

# SCALABLE MUTUAL INFORMATION ESTIMATION USING DEPENDENCE GRAPHS

Morteza Noshad, Yu Zeng, Alfred O. Hero III\*

University of Michigan, Electrical Engineering and Computer Science, Ann Arbor, Michigan, U.S.A

## ABSTRACT

The Mutual Information (MI) is an often used measure of dependency between two random variables utilized in information theory, statistics and machine learning. Recently several MI estimators have been proposed that can achieve parametric MSE convergence rate. However, most of the previously proposed estimators have high computational complexity of at least  $O(N^2)$ . We propose a unified method for empirical non-parametric estimation of general MI function between random vectors in  $\mathbb{R}^d$  based on  $N$  i.i.d. samples. The reduced complexity MI estimator, called the ensemble dependency graph estimator (EDGE), combines randomized locality sensitive hashing (LSH), dependency graphs, and ensemble bias-reduction methods. We prove that EDGE achieves optimal computational complexity  $O(N)$ , and can achieve the optimal parametric MSE rate of  $O(1/N)$  if the density is  $d$  times differentiable. To the best of our knowledge EDGE is the first non-parametric MI estimator that can achieve parametric MSE rates with linear time complexity. We illustrate the utility of EDGE for the analysis of the information plane (IP) in deep learning. Using EDGE we shed light on a controversy on whether or not the compression property of information bottleneck (IB) in fact holds for ReLu and other rectification functions in deep neural networks (DNN).

## 1. INTRODUCTION

The Mutual Information (MI) is an often used measure of dependency between two random variables or vectors [1], and it has a wide range of applications in information theory [1] and machine learning [2, 3]. Non-parametric MI estimation methods have been studied that use estimation strategies including KSG [4], KDE [5] and Parzen window density estimation [6]. The performance of these estimators has been evaluated and compared based on both empirical studies [7] and asymptotic analysis [8]. Recently several MI estimators have been proposed that can achieve parametric MSE rate of convergence. For example, in [9] a KDE plug-in estimator for Rényi divergence and mutual information achieves the MSE rate of  $O(1/N)$  when the densities are at least  $d$  times differentiable. Another KDE based mutual information estimator was proposed in [8] that can achieve the MSE rate of  $O(1/N)$  when

the densities are  $d/2$  times differentiable. Recently Moon et al [10] and Gao et al [11] respectively proposed KDE and KNN based MI estimators for random variables with mixtures of continuous and discrete components. Most of these estimators, however, have high computational cost and require knowledge of the density support boundary.

In this paper we propose a reduced complexity MI estimator called the ensemble dependency graph estimator (EDGE). The estimator combines randomized locality sensitive hashing (LSH), dependency graphs, and ensemble bias-reduction methods. A dependence graph is a bipartite directed graph consisting of two sets of nodes  $V$  and  $U$ . The data points are mapped to the sets  $V$  and  $U$  using a randomized LSH function  $H$  that depends on a hash parameter  $\epsilon$ . Each node is assigned a weight that is proportional to the number of hash collisions. Likewise, each edge between the vertices  $v_i$  and  $u_j$  has a weight proportional to the number of  $(X_k, Y_k)$  pairs mapped to the node pairs  $(v_i, u_j)$ . For a given value of the hash parameter  $\epsilon$ , a base estimator of MI is proposed as a weighted average of non-linearly transformed of the edge weights. The proposed EDGE estimator of MI is obtained by applying the method of weighted ensemble bias reduction [10, 12] to a set of base estimators with different hash parameters. This estimator is a non-trivial extension of the LSH divergence estimator defined in [13]. LSH-based methods have previously been used for KNN search and graph constructions problems [14, 15], and they result in fast and low complexity algorithms.

Recently, Shwartz-Ziv and Tishby utilized MI to study the training process in Deep Neural Networks (DNN) [16]. Let  $X$ ,  $T$  and  $Y$  respectively denote the input, hidden and output layers. The authors of [16] introduced the information bottleneck (IB) that represents the tradeoff between two mutual information measures:  $I(X, T)$  and  $I(T, Y)$ . They observed that the training process of a DNN consists of two distinct phases; 1) an initial fitting phase in which  $I(T, Y)$  increases, and 2) a subsequent compression phase in which  $I(X, T)$  decreases. Saxe *et al* in [17] countered the claim of [16], asserting that this compression property is not universal, rather it depends on the specific activation function. Specifically, they claimed that the compression property does not hold for ReLu activation functions. The authors of [16] challenged these claims, arguing that the authors of [17] had not observed compression due to poor estimates of the MI. We use our proposed rate-optimal ensemble MI estimator to explore this

\*This research was partially supported by ARO grant W911NF-15-1-0479.

controversy, observing that our estimator of MI does exhibit the compression phenomenon in the ReLU network studied by [17].

Our contributions are as follows:

- To the best of our knowledge the proposed MI estimator is the first estimator to have linear complexity and can achieve the optimal MSE rate of  $O(1/N)$ .
- The proposed MI estimator provides a simplified and unified treatment of mixed continuous-discrete variables. This is due to the hash function approach that is adopted.
- EDGE is applied to IB theory of deep learning, and provides evidence that the compression property does indeed occur in ReLU DNNs, contrary to the claims of [17].

The rest of the paper is organized as follows. In Section 2, we introduce the general definition of MI and define the dependence graph. In Section 3, we introduce the hash based MI estimator and give theory for the bias and variance. In section 4 we introduce the ensemble dependence graph MI estimator (EDGE) and show how the ensemble estimation method can be used to improve the convergence rates. Finally, in Section 5 we provide numerical results as well as study the IP in DNNs.

## 2. MUTUAL INFORMATION

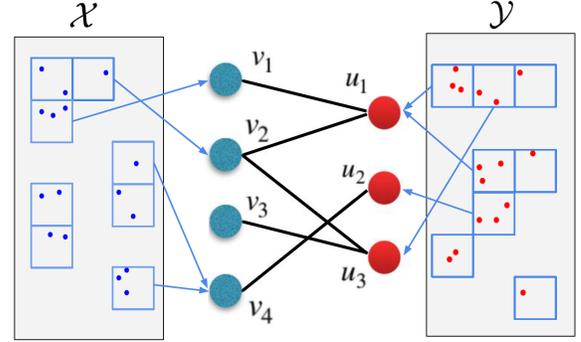
In this section, we introduce the general mutual information function based on the f-divergence measure. Then, we define a consistent estimator for the mutual information function. Consider the probability measures  $P$  and  $Q$  on a Euclidean space  $\mathcal{X}$ . Let  $g : (0, \infty) \rightarrow \mathbb{R}$  be a convex function with  $g(1) = 0$ . The f-divergence between  $P$  and  $Q$  can be defined as follows [18, 19].

$$D(P\|Q) := \mathbb{E}_Q \left[ g \left( \frac{dP}{dQ} \right) \right]. \quad (1)$$

**Mutual Information:** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Euclidean spaces and let  $P_{XY}$  be a probability measure on the space  $\mathcal{X} \times \mathcal{Y}$ . For any measurable sets  $A \subseteq \mathcal{X}$  and  $B \subseteq \mathcal{Y}$ , we define the marginal probability measures  $P_X(A) := P_{XY}(A \times \mathcal{Y})$  and  $P_Y(B) := P_{XY}(\mathcal{X} \times B)$ . Similar to [11, 18], the general MI denoted by  $I(X, Y)$  is defined as

$$D(P_{XY}\|P_X P_Y) = \mathbb{E}_{P_X P_Y} \left[ g \left( \frac{dP_{XY}}{dP_X P_Y} \right) \right], \quad (2)$$

where  $\frac{dP_{XY}}{dP_X P_Y}$  is the Radon-Nikodym derivative, and  $g : (0, \infty) \rightarrow \mathbb{R}$  is, as in (1) a convex function with  $g(1) = 0$ . Shannon mutual information is a particular cases of (1) for which  $g(x) = x \log x$ .



**Fig. 1.** Sample dependence graph with 4 and 3 respective distinct hash values of  $\mathbf{X}$  and  $\mathbf{Y}$  data jointly encoded with LSH, and the corresponding dependency edges.

### 2.1. Dependence Graphs

Consider  $N$  i.i.d samples  $(X_i, Y_i)$ ,  $1 \leq i \leq N$  drawn from the probability measure  $P_{XY}$ , defined on the space  $\mathcal{X} \times \mathcal{Y}$ . Define the sets  $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$  and  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$ . The dependence graph  $G(\mathbf{X}, \mathbf{Y})$  is a directed bipartite graph, consisting of two sets of nodes  $V$  and  $U$  with cardinalities denoted as  $|V|$  and  $|U|$ , and the set of edges  $E_G$ . Each point in the sets  $\mathbf{X}$  and  $\mathbf{Y}$  is mapped to the nodes in the sets  $U$  and  $V$ , respectively, using the hash function  $H$ , described as follows.

A vector valued hash function  $H$  is defined in a similar way as defined in [13]. First, define the vector valued hash function  $H_1 : \mathbb{R}^d \rightarrow \mathbb{Z}^d$  as

$$H_1(x) = [h_1(x_1), h_1(x_2), \dots, h_1(x_d)], \quad (3)$$

where  $x_i$  denotes the  $i$ th component of the vector  $x$ . In (3), each scalar hash function  $h_1(x_i) : \mathbb{R} \rightarrow \mathbb{Z}$  is given by

$$h_1(x_i) = \left\lfloor \frac{x_i + b}{\epsilon} \right\rfloor, \quad (4)$$

for a fixed  $\epsilon > 0$ , where  $\lfloor y \rfloor$  denotes the floor function (the smallest integer value less than or equal to  $y$ ), and  $b$  is a fixed random variable in  $[0, \epsilon]$ . Let  $\mathcal{F} := \{1, 2, \dots, F\}$ , where  $F := c_H N$  and  $c_H$  is a fixed tunable integer. We define a random hash function  $H_2 : \mathbb{Z}^d \rightarrow \mathcal{F}$  with a uniform density on the output and consider the combined hashing function

$$H(x) := H_2(H_1(x)), \quad (5)$$

which maps the points in  $\mathbb{R}^d$  to  $\mathcal{F}$ .

$H(x)$  reveals the index of the mapped vertex in  $G(\mathbf{X}, \mathbf{Y})$ . The weights  $\omega_i$  and  $\omega'_j$  corresponding to the nodes  $v_i$  and  $u_j$ , and  $\omega_{ij}$ , the weight of the edge  $(v_i, u_j)$ , are defined as follows.

$$\omega_i = \frac{N_i}{N}, \quad \omega'_j = \frac{M_j}{N}, \quad \omega_{ij} = \frac{N_{ij} N}{N_i M_j}, \quad (6)$$

where  $N_i$  and  $M_j$  respectively are the the number of hash collisions at the vertices  $v_i$  and  $u_j$ , and  $N_{ij}$  is the number of joint

collisions of the nodes  $(X_k, Y_k)$  at the vertex pairs  $(v_i, u_j)$ . The number of hash collisions is defined as the number of instances of the input variables map to the same output value. In particular,

$$N_{ij} := \#\{(X_k, Y_k) \text{ s.t. } H(X_k) = i \text{ and } H(Y_k) = j\}. \quad (7)$$

Fig. 1 represents a sample dependence graph. Note that the nodes and edges with zero collisions do not show up in the dependence graph.

### 3. THE BASE ESTIMATOR OF MI

#### 3.1. Assumptions

The following are the assumptions we make on the probability measures and  $g$ :

- A1.** The support sets  $\mathcal{X}$  and  $\mathcal{Y}$  are bounded.
- A2.** The following supremum exists and is bounded:

$$\sup_{P_X P_Y} g\left(\frac{dP_{XY}}{dP_X P_Y}\right) \leq U.$$

**A3.** Let  $x_D$  and  $x_C$  respectively denote the discrete and continuous components of the vector  $x$ . Also let  $f_{X_C}(x_C)$  and  $p_{X_D}(x_D)$  respectively denote density and pmf functions of these components associated with the probability measure  $P_X$ . The density functions  $f_{X_C}(x_C)$ ,  $f_{Y_C}(y_C)$ ,  $f_{X_C Y_C}(x_C, y_C)$ , and the conditional densities  $f_{X_C|X_D}(x_C|x_D)$ ,  $f_{Y_C|Y_D}(y_C|y_D)$ ,  $f_{X_C Y_C|X_D Y_D}(x_C, y_C|x_D, y_D)$  are Hölder continuous.

**Hölder continuous functions:** Given a support set  $\mathcal{X}$ , a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is called Hölder continuous with parameter  $0 < \gamma \leq 1$ , if there exists a positive constant  $G_f$ , possibly depending on  $f$ , such that for every  $x \neq y \in \mathcal{X}$ ,

$$|f(y) - f(x)| \leq G_f \|y - x\|^\gamma. \quad (8)$$

**A4.** Assume that the function  $g$  in (2) is Lipschitz continuous; i.e.  $g$  is Hölder continuous with  $\gamma = 1$ .

#### 3.2. The Base Estimator of MI

For a fixed value of the hash parameter  $\epsilon$ , we propose the following base estimator of MI (2) function based on the dependence graph:

$$\hat{I}(X, Y) := \sum_{e_{ij} \in E_G} \omega_i \omega_j' \tilde{g}(\omega_{ij}), \quad (9)$$

where the summation is over all edges  $e_{ij} : (v_i \rightarrow u_j)$  of  $G(X, Y)$  having non-zero weight and  $\tilde{g}(x) := \max\{g(x), U\}$ .

When  $X$  and  $Y$  are strongly dependent, each point  $X_k$  hashed into the bucket (vertex)  $v_i$  corresponds to a unique hash value for  $Y_k$  in  $U$ . Therefore, asymptotically  $\omega_{ij} \rightarrow 1$  and the mutual information estimation in (9) takes its maximum value. On the other hand, when  $X$  and  $Y$  are independent, each point  $X_k$  hashed into the bucket (vertex)  $v_i$  may be associated with different values of  $Y_k$ , and therefore asymptotically  $\omega_{ij} \rightarrow \omega_j$  and the Shannon MI estimation tends to 0.

#### 3.3. Various LSH Functions

There are various types of LSH functions [20–22], and all of them share the common property that they map similar items to the same bins with high probability.

In equations (3) and (4) we considered a simple floor function on the scaled input, however in general, any other type of LSH might be used for our estimation method. In particular, the hash functions based on random projections can reduce the dimensionality of data. SimHash [20], which is based on cosine distance, and the LSH based on p-stable distributions [21] are among well known LSH functions that reduce the dimension of data. For example, the LSH based on p-stable distribution is defined similarly to the floor hash function in (3) and (4), except that the input vector is projected on random hyperplanes with p-stable distributions. The formal definition is  $H_{p\text{-stable}} : \mathbb{R}^d \rightarrow \mathbb{Z}^r$ ,

$$H_{p\text{-stable}}(x) = H_1(XW), \quad (10)$$

where  $H_1$  is defined in (3), and  $W$  is a  $d \times r$  matrix with entries chosen independently from a stable distribution. For high-dimensional datasets one can choose  $r \ll d$  in order to reduce the dimensionality. Finally, note that for theoretical analysis, we only focus on performance of the simple floor hash function defined in (3) and (4).

#### 3.4. Convergence Rates

In the following theorems we state upper bounds on the bias and variance rates of the proposed MI estimator (9). The proofs are given in appendices A and B. We define the notations  $\mathbb{B}[\hat{T}] = \mathbb{E}[\hat{T}] - T$  for bias and  $\mathbb{V}[\hat{T}] = \mathbb{E}[\hat{T}^2] - \mathbb{E}[\hat{T}]^2$  for variance of  $\hat{T}$ . The following theorem states an upper bound on the bias.

**Theorem 3.1.** *Let  $d = d_X + d_Y$  be the dimension of the joint random variable  $(X, Y)$ . Under the aforementioned assumptions **A1-A4**, and assuming that the density functions in **A3** have bounded derivatives up to order  $q \geq 0$ , the following upper bound on the bias of the estimator in (9) holds*

$$\mathbb{B}[\hat{I}(X, Y)] = \begin{cases} O(\epsilon^\gamma) + O\left(\frac{1}{N\epsilon^d}\right), & q = 0 \\ \sum_{i=1}^q C_i \epsilon^i + O(\epsilon^q) + O\left(\frac{1}{N\epsilon^d}\right) & q \geq 1, \end{cases} \quad (11)$$

where  $\epsilon$  is the hash parameter in (4),  $\gamma$  is the smoothness parameter in (8), and  $C_i$  are real constants.

In (11), the hash parameter,  $\epsilon$  needs to be a function of  $N$  to ensure that the bias converges to zero. For the case of  $q = 0$ , the optimum bias is achieved when  $\epsilon = \left(\frac{1}{N}\right)^{\gamma/(\gamma+d)}$ . When  $q \geq 1$ , the optimum bias is achieved for  $\epsilon = \left(\frac{1}{N}\right)^{1/(1+d)}$ .

**Theorem 3.2.** *Under the assumptions A1-A4 the variance of the proposed estimator can be bounded as  $\mathbb{V}[\hat{I}(X, Y)] \leq O(\frac{1}{N})$ . Further, the variance of the variable  $\omega_{ij}$  is also upper bounded by  $O(1/N)$ .*

### 3.5. Computational Complexity

We analyze the computational complexity of the proposed estimator. The procedure for estimation of MI equivalent to the proposed estimator in (9) is given in Algorithm 1.

---

**Algorithm 1:** MI Dependence Graph Estimator

---

**Input :**  $N$  i.i.d samples  $(X_k, Y_k), 1 \leq k \leq N$ .

```

1 for each  $k \in 1 : N$  do
2    $i \leftarrow H(\mathbf{X}_k)$ 
3    $j \leftarrow H(\mathbf{Y}_k)$ 
4    $N_i \leftarrow N_i + 1$ 
5    $M_j \leftarrow M_j + 1$ 
6    $N_{ij} \leftarrow N_{ij} + 1$ 
7  $\omega_i \leftarrow N_i/N; \omega'_j \leftarrow M_j/N; \omega_{ij} \leftarrow N_{ij}N/N_iM_j;$ 
    $\hat{I} \leftarrow \sum_{e_{ij}} \omega_i \omega'_j \tilde{g}(\omega_{ij})$ 

```

**Output :**  $\hat{I}$

---

We go over all of the data points and map them using the hash function  $H$ . We compute the number of marginal and joint collisions  $(N_i, M'_j, N_{ij})$  and based on these we compute the vertex and edge weights. Note that computing the hashing of all of the data points takes about  $O(N)$  time. Finally in the last line we compute the MI estimate by going over all of the edges  $e_{ij}$ . The number the edges is upper bounded by  $O(N)$ , since each edge correspond to at least one pair of  $(X_i, Y_i)$ . Finally, note that for high-dimensional data sets, the computational of computing the hash function of each input may depend on  $d$ , however, is would not be greater than  $O(d)$ . Hence, overall the the computational complexity of computing the proposed MI estimate is linear with respect to both  $N$  and  $d$ .

### 3.6. Comparison to the Other Estimation Methods

So far, various estimation methods for information measures (entropy, divergence and mutual information) have been proposed, most of which are based on kernel density estimates (KDE) [12], k-nearest neighbors (KNN) [23] or histogram binning [24]. Certain KDE and KNN based estimators can achieve the optimal parametric MSE rate [12, 23, 25], however, implementation of the KDE and KNN methods respectively have  $O(N^2)$  and  $O(kN \log N)$  computational complexity, where  $N$  is the number of samples. One could probably could approximation methods for KDE and KNN, however, there would be no theoretical guarantees for the estimation based on these approximations. Empirical histograms, on the other hand, are

simpler and easier to implement, however, their convergence rate is not as good as the KNN and KDE based methods [24].

LSH methods have previously been applied to the approximate nearest neighbor search, however, it has never been directly utilized for estimation of densities or information theoretic quantities. The simplest LSH function considered in (4) has similarities with the histogram estimator in terms of binning, however, there are also certain differences. As opposed to the LSH based method, the histogram binning requires a pre-knowledge of the support set or needs an extra computation to estimate the support set. The number of the bins in the histogram estimator gets exponentially large with increasing dimension which results in a huge computational complexity for high-dimensional datasets. In addition, most of the bins would be empty. In contrast, the LSH based method results in no empty hash bins, the number of the bins is upper bounded by  $O(N)$ , and the computational complexity is linear in dimension and the number of samples. The plug-in methods including histograms, KDE and KNN require the estimation of the densities  $p_X, p_Y$  and  $p_{XY}$  for computation of mutual information, while the proposed LSH-based method finds the mapping of  $\mathbf{X}$  and  $\mathbf{Y}$  data points, and then estimate the mutual information, based on the hash collisions. Finally, by giving an accurate bias and variance rates for the base LSH estimator, we use an ensemble estimation technique to achieve the optimal parametric convergence rate, discussed in the following section.

## 4. ENSEMBLE DEPENDENCE GRAPH ESTIMATOR (EDGE)

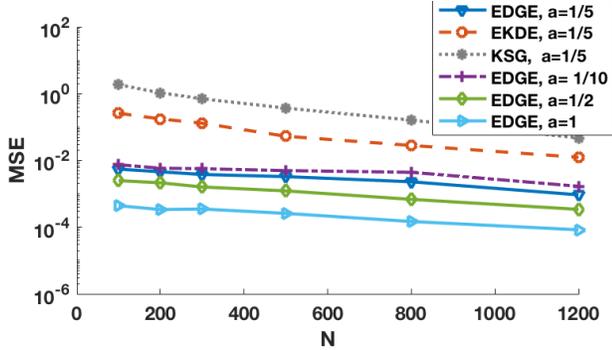
Given the expression for the bias in Theorem 3.1, the ensemble estimation technique proposed in [12] can be applied to improve the convergence rate of the MI estimator (9). Assume that the densities in A3 have continuous bounded derivatives up to the order  $q$ , where  $q \geq d$ . Let  $\mathcal{T} := \{t_1, \dots, t_T\}$  be a set of index values with  $t_i < c$ , where  $c > 0$  is a constant. Let  $\epsilon(t) := tN^{-1/2d}$ . For a given set of weights  $w(t)$  the weighted ensemble estimator is then defined as

$$\hat{I}_w := \sum_{t \in \mathcal{T}} w(t) \hat{I}_{\epsilon(t)}, \quad (12)$$

where  $\hat{I}_{\epsilon(t)}$  is the mutual information estimator with the parameter  $\epsilon(t)$ . Using (11), for  $q > 0$  the bias of the weighted ensemble estimator (12) takes the form

$$\mathbb{B}(\hat{I}_w) = \sum_{i=1}^q C_i N^{-\frac{i}{2d}} \sum_{t \in \mathcal{T}} w(t) t^i + O\left(\frac{t^d}{N^{1/2}}\right) + O\left(\frac{1}{N\epsilon^d}\right) \quad (13)$$

Given the form (13), as long as  $T \geq q$ , we can select the weights  $w(t)$  to force to zero the slowly decaying terms in (13), i.e.  $\sum_{t \in \mathcal{T}} w(t) t^i/d = 0$  subject to the constraint that  $\sum_{t \in \mathcal{T}} w(t) = 1$ . However,  $T$  should be strictly greater than  $q$  in order to control the variance, which is upper bounded



**Fig. 2.** MSE comparison of EDGE, EDKE and KSG Shannon MI estimators.  $X$  is a  $2D$  Gaussian random variable with unit covariance matrix.  $Y = X + aN_U$ , where  $N_U$  is a uniform noise. The MSE rates of EDGE, EKDE and KSG are compared for various values of  $a$ .

by the euclidean norm squared of the weights  $\omega$ . In particular we have the following theorem (the proof is given in Appendix C):

**Theorem 4.1.** For  $T > d$  let  $w_0$  be the solution to:

$$\begin{aligned} \min_w \quad & \|w\|_2 \\ \text{subject to} \quad & \sum_{t \in \mathcal{T}} w(t) = 1, \\ & \sum_{t \in \mathcal{T}} w(t)t^i = 0, i \in \mathbb{N}, i \leq d. \end{aligned} \quad (14)$$

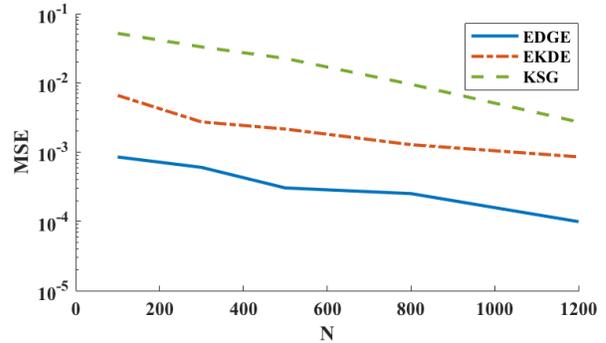
Then the MSE rate of the ensemble estimator  $\hat{I}_{w_0}$  is  $O(1/N)$ .

## 5. EXPERIMENTS

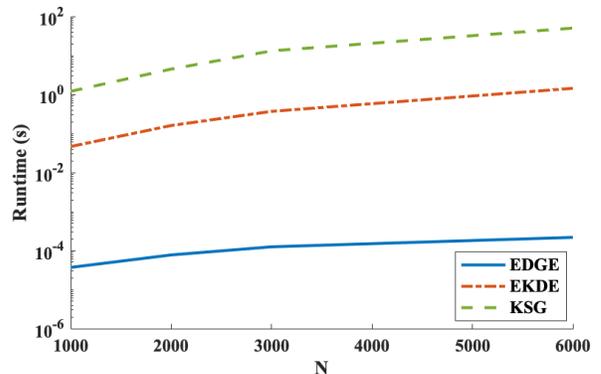
We first use simulated data to compare the proposed estimator to the competing MI estimators Ensemble KDE (EKDE) [10], and generalized KSG [11]. Both of these estimators work on mixed continuous-discrete variables. We also apply EDGE to study the information bottleneck in different networks trained on MNIST hand-written dataset.

Fig. 2, shows the MSE estimation rate of Shannon MI between the continuous random variables  $X$  and  $Y$  having the relation  $Y = X + aN_U$ , where  $X$  is a  $2D$  Gaussian random variable with the mean  $[0, 0]$  and covariance matrix  $C = I_2$ . Here  $I_d$  denote the  $d$ -dimensional identity matrix.  $N_U$  is a uniform random vector with the support  $\mathcal{N}_U = [0, 1] \times [0, 1]$ . We compute the MSE of each estimator for different sample sizes. The MSE rates of EDGE, EKDE and KSG are compared for  $a = 1/5$ . Further, the MSE rate of EDGE is investigated for noise levels of  $a = \{1/10, 1/5, 1/2, 1\}$ . As the dependency between  $X$  and  $Y$  increases the MSE rate becomes slower.

Fig. 3, shows the MSE estimation rate of Shannon MI between a discrete random variables  $X$  and a continuous random variable  $Y$ . We have  $X \in \{1, 2, 3, 4\}$ , and each  $X = x$  is associated with multivariate Gaussian random vector  $Y$ , with



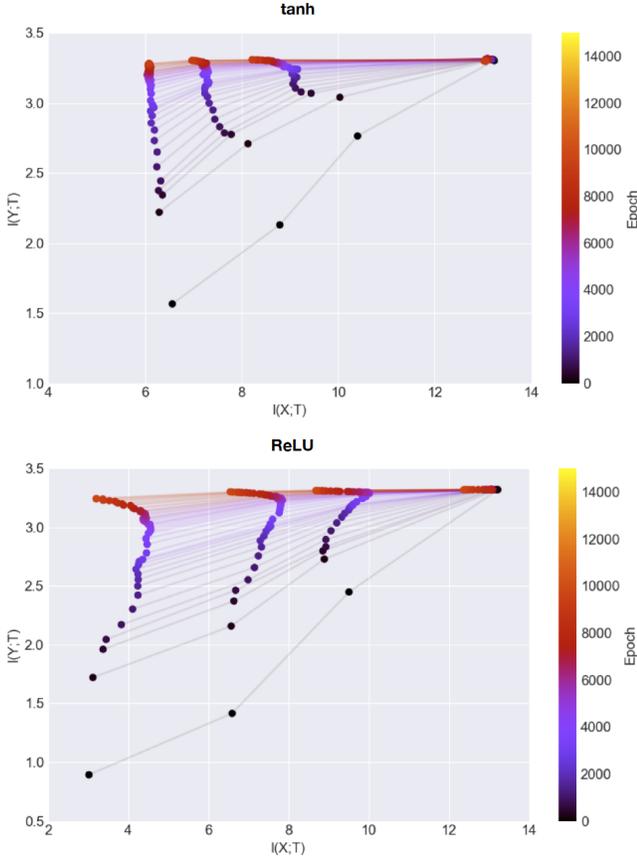
**Fig. 3.** MSE comparison of EDGE, EDKE and KSG Shannon MI estimators.  $X \in \{1, 2, 3, 4\}$ , and each  $X = x$  is associated with multivariate Gaussian random vector  $Y$ , with  $d = 4$ , the mean  $[x/2, 0, 0, 0]$  and covariance matrix  $C = I_4$ .



**Fig. 4.** Runtime comparison of EDGE, EDKE and KSG Shannon MI estimators.  $X \in \{1, 2, 3, 4\}$ , and each  $X = x$  is associated with multivariate Gaussian random vector  $Y$ , with  $d = 4$ , the mean  $[x/2, 0, 0, 0]$  and covariance matrix  $C = I_4$ .

$d = 4$ , the expectation  $[x/2, 0, 0, 0]$  and covariance matrix  $C = I_4$ . In general in Figures 2 and 3, EDGE has better convergence rate than EKDE and KSG estimators. Fig. 4 represents the runtime comparison for the same experiment as in Fig. 3. It can be seen from this graph how fast our proposed estimator performs compared to the other other methods.

Next, we use EDGE to study the information bottleneck [16] in DNNs. Fig. 5 represents the information plane of a DNN with 4 fully connected hidden layers of width  $784 - 1024 - 20 - 20 - 20 - 10$  with tanh and ReLU activations. The sequence of colored points shows different iterations of the training process. Each gray line connects the points with the same iterations for different layers. The left most sequence of points corresponds to the last hidden layer and the right most sequence of points corresponds to the first hidden layer. The network is trained with Adam optimization with a learning rate of 0.003 and cross-entropy loss functions to classify the MNIST handwritten-digits dataset. We repeat the experiment for 20 iterations with different randomized initializations and take the average over all experiments. In both cases of ReLU

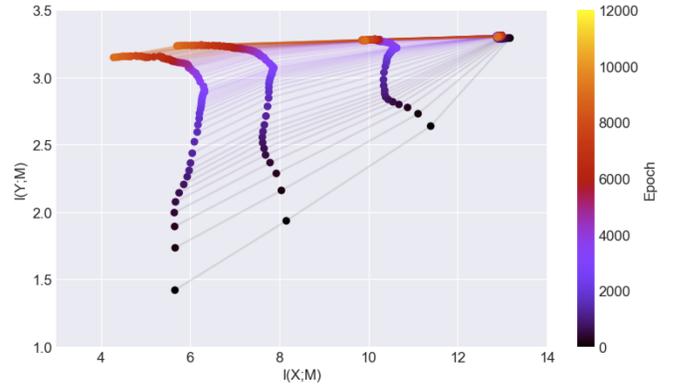


**Fig. 5.** Information plane estimated using EDGE for a neural network of size  $784 - 1024 - 20 - 20 - 20 - 10$  trained on the MNIST dataset with tanh (top) and ReLU (bottom) activations.

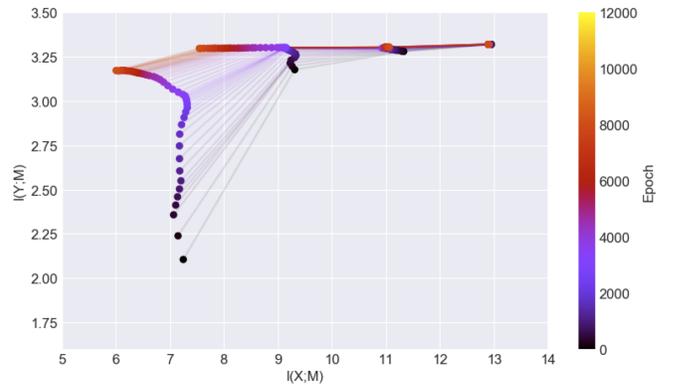
and tanh activations we observe some degree of compression in all of the hidden layers. However, the amount of compressions is different for ReLU and tanh activations. The average test accuracy in both of these networks are around 0.98. This network is the same as the one studied in [17], for which it is claimed that no compression happens with a ReLU activation. The base estimator used in [17] provides KDE-based lower and upper bounds on the true MI [26]. According to our experiments (not shown) the upper bound is in some cases twice as large as the lower bound. In contrast, our proposed ensemble method estimates the exact mutual information with significantly higher accuracy.

Fig. 6 represents the information plane for another network with 4 fully connected hidden layers of width  $784 - 200 - 100 - 60 - 30 - 10$  with ReLU activation. The network is trained with Adam optimization with a learning rate of 0.003 and cross-entropy loss functions to classify the MNIST handwritten-digits dataset. Again, we observe compression for this network with ReLU activation.

Finally, we study the information plane curves in a CNN with three convolutional ReLU layers and a dense ReLU layer. The convolutional layers respectively have depths of 4, 8, 16



**Fig. 6.** Information plane estimated using EDGE for a neural network of size  $784 - 200 - 100 - 60 - 30 - 10$  trained on the MNIST dataset with ReLU activation.



**Fig. 7.** Information plane estimated using EDGE for a CNN consisting of three convolutional ReLU layers with the respective depths of 4, 8, 16 and a dense ReLU layer with the size of 256.

and the dense layer has the dimension 256. Max-pooling functions are used in the second and third layers. Note although for a certain initialization of the weights this model can achieve the test accuracy of 0.99, the average test accuracy (over different weight initializations) is around 0.95. That’s why the converged point of the last layer has smaller  $I(T, Y)$  compared to the examples in Fig. 5, which achieves the average test accuracy of 0.98. Another interesting point about the information plane in CNN is that the convolutional layers have larger  $I(T, Y)$  compared to the hidden layers in the fully connected models in 5 and 6, which implies that the convolutional layers can extract almost all of the useful information about the labels after small number of iterations.

## 6. CONCLUSION

In this paper we proposed a fast non-parametric estimation method for MI based on random hashing, dependence graphs, and ensemble estimation. Remarkably, the proposed estimator has linear computational complexity and attains optimal (parametric) rates of MSE convergence. We provided bias

and variance convergence rate, and validated our results by numerical experiments.

## 7. REFERENCES

- [1] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.
- [2] H. Liu, L. Liu, and H. Zhang, “Ensemble gene selection for cancer classification,” *Pattern Recognition*, vol. 43, no. 8, pp. 2763–2772, 2010.
- [3] A. Hyvärinen and E. Oja, “Independent component analysis: algorithms and applications,” *Neural networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [4] A. Kraskov, H. Stögbauer, and P. Grassberger, “Estimating mutual information,” *Physical Review E*, vol. 69, no. 6, p. 066138, 2004.
- [5] Y. Moon, B. Rajagopalan, and U. Lall, “Estimation of mutual information using kernel density estimators,” *Physical Review E*, vol. 52, no. 3, p. 2318, 1995.
- [6] N. Kwak and C.-H. Choi, “Input feature selection by mutual information based on parzen window,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 12, pp. 1667–1671, 2002.
- [7] S. Khan, S. Bandyopadhyay, A. R. Ganguly, S. Saigal, D. J. Erickson III, V. Protopopescu, and G. Ostrouchov, “Relative performance of mutual information estimation methods for quantifying the dependence among short and noisy data,” *Physical Review E*, vol. 76, no. 2, p. 026209, 2007.
- [8] K. Kandasamy, A. Krishnamurthy, B. Poczos, L. Wasserman, *et al.*, “Nonparametric von mises estimators for entropies, divergences and mutual informations,” in *Advances in Neural Information Processing Systems*, pp. 397–405, 2015.
- [9] S. Singh and B. Póczos, “Exponential concentration of a density functional estimator,” in *Advances in Neural Information Processing Systems*, pp. 3032–3040, 2014.
- [10] K. R. Moon, K. Sricharan, and A. O. Hero III, “Ensemble estimation of mutual information,” *Proceedings of the IEEE Intl Symp. on Information Theory (ISIT), Aachen, June 2017*.
- [11] W. Gao, S. Kannan, S. Oh, and P. Viswanath, “Estimating mutual information for discrete-continuous mixtures,” in *Advances in Neural Information Processing Systems*, pp. 5988–5999, 2017.
- [12] K. R. Moon, K. Sricharan, K. Greenewald, and A. O. Hero, “Improving convergence of divergence functional ensemble estimators,” in *IEEE International Symposium Inf Theory*, pp. 1133–1137, IEEE, 2016.
- [13] M. Noshad and A. O. Hero III, “Scalable hash-based estimation of divergence measures,” *Proceedings of the 22nd Conference on Artificial Intelligence and Statistics, Canary Islands, March 2018, arXiv:1801.00398*.
- [14] Y. Zhang, K. Huang, G. Geng, and C. Liu, “Fast kNN graph construction with locality sensitive hashing,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 660–674, 2013.
- [15] Q. Lv, W. Josephson, Z. Wang, M. Charikar, and K. Li, “Multi-probe LSH: efficient indexing for high-dimensional similarity search,” in *Proceedings of the 33rd International Conference on Very Large Data Bases*, pp. 950–961, VLDB Endowment, 2007.
- [16] R. Shwartz-Ziv and N. Tishby, “Opening the black box of deep neural networks via information,” *arXiv preprint arXiv:1703.00810*, 2017.
- [17] A. M. Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, B. D. Tracey, and D. D. Cox, “On the information bottleneck theory of deep learning,” *ICLR*, 2018.
- [18] Y. Polyanskiy and Y. Wu, “Strong data-processing inequalities for channels and bayesian networks,” in *Convexity and Concentration*, pp. 211–249, Springer, 2017.
- [19] I. Csiszár, “Generalized cutoff rates and renyi’s information measures,” *IEEE Tran on Information Theory*, vol. 41, no. 1, pp. 26–34, 1995.
- [20] A. Andoni and P. Indyk, “Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions,” in *Foundations of Computer Science, 2006. FOCS’06. 47th Annual IEEE Symposium on*, pp. 459–468, IEEE, 2006.
- [21] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni, “Locality-sensitive hashing scheme based on p-stable distributions,” in *Proceedings of the twentieth annual symposium on Computational geometry*, pp. 253–262, ACM, 2004.
- [22] L. Paulevé, H. Jégou, and L. Amsaleg, “Locality sensitive hashing: A comparison of hash function types and querying mechanisms,” *Pattern Recognition Letters*, vol. 31, no. 11, pp. 1348–1358, 2010.
- [23] M. Noshad, K. R. Moon, S. Y. Sekeh, and A. O. Hero III, “Direct estimation of information divergence using nearest neighbor ratios,” *Proc of the IEEE Intl Symp. on Information Theory (ISIT), Aachen, June 2017*.

- [24] L. Paninski, “Estimation of entropy and mutual information,” *Neural computation*, vol. 15, no. 6, pp. 1191–1253, 2003.
- [25] A. Krishnamurthy, K. Kandasamy, B. Póczos, and L. Wasserman, “Nonparametric estimation of Rényi divergence and friends,” *arXiv preprint arXiv:1402.2966*, 2014.
- [26] A. Kolchinsky and B. D. Tracey, “Estimating mixture entropy with pairwise distances,” *Entropy*, vol. 19, no. 7, p. 361, 2017.

## A. Bias Proof

We first prove a theorem that establishes an upper bound on the number of vertices in  $V$  and  $U$ .

**Lemma 7.1.** *Cardinality of the sets  $U$  and  $V$  are upper bounded as  $|V| \leq O(\epsilon^{-d})$  and  $|U| \leq O(\epsilon^{-d})$ , respectively.*

*Proof.* Let  $\{\tilde{X}_i\}_{i=1}^{L_X}$  and  $\{\tilde{Y}_i\}_{i=1}^{L_Y}$  respectively denote distinct outputs of  $H_1$  with the  $N$  i.i.d points  $X_k$  and  $Y_k$  as input. Then according to [13] (Lemma 4.1), we have

$$L_X \leq O(\epsilon^{-d}), \quad L_Y \leq O(\epsilon^{-d}). \quad (15)$$

Simply, because of the deterministic feature of  $H_2$ , the number of its distinct inputs is greater than or equal to the number of its outputs. So,  $|V| \leq L_X$  and  $|U| \leq L_Y$ . Using the bounds in (15) completes the proof.  $\blacksquare$

The bias proof is based on analyzing the hash function defined in (5). The proof consists of two main steps: 1) Finding the expectation of hash collisions of  $H_1$ ; and 2) Analyzing the collision error of  $H_2$ . An important point about  $H_1$  and  $H_2$  is that collision of  $H_1$  plays a crucial role in our estimator, while the collision of  $H_2$  adds extra bias to the estimator. We introduce the following events to formally define these two biases:

$E_{ij}$ : The event that there is an edge between the vertices  $v_i$  and  $u_j$ .

$E_{\mathcal{E}}$ : The event that  $\mathcal{E}$  is the set of all edges in  $G$ , i.e.  $\mathcal{E} = E_G$ .

$E_{v_i}^{>0}$ : The event that there is at least one vector from  $\{\tilde{X}_i\}_{i=1}^{L_X}$  that maps to  $v_i$  using  $H_2$ .

$E_{v_i}^{=1}$ : The event that there is exactly one vector from  $\{\tilde{X}_i\}_{i=1}^{L_X}$  that maps to  $v_i$  using  $H_2$ .

$E_{v_i}^{>1}$ : The event that there are at least two vectors from  $\{\tilde{X}_i\}_{i=1}^{L_X}$  that map to  $v_i$  using  $H_2$ . (16)

$E_{u_i}^{>0}$ ,  $E_{u_i}^{=1}$  and  $E_{u_i}^{>1}$  are defined similarly. Further, let for any event  $E$ ,  $\bar{E}$  denote the complementary event. Let  $E_{ij}^{\bar{=1}} := E_{v_i}^{=1} \cap E_{u_j}^{=1}$ . Finally, we define  $E^{\bar{=1}} := \left(\cap_{i=1}^{L_X} E_{v_i}^{=1}\right) \cap \left(\cap_{j=1}^{L_Y} E_{u_j}^{=1}\right)$ , which represent the event of no collision.

Consider the notation  $\tilde{I}(X, Y) := \sum_{e_{ij} \in E_G} \omega_i \omega'_j \tilde{g}(\omega_{ij})$  (Notice the difference from the definition in (9)). We can derive its expectation as

$$\begin{aligned} \mathbb{E} \left[ \tilde{I}(X, Y) \right] &= \mathbb{E} \left[ \sum_{e_{ij} \in E_G} \omega_i \omega'_j \tilde{g}(\omega_{ij}) \middle| E_G \right] \\ &= \sum_{e_{ij} \in E_G} \mathbb{E} \left[ \omega_i \omega'_j \tilde{g}(\omega_{ij}) \middle| E_{ij} \right] \\ &= \sum_{e_{ij} \in E_G} P(E_{ij}^{\bar{=1}} | E_{ij}) \mathbb{E} \left[ \omega_i \omega'_j \tilde{g}(\omega_{ij}) \middle| E_{ij}^{\bar{=1}}, E_{ij} \right] \\ &\quad + \sum_{e_{ij} \in E_G} P(E_{ij}^{\bar{=1}} | E_{ij}) \mathbb{E} \left[ \omega_i \omega'_j \tilde{g}(\omega_{ij}) \middle| E_{ij}^{\bar{=1}}, E_{ij} \right]. \end{aligned} \quad (17)$$

Note that the second term in (17) is the bias due to collision of  $H_2$  and we denote this term by  $\mathbb{B}_H$ .

### 7.1. Bias Due to Collision

The following lemma states an upper bound on the bias error caused by  $H_2$ .

**Lemma 7.2.** *The bias error due to collision of  $H_2$  is upper bounded as*

$$\mathbb{B}_H \leq O\left(\frac{1}{\epsilon^d N}\right). \quad (18)$$

Before proving this lemma, we provide the following lemma.

**Lemma 7.3.**  $P(E_{ij}^{-1}|E_{ij})$  is given by

$$P(E_{ij}^{-1}|E_{ij}) = 1 - O\left(\frac{1}{\epsilon^d N}\right). \quad (19)$$

*Proof.* Let  $\tilde{\mathbf{X}} = \tilde{\mathbf{x}}$  and  $\tilde{\mathbf{Y}} = \tilde{\mathbf{y}}$  respectively abbreviate the equations  $\tilde{X}_1 = \tilde{x}_1, \dots, \tilde{X}_{L_X} = \tilde{x}_{L_X}$  and  $\tilde{Y}_1 = \tilde{y}_1, \dots, \tilde{Y}_{L_Y} = \tilde{y}_{L_Y}$ . Let  $\tilde{\mathbf{x}} := \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_{L_X}\}$  and  $\tilde{\mathbf{y}} := \{\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_{L_Y}\}$ . Define  $\tilde{\mathbf{z}} := \tilde{\mathbf{x}} \cup \tilde{\mathbf{y}}$  and  $L_Z := |\tilde{\mathbf{z}}|$ .

$$P(E_{ij}^{-1}|E_{ij}) = \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} P(\tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}|E_{ij}) P(E_{ij}^{-1}|E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}). \quad (20)$$

Define  $a = 2$  for the case  $i \neq j$  and  $a = 1$  for the case  $i = j$ . Then we have

$$\begin{aligned} P(E_{ij}^{-1}|E_{ij}) &= \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} P(\tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}|E_{ij}) O\left(\left(\frac{F-a}{F}\right)^{L_Z-a}\right) \\ &= \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} P(\tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}|E_{ij}) \left(1 - O\left(\frac{L_Z}{F}\right)\right) \\ &\leq \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} P(\tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}|E_{ij}) \left(1 - O\left(\frac{L_X + L_Y}{F}\right)\right) \\ &= \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} P(\tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}|E_{ij}) \left(1 - O\left(\frac{1}{\epsilon^d N}\right)\right) \\ &= \left(1 - O\left(\frac{1}{\epsilon^d N}\right)\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} P(\tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}|E_{ij}) \\ &= \left(1 - O\left(\frac{1}{\epsilon^d N}\right)\right), \end{aligned} \quad (21)$$

where in the fourth line we have used (15). ■

**Proof of 7.2.**  $N'_i$  and  $M'_j$  respectively are defined as the number of the input points  $\mathbf{X}$  and  $\mathbf{Y}$  mapped to the buckets  $\tilde{X}_i$  and  $\tilde{Y}_j$  using  $H_1$ . Define  $\mathcal{A}_i := \{j : H_2(\tilde{X}_j) = i\}$  and  $\mathcal{B}_i := \{j : H_2(\tilde{Y}_j) = i\}$ . For each  $i$  we can rewrite  $N_i$  and  $M_i$  as

$$N_i = \sum_{j=1}^{L_X} \mathbb{1}_{\mathcal{A}_i}(j) N'_j, \quad M_i = \sum_{j=1}^{L_Y} \mathbb{1}_{\mathcal{B}_i}(j) M'_j. \quad (22)$$

Thus,

$$\begin{aligned}
\mathbb{B}_H &\leq \sum_{i,j \in \mathcal{F}} P(E_{ij}^{>1}) \mathbb{E} [\mathbb{1}_{E_{ij}} \omega_i \omega'_j \tilde{g}(\omega_{ij}) | E_{ij}^{>1}] \\
&= \sum_{i,j \in \mathcal{F}} P(E_{ij}^{>1}) (P(E_{ij} | E_{ij}^{>1}) \mathbb{E} [\omega_i \omega'_j \tilde{g}(\omega_{ij}) | E_{ij}^{>1}, E_{ij}] + P(\overline{E_{ij}} | E_{ij}^{>1}) \mathbb{E} [\omega_i \omega'_j \tilde{g}(\omega_{ij}) | E_{ij}^{>1}, \overline{E_{ij}}]) \\
&= \sum_{i,j \in \mathcal{F}} P(E_{ij}) P(E_{ij}^{>1} | E_{ij}) \mathbb{E} [\omega_i \omega'_j \tilde{g}(\omega_{ij}) | E_{ij}^{>1}, E_{ij}] \tag{23}
\end{aligned}$$

$$\leq O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{i,j \in \mathcal{F}} P(E_{ij}) \mathbb{E} [\omega_i \omega'_j | E_{ij}^{>1}, E_{ij}] \tag{24}$$

$$= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{i,j \in \mathcal{F}} P(E_{ij}) \mathbb{E} [N_i M_j | E_{ij}^{>1}, E_{ij}]$$

$$= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \sum_{i,j \in \mathcal{F}} P(E_{ij}) \mathbb{E} [N_i M_j | E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}]$$

$$= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}} \sum_{i,j \in \mathcal{F}} P(E_{ij}) \mathbb{E} \left[ \left( \sum_{r=1}^{L_X} \mathbb{1}_{\mathcal{A}_i}(r) N'_r \right) \left( \sum_{s=1}^{L_Y} \mathbb{1}_{\mathcal{B}_j}(s) M'_s \right) \middle| E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}} \right] \tag{25}$$

$$= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}} \sum_{i,j \in \mathcal{F}} P(E_{ij}) \sum_{r=1}^{L_X} \sum_{s=1}^{L_Y} \mathbb{E} \left[ (\mathbb{1}_{\mathcal{A}_i}(r)) (\mathbb{1}_{\mathcal{B}_j}(s)) \middle| E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}} \right] \mathbb{E} [N'_r M'_s | E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}]$$

$$= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}} \sum_{i,j \in \mathcal{F}} P(E_{ij}) \sum_{r=1}^{L_X} \sum_{s=1}^{L_Y} P(r \in \mathcal{A}_i, s \in \mathcal{B}_j | E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}) \mathbb{E} [N'_r M'_s | E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}], \tag{26}$$

where in (23) we have used the Bayes rule, and the fact that  $\tilde{g}(\omega_{ij}) = 0$  conditioned on the event  $\overline{E_{ij}}$ . In (24) we have used the bound in Lemma 7.3, and the upper bound on  $\tilde{g}(\omega_{ij})$ . Equation (25) is due to (22). Now we simplify  $P(r \in \mathcal{A}_i, s \in \mathcal{B}_j | E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}})$  in (26) as follows. First assume that  $\tilde{X}_r \neq \tilde{Y}_s$ .

$$\begin{aligned}
P(r \in \mathcal{A}_i, s \in \mathcal{B}_j | E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}) &\leq P(r \in \mathcal{A}_i, s \in \mathcal{B}_j | E_{v_i}^{>1}, E_{u_j}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}) \\
&= P(r \in \mathcal{A}_i | E_{v_i}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}) P(s \in \mathcal{B}_j | E_{u_j}^{>1}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}), \tag{27}
\end{aligned}$$

where the second line is because the hash function  $H_2$  is random and independent for different inputs.  $P(r \in \mathcal{A}_i | E_{v_i}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}})$  in (27) can be written as

$$P(r \in \mathcal{A}_i | E_{v_i}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}) = \frac{P(r \in \mathcal{A}_i, E_{v_i}^{>1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}})}{P(E_{v_i}^{>1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}})}. \tag{28}$$

We first find  $P(E_{v_i}^{>1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}})$ :

$$\begin{aligned}
P(E_{v_i}^{>1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}) &= 1 - P(E_{v_i}^{=0} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}) - P(E_{v_i}^{=1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}) \\
&= 1 - \left(\frac{F-1}{F}\right)^{L_X} - \left(\frac{L_X}{F} \left(\frac{F-1}{F}\right)^{L_X-1}\right) \\
&= \frac{L_X^2}{2F^2} + o\left(\frac{L_X^2}{2F^2}\right). \tag{29}
\end{aligned}$$

Next, we find  $P\left(r \in \mathcal{A}_i, E_{v_i}^{>1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}\right)$  in (28) as follows.

$$\begin{aligned} P\left(r \in \mathcal{A}_i, E_{v_i}^{>1} | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}\right) &= P\left(E_{v_i}^{>1} | r \in \mathcal{A}_i, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}\right) P\left(r \in \mathcal{A}_i | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}\right) \\ &= \left(1 - \left(\frac{F-1}{F}\right)^{L_X-1}\right) \left(\frac{1}{F}\right) = O\left(\frac{L_X}{F^2}\right) \end{aligned} \quad (30)$$

Thus, using (29) and (30) yields

$$P\left(r \in \mathcal{A}_i | E_{v_i}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}\right) = O\left(\frac{1}{L_X}\right). \quad (31)$$

Similarly, we have

$$P\left(s \in \mathcal{B}_j | E_{u_j}^{>1}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right) = O\left(\frac{1}{L_Y}\right). \quad (32)$$

Now assume the case  $\tilde{X}_r = \tilde{Y}_s$ . Then since  $H_2(\tilde{X}_r) = H_2(\tilde{Y}_s)$ , we can simplify  $P\left(r \in \mathcal{A}_i, s \in \mathcal{B}_j | E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right)$  in (26) as

$$P\left(r \in \mathcal{A}_i, s \in \mathcal{B}_j | E_{ij}^{>1}, E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right) = \delta_{ij} P\left(r \in \mathcal{A}_i | E_{v_i}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right). \quad (33)$$

Recalling the definition  $\tilde{\mathbf{z}} := \tilde{\mathbf{x}} \cup \tilde{\mathbf{y}}$  and  $L_Z := |\tilde{\mathbf{z}}|$ , similar to

$$P\left(r \in \mathcal{A}_i | E_{v_i}^{>1}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right) = O\left(\frac{1}{L_Z}\right). \quad (34)$$

By using equations (27), (31), (32) and (34) in (26), we can write the following upper bound for the bias estimator due to collision.

$$\begin{aligned} \mathbb{B}_H &\leq O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} \sum_{i,j \in \mathcal{F}} P(E_{ij}) \sum_{r=1}^{L_X} \sum_{s=1}^{L_Y} \mathbb{E}\left[N'_r M'_s | E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right] \left(O\left(\frac{1}{L_X L_Y}\right) + \delta_{ij} O\left(\frac{1}{L_Z}\right)\right) \\ &= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} \sum_{i,j \in \mathcal{F}} P(E_{ij}) \mathbb{E}\left[\sum_{r=1}^{L_X} N'_r \sum_{s=1}^{L_Y} M'_s | E_{ij}, \tilde{\mathbf{X}} = \tilde{\mathbf{x}}, \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}\right] \left(O\left(\frac{1}{L_X L_Y}\right) + \delta_{ij} O\left(\frac{1}{L_Z}\right)\right) \\ &= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} \sum_{i,j \in \mathcal{F}} P(E_{ij}) N^2 \left(O\left(\frac{1}{L_X L_Y}\right) + \delta_{ij} O\left(\frac{1}{L_Z}\right)\right) \\ &= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} \left(O\left(\frac{N^2}{L_X L_Y}\right) + O\left(\frac{N}{L_Z}\right)\right) \sum_{i,j \in \mathcal{F}} P(E_{ij}) \\ &= O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} \left(O\left(\frac{N^2}{L_X L_Y}\right) + O\left(\frac{N}{L_Z}\right)\right) \mathbb{E}\left[\sum_{i,j \in \mathcal{F}} \mathbb{1}_{E_{ij}}\right] \\ &\leq O\left(\frac{U}{\epsilon^d N^3}\right) \sum_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} p_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}} \left(O\left(\frac{N^2}{L_X L_Y}\right) + O\left(\frac{N}{L_Z}\right)\right) (L_X L_Y) \\ &\leq O\left(\frac{1}{\epsilon^d N}\right). \end{aligned} \quad (35)$$

■

## 7.2. Bias without Collision

A key idea in proving Theorem 3.1 is to show that the expectation of the edge weights  $\omega_{ij}$  are proportional to the Radon-Nikodym derivative  $dP_{XY}/dP_X P_Y$  at the points that correspond to the vertices  $v_i$  and  $u_j$ . This fact is stated in the following lemma:

**Lemma 7.4.** *Under the assumptions A1-A4, and assuming that the density functions in A3 have bounded derivatives up to order  $q \geq 0$  we have:*

$$\mathbb{E}[\omega_{ij}] = \frac{dP_{XY}}{dP_X P_Y} + \mathbb{B}(N, \epsilon, q, \gamma), \quad (36)$$

where

$$\mathbb{B}(N, \epsilon, q, \gamma) := \begin{cases} O(\epsilon^\gamma) + O\left(\frac{1}{N\epsilon^d}\right), & q = 0 \\ \sum_{i=1}^q C_i \epsilon^i + O(\epsilon^q) + O\left(\frac{1}{N\epsilon^d}\right), & q \geq 1, \end{cases} \quad (37)$$

and  $C_i$  are real constants.

Note that since  $\omega_{ij} = N_{ij}N/N_i M_j$ , and  $N_{ij}$ ,  $N_i$  and  $N_j$  are not independent variables, deriving the expectation is not trivial. In the following we give a lemma that provides conditions under which the expectation of a function of random variables is close to the function of expectations of the random variables. We will use the following lemma to simplify  $\mathbb{E}[\omega_{ij}]$ .

**Lemma 7.5.** *Assume that  $g(Z_1, Z_2, \dots, Z_k) : \mathcal{Z}_1 \times \dots \times \mathcal{Z}_k \rightarrow R$  is a Lipschitz continuous function with constant  $H_g > 0$  with respect to each of variables  $Z_i$ ,  $1 \leq i \leq k$ . Let  $\mathbb{V}[Z_i]$  and  $\mathbb{V}[Z_i|X]$  respectively denote the variance and the conditional variance of each variable  $Z_i$  for a given variable  $X$ . Then we have*

$$\mathbf{a)} \quad |\mathbb{E}[g(Z_1, Z_2, \dots, Z_k)] - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, \mathbb{E}[Z_k])| \leq H_g \sum_{i=1}^k \sqrt{\mathbb{V}[Z_i]}, \quad (38)$$

$$\mathbf{b)} \quad |\mathbb{E}[g(Z_1, Z_2, \dots, Z_k) | X] - g(\mathbb{E}[Z_1|X], \mathbb{E}[Z_2|X], \dots, \mathbb{E}[Z_k|X])| \leq H_g \sum_{i=1}^k \sqrt{\mathbb{V}[Z_i|X]}. \quad (39)$$

*Proof.*

$$\begin{aligned} |\mathbb{E}[g(Z_1, Z_2, \dots, Z_k)] - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, \mathbb{E}[Z_k])| &= |\mathbb{E}[g(Z_1, Z_2, \dots, Z_k) - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, \mathbb{E}[Z_k])]| \\ &\leq \mathbb{E}[|g(Z_1, Z_2, \dots, Z_k) - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, \mathbb{E}[Z_k])|] \\ &\leq \mathbb{E}[|g(Z_1, Z_2, \dots, Z_k) - g(\mathbb{E}[Z_1], Z_2, \dots, Z_k)| + \\ &\quad + |g(\mathbb{E}[Z_1], Z_2, \dots, Z_k) - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, Z_k)| \\ &\quad + \dots \\ &\quad + |g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, \mathbb{E}[Z_{k-1}], Z_k) - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, \mathbb{E}[Z_k])|] \\ &\leq \mathbb{E}\left[|g(Z_1, Z_2, \dots, Z_k) - g(\mathbb{E}[Z_1], Z_2, \dots, Z_k)|\right] \\ &\quad + \mathbb{E}\left[|g(\mathbb{E}[Z_1], Z_2, \dots, Z_k) - g(\mathbb{E}[Z_1], \mathbb{E}[Z_2], \dots, Z_k)|\right] \\ &\quad + \dots \\ &\quad + \mathbb{E}\left[|g(\mathbb{E}[Z_1], \dots, \mathbb{E}[Z_{k-1}], Z_k) - g(\mathbb{E}[Z_1], \dots, \mathbb{E}[Z_k])|\right] \\ &\leq H_g \mathbb{E}[|Z_1 - \mathbb{E}[Z_1]|] + H_g \mathbb{E}[|Z_2 - \mathbb{E}[Z_2]|] + \dots + H_g \mathbb{E}[|Z_k - \mathbb{E}[Z_k]|] \\ &\leq H_g \sum_{i=1}^k \sqrt{\mathbb{V}[Z_i]}. \end{aligned} \quad (40)$$

$$\begin{aligned} &\leq H_g \mathbb{E}[|Z_1 - \mathbb{E}[Z_1]|] + H_g \mathbb{E}[|Z_2 - \mathbb{E}[Z_2]|] + \dots + H_g \mathbb{E}[|Z_k - \mathbb{E}[Z_k]|] \\ &\leq H_g \sum_{i=1}^k \sqrt{\mathbb{V}[Z_i]}. \end{aligned} \quad (41)$$

$$\begin{aligned} &\leq H_g \mathbb{E}[|Z_1 - \mathbb{E}[Z_1]|] + H_g \mathbb{E}[|Z_2 - \mathbb{E}[Z_2]|] + \dots + H_g \mathbb{E}[|Z_k - \mathbb{E}[Z_k]|] \\ &\leq H_g \sum_{i=1}^k \sqrt{\mathbb{V}[Z_i]}. \end{aligned} \quad (42)$$

$$\begin{aligned} &\leq H_g \sum_{i=1}^k \sqrt{\mathbb{V}[Z_i]}. \end{aligned} \quad (43)$$

In (40) and (41) we have used triangle inequalities. In (42) we have applied Lipschitz condition, and finally in (43) we have used Cauchy-Schwarz inequality. Since the proofs of parts (a) and (b) are similar, we omit the proof of part (b).  $\blacksquare$

**Lemma 7.6.** Define  $\nu_{ij} = N_{ij}/N$ , and recall the definitions  $\omega_{ij} = N_{ij}N/N_iN_j$ ,  $\omega_i = N_i/N$ , and  $\omega'_j = N_j/N$ . Then we can write

$$\mathbb{E}[\omega_{ij}] = \frac{\mathbb{E}[\nu_{ij}]}{\mathbb{E}[\omega_i]\mathbb{E}[\omega'_j]} + O\left(\sqrt{\frac{1}{N}}\right) \quad (44)$$

*Proof.* The proof follows by Lemma 7.5 and the fact that  $\mathbb{V}[\omega_{ij}] \leq O(1/N)$  (proved in Lemma 7.10).  $\blacksquare$

Let  $x_D$  and  $x_C$  respectively denote the discrete and continuous components of the vector  $x$ , with dimensions  $d_D$  and  $d_C$ . Also let  $f_{X_C}(x_C)$  and  $p_{X_D}(x_D)$  respectively denote density and pmf functions of these components associated with the probability measure  $P_X$ . Let  $S(x, r)$  be the set of all points that are within the distance  $r/2$  of  $x$  in each dimension  $i$ , i.e.

$$S(x, r) : \{x | \forall i \leq d, |X_i - x_i| < r/2\}. \quad (45)$$

Denote  $P_r(x) := P(x \in S(x, r))$ . Then we have the following lemma.

**Lemma 7.7.** Let  $r < s_X$ , where  $s_X$  is the smallest possible distance in the discrete components of the support set,  $\mathcal{X}$ . Under the assumption **A3**, and assuming that the density functions in **A3** have bounded derivatives up to the order  $q \geq 0$ , we have

$$P_r(x) = P(X_D = x_D)r^{d_C}(f(x_C|x_D) + \mu(r, \gamma, q, \mathbf{C}_X)), \quad (46)$$

where

$$\mu(r, \gamma, q, \mathbf{C}_X) := \begin{cases} O(r^\gamma), & q = 0 \\ \sum_{i=1}^q C_i r^i + O(r^q), & q \geq 1. \end{cases} \quad (47)$$

In the above equation,  $\mathbf{C}_X := (C_1, C_2, \dots, C_q)$ , and  $C_i$  are real constants depending on the probability measure  $P_X$ .

*Proof.* The proof is straightforward by using (8) for the case  $q = 0$  (similar to (27)-(29) in [13]), and using the Taylor expansion of  $f(x_C|x_D)$  for the case  $q \geq 1$  (similar to (36)-(37) in [13]).  $\blacksquare$

**Lemma 7.8.** Let  $H(x) = i, H(y) = j$ . Under the assumptions **A1-A3**, and assuming that the density functions in **A3** have bounded derivatives up to the order  $q \geq 0$ , we have

$$\mathbb{E}[\omega_{ij}|E_{ij}^{\leq 1}] = \frac{dP_{XY}}{dP_X P_Y}(x, y) + \mu(\epsilon, \gamma, q, \mathbf{C}'_{XY}) + O\left(\frac{1}{\sqrt{N}}\right), \quad (48)$$

where  $\mu(\epsilon, \gamma, q, \mathbf{C}'_{XY})$  is defined in (47).

*Proof.* Define  $\nu_{ij} = N_{ij}/N$ , and recall the definitions  $\omega_{ij} = N_{ij}N/N_iN_j$ ,  $\omega_i = N_i/N$ , and  $\omega'_j = N_j/N$ . Using Lemma 7.5 we have

$$\mathbb{E}[\omega_{ij}|E_{ij}^{\leq 1}] = \frac{\mathbb{E}[\nu_{ij}|E_{ij}^{\leq 1}]}{\mathbb{E}[\omega_i|E_{ij}^{\leq 1}]\mathbb{E}[\omega'_j|E_{ij}^{\leq 1}]} + O\left(\frac{1}{\sqrt{N}}\right) \quad (49)$$

Assume that  $H(x) = i$ . Let  $\mathcal{X}$  have  $d_C$  and  $d_D$  continuous and discrete components, respectively. Also let  $\mathcal{Y}$  have  $d'_C$  and  $d'_D$  continuous and discrete components, respectively. Then we can write

$$\begin{aligned} \mathbb{E}[\omega_i|E_{ij}^{\leq 1}] &= \frac{1}{N}\mathbb{E}[N_i|E_{ij}^{\leq 1}] \\ &= P(X \in S(x, \epsilon)) \\ &= P(X_D = x_D)\epsilon^{d_C}(f(x_C|x_D) + \mu(\epsilon, \gamma, q, \mathbf{C}_X)), \end{aligned} \quad (50)$$

where in the third line we have used Lemma 7.7. Similarly we can write

$$\begin{aligned}\mathbb{E} \left[ \omega'_j | E_{ij}^{\leq 1} \right] &= P(Y_D = y_D) \epsilon^{d'_C} (f(y_C | y_D) + \mu(\epsilon, \gamma, q, \mathbf{C}_X)), \\ \mathbb{E} \left[ \nu_{ij} | E_{ij}^{\leq 1} \right] &= P(X_D = x_D, Y_D = y_D) \epsilon^{(d_C + d'_C)} (f(x_C, y_C | x_D, y_D) + \mu(\epsilon, \gamma, q, \mathbf{C}_{XY})).\end{aligned}\quad (51)$$

Using (50) and (51) in (49) results in

$$\mathbb{E} \left[ \omega_{ij} | E_{ij}^{\leq 1} \right] = \frac{P(X_D = x_D) P(Y_D = y_D) f(x_C | x_D) f(y_C | y_D)}{P(X_D = x_D, Y_D = y_D) f(x_C, y_C | x_D, y_D)} + \mu(\epsilon, \gamma, q, \mathbf{C}'_{XY}) + O\left(\frac{1}{\sqrt{N}}\right), \quad (52)$$

where  $\mathbf{C}'_{XY}$  depends only on  $P_{XY}$ . Now note that using Lemma 7.7,  $\frac{dP_{XY}}{dP_X P_Y}(x, y)$  can be simplified as

$$\frac{dP_{XY}}{dP_X P_Y}(x, y) = \frac{\frac{dP_{XY,r}(x, y)}{dr}}{\frac{dP_{X,r} P_{Y,r}(x, y)}{dr}} = \frac{P(X_D = x_D) P(Y_D = y_D) f(x_C | x_D) f(y_C | y_D)}{P(X_D = x_D, Y_D = y_D) f(x_C, y_C | x_D, y_D)} + \mu(\epsilon, \gamma, q, \mathbf{C}''_{XY}). \quad (53)$$

Finally, using (53) in (52) gives

$$\mathbb{E} \left[ \omega_{ij} | E_{ij}^{\leq 1} \right] = \frac{dP_{XY}}{dP_X P_Y}(x, y) + \mu(\epsilon, \gamma, q, \tilde{\mathbf{C}}_{XY}) + O\left(\frac{1}{\sqrt{N}}\right), \quad (54)$$

where  $H(x) = i, H(y) = j$ . ■

**Proof of Lemma 7.4.** Lemma 7.4 is a simple consequence of Lemma 7.8. We have

$$\mathbb{E} [\omega_{ij}] = P\left(E_{ij}^{\leq 1}\right) \mathbb{E} \left[ \omega_{ij} | E_{ij}^{\leq 1} \right] + P\left(E_{ij}^{> 1}\right) \mathbb{E} \left[ \omega_{ij} | E_{ij}^{> 1} \right]. \quad (55)$$

Recall the definitions  $\tilde{\mathbf{X}} := (\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_{L_X})$  and  $\tilde{\mathbf{Y}} := (\tilde{Y}_1, \tilde{Y}_2, \dots, \tilde{Y}_{L_Y})$  as the mapped  $\mathbf{X}$  and  $\mathbf{Y}$  points through  $H_1$ . Let  $\tilde{\mathbf{Z}} := \tilde{\mathbf{X}} \cup \tilde{\mathbf{Y}}$  and  $L_Z := |\tilde{\mathbf{Z}}|$ . We first find  $P\left(E_{ij}^{\leq 1}\right)$  as follows. For a fixed set  $\tilde{\mathbf{Z}}$  we have

$$\begin{aligned}P\left(E_{ij}^{\leq 1}\right) &= P\left(E_{v_i}^{=0} \cap E_{u_j}^{=0}\right) + P\left(E_{v_i}^{=0} \cap E_{u_j}^{=1}\right) + P\left(E_{v_i}^{=1} \cap E_{u_j}^{=0}\right) + P\left(E_{v_i}^{=1} \cap E_{u_j}^{=1}\right) \\ &= \frac{(F-2)^{L_Z}}{F^{L_Z}} + \frac{L_Y (F-2)^{L_Z-1}}{F^{L_Z}} + \frac{L_X (F-2)^{L_Z-1}}{F^{L_Z}} + \frac{L_Y L_X (F-2)^{L_Z-2}}{F^{L_Z}} \\ &= 1 - O\left(\frac{L_Z}{F}\right) \\ &\leq 1 - O\left(\frac{L_X + L_Y}{F}\right) \\ &= 1 - O\left(\frac{1}{\epsilon^d N}\right).\end{aligned}\quad (56)$$

Now note that the second term in (55) is the bias due to collision of  $H_2$ , and similar to (35) it is upper bounded by  $O\left(\frac{1}{\epsilon^d N}\right)$ . Thus, (56) and (55) give rise to

$$\mathbb{E} [\omega_{ij}] = \frac{dP_{XY}}{dP_X P_Y}(x, y) + \mu(\epsilon, \gamma, q, \tilde{\mathbf{C}}_{XY}) + O\left(\frac{1}{\sqrt{N}}\right) + O\left(\frac{1}{\epsilon^d N}\right). \quad (57)$$

which completes the proof. ■

In the following lemma we make a relation between the bias of an estimator and the bias of a function of that estimator.

**Lemma 7.9.** Assume that  $g(x) : \mathcal{X} \rightarrow \mathbb{R}$  is infinitely differentiable. If  $\widehat{Z}$  is a random variable estimating a constant  $Z$  with the bias  $\mathbb{B}[\widehat{Z}]$  and the variance  $\mathbb{V}[\widehat{Z}]$ , then the bias of  $g(\widehat{Z})$  can be written as

$$\mathbb{E} \left[ g(\widehat{Z}) - g(Z) \right] = \sum_{i=1}^{\infty} \xi_i \left( \mathbb{B}[\widehat{Z}] \right)^i + O \left( \sqrt{\mathbb{V}[\widehat{Z}]} \right), \quad (58)$$

where  $\xi_i$  are real constants.

*Proof.*

$$\begin{aligned} \mathbb{E} \left[ g(\widehat{Z}) - g(Z) \right] &= g \left( \mathbb{E}[\widehat{Z}] \right) - g(Z) + \mathbb{E} \left[ g(\widehat{Z}) - g \left( \mathbb{E}[\widehat{Z}] \right) \right] \\ &= \sum_{i=1}^{\infty} \left( \mathbb{E}[\widehat{Z}] - Z \right)^i \frac{g^{(i)}(Z)}{i!} + O \left( \mathbb{E} \left[ \left| g(\widehat{Z}) - g \left( \mathbb{E}[\widehat{Z}] \right) \right| \right] \right) \\ &= \sum_{i=1}^{\infty} \xi_i \left( \mathbb{B}[\widehat{Z}] \right)^i + O \left( \sqrt{\mathbb{V}[\widehat{Z}]} \right). \end{aligned} \quad (59)$$

In the second line we have used Taylor expansion for the first term, and triangle inequality for the second term. In the third line we have used the definition  $\xi_i := g^{(i)}(Z)/i!$ , and the Cauchy-Schwarz inequality for the second term. ■

In the following we compute the expectation of the first term in (17) and prove Theorem 3.1.

**Proof of Theorem 3.1.** Recall that  $N'_i$  and  $M'_j$  respectively are defined as the number of the input points  $\mathbf{X}$  and  $\mathbf{Y}$  mapped to the buckets  $\widetilde{X}_i$  and  $\widetilde{Y}_j$  using  $H_1$ . Similarly,  $N'_{ij}$  is defined as the number of input pairs  $(\mathbf{X}, \mathbf{Y})$  mapped to the bucket pair  $(\widetilde{X}_i, \widetilde{Y}_j)$  using  $H_1$ . Define the notations  $r(i) := H_2^{-1}(i)$  for  $i \in \mathcal{F}$  and  $s(x) := H_1(x)$  for  $x \in \mathcal{X} \cup \mathcal{Y}$ . Then from (52) since there is no collision of mapping with  $H_2$  into  $v_i$  and  $u_j$  we have

$$\mathbb{E} \left[ \frac{N'_{s(x)s(y)} N}{N'_{s(x)} N'_{s(y)}} \right] = \frac{dP_{XY}}{dP_X P_Y}(x, y) + \mu(\epsilon, \gamma, q, \widetilde{\mathbf{C}}_{XY}) + O \left( \frac{1}{\sqrt{N}} \right), \quad (60)$$

By using (56) and defining  $\tilde{h}(x) = \tilde{g}(x)/x$  we can simplify the first term of (17) as

$$\begin{aligned}
\sum_{i,j \in \mathcal{F}} P(E_{ij}^{\leq 1}) \mathbb{E} \left[ \mathbb{1}_{E_{ij}} \omega_i \omega'_j \tilde{g}(\omega_{ij}) \middle| E_{ij}^{\leq 1} \right] &= \left( 1 - O\left(\frac{1}{\epsilon^d N}\right) \right) \sum_{i,j \in \mathcal{F}} \mathbb{E} \left[ \mathbb{1}_{E_{ij}} \omega_i \omega'_j \tilde{g}(\omega_{ij}) \middle| E_{ij}^{\leq 1} \right] \\
&= \sum_{i,j \in \mathcal{F}} \mathbb{E} \left[ \mathbb{1}_{E_{ij}} \frac{N_i M_j}{N^2} \tilde{g}\left(\frac{N_{ij} N}{N_i M_j}\right) \middle| E_{ij}^{\leq 1} \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \sum_{i,j \in \mathcal{F}} \mathbb{E} \left[ \mathbb{1}_{E_{ij}} \frac{N'_{r(i)} M'_{r(j)}}{N^2} \tilde{g}\left(\frac{N'_{r(i)} r(j) N}{N'_{r(i)} M'_{r(j)}}\right) \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \sum_{i,j \in \mathcal{F}} \mathbb{E} \left[ \mathbb{1}_{E_{ij}} \frac{N'_{r(i)} r(j)}{N} \tilde{h}\left(\frac{N'_{r(i)} r(j) N}{N'_{r(i)} M'_{r(j)}}\right) \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \frac{1}{N} \sum_{i,j \in \mathcal{F}} \mathbb{E} \left[ N'_{r(i)} r(j) \tilde{h}\left(\frac{N'_{r(i)} r(j) N}{N'_{r(i)} M'_{r(j)}}\right) \right] + O\left(\frac{1}{\epsilon^d N}\right) \tag{61} \\
&= \frac{1}{N} \mathbb{E} \left[ \sum_{i,j \in \mathcal{F}} N'_{r(i)} r(j) \tilde{h}\left(\frac{N'_{r(i)} r(j) N}{N'_{r(i)} M'_{r(j)}}\right) \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \frac{1}{N} \mathbb{E} \left[ \sum_{i=1}^N \tilde{h}\left(\frac{N'_{s(X)} s(Y) N}{N'_{s(X)} M'_{s(Y)}}\right) \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \tilde{h}\left(\frac{N'_{s(X)} s(Y) N}{N'_{s(X)} M'_{s(Y)}}\right) \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \mathbb{E}_{(X,Y) \sim P_{XY}} \left[ \mathbb{E} \left[ \tilde{h}\left(\frac{N'_{s(X)} s(Y) N}{N'_{s(X)} M'_{s(Y)}}\right) \middle| X = x, Y = y \right] \right] + O\left(\frac{1}{\epsilon^d N}\right) \\
&= \mathbb{E}_{(X,Y) \sim P_{XY}} \left[ \frac{dP_{XY}}{dP_X P_Y} \right] + \mu(\epsilon, \gamma, q, \bar{\mathbf{C}}_{XY}) + O\left(\frac{1}{\sqrt{N}}\right) + O\left(\frac{1}{\epsilon^d N}\right). \tag{62} \\
&\tag{63}
\end{aligned}$$

(61) is due to the fact that  $N'_{r(i)} r(j) = 0$  if there is no edge between  $v_i$  and  $u_j$ . Also, (62) is due to (60).

From (62) and (17) we obtain

$$\mathbb{E} \left[ \tilde{I}(X, Y) \right] = \mathbb{E} \left[ \sum_{e_{ij} \in E_G} \omega_i \omega'_j \tilde{g}(\omega_{ij}) \right] = \mathbb{E}_{(X,Y) \sim P_{XY}} \left[ \frac{dP_{XY}}{dP_X P_Y} \right] + \mu(\epsilon, \gamma, q, \bar{\mathbf{C}}_{XY}) + O\left(\frac{1}{\sqrt{N}}\right) + O\left(\frac{1}{\epsilon^d N}\right). \tag{64}$$

Finally using Lemma 7.9 results in (11). ■

## B. Variance Proof

In this section we first prove bounds on the variances of the edge and vertex weights and then we provide the proof of Theorem 3.2.

**Lemma 7.10.** *Under the assumptions A1-A4, the following variance bounds hold true.*

$$\mathbb{V}[\omega_i] \leq O\left(\frac{1}{N}\right), \quad \mathbb{V}[\omega'_j] \leq O\left(\frac{1}{N}\right), \quad \mathbb{V}[\omega_{ij}] \leq O\left(\frac{1}{N}\right), \quad \mathbb{V}[\nu_{ij}] \leq O\left(\frac{1}{N}\right). \tag{65}$$

*Proof.* Here we only provide the variance proof of  $\omega_i$ . The variance bounds of  $\omega'_j$ ,  $\omega_{ij}$  and  $\nu_{ij}$  can be proved in the same way. The proof is based on Efron-Stein inequality. Define  $Z_i := (X_i, Y_i)$ . For using the Efron-Stein inequality on  $\mathbf{Z} := (Z_1, \dots, Z_N)$ ,

we consider another independent copy of  $\mathbf{Z}$  as  $\mathbf{Z}' := (Z'_1, \dots, Z'_N)$  and define  $\mathbf{Z}^{(i)} := (Z_1, \dots, Z_{i-1}, Z'_i, Z_{i+1}, \dots, Z_N)$ . Define  $\omega_i(\mathbf{Z})$  as the weight of vertex  $v_i$  in the dependence graph constructed by the set  $\mathbf{Z}$ . By applying Efron-Stein inequality [23] we have

$$\begin{aligned}
\mathbb{V}[\omega_i] &\leq \frac{1}{2} \sum_{i=1}^N \mathbb{E} \left[ \left( \omega_i(\mathbf{Z}) - \omega_i(\mathbf{Z}^{(i)}) \right)^2 \right] \\
&= \frac{1}{2N^2} \sum_{i=1}^N \mathbb{E} \left[ \left( N_i(\mathbf{Z}) - N_i(\mathbf{Z}^{(i)}) \right)^2 \right] \\
&\leq \frac{1}{2N^2} O(N) \\
&\leq O\left(\frac{1}{N}\right).
\end{aligned} \tag{66}$$

In the third line we have used the fact that the absolute value of  $N_i(\mathbf{Z}) - N_i(\mathbf{Z}^{(i)})$  is at most 1. ■

**Proof of Theorem 3.2.** We follow similar steps as the proof of Lemma 7.10. Define  $\widehat{I}_g(\mathbf{Z})$  as the mutual information estimation using the set  $\mathbf{Z}$ . By applying Efron-Stein inequality we have

$$\begin{aligned}
\mathbb{V}[\widehat{I}(X, Y)] &\leq \frac{1}{2} \sum_{k=1}^N \mathbb{E} \left[ \left( \widehat{I}(\mathbf{Z}) - \widehat{I}(\mathbf{Z}^{(k)}) \right)^2 \right] \\
&\leq \frac{N}{2} \mathbb{E} \left[ \left( \sum_{e_{ij} \in E_G} \omega_i(\mathbf{Z}) \omega'_j(\mathbf{Z}) \widetilde{g}(\omega_{ij}(\mathbf{Z})) - \sum_{e_{ij} \in E_G} \omega_i(\mathbf{Z}^{(k)}) \omega'_j(\mathbf{Z}^{(k)}) \widetilde{g}(\omega_{ij}(\mathbf{Z}^{(k)})) \right)^2 \right] \\
&= \frac{N}{2N^4} \mathbb{E} \left[ \left( \sum_{e_{ij} \in E_G} N_i(\mathbf{Z}) M_j(\mathbf{Z}) \widetilde{g}\left(\frac{N_{ij}(\mathbf{Z}) N}{N_i(\mathbf{Z}) M_j(\mathbf{Z})}\right) - \sum_{e_{ij} \in E_G} N_i(\mathbf{Z}^{(k)}) M_j(\mathbf{Z}^{(k)}) \widetilde{g}\left(\frac{N_{ij}(\mathbf{Z}^{(k)}) N}{N_i(\mathbf{Z}^{(k)}) M_j(\mathbf{Z}^{(k)})}\right) \right)^2 \right]
\end{aligned} \tag{67}$$

$$\leq \frac{1}{2N^3} \mathbb{E} \left[ (\Sigma_{n_1} + \Sigma_{n_2} + \Sigma_{m_1} + \Sigma_{m_2} + D_{n_1 m_1} + D_{n_2 m_2})^2 \right]. \tag{68}$$

Note that in equation (68), when  $(X_k, Y_k)$  is resampled, at most two of  $N_i$  for  $i \in \mathcal{F}$  are changed exactly by one (one decrease and the other increase). The same statement holds true for  $M_j$ . Let these vertices be  $v_{n_1}, v_{n_2}, v_{m_1}$  and  $v_{m_2}$ . Also the pair collision counts  $N_{ij}$  are fixed except possibly  $N_{n_1 m_1}$  and  $N_{n_2 m_2}$  that may change by one. So, in the fourth line  $\Sigma_{n_1}$  and  $\Sigma_{n_2}$  account for the changes in MI estimation due to the changes in  $N_{n_1}$  and  $N_{n_2}$ , and  $\Sigma_{m_1}$  and  $\Sigma_{m_2}$  account for the changes in  $M_{m_1}$  and  $M_{m_2}$ , respectively. Finally  $D_{n_1 m_1}$  and  $D_{n_2 m_2}$  account for the changes in MI estimation due to the changes in  $N_{n_1 m_1}$  and  $N_{n_2 m_2}$ . For example,  $\Sigma_{n_1}$  is precisely defined as follows:

$$\Sigma_{n_1} := \sum_{j: e_{m_j} \in E_G} N_m M_j \widetilde{g}\left(\frac{N_{mj} N}{N_m N_j}\right) - (N_m + 1) M_j \widetilde{g}\left(\frac{N_{mj} N}{(N_m + 1) M_j}\right) \tag{69}$$

where we have used the notations  $N_i$  and  $N_i^{(k)}$  instead of  $N_i(\mathbf{Z})$  and  $N_i(\mathbf{Z}^{(k)})$  for simplicity. Now note that by assumption **A4** we have

$$\begin{aligned}
\left| \widetilde{g}\left(\frac{N_{mj} N}{N_m M_j}\right) - \widetilde{g}\left(\frac{N_{mj} N}{(N_m + 1) M_j}\right) \right| &\leq G_g \left| \frac{N_{mj} N}{N_m M_j} - \frac{N_{mj} N}{(N_m + 1) M_j} \right| \\
&\leq O\left(\frac{N_{mj} N}{N_m^2 M_j}\right).
\end{aligned} \tag{70}$$

Thus, using (70),  $\Sigma_{n_1}$  can be upper bounded as follows

$$\Sigma_{n_1} \leq \sum_{j: e_{m_j} \in E_G} O\left(\frac{N_{m_j} N}{N_m^2}\right) = O\left(\frac{N}{N_m}\right) \leq O(N). \quad (71)$$

It can similarly be shown that  $N_{n_2}$ ,  $\Sigma_{m_1}$ ,  $\Sigma_{m_2}$ ,  $D_{n_1 m_1}$  and  $D_{n_2 m_2}$  are upper bounded by  $O(N)$ . Thus, (68) simplifies as follows

$$\mathbb{V}[\widehat{I}(X, Y)] \leq \frac{36O(N^2)}{2N^3} = O\left(\frac{1}{N}\right). \quad (72)$$

■

## C. Optimum MSE Rates of EDGE

In this short section we prove Theorem 4.1.

**Proof of Theorem 4.1.** The proof simply follows by using the ensemble theorem in ([12], Theorem 4) with the parameters  $\psi_i(t) = t^i$  and  $\phi_{i,d}(N) = N^{-i/2d}$  for the bias result in Theorem 3.1. Thus, the following weighted ensemble estimator (EDGE) can achieve the optimum parametric MSE convergence rate of  $O(1/N)$  for  $q \geq d$ .

$$\widehat{I}_w := \sum_{t \in \mathcal{T}} w(t) \widehat{I}_{\epsilon(t)}, \quad (73)$$

■