

CURSE OF SLICING: WHY SLICED MUTUAL INFORMATION IS A DECEPTIVE MEASURE OF STATISTICAL DEPENDENCE

The Fourteenth International Conference on Learning Representations

Alexander Semenenko, Ivan Butakov, Ivan Oseledets, Alexey Frolov

semenenko@applied-ai.ru, ivan.butakov@applied-ai.ru

The **mutual information (MI)** between random vectors $X : \Omega \rightarrow \mathcal{X}$ and $Y : \Omega \rightarrow \mathcal{Y}$ with joint distribution $\mathbb{P}_{X,Y}$ and marginals \mathbb{P}_X and \mathbb{P}_Y is defined as

$$I(X; Y) = \mathbb{E} \log \frac{d\mathbb{P}_{X,Y}}{d\mathbb{P}_X \otimes \mathbb{P}_Y} = \text{KL}[\mathbb{P}_{X,Y} \parallel \mathbb{P}_X \otimes \mathbb{P}_Y].$$

Properties:

- Nullification Property: $I(X; Y) = 0$ iff X and Y are independent
- Invariance to bijections, chain rule, etc.

Applications:

- Independence testing
- Feature selection
- Representation learning/disentanglement
- Analysis of DNNs

The k -sliced mutual information (k -SMI) (Goldfeld et al., 2022) between X and Y is defined as

$$SI_k(X; Y) = \int_{\text{St}(k, d_x)} \int_{\text{St}(k, d_y)} I(\Theta^\top X; \Phi^\top Y) d\mu_{\text{St}(k, d_x)}(\Theta) d\mu_{\text{St}(k, d_y)}(\Phi),$$

where $\mu_{\text{St}(k, d_x)}$ denotes normalized Haar (uniform) probability measure on a Stiefel manifold $\text{St}(k, d) = \{Q \in \mathbb{R}^{d \times k} : Q^\top Q = I\}$. Setting $k = 1$ recovers the SMI

Pros:

- Scalability
- Nullification Property: $SI_k(X; Y) = 0$ iff X and Y are independent
- Monotonic: $SI_{k_1}(X; Y) \leq SI_{k_2}(X; Y)$, $k_1 < k_2$
- Chain rule

Cons (not well-covered in the literature):

- Asymptotics in high-dimensional regime
- Data Processing Inequality violation
- Suboptimality of random slicing

Lemma 1. Consider the following pair of jointly Gaussian d -dimensional random vectors: $(X, Y) \sim \mathcal{N}\left(0, \begin{pmatrix} \mathbf{I} & \rho\mathbf{I} \\ \rho\mathbf{I} & \mathbf{I} \end{pmatrix}\right)$, $\rho \in (-1; 1)$. In this setup, MI and SMI can be calculated analytically:

$$I(X; Y) = -\frac{d}{2} \log(1 - \rho^2), \quad \text{SI}(X; Y) = \frac{\rho^2}{2d} {}_3F_2\left(1, 1, \frac{3}{2}; \frac{d}{2} + 1, 2; \rho^2\right),$$

where ${}_3F_2$ is the *generalized hypergeometric function*. Additionally, the following limits with the *digamma function* ψ hold:

$$\lim_{d \rightarrow \infty} I(X; Y) = +\infty \quad \lim_{d \rightarrow \infty} \text{SI}(X; Y) = 0$$

$$\lim_{\rho^2 \rightarrow 1} I(X; Y) = +\infty \quad \lim_{\rho^2 \rightarrow 1} \text{SI}(X; Y) = \psi(d - 1) - \psi\left(\frac{d - 1}{2}\right) - \log 2 \leq \frac{1}{d - 1}.$$

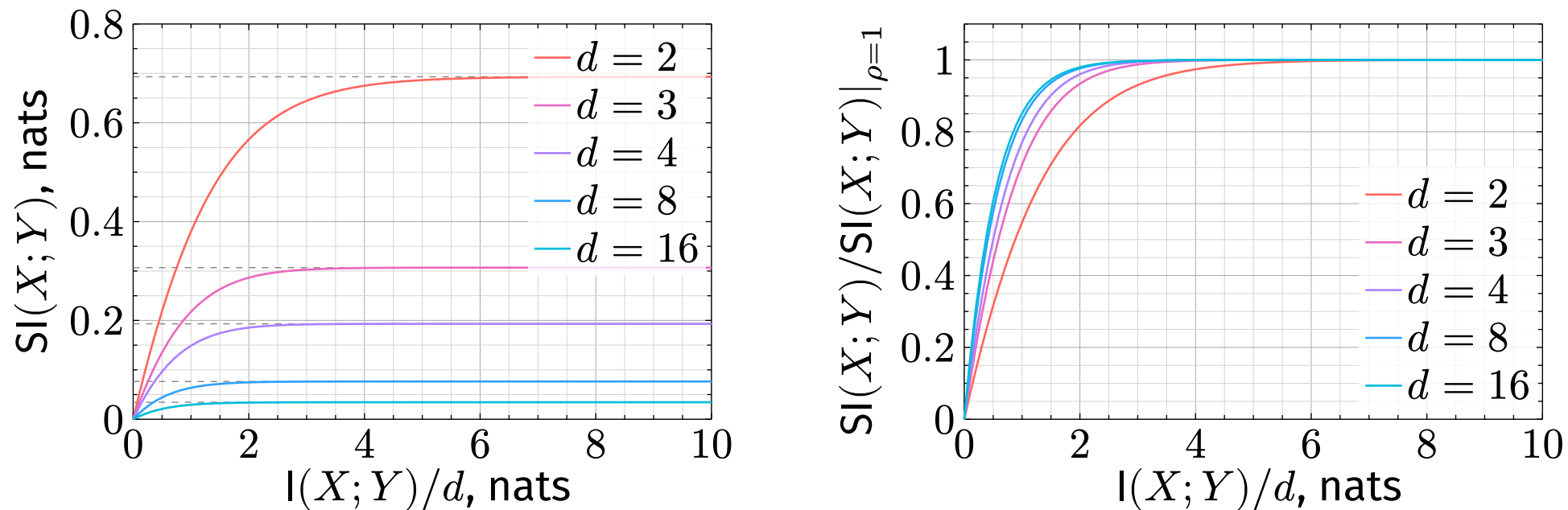
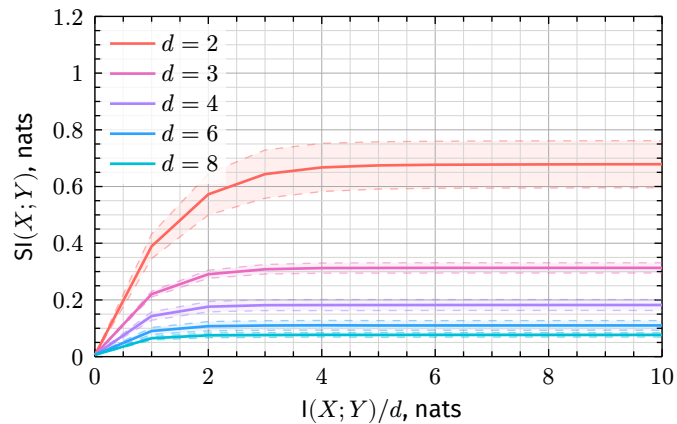
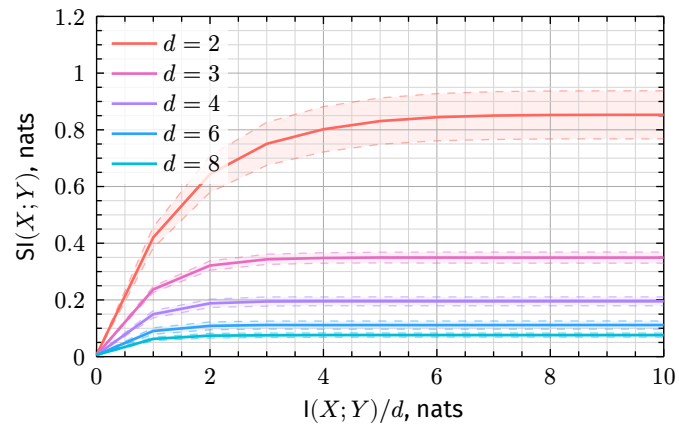


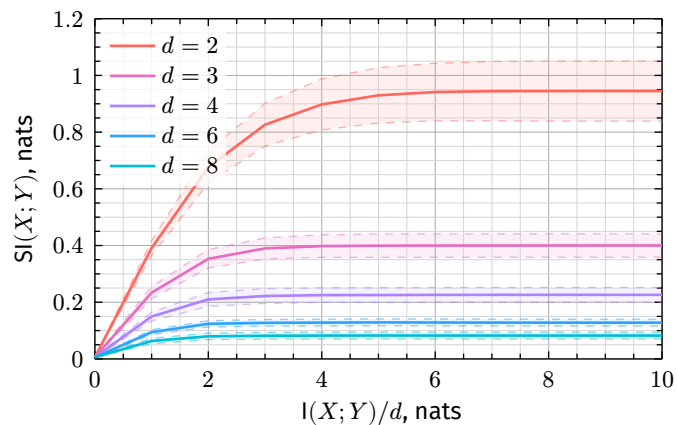
Figure 1: Saturation of $SI(X; Y)$ as function of $I(X; Y)/d$ for the example from [Lemma 1](#), non-normalized (left) and normalized (right) versions. Note that the problem becomes more prominent in higher dimensions, both because of lower plateau and faster saturation.



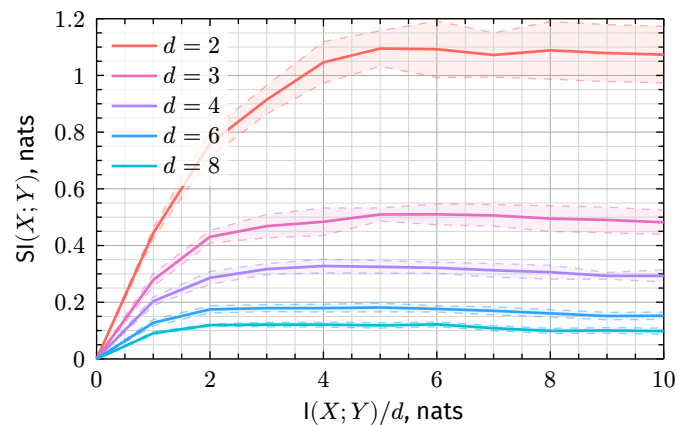
(a) Correlated Normal



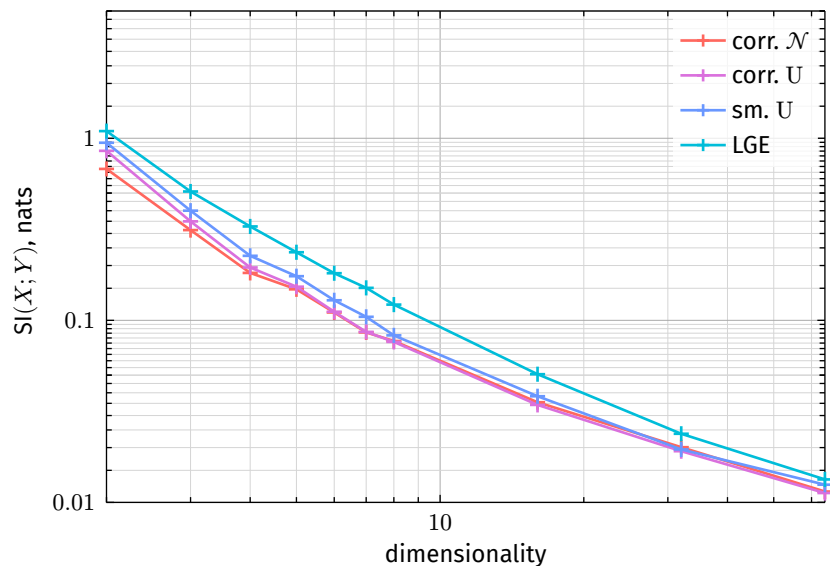
(b) Correlated Uniform



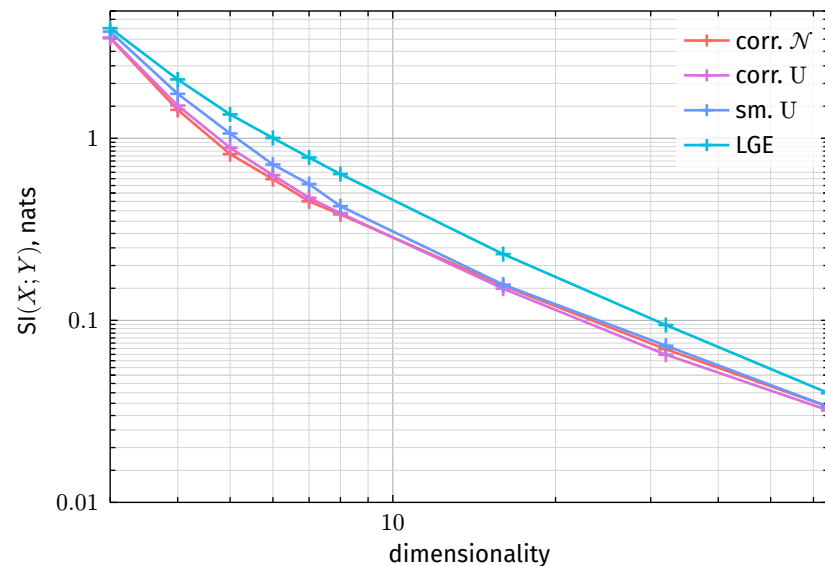
(c) Smoothed Uniform



(d) Log-Gamma-Exponential



(a) $k = 1$



(b) $k = 2$

Figure 3: Decaying trends of k -SMI for *correlated normal* (corr. \mathcal{N}), *correlated uniform* (corr. U), *smoothed uniform* (sm. U) and *log-gamma-exponential* (LGE).

Proposition 1. Let X and Y be d_x, d_y -dimensional random vectors correspondingly, with $d_x, d_y < k$. Let $A \in \mathbb{R}^{m_x \times d_x}$ and $B \in \mathbb{R}^{m_y \times d_y}$ be matrices of ranks d_x, d_y . Then $SI_k(AX; BY) = I(X; Y)$.

Corollary 1. Consider the following pair of jointly Gaussian d -dimensional random vectors:

$$(X, Y) \sim \mathcal{N}\left(0, \begin{pmatrix} J & \rho J \\ \rho J & J \end{pmatrix}\right), \quad \rho \in (-1; 1),$$

where $J = \mathbf{1} \cdot \mathbf{1}^\top$ with $\mathbf{1}^\top = (1, \dots, 1)$. Then $SI_k(X; Y) = I(X; Y) = -\frac{1}{2} \log(1 - \rho^2)$.

Remark 1. Applying $V = \mathbf{1} \cdot e_1^\top$ to the random vectors from [Lemma 1](#) individually yields the example from [Corollary 1](#). Therefore,

$$SI_k(VX; VY) > SI_k(X; Y), \quad \text{but} \quad I(VX; VY) < I(X; Y).$$

- Let X be a high-dimensional random vector and f be an encoder (approximated by a neural network). We aim to obtain low-dimensional representation $f(X)$ by maximizing MI (Hjelm et al., 2019):

$$I(X; f(X)) \rightarrow \max.$$

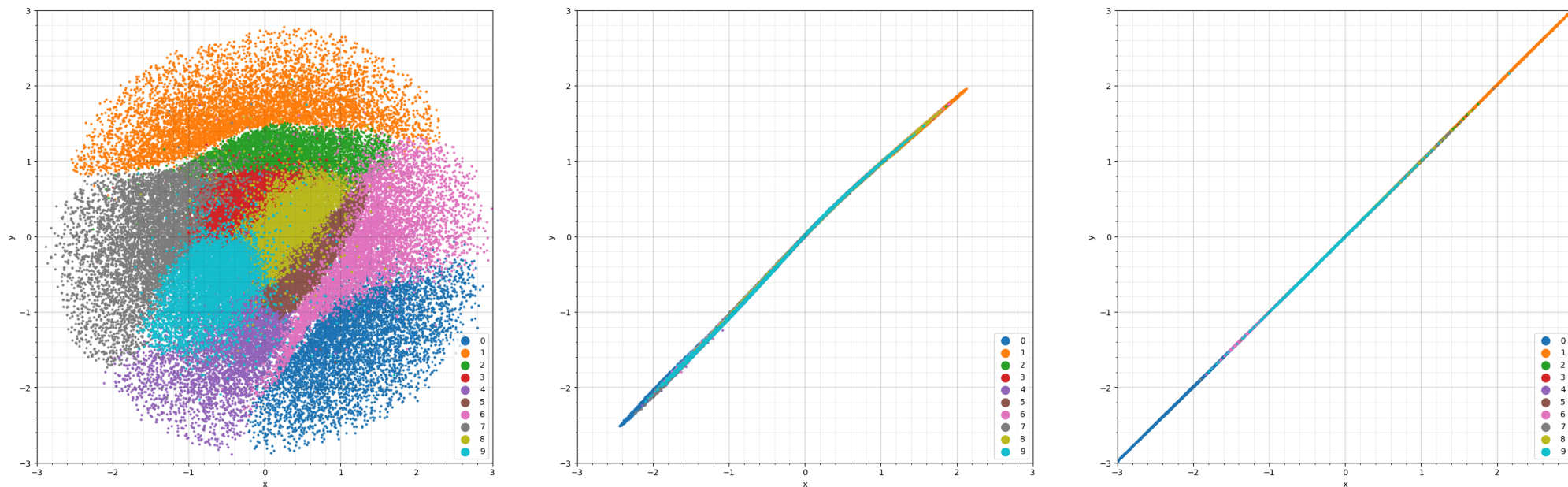
- We replace MI with SMI in Deep InfoMax framework

$$SI_k(X; f(X)) \rightarrow \max.$$

This substitution is easy since both MI and SMI admit Donsker-Varadhan bounds:

$$I(X; Y) = \sup_{T: \Omega \rightarrow \mathbb{R}} \left[\mathbb{E}_{\mathbb{P}_{X, Y}} T(X, Y) - \log \left(\mathbb{E}_{\mathbb{P}_X \otimes \mathbb{P}_Y} e^{T(X, Y)} \right) \right],$$

$$SI_k(X; Y) = \sup_{T: \Omega \rightarrow \mathbb{R}} \mathbb{E}_{\Theta, \Phi} \left[\mathbb{E}_{\mathbb{P}_{X, Y}} T(\Theta^\top X, \Phi^\top Y, \Theta, \Phi) - \log \left(\mathbb{E}_{\mathbb{P}_X \otimes \mathbb{P}_Y} e^{T(\Theta^\top X, \Phi^\top Y, \Theta, \Phi)} \right) \right].$$



(a) MI \rightarrow max, 2000 epochs. (b) SMI \rightarrow max, 10 epochs. (c) SMI \rightarrow max, 2000 epochs.

Figure 4: Visualizations of embeddings from the representation learning experiments, with points colored by class. Note that MI maximization (left) produces clustered low-redundancy representations, while SMI maximization results in immediate (after 10 epochs) collapse.

- **Replication of original SMI experiments:** feature extraction and independence testing. With minor modifications SMI fails where MI succeeds, contradicting earlier claims.
- **Theoretical analysis beyond Gaussians:** saturation and decay hold for general distributions with independent components.
- **k -SMI and optimal slicing (mSMI, oSMI)** do not escape the redundancy bias.
- **Gaussian channel & feature extraction:** SMI prefers ill-conditioned transformations (redundancy bias) and collapses to degenerate solutions in InfoMax tasks.

Thank you for your attention!



REFERENCES

- [1] Z. Goldfeld, K. Greenewald, T. Nuradha, and G. Reeves, “ k -Sliced Mutual Information: A Quantitative Study of Scalability with Dimension,” in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available: <https://openreview.net/forum?id=L-ceBdl2DPb>
- [2] R. D. Hjelm *et al.*, “Learning deep representations by mutual information estimation and maximization,” in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=Bklr3j0cKX>