf{y}|\mathbf{x}_j) \right\| \right) \right]. \end{aligned}$$

Apply the condition in Theorem C.2 to each $N - 1$ combination of $\{\mathbf{x}_j\}_{j \neq i}$, we conclude $\text{I}(\mathbf{x}; \mathbf{y}) \leq \text{I}_{\text{vL1Out}}(\mathbf{x}; \mathbf{y})$. □

Theorem C.1 and Theorem C.2 indicate that if the approximation $q_\theta(\mathbf{y}|\mathbf{x})$ is good enough, the estimators I_{vVUB} and I_{vL1Out} can remain as MI upper bounds. Based on the analysis in Section B, when implemented with neural networks, the approximation can be far more accurate to preserve the variational estimators as MI upper bounds.

D. Implementation Details

vCLUB with Gaussian Approximation When $q_\theta(\mathbf{y}|\mathbf{x})$ is parameterized by $\mathcal{N}(\mathbf{y}|\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\sigma}^2(\mathbf{x}) \cdot \mathbf{I})$, then given samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, we denote $\boldsymbol{\mu}_i = \boldsymbol{\mu}(\mathbf{x}_i)$ and $\boldsymbol{\sigma}_i = \boldsymbol{\sigma}(\mathbf{x}_i)$. Moreover, $\boldsymbol{\mu}_i = [\mu_i^{(1)}, \mu_i^{(2)}, \dots, \mu_i^{(D)}]^\text{T}$, $\boldsymbol{\sigma}_i =$

$[\sigma_i^{(1)}, \sigma_i^{(2)}, \dots, \sigma_i^{(D)}]^\top$, are D -dimensional vectors as $\mathbf{y}_i = [y_i^{(1)}, y_i^{(2)}, \dots, y_i^{(D)}]^\top$. Then the conditional distribution

$$q_\theta(\mathbf{y}_j|\mathbf{x}_i) = \prod_{d=1}^D (2\pi(\sigma_i^{(d)})^2)^{-1/2} \exp \left\{ -\frac{(y_j^{(d)} - \mu_i^{(d)})^2}{2(\sigma_i^{(d)})^2} \right\}. \quad (19)$$

Therefore, the log-ratio

$$\begin{aligned} \log q_\theta(\mathbf{y}_i|\mathbf{x}_i) - \log q_\theta(\mathbf{y}_j|\mathbf{x}_i) &= \log \left(\prod_{d=1}^D (2\pi(\sigma_i^{(d)})^2)^{-1/2} \right) + \log \left(\prod_{d=1}^D \exp \left\{ -\frac{(y_i^{(d)} - \mu_i^{(d)})^2}{2(\sigma_i^{(d)})^2} \right\} \right) \\ &\quad - \log \left(\prod_{d=1}^D (2\pi(\sigma_i^{(d)})^2)^{-1/2} \right) - \log \left(\prod_{d=1}^D \exp \left\{ -\frac{(y_j^{(d)} - \mu_i^{(d)})^2}{2(\sigma_i^{(d)})^2} \right\} \right) \\ &= \sum_{d=1}^D \left\{ -\frac{(y_i^{(d)} - \mu_i^{(d)})^2}{2(\sigma_i^{(d)})^2} \right\} - \sum_{d=1}^D \left\{ -\frac{(y_j^{(d)} - \mu_i^{(d)})^2}{2(\sigma_i^{(d)})^2} \right\} \\ &= -\frac{1}{2}(\mathbf{y}_i - \boldsymbol{\mu}_i)^\top \text{Diag}[\boldsymbol{\sigma}_i^{-2}](\mathbf{y}_i - \boldsymbol{\mu}_i) + \frac{1}{2}(\mathbf{y}_j - \boldsymbol{\mu}_i)^\top \text{Diag}[\boldsymbol{\sigma}_i^{-2}](\mathbf{y}_j - \boldsymbol{\mu}_i), \end{aligned}$$

where $\text{Diag}[\boldsymbol{\sigma}_i^{-2}]$ is a $D \times D$ diagonal matrix with $(\text{Diag}[\boldsymbol{\sigma}_i^{-2}])_{d,d} = (\sigma_i^{(d)})^{-2}$, $d = 1, 2, \dots, D$. The vCLUB estimator can be calculated by

$$\begin{aligned} \hat{\mathbb{I}}_{\text{vCLUB}} &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N [\log q_\theta(\mathbf{y}_i|\mathbf{x}_i) - \log q_\theta(\mathbf{y}_j|\mathbf{x}_i)] \\ &= -\frac{1}{2} \left\{ \frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \boldsymbol{\mu}_i)^\top \text{Diag}[\boldsymbol{\sigma}_i^{-2}](\mathbf{y}_i - \boldsymbol{\mu}_i) \right\} + \frac{1}{2} \left\{ \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\mathbf{y}_j - \boldsymbol{\mu}_i)^\top \text{Diag}[\boldsymbol{\sigma}_i^{-2}](\mathbf{y}_j - \boldsymbol{\mu}_i) \right\}. \end{aligned}$$

E. Detailed Experimental Setups

Information Bottleneck: For the experiment on information bottleneck, we follow the setup from [Aleml et al. \(2016\)](#). The parameters $\mu_\sigma(\mathbf{x})$ and $\Sigma_\sigma(\mathbf{x})$ are the output from a MLP with layers $784 \rightarrow 1024 \rightarrow 1024 \rightarrow 2K$, where K is the size of the bottleneck. We set $K = 256$. For the variational classifier to implement the Barber-Agakov MI lower bound, the structure is set to a one-layer MLP. The batch size is 100. We set our learning rate to 10^{-4} , with an exponential decay rate of 0.97 and a decay step of 1200.

Domain Adaptation: The network is constructed as follows. Both feature extractors (*i.e.*, E_c and E_d) are nine-layer convolutional neural network with leaky ReLU non-linearities. The content classifier C and the domain discriminator D are a one-layer and a two-layer MLPs, respectively. Images from each domain are normalized using Gaussian normalization.

Classifier C	Discriminator D	Extractor (both E_c and E_d)
Content feature \mathbf{z}_c^s	Domain feature \mathbf{z}_d	Input data \mathbf{x}
MLP output $C(\mathbf{z}_c^s)$ with shape 10	MLP output $D(\mathbf{z}_d)$ with shape 2	3×3 conv. 64 lReLU, stride 1 3×3 conv. 64 lReLU, stride 1 3×3 conv. 64 lReLU, stride 1 2×2 max pool, stride 2, dropout, $p = 0.5$, Gaussian noise, $\sigma = 1$ 3×3 conv. 64 lReLU, stride 1 3×3 conv. 64 lReLU, stride 1 3×3 conv. 64 lReLU, stride 1 2×2 max pool, stride 2, dropout, $p = 0.5$, Gaussian noise, $\sigma = 1$ 3×3 conv. 64 lReLU, stride 1 3×3 conv. 64 lReLU, stride 1 3×3 conv. 64 lReLU, stride 1 global average pool, output feature with shape 64

F. Numerical Results of MI Estimation

We report the numerical results of MI estimation quality in Table 3. The detailed setups are provided in Section 4.1. Our CLUB estimator has the lowest estimation error when the ground-truth MI value goes larger.

	Gaussian					Cubic				
MI true value	2	4	6	8	10	2	4	6	8	10
VUB	3.85	15.33	34.37	61.25	95.70	2.09	10.38	25.56	47.84	77.59
NWJ	1.67	7.20	17.46	33.26	55.34	1.10	5.54	14.68	30.25	51.07
MINE	1.61	6.66	16.01	29.60	49.87	1.53	6.58	17.4	34.20	59.46
NCE	0.59	2.85	8.56	19.66	37.79	0.45	1.89	6.70	17.48	35.86
L1Out	0.13	0.11	0.75	4.65	17.08	2.30	5.58	8.92	8.27	7.19
CLUB	0.15	0.12	0.70	4.53	16.57	2.22	5.89	8.25	8.23	6.93
CLUBSample	0.38	0.44	1.31	5.30	17.63	2.37	5.89	8.07	8.87	7.54

Table 3. MSE of MI estimation