

# Estimating Mutual Information via Geodesic $k$ NN

Alexander Marx \*

Jonas Fischer †

## Abstract

Estimating mutual information (MI) between two continuous random variables  $X$  and  $Y$  allows to capture non-linear dependencies between them, non-parametrically. As such, MI estimation lies at the core of many data science applications. Yet, robustly estimating MI for high-dimensional  $X$  and  $Y$  is still an open research question.

In this paper, we formulate this problem through the lens of manifold learning. That is, we leverage the common assumption that the information of  $X$  and  $Y$  is captured by a low-dimensional manifold embedded in the observed high-dimensional space and transfer it to MI estimation. As an extension to state-of-the-art  $k$ NN estimators, we propose to determine the  $k$ -nearest neighbors via geodesic distances on this manifold rather than from the ambient space, which allows us to estimate MI even in the high-dimensional setting. An empirical evaluation of our method, G-KSG, against the state-of-the-art shows that it yields good estimations of MI in classical benchmark and manifold tasks, even for high dimensional datasets, which none of the existing methods can provide.

## 1 Introduction

Quantifying the strength of a dependence between two continuous random variables is an essential task in data science [32]. Due to its non-parametric nature, and hence its ability to measure complex non-linear dependencies, mutual information is ideal for this task [3], which is why it is routinely applied for challenging settings such as gene network inference.

Given two multidimensional continuous random variables  $\mathbf{X} \in \mathbb{R}^{d_{\mathbf{X}}}$  and  $\mathbf{Y} \in \mathbb{R}^{d_{\mathbf{Y}}}$ , mutual information

$$I(\mathbf{X}; \mathbf{Y}) = h(\mathbf{X}) + h(\mathbf{Y}) - h(\mathbf{X}, \mathbf{Y})$$

can be expressed as sum of differential entropies

$$h(\mathbf{X}) = - \int_{\mathbb{R}^{d_{\mathbf{X}}}} f_{\mathbf{X}}(\mathbf{x}) \log f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x},$$

where  $\log$  refers to the natural logarithm. Although the differential entropy of a random variable can be negative, the chain rule does still apply for differential entropy, and hence  $I(\mathbf{X}; \mathbf{Y}) \geq 0$  with equality if and only if  $\mathbf{X}$  is independent of  $\mathbf{Y}$  [3].

In an ideal scenario, where we are given an iid sample  $(\mathbf{x}_i, \mathbf{y}_i)_{i=1, \dots, n} \sim f_{\mathbf{X}\mathbf{Y}}$ , and an unbiased estimator  $\hat{f}_{\mathbf{X}\mathbf{Y}}$ , we could simply estimate  $I(\mathbf{X}; \mathbf{Y})$  by individually estimating the differential entropies as  $\hat{h}(\mathbf{X}) = -\frac{1}{n} \sum_{i=1}^n \log \hat{f}_{\mathbf{X}}(\mathbf{x}_i)$ . In practice, such an unbiased density estimate is usually not available and state-of-the-art approaches instead resort to  $k$ NN based estimates of  $\hat{h}(\mathbf{X})$  [6, 14]. Simply put, those methods estimate the log-density locally around each point, e.g. by enclosing all its  $k$ NNs into a unit ball and computing its volume [12, 14].

This approach, i.e. estimating each of the entropy terms individually with the same  $k$ , was, however, shown to induce a bias since the volume-related correction terms do not cancel [15]. To correct for this bias, Kraskov, Stögbauer, and Grassberger (KSG) [15] suggested to determine the distance to the  $k$ th neighbor only on the joint space and retrospectively count the data points falling within this region in  $\mathbf{X}$  and  $\mathbf{Y}$ . Other approaches try to reduce this bias by using different geometries that more tightly enclose the  $k$ NNs to better model the local densities [7, 9, 16, 17]. Yet, none of these approaches has been successfully applied to high-dimensional data, which is exactly the setting we are interested in.

To estimate MI on high-dimensional data, we build upon the manifold assumption [2]. This common assumption in machine learning states that high-dimensional data resides on a low-dimensional manifold embedded in the ambient space. Under this assumption, we propose G-KSG, which instantiates the KSG estimator by determining the nearest neighbors via geodesic distances, i.e. the shortest path between two points on the manifold they reside on. To estimate geodesic distances, we make use of and extend a recent proposal for manifold learning, called Geodesic Forests [18], which is an unsupervised random forest based on sparse linear projections. As such, our method is well suited to estimate local densities and therewith mutual information on high-dimensional data implementing the manifold assumption. Our main contributions are, we

- establish a formal connection between manifold learning and mutual information estimation, for which we derive identifiability results in Sec. 4,

\*Department of Computer Science, ETH Zurich and ETH AI Center, Zurich, Switzerland. alexander.marx@ai.ethz.ch

†Max Planck Institute for Informatics, Saarbrücken, Germany. fischer@mpi-inf.mpg.de

- propose G-KSG, an instantiation of KSG using geodesic distances, which we approximate via Geodesic Forests [2], as explained in Sec 5,
- derive a locally adjusted dissimilarity measure, as well as a more efficient,  $\mathcal{O}(n)$ , split criterium for unsupervised forests in Sec 5, and
- provide an extensive empirical evaluation of G-KSG in Sec. 6.

Before that, we discuss related work.

## 2 Related Work

Mutual information estimation is a well studied problem for discrete, continuous and even discrete-continuous mixture data [20, 29, 8, 14, 15, 24, 19, 21]. Here, we focus on continuous data, for which a broad spectrum of MI estimators exists. Most common are estimators based on discretization [4, 21, 13], kernel density estimation [24, 10], and  $k$ -nearest neighbor estimation [6, 14, 15]. Recently,  $k$ NN-based estimators have been established as state-of-the-art and can be computed efficiently, e.g. via the  $k$ -D trie method [31].

Simply put,  $k$ NN-based methods estimate the local density around each point  $i$  via its  $k$ -nearest neighbors [15, 16]. Critical for the performance of these estimators are assumptions about the shape of the local volumes used to calculate the densities. One group of estimators measure the local distances via  $L_2$  or  $L_\infty$ -norm [6, 14]. Other approaches try to estimate the volumes via locally computing an SVD [16] or PCA [17] transformation, or use a local Gaussian kernel [10]. Alternatively, the KSG [15] estimator avoids estimating the volumes all along by computing the distance to the  $k$ th neighbor on the joint space, while simply counting the neighbors falling within this region in  $\mathbf{X}$  and  $\mathbf{Y}$ . Gao et al. [11] proved that the KSG estimator is consistent and proposed a bias-corrected alternative, which focuses on low-dimensional data. Closest to our approach is  $k$ NN-based estimation using random forests to estimate the nearest neighbors, which requires either  $\mathbf{X}$  or  $\mathbf{Y}$  to be discrete [23]. None of these estimators has been evaluated on more than 20 dimensions.

To efficiently estimate MI in a high-dimensional setting, we build upon the manifold assumption [2], based on which embedding techniques were developed that successfully capture the most relevant information in few dimensions by focussing on preserving local Euclidean distances and estimating geodesics that resemble the data location on the manifold [22, 30, 26].

In particular, we suggest a novel approach that estimates mutual information considering geodesic distances, combining ideas from manifold learning and MI

estimation. We leverage recent advances of Madhyasta et al. [18] in approximating geodesic distances based on tree estimates on sparse linear projections of the original space [1, 5], which is suitable for  $k$ NN estimation on high-dimensional data.

## 3 Preliminaries

Next, we briefly introduce the line of  $k$ NN based MI estimators more formally, and then shortly discuss the KSG [15] estimator and its limitations.

To estimate the differential entropy  $\hat{h}(\mathbf{X})$  of a random variable  $\mathbf{X}$ ,  $k$ NN based estimators [6, 14] estimate the log-density locally around each point  $\mathbf{x}_i$ , e.g. by computing the volume  $V_{\mathbf{x}_i}$  of a ball enclosing all its  $k$ NNs [12, 14]. That is,

$$(3.1) \quad \log \hat{f}_{\mathbf{X}}(\mathbf{x}_i) = \psi(k) - \psi(n) - \log V_{\mathbf{x}_i},$$

where  $\psi(k)$  and  $\psi(n)$  are the correction terms, with  $\psi$  being the digamma function.<sup>1</sup> A straight-forward approach to estimate  $I(\mathbf{X}; \mathbf{Y})$  would be to estimate each of the involved entropy terms individually using Eq. (3.1) based on the same  $k$ . Kraskov et al. [15] showed that this will, however, induce a bias. Instead, they compute the distance to the  $k$ th neighbor only on the joint space and retrospectively count the data points falling within this region in  $\mathbf{X}$  and  $\mathbf{Y}$ .

**KSG Estimator** Let  $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$  be the joint space spanned by  $\mathbf{X}$  and  $\mathbf{Y}$ . For any two data points  $\mathbf{z}_i$  and  $\mathbf{z}_j$ , we define the distance between them as the maximum distance from their projections in  $\mathbf{X}$  resp.  $\mathbf{Y}$ ,

$$d(\mathbf{z}_i, \mathbf{z}_j)_{\max} = \max\{d(\mathbf{x}_i, \mathbf{x}_j), d(\mathbf{y}_i, \mathbf{y}_j)\},$$

where the distances measured on the subspaces, i.e.  $d(\mathbf{x}_i, \mathbf{x}_j)$  and  $d(\mathbf{y}_i, \mathbf{y}_j)$ , can be instantiated with any norm. Furthermore, it is possible to use different norms for the subspaces  $\mathbf{X}$  and  $\mathbf{Y}$  [15]. In practice, both  $d(\mathbf{x}_i, \mathbf{x}_j)$  and  $d(\mathbf{y}_i, \mathbf{y}_j)$  are instantiated with the  $L_\infty$ -norm, and hence  $d(\mathbf{z}_i, \mathbf{z}_j)_{\max}$  reduces to the  $L_\infty$ -norm over the joint space  $(\mathbf{X}, \mathbf{Y})$ .

Next, let  $\frac{1}{2}\rho_{i,k}$  be the distance to the  $k$ -th neighbor on the  $\mathbf{Z}$  space using the maximum norm as defined above. We define the number of data points  $\mathbf{x}_j$  with a distance smaller than  $\frac{1}{2}\rho_{i,k}$  to point the  $\mathbf{x}_i$  in the  $\mathbf{X}$  subspace as  $n_{\mathbf{x},i}$ , i.e.

$$n_{\mathbf{x},i} = \left| \left\{ \mathbf{x}_j : d(\mathbf{x}_i, \mathbf{x}_j) < \frac{1}{2}\rho_{i,k}, i \neq j \right\} \right|,$$

and similarly, we define  $n_{\mathbf{y},i}$  as the number of data points with a smaller distance to point  $\mathbf{y}_i$  than  $\frac{1}{2}\rho_{i,k}$

<sup>1</sup>The digamma function is defined as  $\psi(x) = \Gamma(x)^{-1}d\Gamma(x)/dx$ . It satisfies the recursion  $\psi(x+1) = \psi(x) + \frac{1}{x}$  with  $\psi(1) = -C$ , where  $C = 0.577215\dots$  is the Euler-Mascheroni constant.

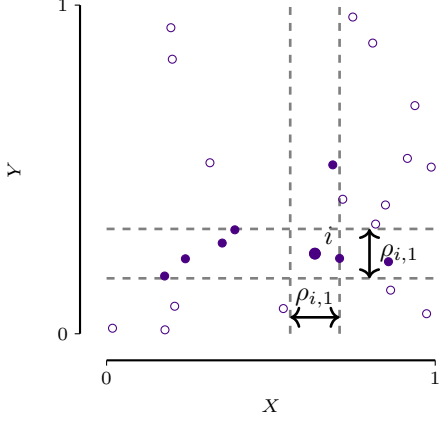


Figure 1: Example calculation for  $\frac{1}{2}\rho_{i,k}$ , where  $k = 1$ , and we use the  $L_\infty$ -norm as distance measure for  $\mathbf{X}$  and  $\mathbf{Y}$ . In this case,  $n_{\mathbf{x},i} = 1$  and  $n_{\mathbf{y},i} = 6$ .

on the  $\mathbf{Y}$  subspace. As an example consider the two-dimensional plot in Figure 1. In this example  $n_{\mathbf{x},i} = 1$ , that is, except  $\mathbf{x}_i$  itself, there exists only one further point with a distance  $< \frac{1}{2}\rho_{i,k}$  to  $\mathbf{x}_i$  for  $k = 1$ . On the other hand, there exist 6 data points, which fulfill this criterium for the  $\mathbf{Y}$  subspace. In general, it holds that  $n_{\mathbf{x},i} + 1 \geq k$ , as well as  $n_{\mathbf{y},i} + 1 \geq k$ .

To compute the entropies for the subspaces  $\mathbf{X}$  and  $\mathbf{Y}$ , we consider the volumes related to  $\rho_{i,k}$  computed on  $\mathbf{Z}$  and determine the corresponding  $k$  retrospectively, i.e. we define  $h^{\text{KSG}}(\mathbf{X})$  as

$$h^{\text{KSG}}(\mathbf{X}) = -\frac{1}{n} \sum_{i=1}^n \psi(n_{\mathbf{x},i} + 1) - \psi(n) - \log V_{\mathbf{x}_i},$$

where  $\log V_{\mathbf{x}_i}$  is computed as  $\log c_{d_{\mathbf{X}}} - d_{\mathbf{X}} \rho_{i,k}$ , with  $c_{d_{\mathbf{X}}}$  being the volume of a  $d_{\mathbf{X}}$ -dimensional unit ball. By adding up the individual entropy terms, we arrive at the KSG estimator, defined as

$$I^{\text{KSG}}(\mathbf{X}; \mathbf{Y}) = \psi(k) + \psi(n) - \frac{1}{n} \sum_{i=1}^n \psi(n_{\mathbf{x},i} + 1) + \psi(n_{\mathbf{y},i} + 1),$$

where we can see that all volume terms cancel.

Although, it was shown that the KSG estimator is consistent [11], its bias increases with the number of dimensions  $d$  as  $\mathcal{O}(n^{-1/d})$  [11], which is problematic for the high-dimensional setting where  $d > n$ . In practice, we observe that the more dimensions we consider, the larger  $n_{\mathbf{x},i}$  and  $n_{\mathbf{y},i}$  become on average. This phenomenon occurs naturally, since we consider the maximum norm in the joint space. In extreme cases,  $\psi(n_{\mathbf{x},i} + 1) + \psi(n_{\mathbf{y},i} + 1)$  is on average larger than  $\psi(k) + \psi(n)$  and hence the estimate can be negative.

To address the limitations of KSG in high-dimensional data, we leverage insights from manifold learning and mutual information estimation. In the next section, we will formally define the assumed data generative model and provide identifiability results.

#### 4 Mutual Information & Manifold Learning

The classical objective of MI estimation is to estimate  $I(\mathbf{X}; \mathbf{Y})$  given an iid sample of the joint distribution  $P_{\mathbf{X}, \mathbf{Y}}$ . Especially for high-dimensional and possibly noisy  $\mathbf{X}$  and  $\mathbf{Y}$ , estimating mutual information is challenging [17, 16, 10]. Here, we view this problem from a manifold learning perspective [2], where we assume that shared information between  $\mathbf{X}$  and  $\mathbf{Y}$  is encoded in an intrinsic low-dimensional space  $(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}})$ , whereas the majority of the dimensions of  $\mathbf{X}$  and  $\mathbf{Y}$  are independent of each other, or correspond to noise dimensions, i.e.  $I(\mathbf{X}; \mathbf{Y})$  is upper-bounded by  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ .

To rigorously define the problem setting, we write down our assumptions about the data generative process as a structural causal model [25]. Simply put, we assume that our observed variables  $\mathbf{X}$  and  $\mathbf{Y}$  are both generated from a shared variable  $\tilde{\mathbf{Z}}$  and individual variables  $\mathbf{E}_X$  and  $\mathbf{E}_Y$  that are independent of  $\tilde{\mathbf{Z}}$  and independent of each other (see Figure 2). Further, we assume that the information that  $\mathbf{X}$  contains about  $\tilde{\mathbf{Z}}$  is first passed through  $\tilde{\mathbf{X}}$ . Accordingly, the information that  $\mathbf{Y}$  has about  $\tilde{\mathbf{Z}}$  is processed through the path  $\tilde{\mathbf{Z}} \rightarrow \tilde{\mathbf{Y}} \rightarrow \mathbf{Y}$ . We chose this model to reason about the shared information of  $\mathbf{X}$  and  $\mathbf{Y}$  in terms of  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ , which is assumed to be low-dimensional.

We formally define the generative model below.

MODEL 1. *Given multi-dimensional random vectors  $\tilde{\mathbf{Z}}, \mathbf{E}_X, \mathbf{E}_Y$ , which are pairwise independent. We generate  $\mathbf{X}$  and  $\mathbf{Y}$  according to*

$$\begin{aligned} \tilde{\mathbf{X}} &= g_{\tilde{\mathbf{X}}}(\tilde{\mathbf{Z}}) & \tilde{\mathbf{Y}} &= g_{\tilde{\mathbf{Y}}}(\tilde{\mathbf{Z}}) \\ \mathbf{X} &= g_{\mathbf{X}}(\tilde{\mathbf{X}}, \mathbf{E}_X) & \mathbf{Y} &= g_{\mathbf{Y}}(\tilde{\mathbf{Y}}, \mathbf{E}_Y), \end{aligned}$$

where we require that

1.  $g_{\tilde{\mathbf{X}}}$  and  $g_{\tilde{\mathbf{Y}}}$  preserve some information about  $\tilde{\mathbf{Z}}$ , s.t.  $I(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}) > 0$ ,
2.  $g_{\mathbf{X}}$ , as well as  $g_{\mathbf{Y}}$  are required to be homeomorphisms (smooth and uniquely invertible maps), and
3. for all sub-spaces  $\tilde{\mathbf{X}}'$  or  $\tilde{\mathbf{Y}}'$  containing only a proper subset of rows of  $\tilde{\mathbf{X}}$  resp.  $\tilde{\mathbf{Y}}$ , it holds that  $I(\tilde{\mathbf{X}}'; \tilde{\mathbf{Y}}) < I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ , resp.  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}}') < I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ .

Conditions 1-3 in Model 1 are very light requirements. If Cond. 1 would be violated, our estimator

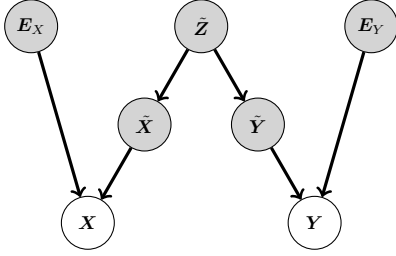


Figure 2: [Generative Model] We are interested in estimating  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ , from  $\mathbf{X}$  and  $\mathbf{Y}$ . The random vector  $\mathbf{X}$  is generated as a function of two unobserved variables (denoted as shaded nodes)  $\tilde{\mathbf{X}}$  and  $\mathbf{E}_X$ , where  $\tilde{\mathbf{X}} \perp\!\!\!\perp \mathbf{E}_X$ ;  $\mathbf{Y}$  is generated accordingly. Further,  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$  are generated from a shared latent factor  $\tilde{\mathbf{Z}}$ .

could still detect that there is no shared information between  $\mathbf{X}$  and  $\mathbf{Y}$ , but there would be little point in modeling the generative process from a manifold learning perspective. Similarly, Cond. 3 requires that all variables in the low-dimensional representation are contributing to the shared information about  $\mathbf{X}$  and  $\mathbf{Y}$ . Any subsets of variables not fulfilling this requirement would be modeled within the individual factor. Last, Cond. 2 assures that  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$  can be recovered from  $I(\mathbf{X}; \mathbf{Y})$ , as shown in Proposition 4.1. This is a standard assumption in manifold learning: the low-dimensional manifold can be modeled as a homeomorphism of the ambient space or a subspace of it [2].

In the following we will show that for a data generative process as defined in Model 1, we can compute  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$  given only  $I(\mathbf{X}; \mathbf{Y})$ . For arbitrary functions  $g$  and  $g'$ , where the transformed variables  $g(\mathbf{X})$  and  $g'(\mathbf{Y})$  are differentiable almost everywhere, the data processing inequality states that  $I(g(\mathbf{X}); g'(\mathbf{Y})) \leq I(\mathbf{X}; \mathbf{Y})$  [3]. Applied to our scenario,  $I(\mathbf{X}; \mathbf{Y})$  would at most provide a lower bound for  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ . When requiring that  $g$  and  $g'$  are homeomorphisms, as stated in Condition 2 in Model 1, we can make a stronger statement.

**PROPOSITION 4.1.** *Given a data generative process as defined in Model 1, then  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}}) = I(\mathbf{X}; \mathbf{Y})$ .*

*Proof.* Due to Condition 2 in Model 1, we know that  $g_{\mathbf{X}}$  and  $g_{\mathbf{Y}}$  are homeomorphisms. Kraskov et al. [15] showed for any transformations  $\mathbf{X}^t = g(\mathbf{X})$ ,  $\mathbf{Y}^t = g'(\mathbf{Y})$ , where  $g$  and  $g'$  are homeomorphisms, it holds that  $I(\mathbf{X}; \mathbf{Y}) = I(\mathbf{X}^t; \mathbf{Y}^t)$ . As an immediate consequence, we observe that in our case  $I(\mathbf{X}; \mathbf{Y}) =$

$I(\tilde{\mathbf{X}}, \mathbf{E}_X; \tilde{\mathbf{Y}}, \mathbf{E}_Y)$ . Hence, we can derive that

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= h(\tilde{\mathbf{X}}, \mathbf{E}_X) - h(\tilde{\mathbf{X}}, \mathbf{E}_X | \tilde{\mathbf{Y}}, \mathbf{E}_Y) \\ &= h(\tilde{\mathbf{X}}, \mathbf{E}_X) - h(\tilde{\mathbf{X}}, \mathbf{E}_X | \tilde{\mathbf{Y}}) \\ &= h(\tilde{\mathbf{X}}) + h(\mathbf{E}_X) - (h(\tilde{\mathbf{X}} | \tilde{\mathbf{Y}}) + h(\mathbf{E}_X)) \\ &= h(\tilde{\mathbf{X}}) - h(\tilde{\mathbf{X}} | \tilde{\mathbf{Y}}) = I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}}). \end{aligned}$$

In Line 2, we can omit  $\mathbf{E}_Y$  since  $\mathbf{E}_Y \perp\!\!\!\perp \mathbf{E}_X$ ,  $\mathbf{E}_Y \perp\!\!\!\perp \tilde{\mathbf{Z}}$  by definition. Due to the Markov chain structure,  $\mathbf{E}_Y \perp\!\!\!\perp \tilde{\mathbf{X}}$ , as well. Similarly, in Line 3 we can exploit that  $\tilde{\mathbf{X}} \perp\!\!\!\perp \mathbf{E}_X$  and  $\tilde{\mathbf{Y}} \perp\!\!\!\perp \mathbf{E}_X$ .  $\square$

Simply put, Proposition 4.1 shows that, for a data generative process as defined in Model 1, estimating  $I(\mathbf{X}; \mathbf{Y})$  is equivalent to estimating  $I(\tilde{\mathbf{X}}; \tilde{\mathbf{Y}})$ . If we are provided only a finite number of iid samples, which is the case in practice, it allows us to significantly improve the sample efficiency for estimating  $I(\mathbf{X}; \mathbf{Y})$ . In particular, assume that we could recover the low-dimensional manifold  $(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}})$  with  $\tilde{d}$  dimensions, which preserves the shared information on the high-dimensional space  $(\mathbf{X}, \mathbf{Y})$  with  $d \gg \tilde{d}$  dimensions. This could generally lead to a significant improvement of sample complexity for MI estimators. For the KSG estimator, for example, the bias would decrease from  $\mathcal{O}(n^{-1/d})$  to  $\mathcal{O}(n^{-1/\tilde{d}})$ , and hence only depend on the relevant dimensions.

In practice, this insight suggests a two-step procedure. In the first step, we aim to learn the low-dimensional manifold  $(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}})$  from  $(\mathbf{X}, \mathbf{Y})$ , and in the second step, we estimate the mutual information between  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$  from the distances on the learned manifold  $(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}})$ . In the next section, we propose such an approach based on Geodesic Forests.

## 5 Geodesic Mutual Information Estimation

In this section, we explain how to estimate mutual information via geodesic distances, where we first introduce our new estimator and then explain how we compute the corresponding quantities via Geodesic Forests.

**5.1 Geodesic KSG** To efficiently estimate mutual information from data generated according to Model 1, we propose to first learn the embedded low-dimensional manifold of the joint space  $(\mathbf{X}, \mathbf{Y})$  via Geodesic Forests [18], as described subsequently in Sec 5.2, and then compute the local  $k$ NN distances from this representation. In other words, we aim to approximate the distance between two points via the length of its shortest path on the manifold, i.e. its geodesic distance.

More specifically, we approximate  $d(\mathbf{x}_i, \mathbf{x}_j)$  and  $d(\mathbf{y}_i, \mathbf{y}_j)$  in Eq. 3 with geodesic distances  $d_G(\mathbf{x}_i, \mathbf{x}_j)$  and

$d_G(\mathbf{y}_i, \mathbf{y}_j)$  obtained via marginalization (see Sec. 5.3) from the Geodesic Forest trained on  $(\mathbf{X}, \mathbf{Y})$ . To compute the distances on the joint space  $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ , we follow the KSG approach [15] and stick to the maximum between the distances on  $\mathbf{X}$  and  $\mathbf{Y}$ . Thus, we define  $\frac{1}{2}\rho_{i,k}^G$  as the distance to the  $k$ th neighbor on the joint space  $\mathbf{Z}$  and obtain  $n_{\mathbf{x},i}^G$  as

$$n_{\mathbf{x},i}^G = \left| \left\{ \mathbf{x}_j : d_G(\mathbf{x}_i, \mathbf{x}_j) < \frac{1}{2}\rho_{i,k}^G, i \neq j \right\} \right|,$$

and compute  $n_{\mathbf{y},i}^G$  accordingly. Finally, we derive our proposed G-KSG estimator as

$$I^{\text{G-KSG}}(\mathbf{X}; \mathbf{Y}) = \psi(k) + \psi(n) - \frac{1}{n} \sum_{i=1}^n \psi(n_{\mathbf{x},i}^G + 1) + \psi(n_{\mathbf{y},i}^G + 1).$$

Next, we explain how to estimate geodesic distances  $d_G$  from Geodesic Forests.

**5.2 Geodesic Forests** The term *Geodesic Forest* (GF) has been introduced by Madhyastha et al. [18] and describes an unsupervised version of sparse projection oblique random forests [28]. In a nutshell, each node of a tree in a GF is split based on a sparse linear projection of each data point onto a one-dimensional feature. Classical splitting criteria allow to compute binary splits of the projected samples in this 1D space efficiently. For a collection of trees, the relative geodesic similarity of two data points  $\mathbf{x}_i, \mathbf{x}_j$  is estimated by the fraction of leaves they occur in *together*, which has been proven successful to estimate geodesic distances even in the presence of many noise dimensions [18].

More formally, given a sample  $\mathbf{x}^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of a  $d_{\mathbf{X}}$ -dimensional random vector  $\mathbf{X}$ , GF builds a set of  $T$  trees. Each tree  $T_i$  is trained on a bootstrapped subsample of size  $m < n$ , as typical for learning random forests [1]. To grow a tree, we recursively split each parent node into its two child nodes until a certain stopping criterium is met. The two critical features, in which Geodesic Forests are different from classical random forests are the node splitting and the stopping criterium, which we describe in more detail below.

**Node Splitting** Instead of splitting on a random feature, GF computes  $p$  sparse random projections of the feature space and splits on that 1D projection, which minimizes the fast-BIC criterium (see below). To generate sparse projections, GF samples a random projection matrix  $\mathbf{A} \in \{-1, 0, 1\}^{d_{\mathbf{X}} \times p}$ , where an entry  $a_{ij}$  is non-zero with probability  $\lambda$ , i.e.  $P(a_{ij} = 1) = P(a_{ij} = -1) = \frac{\lambda}{2}$ , and zero otherwise. The sparsity parameter  $\lambda$  is typically set to  $\lambda = \frac{1}{d_{\mathbf{X}}}$ . Given projection

matrix  $\mathbf{A}$ , the projected feature matrix is  $\mathbf{X}' = \mathbf{A}^T \mathbf{X}$ , from which we can extract  $p$  one-dimensional features, which are each evaluated by the splitting criterion.

**Fast-Bic** To find a cut-point in a one-dimensional vector, the authors of GF [18] introduce fast-BIC, which is a regularized version of the classical two-means criterium [5]. Both criteria induce a hard cluster assignment to either the left or right cluster.

The general Bayesian Information Coefficient (BIC) for a model  $M$  with parameter vector  $\theta_M$  of length  $|\theta_M|$  can be written as  $BIC(M) = -2 \log \hat{L} + \log(n)|\theta_M|$ , where  $\log \hat{L}$  is the empirical log-likelihood of the data given model  $M$ . In our case, the model consists of five parameters, the cluster assignment and the parameters  $\hat{\mu}_i$  and  $\hat{\sigma}_i^2$ , for  $i \in \{1, 2\}$ , which parameterize the assumed Gaussian distribution for cluster  $i$ . Accordingly, the empirical negative log likelihood is defined as

$$-\log \hat{L} = \sum_{i=1}^2 \frac{n_i}{2} (2 \log w_i - \log 2\pi \hat{\sigma}_i^2),$$

where  $w_i = n_i/n$  is the probability of a data point being assigned to cluster  $i$ . Given an ordered one-dimensional sequence  $x_1, \dots, x_n$ , such that  $x_i \leq x_{i+1}$ , we can compute the optimal split point in  $\mathcal{O}(n^2)$  time, since we need to obtain for each of the potential  $n - 1$  cut-points the variances for both clusters. We can, however, compute it even faster.

**Faster Fast-Bic** For an ordered sequence, we can reduce the runtime complexity, to determine the best split point, to  $\mathcal{O}(n)$ . Hence, even for an unordered sequence, we obtain a runtime in  $\mathcal{O}(n \log n)$  by sorting, which is still faster than the original fast-BIC computation. To achieve this speed-up, we use an elegant trick developed for segmentation. As derived by Terzi [27, Ch. 2], we can compute the empirical variance  $\hat{\sigma}_{i,j}^2$  for an arbitrary segment  $x_i, \dots, x_j$ , with  $1 \leq i \leq j \leq n$  as

$$\hat{\sigma}_{i,j}^2 = \frac{1}{j-i+1} \left( (css_j - css_{i-1}) - \frac{1}{j-i+1} (cs_j - cs_{i-1})^2 \right),$$

where  $cs_i = \sum_1^i x_i$  is the cumulative sum of the first  $i$  entries of the ordered sequence  $x_1, \dots, x_n$  and  $css_i = \sum_1^i x_i^2$  the corresponding sum of squares, with  $css_0 = cs_0 = 0$ . In other words, after precomputing  $cs$  and  $css$  in linear time, we can compute the variance for an arbitrary segment in constant time.

**5.3 Approximate Geodesic Distances** To obtain a dissimilarity measure from a random forest, we can utilize the proximity score for random forests proposed by Breiman [1]. That is, let  $L_{ij}$  denote the number of trees for which data points  $i$  and  $j$  end up in the

same leaf and let  $T$  be the number of trees, the proximity score between two data points  $i$  and  $j$  is defined as  $p^F(\mathbf{x}_i, \mathbf{x}_j) = L_{ij}/T$ . Since  $p^F(\mathbf{x}_i, \mathbf{x}_j) \in [0, 1]$ , we can compute a dissimilarity between two points as  $d_F(\mathbf{x}_i, \mathbf{x}_j) = 1 - p^F(\mathbf{x}_i, \mathbf{x}_j)$ . In their empirical evaluation, Madhyastha et al. [18] demonstrate that  $d_F$  robustly recovers geodesic neighborhoods from high-dimensional data with many independent or noise dimensions. As geodesic nearest neighbors, they define points that lie close on the low-dimensional manifold.

In theory, estimating the geodesic  $k$ -nearest neighbors is exactly what we are after, however,  $d_F$  is not a proper distance metric. In particular,  $d_F$  satisfies the reflexivity property ( $d_F(\mathbf{x}_i, \mathbf{x}_i) = 0$ ), the non-negativity property, and the symmetry property ( $d_F(\mathbf{x}_i, \mathbf{x}_j) = d_F(\mathbf{x}_j, \mathbf{x}_i)$ ). Despite those,  $d_F(\mathbf{x}_i, \mathbf{x}_j) = 0$  does not imply that  $\mathbf{x}_i = \mathbf{x}_j$ , since two points could always end up in the same leaf, especially for small forests. Thus, the definiteness property is violated. In addition,  $d_F$  does not satisfy the triangle inequality, i.e.  $d_F(\mathbf{x}_i, \mathbf{x}_k) \leq d_F(\mathbf{x}_i, \mathbf{x}_j) + d_F(\mathbf{x}_j, \mathbf{x}_k)$  cannot be guaranteed.

Hence, we propose an adjusted dissimilarity measure. The key idea of this adjusted measure builds upon the fact that locally a manifold resembles a Euclidean space, thus the geodesic distances locally become the  $L_2$ -norm. In our context, we assume that two points are close on the manifold, if  $d_F(\mathbf{x}_i, \mathbf{x}_j) = 0$ . Under this premise, we propose the distance measure  $d_G$ , i.e.

$$d_G(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} d_F(\mathbf{x}_i, \mathbf{x}_j) & \text{if } d_F(\mathbf{x}_i, \mathbf{x}_j) > 0, \\ \frac{d_2(\mathbf{x}_i, \mathbf{x}_j)}{c(T+\epsilon)} & \text{otherwise.} \end{cases}$$

In short, for all pairs  $i, j$  for which  $d_F(\mathbf{x}_i, \mathbf{x}_j) = 0$ , we approximate their local geodesic distance via their Euclidian distance normalized by a constant factor. The normalization factor ensures that the normalized distances are always smaller than  $\frac{1}{T}$ , i.e. the smallest non-zero value that  $d_F$  can attain. It consists of  $c$ , the maximum  $L_2$ -norm between any two pairs  $i, j$ , the number of trees  $T$  and a small constant  $\epsilon > 0$ .

As a result,  $d_G$  satisfies reflexivity, non-negativity and symmetry, and in addition, satisfies definiteness and locally (for those points, for which  $d_F(\mathbf{x}_i, \mathbf{x}_j) = 0$ ) also satisfies the triangle inequality. We argue that possible violations of the triangle inequality for data points, for which  $d_F(\mathbf{x}_i, \mathbf{x}_j) > 0$ , are on average not relevant for  $k$ NN estimation with small  $k \leq 10$  in a high-dimensional setting, which is our main use-case.

**Marginal Distances** Next, we briefly outline how we can compute the distance between two points  $i, j$  on a subspace  $\mathcal{S}$  of  $\mathbf{X}$  given the forest learned on  $\mathbf{X}$ .

In essence, we can compute  $d_G(\mathbf{s}_i, \mathbf{s}_j)$  in a straight forward manner. To compute the local distances, i.e.  $\frac{d_2(\mathbf{s}_i, \mathbf{s}_j)}{c(T+\epsilon)}$ , we need to set  $c$  to refer to the maximum

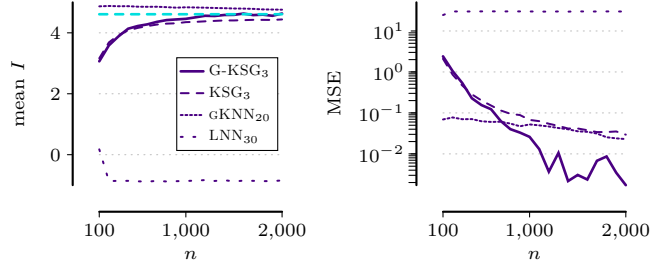


Figure 3: MI estimates with increasing sample size on uniform data (no noise), true MI is dashed cyan line.

distance between two points in the subspace  $\mathcal{S}$ , and recompute  $d_F(\mathbf{s}_i, \mathbf{s}_j)$  for each pair  $i, j$ . To compute  $d_F$  for a subspace  $\mathcal{S}$ , we first need to recompute all leaf assignments. That is, given a tree  $T_j$ , we assign each point  $i$  to that leaf in  $T_j$ , to which it would be assigned, if projected onto  $\mathcal{S}$ . After reassigning the leaves, we can compute  $d_F(\mathbf{s}_i, \mathbf{s}_j)$  as above. Based on the above procedure, we can compute  $d_G(\mathbf{x}_i, \mathbf{x}_j)$  and  $d_G(\mathbf{y}_i, \mathbf{y}_j)$  for each pair  $i, j$  and use these distances to compute our G-KSG estimator as described in Sec. 5.1.

Next, we will empirically evaluate G-KSG.

## 6 Experiments

We extensively evaluate G-KSG against state-of-the-art MI estimators KSG, gKNN, and LNN [15, 16, 10]. gKNN uses ellipsoids computed from principal components to better fit the local data distribution, LNN uses KDE with bandwidths automatically determined from the nearest neighbors. In particular, we compare on standard synthetic benchmark data as well as two simulated manifolds with known baseline MI for varying sample sizes, dimensionality of the data, and neighborhood size for each estimator. For all experiments, we train geodesic forests with original parameters [18], i.e.  $\lambda = 1/d$  number of dimensions,  $T = 300$  trees, and  $\sqrt{2n}$  minimum number of points to split a node. We use hyperparameters as suggested by the respective methods, details on which can be found in Supplementary Material S.2, and report the average across MI estimates of 20 repetitions for all experiments. For reproducibility, we make code and data publicly available.<sup>2</sup>

**6.1 Synthetic Data** We first evaluate the performance for estimating the MI between two variables generated from simple distributions, uniform and Gaussian.

**Sample Efficiency** To measure how well MI can be estimated with respect to the sample size, we draw datasets of size 100 to 2000 of  $X$  and  $Y$ , where  $X$  is

<sup>2</sup><https://github.com/a-marx/geodesic-mi>



uniformly distributed between 0 and 1 and  $Y = X + N$  with  $N \sim \text{Unif}(-\alpha/2, \alpha/2)$  and  $\alpha = 0.01$ . The ground truth MI is given by  $I(X; Y) = h(Y) - h(Z) = \frac{\alpha}{2} - \log \alpha$  [3, Ex. 8.3]. Note that no independent variables  $Z$  are added to the data, yet. We observe that even without independent variables in the data, G-KSG is able to more efficiently estimate the MI, with an order of magnitude lower mean squared error (MSE) than classical KSG (see Fig. 3). We further see that LNN greatly underestimates the true mutual information on this uniform data. While efficient even for as few as 100 samples, gKNN constantly overestimates the true MI slightly, showing an order of magnitude larger mean squared error for  $n \geq 1000$  samples compared to G-KSG. Note that, despite its complexity, G-KSG is only a factor 10 slower than classical KSG, regardless of samples size (see Supplementary Material S.3). For the rest of the experiments, we will use  $n = 500$  samples, which is the largest sample size where the original KSG could keep up with G-KSG.

**Uniform** For the same data as above, we now add an increasing number of  $\{0, 2, \dots, 20\}$  dimensions each sampled from a standard normal distribution to the original data. We report the results in Fig. 5 (top right), where we can observe that while gKNN, KSG, and G-KSG yield good estimates of the true MI without additional dimensions, LNN drastically overestimates the true MI, exponentially increasing with the number of dimensions. The predicted MI of classical KSG quickly falls to 0 for as few as 4 dimensions and is hence useless for this data. While also decreasing slightly, G-KSG yields the most stable prediction of MI with respect to dimensionality, and with the lowest error on 15 or more dimensions. gKNN underestimates even the original data without additional dimensions, and overestimates for higher-dimensional data, having numerical issues already for 10 additional dimensions.

**High-Dimensional Data** For the same uniform distribution, we generate data and add  $\{50, 100, \dots, 600\}$  independent dimensions, adding half of the dimensions to  $X$  and the other half to  $Y$ . As expected, while KSG immediately estimates a MI of 0, and is therefore not useful at all, the MI estimate of G-KSG decreases slightly but then remains stable with an increasing number of (independent) dimensions in the data. Both gKNN as well as LNN fail to yield any result even for 50 dimensions due to numerical errors. We provide results on this experiment in Supplementary Material S.3. Next, we investigate the behavior for variables from a different distribution.

**Gaussian** We generate synthetic data of Gaussian distributed random variables  $X$  and  $Y$  with zero mean, unit variance and covariance of 0.9. Consequently, with

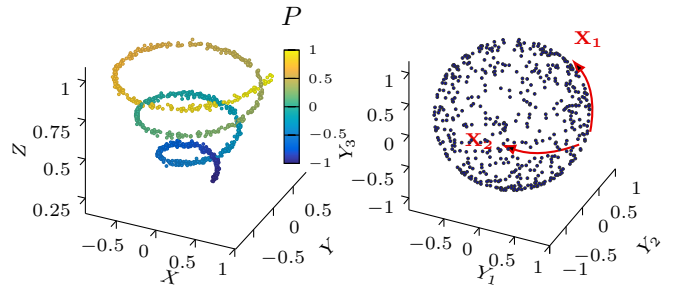


Figure 4: [Simulation Data] Shown are samples from the **Sphere** with input  $(X_1, X_2)$  corresponding to (longitude, latitude) and output  $(Y_1, Y_2, Y_3)$  cartesian coordinates (left), and **Helix** with input  $P$  indicated by color, and cartesian output coordinates  $(X, Y, Z)$  (right).

correlation  $\rho$  between  $X$  and  $Y$  being 0.9, the true MI can be calculated as  $I(X; Y) = -\frac{1}{2} \log(1 - \rho^2)$ . For varying neighborhood sizes for density estimation, we report the results in Fig. 5 (top left) for an increasing number of independent dimensions added to the original data. We observe that both gKNN as well as LNN provide a decent estimate on the simple Gaussian data without additional dimensions, but then greatly overestimate the true MI for as few as 5 independent or noise dimensions. gKNN does not yield meaningful MI estimates beyond 10 additional dimensions. Classical KSG provides a robust estimate on the simple Gaussian data, but quickly deteriorates with increasing dimensions, underestimating the true MI capturing almost zero mutual information. While G-KSG also experiences the same effect, it does so much more slowly and consistently yields the lowest error in terms of the true MI.

**6.2 Simulation Study** To test our approach for more complex data, and in particular data distributions resembling a manifold, we consider two simulated data sets. The first dataset we study is a sphere, resembling e.g. a planet (see Fig. 4 left). Data is distributed on the sphere, for which  $\mathbf{X}$  is given as longitude and latitude, and  $\mathbf{Y}$  is given as the 3D coordinates. The second dataset is a helix (see Fig. 4 right), for which  $\mathbf{X}$  specifies the distance along the helix (1D), and  $\mathbf{Y}$  is given as the 3D coordinates, a typical problem from manifold learning. We give details on how these datasets are sampled and how to compute a lower bound on the ground truth MI in Supplementary Material S.1.

**6.3 Helix** First, we consider the **Helix** data. Varying the number of independent dimensions for this data, we see that KSG estimates that the data contains 0 mutual information, as soon as any independent dimensions

are added, whereas our method maintains a stable MI estimate even for many independent dimensions (see Fig. 5 bottom left). For LNN we observe a sharp linear increase in predicted mutual information as a function of number of dimensions, first under and then over-predicting the true MI by a wide margin. For gKNN, we observe first an overestimation of the MI which then comes close to the true MI for 2-6 additional dimensions. However, gKNN is not able to yield an estimate for more than 2 additional dimensions when using  $k = 20$  respectively more than 6 additional dimensions for  $k = 30$ . Increasing the neighborhood size further does not make sense, as gKNN would then miss on the locality of the data when computing the volume, thus estimating global rather than local density.

**6.4 Sphere** Next, we consider a simulation using the **Sphere** data, a simple dataset which, however, shows to be an astonishingly hard challenge for the state-of-the-art MI estimators. We show the results in Fig. 5 bottom right. Overall, we observe similar trends as for **Helix**, giving further evidence that classical MI estimates are not able to capture manifolds within data with a large number of independent or noise dimensions. As before, KSG quickly deteriorates to estimate 0 mutual information. Similarly, LNN shows a steep, near linear dependence between predicted MI and independent dimensions, which has little to do with the true MI. Again, gKNN shows to perform poorly on the original task without added dimensions and fails to compute a result due to numerical issues on data of more than a handful of dimensions. On this dataset, G-KSG consistently predicts MI close to the ground truth, regardless dimensions. Even for hundreds of added dimensions, G-KSG remains stable estimates, whereas all other methods fail to do so (see Supplementary Material S.3).

## 7 Discussion & Conclusion

Mutual information, due to its non-parametric nature, captures non-linear dependencies and, hence, lends itself for measuring complex dependencies between random variables, and is routinely applied for challenging tasks such as gene network inference. Yet, estimation of MI on high-dimensional data remains an open problem. In this work, we considered the problem of estimating MI for high-dimensional data under the manifold assumption, i.e. that the data has low intrinsic dimensionality.

To tackle this problem, we combined ideas from classical MI estimation and manifold learning. In particular, we proposed to use geodesic distances to estimate local densities on the manifold rather than in the ambient space. We leveraged the recently proposed Geodesic

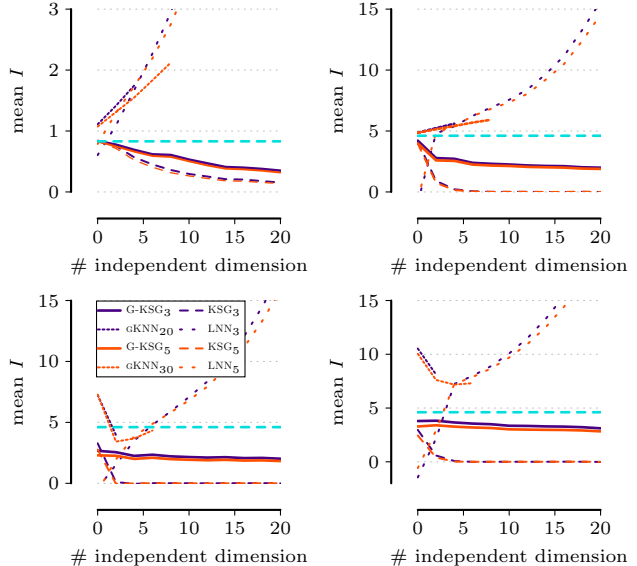


Figure 5: Mutual information estimates of G-KSG, KSG, gKNN, and LNN for different number of  $k$  indicated by method subscript for an increasing number of independent dimensions ( $n=500$ ). Top-left Gaussian correlation, top-right uniform linear, bottom-left **Helix**, bottom-right **Sphere**, true MI is the dashed cyan line.

Forests, which can estimate relative geodesic distances even for high-dimensional data with many independent or noise dimensions. To be able to scale to our setting, we further proposed a faster algorithm to compute the splitting criterion for individual nodes. Based on computed geodesic distances, we then extended the state-of-the-art KSG MI estimator to operate on these distances, leading to the G-KSG estimator.

We evaluated G-KSG against state-of-the-art MI estimators on standard benchmark data and two simulated manifold datasets. The results show that G-KSG outperforms its competitors, better estimating the true MI across tasks and settings. Furthermore, it provides more stable results across data of different dimensionality, whereas the state-of-the-art fails to scale to higher dimensions, or greatly over- or underestimates the true MI. At the same time, G-KSG is sample efficient, and is only slightly slower than classical KSG, independent of sample size.

In summary, our approach allows us to efficiently and robustly estimate MI of high-dimensional data modeling the manifold assumption. While thorough empirical evaluation showed that G-KSG ably estimates the true MI across different datasets and dimensionalities, it would make for engaging future work to study the theoretical guarantees of our estimator. The estima-



tion of geodesic distances renders this theoretical aspect extremely challenging. Besides, we would be interested to study MI estimation for different settings, such as discrete-continuous mixtures or on time-series.

## Acknowledgements

AM is supported by the ETH AI Center. JF is supported by the International Max Planck Research School for Computer Science (IMPRS-CS).

## References

- [1] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [2] L. Cayton, “Algorithms for manifold learning,” *Univ. of California at San Diego Tech. Rep.*, vol. 12, no. 1-17, p. 1, 2005.
- [3] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Wiley-Interscience New York, 2006.
- [4] G. A. Darbellay and I. Vajda, “Estimation of the information by an adaptive partitioning of the observation space,” *IEEE Trans. Inf. Theorie*, vol. 45, no. 4, pp. 1315–1321, 1999.
- [5] S. Dasgupta and Y. Freund, “Random projection trees and low dimensional manifolds,” in *STOC*, 2008, pp. 537–546.
- [6] S. Frenzel and B. Pompe, “Partial mutual information for coupling analysis of multivariate time series,” *Physical Review Letters*, vol. 99, no. 20, p. 204101, 2007.
- [7] S. Gao, G. Ver Steeg, and A. Galstyan, “Efficient estimation of mutual information for strongly dependent variables,” in *AISTATS*. PMLR, 2015, pp. 277–286.
- [8] W. Gao, S. Kannan, S. Oh, and P. Viswanath, “Estimating mutual information for discrete-continuous mixtures,” in *NIPS*, 2017, pp. 5986–5997.
- [9] W. Gao, S. Oh, and P. Viswanath, “Breaking the bandwidth barrier: Geometrical adaptive entropy estimation,” in *NIPS*. Curran Associates, Inc., 2016, pp. 2460–2468.
- [10] —, “Density functional estimators with k-nearest neighbor bandwidths,” in *ISIT*. IEEE, 2017, pp. 1351–1355.
- [11] —, “Demystifying fixed k-nearest neighbor information estimators,” *IEEE Trans. Inf. Theorie*, vol. 64, no. 8, pp. 5629–5661, 2018.
- [12] P. Grassberger, “Generalizations of the hausdorff dimension of fractal measures,” *Physics Letters A*, vol. 107, no. 3, pp. 101–105, 1985.
- [13] M. Koeman and T. Heskes, “Mutual information estimation with random forests,” in *ICONIP*. Springer, 2014, pp. 524–531.
- [14] L. F. Kozachenko and N. N. Leonenko, “Sample estimate of the entropy of a random vector,” *Probl. Peredachi Inf.*, vol. 23, pp. 9–16, 1987.
- [15] A. Kraskov, H. Stögbauer, and P. Grassberger, “Estimating mutual information,” *Phys. Rev. E*, vol. 69, no. 6, p. 066138, 2004.
- [16] W. M. Lord, J. Sun, and E. M. Bollt, “Geometric k-nearest neighbor estimation of entropy and mutual information,” *Chaos: An Interdisciplinary Journal of Nonlinear Science*, vol. 28, no. 3, p. 033114, 2018.
- [17] C. Lu and J. Peltonen, “Enhancing nearest neighbor based entropy estimator for high dimensional distributions via bootstrapping local ellipsoid,” in *AAAI*, vol. 34, no. 04, 2020, pp. 5013–5020.
- [18] M. Madhyastha, G. Li, V. Strnadová-Neeley, J. Browne, J. T. Vogelstein, R. Burns, and C. E. Priebe, “Geodesic forests,” in *KDD*, 2020.
- [19] P. Mandros, D. Kaltenpoth, M. Boley, and J. Vreeken, “Discovering functional dependencies from mixed-type data,” in *KDD*, 2020, pp. 1404–1414.
- [20] A. Marx and J. Vreeken, “Testing conditional independence on discrete data using stochastic complexity,” in *AISTATS*, 2019, pp. 496–505.
- [21] A. Marx, L. Yang, and M. van Leeuwen, “Estimating conditional mutual information for discrete-continuous mixtures using multi-dimensional adaptive histograms,” in *SDM*, 2021, pp. 387–395.
- [22] L. McInnes, J. Healy, and J. Melville, “Umap: Uniform manifold approximation and projection for dimension reduction,” 2020.
- [23] R. Mehta, R. Guo, J. Arroyo, M. Powell, H. Helm, C. Shen, and J. T. Vogelstein, “Estimating information-theoretic quantities with uncertainty forests,” *arXiv preprint arXiv:1907.00325*, 2019.
- [24] L. Paninski and M. Yajima, “Undersmoothed kernel entropy estimators,” *IEEE Trans. Inf. Theorie*, vol. 54, no. 9, pp. 4384–4388, 2008.
- [25] J. Pearl, *Causality: Models, Reasoning and Inference*, 2nd ed. Cambridge University Press, 2009.
- [26] J. B. Tenenbaum, V. De Silva, and J. C. Langford, “A Global Geometric Framework for Nonlinear Dimensionality Reduction,” *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [27] E. Terzi, *Problems and algorithms for sequence segmentations*. Helsingin yliopisto, 2006.
- [28] T. M. Tomita, J. Browne, C. Shen, J. Chung, J. L. Patsolic, B. Falk, C. E. Priebe, J. Yim, R. Burns, M. Maggioni, and J. T. Vogelstein, “Sparse projection oblique randomer forests,” *JMLR*, vol. 21, no. 104, 2020.
- [29] G. Valiant and P. Valiant, “Estimating the unseen: an  $n/\log(n)$ -sample estimator for entropy and support size, shown optimal via new clts,” in *STOC*, 2011, pp. 685–694.
- [30] L. van der Maaten and G. Hinton, “Visualizing data using t-sne,” *JMLR*, vol. 9, no. 86, pp. 2579–2605, 2008.
- [31] M. Vejmelka and K. Hlaváčková-Schindler, “Mutual information estimation in higher dimensions: A speed-up of a k-nearest neighbor based estimator,” in *ICANNGA*. Springer, 2007, pp. 790–797.
- [32] Q. Wang, S. R. Kulkarni, and S. Verdú, “Divergence estimation for multidimensional densities via k-nearest-neighbor distances,” *IEEE Trans. Inf. Theorie*, vol. 55, no. 5, pp. 2392–2405, 2009.

## Supplementary Material

The Supplementary Material, is split into three sections. In the first section, we provide the data generating mechanisms for the **Helix** and **Sphere** data, for which we derive lower bounds on the ground truth mutual information. Subsequently, we list the hyperparameter for G-KSG and for the baselines in Section S.2, and provide additional experiments in Section S.3.

**S.1 Data Generation and MI Bounds** In the following, we will first briefly recap the derivation of the ground truth of the mutual information for the linear uniform scenario and then explain the data generation, as well as, the derivation of the lower bounds for the ground truth for the **Helix** and **Sphere** data.

To derive the ground truth value for the linear uniform example, where  $X \sim Unif(0, 1)$ ,  $Z \sim Unif(-\alpha/2, \alpha/2)$  and

$$Y = X + Z,$$

we follow Exercise 8.3 in [3]. First, note that

$$\begin{aligned} I(X; Y) &= h(Y) - h(Y | X) \\ &= h(Y) - h(Z) = h(Y) - \log \alpha. \end{aligned}$$

Thus, it remains to compute the entropy for  $Y$ , for which the density function can be computed as a convolution of  $X$  and  $Z$ . Since  $X$  and  $Z$  follow a uniform distribution,  $Y$  is a trapezoid and it can be derived that  $h(Y) = \frac{\alpha}{2}$  for  $\alpha \leq 1$  [3].

**Helix** For the **Helix** data, we generate a radius as  $R \sim Unif(0, 1)$ , which we then embed as a spiral in a three-dimensional space. Accordingly, we generate the  $X, Y, Z$  dimensions as

$$\begin{aligned} P &= 5\pi + 3\pi R, \\ X &= \frac{P \cos(P)}{8\pi} + N_1, \\ Y &= \frac{P \sin(P)}{8\pi} + N_2, \\ Z &= \frac{P}{8\pi} + N_3, \end{aligned}$$

where  $N_1, N_2, N_3 \sim Unif(-\alpha/2, \alpha/2)$  are noise variables and  $\alpha = 0.01$ .

To approximate the ground truth value of such a **Helix**, we can build upon the derivation for the uniform linear data to obtain a lower bound. A lower bound is sufficient in our case, since we do not measure how close we can estimate the true value, but how much information we can still recover in a noisy setting.

First, note that  $I(X; X + Z) = I(X; X \sin(X) + Z)$  (see [7]). More generally,  $I(P; \{X, Y, Z\}) =$

$I(P; \{X', Y', Z'\})$ , where  $X' = P + N_1$ ,  $Y' = P + N_2$  and  $Z' = P + N_3$ . We get that

$$\begin{aligned} I(P; \{X', Y', Z'\}) &= h(X', Y', Z') - h(X', Y', Z' | P) \\ &= h(X', Y', Z') - h(N_1, N_2, N_3) \\ &= h(X') + h(Y' | X') + h(Z' | X', Y') \\ &\quad - h(N_1, N_2, N_3), \end{aligned}$$

where again, we can rewrite  $h(X', Y', Z' | P)$  as  $h(N_1, N_2, N_3)$  similar to the linear case. Additionally, all noise terms are independent of each other and thus  $h(N_1, N_2, N_3) = h(N_1) + h(N_2) + h(N_3) = 3 \log \alpha$ . Now, due to the Markov chain structure, we can only approximate  $h(Y' | X') \geq h(Y' | P)$ . However, we conjecture that this approximation is quite close due to the low amount of noise added in the data generation. Similarly,  $h(Z' | X', Y') \geq h(Z' | P)$ . Thus

$$\begin{aligned} I(P; \{X', Y', Z'\}) &\geq h(X') + h(Y' | P) + h(Z' | P) - 3 \log \alpha \\ &= h(X') - \log \alpha \\ &= \frac{\alpha}{2} - \log \alpha. \end{aligned}$$

**Sphere** We generate samples from a **Sphere** by drawing latitude and longitude as  $X_1, X_2 \sim Unif(0, 1)$  and compute the cartesian coordinates as

$$\begin{aligned} Y_1 &= \cos(X_1) \cos(X_2)R + N_1, \\ Y_2 &= \cos(X_1) \sin(X_2)R + N_2, \\ Y_3 &= \sin(X_2)R + N_3. \end{aligned}$$

where  $N_1, N_2, N_3 \sim Unif(-\alpha/2, \alpha/2)$ , and we use a radius of  $R = 1$  and  $\alpha = 0.01$  for our experiments.

To derive a lower bound on the mutual information, we follow a similar procedure as for the **Helix** data. First, we can rewrite  $Y'_3 = X_2 + N_3$  and get that

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= h(Y_1, Y_2, Y'_3) - h(Y_1, Y_2, Y'_3 | X_1, X_2) \\ &= h(Y'_3) + h(Y_2 | Y'_3) + h(Y_1 | Y_2, Y'_3) \\ &\quad - h(Y_1, Y_2, Y'_3 | X_1, X_2), \end{aligned}$$

We can similarly decompose the conditional term to  $h(Y'_3 | X_1, X_2) + h(Y_2 | X_1, X_2, Y'_3) + h(Y_1 | X_1, X_2, Y'_3, Y_2)$ . From the linear case, we know that  $h(Y'_3) - h(Y'_3 | X_1, X_2) = \frac{\alpha}{2} - \log \alpha$ , since  $X_1 \perp\!\!\!\perp Y'_3$ . Additionally, it is clear that  $h(Y_2 | Y'_3) \geq h(Y_2 | X_1, X_2, Y'_3)$  and  $h(Y_1 | Y'_3, Y_2) \geq h(Y_1 | X_1, X_2, Y'_3, Y_2)$  and thus  $I(\mathbf{X}; \mathbf{Y}) \geq \frac{\alpha}{2} - \log \alpha$ .

**S.2 Hyperparameter** For G-KSG, we use projections with  $\lambda = 1/d$  number of dimensions,  $T = 300$

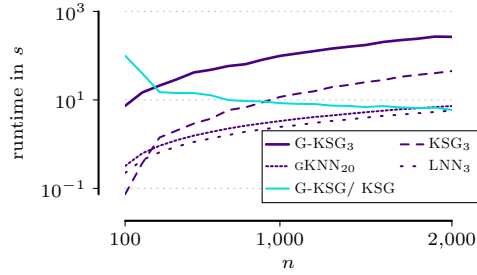


Figure 6: Shown is the runtime in seconds for G-KSG, KSG, gKNN and LNN for the experiment on sample efficiency in Sec 6.1. In addition, we show the quotient of the runtime of G-KSG divided by the runtime of KSG. For larger sample sizes ( $n \geq 1000$ ), KSG is faster by a factor smaller than 10.

trees, and  $\sqrt{2n}$  minimum number of points in a node to consider it for splitting. These parameters are suggested by the authors of geodesic forests [18].

In case of LNN, we use the suggested setting, which is  $k = 30$  neighbours for MI estimation and evaluate  $k' = 3$ , as well as  $k' = 5$ , for bandwidth prediction for the kernel density estimation.

For gKNN and KSG we set the neighbourhood size as specified in the main paper.

**S.3 Additional Experiments** Next, we supplement the runtimes for the experiment on linear uniform data with increasing sample size and provide the results for the experiments on high-dimensional data.

**Time comparison** We provide a time comparison of all methods in Fig. 6 for 10 repetitions. We observe that G-KSG only takes 10 times longer than KSG, independent of number of samples.

**High dimensional results** We provide the results for high dimensional data for the uniform linear data and **Sphere** across 3 repetitions in Fig. 7. We varied the number of independent dimensions in  $\{0, 100, 200, \dots, 600\}$ , splitting equally between  $X$  and  $Y$ . Both gKNN as well as LNN were not able to compute the MI estimate due to numerical issues even for as few as 50 dimensions.

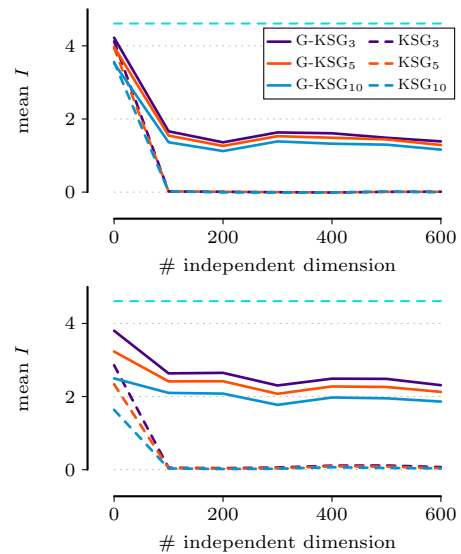


Figure 7: Mutual information estimates of G-KSG and KSG  $k = \{3, 5, 10\}$  on linear data with uniform source (top) and sphere data (bottom) with an increasing number of noise dimensions up to the high-dimensional setting ( $d > n = 500$ ). The dashed cyan line is a lower bound on the ground truth MI.