
Curse of Slicing: Why Sliced Mutual Information is a Deceptive Measure of Statistical Dependence

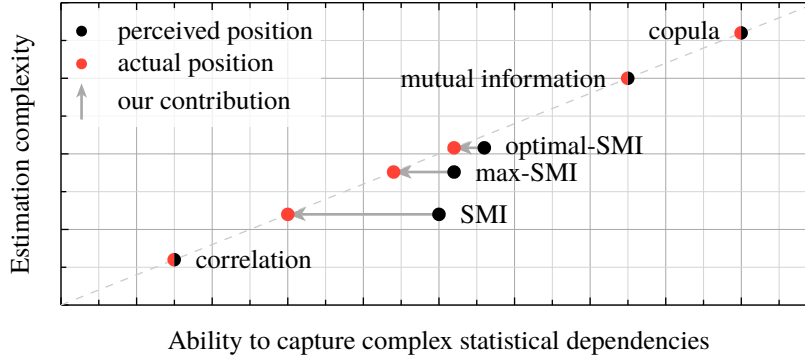
Anonymous Author(s)

Affiliation, Address

anon.email@example.org

Abstract

Sliced Mutual Information (SMI) is widely used as a scalable alternative to mutual information for measuring non-linear statistical dependence. Despite its advantages, such as faster convergence, robustness to high dimensionality, and nullification only under statistical independence, we demonstrate that SMI is highly susceptible to data manipulation and exhibits counterintuitive behavior. Through extensive benchmarking and theoretical analysis, we show that SMI saturates easily, fails to detect increases in statistical dependence (even under linear transformations designed to enhance the extraction of information), prioritizes redundancy over informative content, and in some cases, performs worse than simpler dependence measures like the correlation coefficient.



1 Introduction

Mutual information (MI) is a fundamental and invariant measure of nonlinear statistical dependence between two random vectors, defined as the Kullback-Leibler divergence between the joint distribution and the product of marginals [1]:

$$I(X; Y) = D_{\text{KL}}(\mathbb{P}_{X,Y} \parallel \mathbb{P}_X \otimes \mathbb{P}_Y).$$

Due to several outstanding properties, such as nullification only under statistical independence, invariance to invertible transformations, and ability to capture non-linear dependencies, MI is used extensively for theoretical analysis of overfitting [2], [3], hypothesis testing [4], feature selection [5], [6], [7], representation learning [8], [9], [10], [11], [12], [13], and studying the mechanisms behind generalization in deep neural networks (DNNs) [14], [15], [16], [17].

In practical scenarios, $\mathbb{P}_{X,Y}$ and $\mathbb{P}_X \otimes \mathbb{P}_Y$ are unknown, requiring MI to be estimated from finite samples. Despite all the aforementioned merits, this reliance on empirical estimates leads to the curse

of dimensionality: the sample complexity of MI grows exponentially with the number of dimensions [18], [19]. A common strategy to mitigate this issue is to use alternative measures of statistical dependence that are more stable in high dimensions. However, such measures usually offer only a fraction of MI capabilities. Therefore, it is crucial to maintain a balance between robustness to the curse of dimensionality and the ability to detect complex dependency structures.

To strike this balance, popular techniques often retain MI as a backbone statistical measure but employ dimensionality reduction before estimation. While some studies explore sophisticated nonlinear compression methods [17], [20], others favor more scalable linear projection approaches [21], [22], [23], [24], [25]. Among the latter group, the *Sliced Mutual Information* (SMI) [22], [23] stands out, leveraging random projections to cover all directions uniformly:

$$SI(X; Y) = \frac{1}{\oint_{\mathbb{S}^{d_x-1}} d\theta} \frac{1}{\oint_{\mathbb{S}^{d_y-1}} d\phi} \oint_{\mathbb{S}^{d_x-1}} \oint_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\theta d\phi. \quad (1)$$

Uniform slicing allows SMI to maintain some crucial properties of MI (e.g., being zero if and only if X and Y are independent), while remaining completely free from additional optimization problems (e.g., from finding optimal projections, as in [24], [25]). Combined with fast convergence rates, this has established SMI as a scalable alternative to MI. Consequently, it has been widely adopted for studying DNNs [26], [27], [28], [29], [30], deriving generalization bounds [31], independence testing [32] and auditing differential privacy [33]. It was also proposed to use SMI for feature selection [22] and preventing mode collapse in generative models [23].

Despite its popularity, the research community has largely overlooked potential shortcomings of SMI. Some studies prematurely attribute their results to underlying phenomena without rigorously investigating whether they stem from artifacts introduced by random projections. Furthermore, existing works fail to comprehensively address issues related to random slicing, focusing primarily on suboptimality of random projections for information preservation [24], [25].

Contribution. In this article, we address this gap by systematically analyzing SMI across diverse settings, demonstrating that it frequently exhibits counterintuitive behavior and fails to accurately capture statistical dependence dynamics. Our key contributions are:

1. **Saturation and Sensitivity Analysis.** Through theoretical analysis and extensive benchmarking, we show that SMI saturates prematurely, even for low-dimensional synthetic problems, and fails to detect significant increases in statistical dependence.
2. **Redundancy Bias.** We refute the prevailing assumption that SMI favors linearly extractable information by constructing an explicit example where introducing such structure increases MI and even linear correlation, but decreases SMI. In fact, we show that SMI prioritizes information *redundancy* over information content. We argue that this bias can lead to catastrophic failures in some applications, e.g. collapses in representation learning.
3. **Curse of Dimensionality.** We revisit the dynamics of SMI for increasing dimensionality and argue that SMI is, in fact, cursed, with the curse of dimensionality manifesting itself not through sample complexity, but via asymptotic decay to zero in high-dimensional regimes due to diminishing redundancy.
4. **Reestablishing the Trade-off.** Finally, we discuss to which extent the aforementioned problems can be solved by using non-uniform/non-random slicing strategies, and how they affect the trade-off between scalability and utility of different measures of statistical dependence.

Our paper is structured as follows. In Section 2, we provide the mathematical background that is necessary for our analysis. In Section 3, we discuss previous findings which are related to the research topic of this work. Section 4 consists of our main theoretical results, with the complete proofs being provided in Section B. In Section 5, we employ synthetic benchmarks to show the disconnection between dynamics of MI and SMI. Section 6 illustrates that tasks related to SMI maximization may yield degenerate solutions, contrary to MI maximization. Finally, we discuss our results in Section 7.

77 2 Preliminaries

78 **Elements of Information Theory.** Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space with sample space Ω , σ
 79 -algebra \mathcal{F} , and probability measure \mathbb{P} defined on \mathcal{F} . Consider random vectors $X : \Omega \rightarrow \mathbb{R}^{d_x}$ and
 80 $Y : \Omega \rightarrow \mathbb{R}^{d_y}$ with joint distribution $\mathbb{P}_{X,Y}$ and marginals \mathbb{P}_X and \mathbb{P}_Y , respectively. Wherever it is
 81 needed, we assume the relevant Radon-Nikodym derivatives exist. For any probability measure $\mathbb{Q} \ll$
 82 \mathbb{P} , the Kullback-Leibler (KL) divergence is $D_{\text{KL}}(\mathbb{Q} \parallel \mathbb{P}) = \mathbb{E}_{\mathbb{Q}} \left[\log \frac{d\mathbb{Q}}{d\mathbb{P}} \right]$, which is non-negative and
 83 vanishes if and only if (iff) $\mathbb{P} = \mathbb{Q}$. The mutual information (MI) between X and Y quantifies the
 84 divergence between the joint distribution and the product of marginals:

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

85 When \mathbb{P}_X admits a probability density function (PDF) $p(X)$ with respect to (w.r.t.) the Lebesgue
 86 measure, the differential entropy is defined as $h(X) = -\mathbb{E}[\log p(X)]$, where $\log(\cdot)$ denotes the
 87 natural logarithm. Likewise, the joint entropy $h(X, Y)$ is defined via the joint density $p(X, Y)$,
 88 and conditional entropy is $h(X|Y) = -\mathbb{E}[\log p(X|Y)] = -\mathbb{E}_Y [\mathbb{E}_{X|Y} [\log p(X|Y)]]$. Under the
 89 existence of PDFs, MI satisfies the identities

$$I(X; Y) = h(X) - h(X|Y) = h(Y) - h(Y|X) = h(X) + h(Y) - h(X, Y). \quad (2)$$

90 In this work, we denote by μ_M the normalized Haar (uniform) probability measure on a compact
 91 manifold M , i.e., the unique bi-invariant measure satisfying $\mu_M(M) = 1$. Hence, to sample uniformly
 92 from specific spaces we write $W \sim \mu_{O(d)}$, $\theta \sim \mu_{\mathbb{S}^{d-1}}$, $A \sim \mu_{\text{St}(k,d)}$, indicating draws from the Haar
 93 measures on orthogonal group $O(d) = \{Q \in \mathbb{R}^{d \times d} : Q^T Q = Q Q^T = I\}$, the unit sphere $\mathbb{S}^{d-1} =$
 94 $\{X \in \mathbb{R}^d : \|X\|_2 = 1\}$, and the Stiefel manifold $\text{St}(k, d) = \{Q \in \mathbb{R}^{d \times k} : Q^T Q = I\}$, respectively.

95 **Sliced Mutual Information.** To mitigate the curse of dimensionality, one may average MI over
 96 all k -dimensional projections. The k -sliced mutual information (k -SMI) [23] between X and Y is
 97 defined as

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

98 which can be efficiently estimated. Setting $k = 1$ recovers the standard sliced mutual information (1).

99 3 Background

100 Merits of SMI are straightforward and have been investigated thoroughly in [22], [23]. We remind
 101 the reader of the two most important of them:

- 102 1. **Scalability** (i.e., fast convergence in high dimensions), enabled by low-dimensional projections.
- 103 2. **Nullification Property** (i.e., $\text{Sl}_k(X; Y) = 0$ iff X and Y are independent), which stems from the
 104 projections being random and independent.

105 In contrast, demerits of SMI are not very obvious and not well-covered in the literature. In this
 106 section, we recapitulate and analyze previous works which address the shortcomings of SMI. To
 107 facilitate the analysis, we divide them into three main categories.

108 **Suboptimality of random slicing.** In [24] and [25], it is argued that a uniform slicing strategy can
 109 produce suboptimal projections, impairing SMI's ability to capture dependencies in the presence of
 110 noisy or non-informative components. To address this issue, [24] proposed max-sliced MI (mSMI),
 111 which selects non-random projectors that maximize the MI between projected representations. This
 112 approach is also claimed to improve interpretability and convergence rates.

113 However, deterministic slicing may overlook dependencies captured by non-optimal components.
 114 To mitigate this, [25] extends the max-sliced approach by optimizing SMI over probability distrib-
 115 utions of projectors, with regularization to maintain slice diversity. While the authors emphasize
 116 that optimization should occur over *joint* distributions, their motivation primarily addresses the issue
 117 of non-optimal *marginal* distributions of θ and ϕ — specifically, the presence of non-informative

components in X and Y . We contend that this represents only a partial understanding of the problem, as many SMI artifacts arise from other factors. Needless to say that optimization over probability distributions is also a heavy burden, which does not align with the slicing philosophy.

Data Processing Inequality violation. A fundamental property of MI is that it cannot be increased by deterministic processing or, more generally, by Markov kernels. Furthermore, MI is preserved under invertible transformations. This is formalized by the *data processing inequality* (DPI).

Theorem 3.1. (Theorem 3.7 in [1]) For a Markov chain $X \rightarrow Y \rightarrow Z$, $I(X; Y) \geq I(X; Z)$. Additionally, if $Z = f(Y)$ where f is measurably invertible, then equality holds.

In contrast to MI, SMI violates the DPI (see Section 3.2 in [22] for an example). While the intuition behind DPI is clear (raw data already contains full information, and processing can only destroy it), the implications of DPI violation are less straightforward.

Existing works suggest that SMI’s violation of DPI can reflect a preference for linearly extractable features, framing this as a useful property that aligns with the informal understanding of “practically available” (i.e., easily accessible) information [22], [26], [30]. However, this interpretation can be misleading if the factors behind SMI increases are misidentified. Our analysis reveals that this is indeed the case, as SMI exhibits more inherent biases than previously recognized.

Asymptotics in high-dimensional regime. Convergence analysis suggests that the sample complexity of SMI estimation is far less sensitive to data dimensionality compared to that of MI. In fact, it has been argued that the estimation error may even decrease with dimensionality in some cases (see Remark 4 in [23]). However, an analysis of SMI itself reveals that this behavior may result from the fact that SMI can decrease as dimensionality grows. Specifically, Theorem 3 in [23] provides an asymptotic expression (as $d \rightarrow \infty$) for SMI in the case of jointly normal X and Y , which decays hyperbolically with d under some circumstances.

To date, no explanation for this phenomenon has been provided in the literature. We therefore elaborate on this finding by deriving non-asymptotic expressions, along with experimental results for non-Gaussian data, which reveal further nuances behind the decay.

4 Theoretical analysis

We start our analysis with considering a simple example, which (a) admits closed-form expression for SMI and (b) is capable of illustrating severe problems of the quantity in question.

Lemma 4.1. Consider the following pair of jointly Gaussian d -dimensional random vectors:

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

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 hold:

$$\begin{aligned} \lim_{d \rightarrow \infty} I(X; Y) &= +\infty & \lim_{d \rightarrow \infty} \text{SI}(X; Y) &= 0 \\ \lim_{\rho^2 \rightarrow 1} I(X; Y) &= +\infty & \lim_{\rho^2 \rightarrow 1} \text{SI}(X; Y) &= \psi(d-1) - \psi\left(\frac{d-1}{2}\right) - \log 2 \leq \frac{3}{d-1}, \end{aligned}$$

with ψ being the *digamma function*.

Note that while MI correctly captures the growing statistical dependence as $d \rightarrow \infty$ (since additional components contribute shared information), SMI drops to zero, exposing a fundamental problem. This issue was briefly noted in [23], but only through providing an asymptotic expression without further discussion. We interpret this behavior as a distinct manifestation of the **curse of dimensionality**: as d grows, SMI uniformly decays to zero and becomes ineffective for statistical analysis.

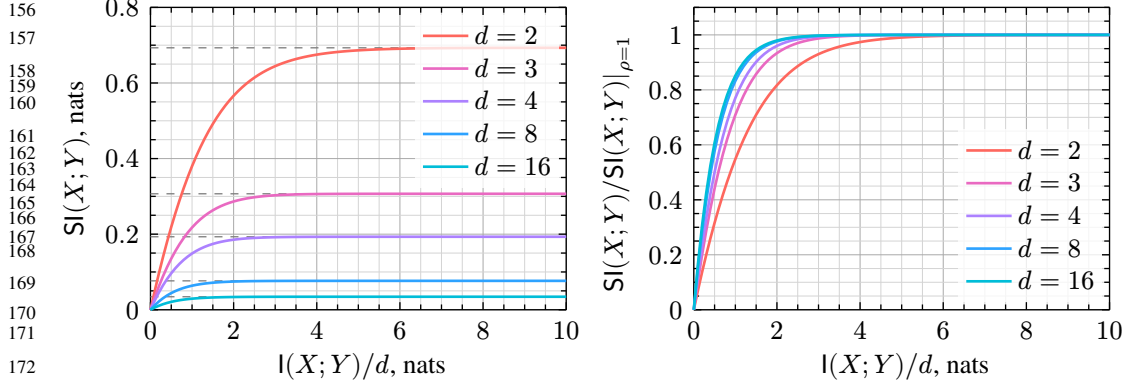


Figure 2: Saturation of $\text{SI}(X; Y)$ as function of $I(X; Y)/d$ for the example from Lemma 4.1, non-normalized (left) and normalized (right) versions. Note that the problem becomes more prominent in higher dimensions, both because of lower plateau and faster saturation.

The second pair of limits reveals another critical flaw of SMI. When $\rho^2 \rightarrow 1$, the X - Y relationship becomes deterministic — a property MI reflects successfully. In stark contrast, SMI remains bounded by a dimension-dependent factor that decays hyperbolically. Furthermore, plotting SMI against MI shows this bound is reached prematurely, demonstrating SMI's **rapid saturation** with increasing dependence (Figure 2). In this saturated regime, SMI becomes effectively insensitive to further growth in shared information. Moreover, this renders estimates of SMI for different dimensionalities fundamentally incomparable, as they are theoretically bounded by factors depending on d .

These phenomena can not be explained by suboptimality of individual projections. In fact, each individual projection is optimal, as $I(\theta^\top X; Y)$ does not depend on θ in this particular example. The proof of Lemma 4.1 suggests that the problem arises from the majority of *pairs* of projectors being suboptimal, yielding near-independent $\theta^\top X$ and $\phi^\top Y$ in the most outcomes, even for $d = 2$. Although similar analysis for k -SMI is extremely challenging, we argue that the problems in question prevail even when employing k -rank projectors.

Proposition 4.2. Under the setup of Lemma 4.1, k -SMI has the following integral representation

$$\text{SI}_k(X; Y) = -\frac{1}{2} \int_{[0,1]^k} \sum_{i=1}^k \log(1 - \rho^2 \lambda_i) p(\lambda) d\lambda,$$

where $p(\lambda) \propto \prod_{i < j} |\lambda_j - \lambda_i| \underbrace{\prod_{i=1}^k (1 - \lambda_i)^{(d-2k-1)/2}}_{(*)}$.

Remark. 4.3. As the dimension d grows, the term $(*)$ asymptotically concentrates the eigenvalues λ_i near zero, leading to the decay of SI_k to zero.

We argue that the limitations we uncovered can be attributed to a strong bias of SMI toward **information redundancy**. That is, SMI favors repetition of information across different axes, and suffers from the curse of dimensionality if X and Y have high entropy. The following proposition and remark present a simple example to clarify this bias.

Proposition 4.4. 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 $\text{SI}_k(AX; BY) = I(X; Y)$.

Corollary 4.5. 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 $\text{SI}_k(X; Y) = I(X; Y) = -\frac{1}{2} \log(1 - \rho^2)$.

Remark. 4.6. Applying $\mathbf{1} \cdot e_1^\top$ to the random vectors from Lemma 4.1 individually yields the example from Corollary 4.5. Therefore, this linear transform increases SMI despite decreasing MI.

203 4.1 Extension to optimal slicing

204 Although our work primarily focuses on conventional (average) sliced mutual information (SMI),
 205 as it is the most widely used variant, we also provide some intuition regarding the limitations of
 206 its “optimal” counterparts: max-sliced MI (mSMI) [24] and *optimal-sliced* MI (oSMI) [25]. Since
 207 mSMI is a special case of oSMI without regularization constraints, we restrict our discussion to
 208 mSMI, though our reasoning extends to oSMI as well. The k -mSMI is defined as:

$$\overline{\text{SI}}_k(X; Y) = \sup_{\substack{\Theta \in \text{St}(d_x, k) \\ \Phi \in \text{St}(d_y, k)}} I(\Theta^\top X; \Phi^\top Y) \quad (3)$$

209 To highlight the shortcomings of linear compression, we revisit a Gaussian example. The following
 210 proposition demonstrates that even in this simple setting, mSMI captures only a subset of depen-
 211 dencies and can exhibit opposite trends to MI. This occurs, for instance, when dependencies become
 212 more evenly distributed across components, which again returns us to the **redundancy bias**.

213 **Proposition 4.7.** (Proposition 2 in [24]) Let $(X, Y) \sim \mathcal{N}(\mu, \Sigma)$, with marginal covariances Σ_X ,
 214 Σ_Y and cross-covariance Σ_{XY} . Suppose the matrix $\Sigma_X^{-\frac{1}{2}} \Sigma_{XY} \Sigma_Y^{-\frac{1}{2}}$ exists, and let $\{\rho_i\}_{i=1}^d$ denote its
 215 singular values in descending order, where $d = \min(d_x, d_y)$. Then

$$I(X; Y) = -\frac{1}{2} \sum_{i=1}^d \log(1 - \rho_i^2), \quad \overline{\text{SI}}_k(X; Y) = -\frac{1}{2} \sum_{i=1}^k \log(1 - \rho_i^2).$$

216 5 Synthetic Experiments

217 To complement the theoretical analysis from the previous section and address complex, non-Gaussian
 218 cases, we conduct an extensive benchmarking of SMI using synthetic tests from [34], based on the
 219 works of [35], [36]. This benchmark suite is used to evaluate MI estimators. However, we do not
 220 assess whether SMI estimates converge to ground-truth MI values. SMI is a *distinct measure of*
 221 *statistical dependence*, and should not be viewed as an approximation of MI. Instead, our analysis
 222 focuses on the relationship between the two measures: since MI captures the true degree of statistical
 223 dependence, opposing trends in MI and SMI reveal problems with the latter quantity.

224 For the experiments, we use *correlated normal*, *correlated uniform*, *smoothed uniform* and *log-*
 225 *gamma-exponential* distributions, for which the ground-truth value of MI is available. To increase
 226 the dimensionality, we use independent components with equally distributed per-component MI.
 227 These setups will be referred to as “randomized” and “non-randomized” correspondingly. For each
 228 distribution, we vary both the data dimensionality (d) and the projection dimensionality ($k < d$).

229 To estimate MI between projections, we use the KSG estimator [35] with the number of neighbors
 230 fixed at 1. For each configuration, we conduct 10 independent runs with different random seeds
 231 to compute means and standard deviations. Our experiments use 10^4 samples for (X, Y) and 128
 232 samples for (Θ, Φ) .

233 To experimentally verify saturation, we plot SMI against MI normalized by dimensionality d in
 234 Figure 3. The plots clearly show that SMI reaches a plateau relatively early for all the featured
 235 distributions. The results for the normal distribution also align well with those from Lemma 4.1. We
 236 further confirm the saturation of k -SMI for $k \in \{2, 3\}$ experimentally in Section C. Finally, we plot
 237 the saturated values against d on a log-log scale, demonstrating that the $1/d$ trend from Lemma 4.1
 238 also holds for non-Gaussian distributions.

239 6 SMI for InfoMax-like tasks

240 Since mutual information is interpretable and captures non-linear dependencies, it is widely used as
 241 a training objective. Many applications involve maximizing MI (InfoMax) for feature selection [5],
 242 [6], [7] and self-supervised representation learning [8], [9], [10], [11], [12], [13]. However, due to
 243 the curse of dimensionality, alternative objectives have been proposed, with some works using sliced
 244 mutual information maximization for feature extraction [22] and disentanglement in InfoGAN [23].

In this section, we argue that SMI is not a suitable alternative to MI for InfoMax tasks. Since SMI exhibits a strong preference for redundancy, SMI maximization may lead to collapsed (high-redundancy) solutions. We demonstrate this through two experiments. Firstly, we revisit the Gaussian noisy channel to demonstrate that SMI favors linear mappings which decrease robustness to noise. Then, we consider a self-supervised representation learning task and show that using SMI immediately leads to collapsed representations.

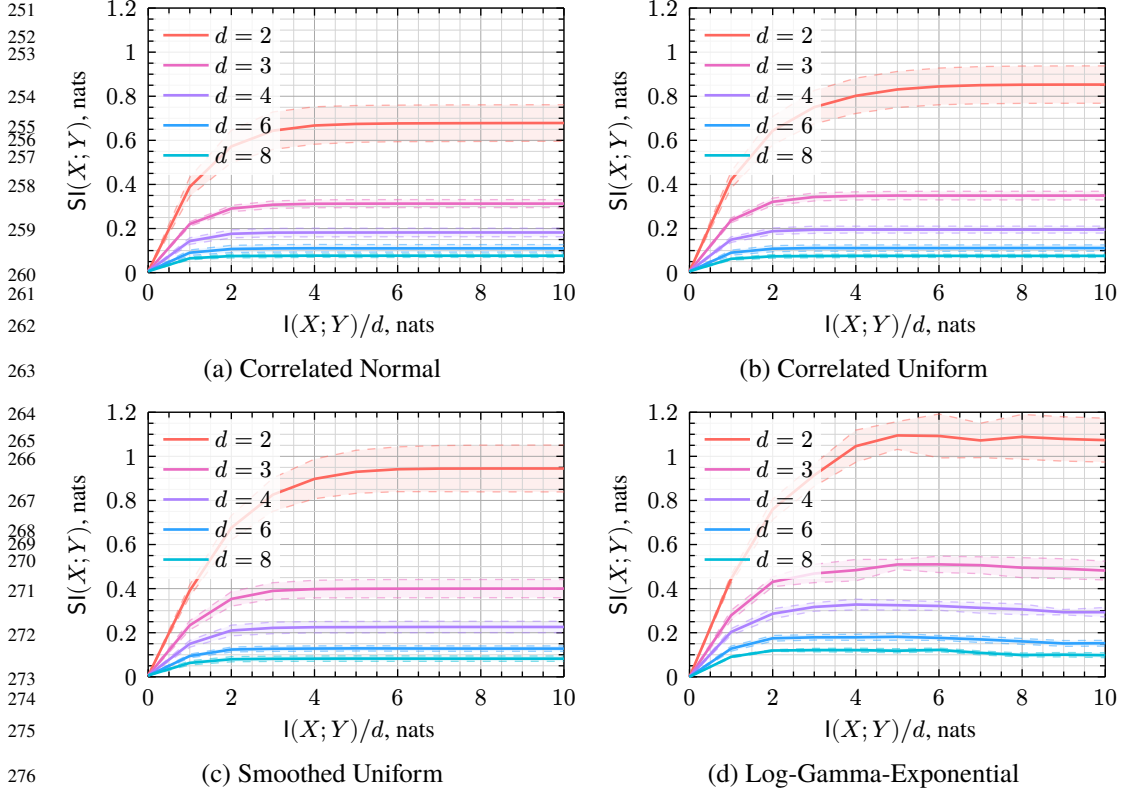


Figure 3: Results of synthetic experiments with different distributions. We report mean values and standard deviations computed across 10 runs, with 10^4 samples used for MI estimation and 128 for averaging across projections.

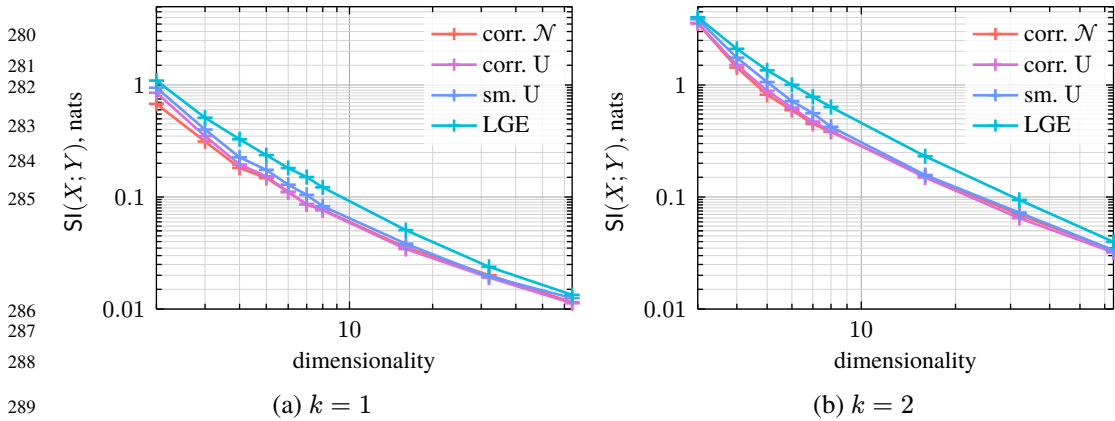


Figure 4: Decaying trends of k -SMI for *correlated normal* (corr. \mathcal{N}), *correlated uniform* (corr. \mathcal{U}), *smoothed uniform* (sm. \mathcal{U}) and *log-gamma-exponential* (LGE). We plot saturated values of k -SMI against data dimensionality d . Log scale is used to illustrate the $1/d$ trend predicted in Lemma 4.1.

293 6.1 Gaussian Channel

294 Let X be a zero-mean d -dimensional random vector, and let $Z \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ be an independent
 295 noise. Additive white noise Gaussian (AWGN) channel is defined as $X \rightarrow X + Z$. Maximization
 296 of $I(X; X + Z)$ w.r.t. the distribution of X is a classical information transmission problem, which
 297 arises in many fields under the Gaussian noise assumption. Given energy constraints, it admits an
 298 analytical solution [37]:

$$\sup_{\mathbb{E} X_i^2=1} I(X; X + Z) = \frac{d}{2} \log \left(1 + \frac{1}{\sigma^2} \right), \quad X_{\text{opt}} \sim \mathcal{N}(0, \mathbf{I}) \quad (4)$$

299 It is somewhat intuitive that unit covariance matrix allows for more information to be transmitted, as
 300 all the components of X are utilized to full extent. However, due to the redundancy bias, SMI prefers
 301 less robust distributions. To demonstrate this, we consider two linear normalization mappings which
 302 impose energy constraints on a vector X with zero mean and covariance Σ :

- 303 1. *Whitening*: $\Sigma^{-1/2} X$;
- 304 2. *Standardization*: $D^{-1/2} X$, where $D = \text{diag}(\Sigma)$.

305 We conduct numerical experiments for $\sigma = 0.1$, $X' \sim A \cdot U([-1; 1]^5)$ and $X'' \sim A \cdot \mathcal{N}(0, \mathbf{I}_5)$,
 306 where $A = 10^{-2} \cdot \mathbf{I} + \mathbf{1} \cdot \mathbf{1}^T$ is an ill-conditioned matrix. We employ the same estimators and
 307 hyperparameters as in Section 5. The results are presented in Table 1.

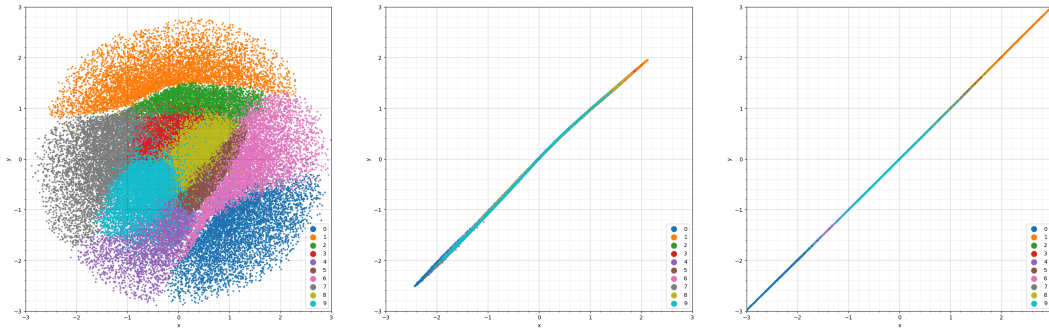
308 Table 1: Results for additive white Gaussian noise channel ($\sigma = 0.1$), mean and std for 10 runs.

	MI		SMI		2-SMI	
	$\Sigma^{-1/2}$	$D^{-1/2}$	$\Sigma^{-1/2}$	$D^{-1/2}$	$\Sigma^{-1/2}$	$D^{-1/2}$
310 X'	7.48 ± 0.01	3.04 ± 0.01	0.17 ± 0.02	1.82 ± 0.04	0.96 ± 0.04	2.46 ± 0.03
312 X''	7.49 ± 0.02	3.04 ± 0.01	0.14 ± 0.02	1.83 ± 0.04	0.82 ± 0.05	2.49 ± 0.05

313 6.2 Representation Learning

314 To further demonstrate SMI's sensitivity to information redundancy, we examine its performance
 315 in learning compressed representations through mutual information maximization (*Deep InfoMax*)
 316 [8]. This approach is known to be equivalent to many popular contrastive self-supervised learning
 317 methods [13].

318 In Deep InfoMax, an encoder network f is trained to maximize a lower bound on $I(X; f(X))$, where
 319 X represents input data and $f(X)$ its compressed representation. This method is theoretically sound,
 320 as maximizing MI ensures the most informative embeddings under the latent space dimensionality



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

322 Figure 5: Visualizations of embeddings from the representation learning experiments, with points
 323 colored by class. Note that mutual information maximization (left) produces clustered low-redundancy
 324 representations, while SMI maximization results in immediate (after 10 epochs) collapse.

constraint. For our study, we replace MI with SMI in this framework. This substitution is straightforward since both MI and SMI admit Donsker-Varadhan variational lower bounds [38]:

$$\begin{aligned} I(X; Y) &= \sup_{T: \Omega \rightarrow \mathbb{R}} \left[\mathbb{E}_{\mathbb{P}_{X,Y}} T(X, Y) - \log \left(\mathbb{E}_{\mathbb{P}_X \otimes \mathbb{P}_Y} e^{T(X,Y)} \right) \right], \\ \text{Sl}_k(X; Y) &= \sup_{T: \Omega \rightarrow \mathbb{R}} \mathbb{E}_{\Theta, \Phi} \left[\mathbb{E}_{\mathbb{P}_{X,Y}} T(\Theta^\top X, \Phi^\top Y, \Theta, \Phi) - \log \left(\mathbb{E}_{\mathbb{P}_X \otimes \mathbb{P}_Y} e^{T(\Theta^\top X, \Phi^\top Y, \Theta, \Phi)} \right) \right], \end{aligned} \quad (5)$$

where T is a critic function, which is also approximated in practice by a neural network. For detailed derivations of these bounds, we refer the reader to [39] (MI) and [22], [23] (SMI).

We strictly follow the experimental protocol from [13]. In particular, we use MNIST handwritten digits dataset [40], employ InfoNCE loss [41] to approximate (5), use convolutional network for f and fully-connected network for T . Latent space dimensionality is fixed at $d = 2$ for visualization purposes. Small Gaussian noise is added to the outlet of the encoder to combat representation collapse [13]. More details are provided in Section D. We focus on this simple setup because our objective is to show that SMI produces degenerate results even in elementary tasks, making more complex configurations unnecessary for this demonstration.

Results are presented in Figure 5. As our theory predicts, maximization of SMI immediately leads to collapsed representations, while conventional InfoMax yields embeddings with low or even zero redundancy (components are close to $\mathcal{N}(0, \mathbf{I})$). This behavior is consistent across different runs.

7 Discussion

Results. Sliced mutual information (SMI) has been proposed as a scalable alternative to Shannon’s mutual information. While SMI enables efficient computation in high-dimensional settings and satisfies the nullification property, our findings reveal critical deficiencies that undermine its reliability for feature extraction and related tasks.

We demonstrate that SMI saturates rapidly, failing to capture variations in statistical dependence. This makes it difficult to distinguish between intrinsic SMI fluctuations and genuine changes in dependence structure. Furthermore, we invalidate the common hypothesis that SMI favors linear features through a counterexample where even correlation coefficients reflect dependence more faithfully than SMI, which exhibits inverted behavior.

In high-dimensional spaces, SMI decays with increasing dimensionality, contrary to MI’s monotonic behavior. This is established analytically for Gaussian cases and validated empirically across diverse synthetic experiments. Consequently, SMI variations may reflect redundancy, dependence changes, or high-dimensional artifacts without a principled way to disentangle these factors.

Impact. Thanks to fast convergence rates and the absence of additional optimization problems, SMI has been widely applied across various fields of statistics and machine learning. Given our findings, it is therefore crucial to recognize how the inherent biases of SMI affect practical applications.

The works [22] and [23] propose using SMI in a Deep InfoMax setting. However, we demonstrate that maximizing SMI can lead to collapsed solutions due to redundancy bias. Meanwhile, [26], [27], [28], [30] study deep neural networks by measuring SMI between intermediate layers. Yet, as our analysis reveals, changes in SMI do not always reflect true shifts in statistical dependence; they may instead result from differences in layer dimensionality, redundancy in intermediate representations, low sensitivity in saturated regimes, or other factors. Finally, [33] suggests using SMI for independence testing in differential privacy tasks. We contend that this approach poses critical issues, as SMI estimates can become statistically indistinguishable from zero in high-dimensional or low-redundancy settings.

Limitations. While we support our claims with both theoretical analysis and experimental evidence, we were able to derive analytical expressions for the Gaussian case only. Furthermore, our synthetic tests do not feature complex, highly non-linear distributions (such as structured image data used in [17]). Nevertheless, we demonstrate that our findings are more than sufficient to expose fundamental limitations of SMI, and to support all the claims we made.

References

- [1] Y. Polyanskiy and Y. Wu, *Information Theory: From Coding to Learning*. Cambridge University Press, 2024. [Online]. Available: <https://books.google.ru/books?id=CySo0AEACAAJ>
- [2] A. Asadi, E. Abbe, and S. Verdu, “Chaining Mutual Information and Tightening Generalization Bounds,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., Curran Associates, Inc., 2018, p. . [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2018/file/8d7628dd7a710c8638dbd22d4421ee46-Paper.pdf
- [3] J. Negrea, M. Haghifam, G. K. Dziugaite, A. Khisti, and D. M. Roy, “Information-Theoretic Generalization Bounds for SGLD via Data-Dependent Estimates,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., Curran Associates, Inc., 2019, p. . [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2019/file/05ae14d7ae387b93370d142d82220f1b-Paper.pdf
- [4] B. Duong and T. Nguyen, “Conditional Independence Testing via Latent Representation Learning,” in *2022 IEEE International Conference on Data Mining (ICDM)*, Los Alamitos, CA, USA: IEEE Computer Society, Dec. 2022, pp. 121–130. doi: [10.1109/ICDM54844.2022.00022](https://doi.org/10.1109/ICDM54844.2022.00022).
- [5] S. Yang and J. Gu, “Feature selection based on mutual information and redundancy-synergy coefficient,” *J. Zhejiang Univ. Sci.*, vol. 5, no. 11, pp. 1382–1391, Nov. 2004.
- [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, doi: [10.1109/TPAMI.2002.1114861](https://doi.org/10.1109/TPAMI.2002.1114861).
- [7] M. A. Sulaiman and J. Labadin, “Feature selection based on mutual information,” in *2015 9th International Conference on IT in Asia (CITA)*, 2015, pp. 1–6. doi: [10.1109/CITA.2015.7349827](https://doi.org/10.1109/CITA.2015.7349827).
- [8] 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>
- [9] P. Bachman, R. D. Hjelm, and W. Buchwalter, “Learning Representations by Maximizing Mutual Information Across Views,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., Curran Associates, Inc., 2019, p. . [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2019/file/ddf354219aac374fd40b7e760ec5bb7-Paper.pdf
- [10] P. Veličković, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, “Deep Graph Infomax,” in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=rklz9iAcKQ>
- [11] M. Tschannen, J. Djolonga, P. K. Rubenstein, S. Gelly, and M. Lucic, “On Mutual Information Maximization for Representation Learning,” in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=rkxoh24FPH>
- [12] X. Yu, “Leveraging Superfluous Information in Contrastive Representation Learning.” [Online]. Available: <https://arxiv.org/abs/2408.10292>
- [13] I. Butakov, A. Semenenko, A. Tolmachev, A. Gladkov, M. Munkhoeva, and A. Frolov, “Efficient Distribution Matching of Representations via Noise-Injected Deep InfoMax,” in *The Thirteenth International Conference on Learning Representations*, 2025. [Online]. Available: <https://openreview.net/forum?id=mAmCdASmJ5>
- [14] N. Tishby and N. Zaslavsky, “Deep learning and the information bottleneck principle,” *2015 IEEE Information Theory Workshop (ITW)*, pp. 1–5, 2015.
- [15] R. Shwartz-Ziv and N. Tishby, “Opening the Black Box of Deep Neural Networks via Information.” 2017.
- [16] Z. Goldfeld *et al.*, “Estimating Information Flow in Deep Neural Networks,” in *Proceedings of the 36th International Conference on Machine Learning*, K. Chaudhuri and R. Salakhutdinov, Eds., in Proceedings of Machine Learning Research, vol. 97. PMLR, 2019, pp. 2299–2308. [Online]. Available: <https://proceedings.mlr.press/v97/goldfeld19a.html>
- [17] I. Butakov, A. Tolmachev, S. Malanchuk, A. Neopryatnaya, A. Frolov, and K. Andreev, “Information Bottleneck Analysis of Deep Neural Networks via Lossy Compression,” in *The Twelfth International*

- 423 *Conference on Learning Representations*, 2024. [Online]. Available: [https://openreview.net/forum?id=](https://openreview.net/forum?id=huGECz8dPp)
424 [huGECz8dPp](https://openreview.net/forum?id=huGECz8dPp)
- 425 [18] Z. Goldfeld, K. Greenewald, J. Niles-Weed, and Y. Polyanskiy, “Convergence of Smoothed Empirical
426 Measures With Applications to Entropy Estimation,” *IEEE Transactions on Information Theory*, vol. 66,
427 no. 7, pp. 4368–4391, 2020, doi: [10.1109/TIT.2020.2975480](https://doi.org/10.1109/TIT.2020.2975480).
- 428 [19] D. McAllester and K. Stratos, “Formal Limitations on the Measurement of Mutual Information,” in
429 *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, S.
430 Chiappa and R. Calandra, Eds., in Proceedings of Machine Learning Research, vol. 108. PMLR, 2020,
431 pp. 875–884. [Online]. Available: <https://proceedings.mlr.press/v108/mcallester20a.html>
- 432 [20] G. Gowri, X. Lun, A. M. Klein, and P. Yin, “Approximating mutual information of high-dimensional
433 variables using learned representations,” in *The Thirty-eighth Annual Conference on Neural Information*
434 *Processing Systems*, 2024. [Online]. Available: <https://openreview.net/forum?id=HN05DQxyLI>
- 435 [21] K. H. Greenewald, B. Kingsbury, and Y. Yu, “High-Dimensional Smoothed Entropy Estimation via
436 Dimensionality Reduction,” in *IEEE International Symposium on Information Theory, ISIT 2023, Taipei,*
437 *Taiwan, June 25-30, 2023*, IEEE, 2023, pp. 2613–2618. doi: [10.1109/ISIT54713.2023.10206641](https://doi.org/10.1109/ISIT54713.2023.10206641).
- 438 [22] Z. Goldfeld and K. Greenewald, “Sliced Mutual Information: A Scalable Measure of Statistical Depen-
439 dence,” in *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang,
440 and J. W. Vaughan, Eds., 2021. [Online]. Available: <https://openreview.net/forum?id=27qon5Ut4PSI>
- 441 [23] Z. Goldfeld, K. Greenewald, T. Nuradha, and G. Reeves, “ $\$k$ -Sliced Mutual Information: A Quantitative
442 Study of Scalability with Dimension,” in *Advances in Neural Information Processing Systems*, A. H. Oh,
443 A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available: [https://openreview.net/forum?id=](https://openreview.net/forum?id=L-ceBdl2DPb)
444 [L-ceBdl2DPb](https://openreview.net/forum?id=L-ceBdl2DPb)
- 445 [24] D. Tsur, Z. Goldfeld, and K. Greenewald, “Max-Sliced Mutual Information,” in *Thirty-seventh Conference*
446 *on Neural Information Processing Systems*, 2023. [Online]. Available: [https://openreview.net/forum?id=](https://openreview.net/forum?id=ce9B2x3zQa)
447 [ce9B2x3zQa](https://openreview.net/forum?id=ce9B2x3zQa)
- 448 [25] A. Fayad and M. Ibrahim, “On Slicing Optimality for Mutual Information,” in *Thirty-seventh Conference*
449 *on Neural Information Processing Systems*, 2023. [Online]. Available: [https://openreview.net/forum?id=](https://openreview.net/forum?id=JMuKfZx2xU)
450 [JMuKfZx2xU](https://openreview.net/forum?id=JMuKfZx2xU)
- 451 [26] S. Wongso, R. Ghosh, and M. Motani, “Understanding Deep Neural Networks Using Sliced Mutual
452 Information,” in *2022 IEEE International Symposium on Information Theory (ISIT)*, 2022, pp. 133–138.
453 doi: [10.1109/ISIT50566.2022.9834357](https://doi.org/10.1109/ISIT50566.2022.9834357).
- 454 [27] S. Wongso, R. Ghosh, and M. Motani, “Using Sliced Mutual Information to Study Memorization and
455 Generalization in Deep Neural Networks,” in *Proceedings of The 26th International Conference on Arti-*
456 *ficial Intelligence and Statistics*, F. Ruiz, J. Dy, and J.-W. van de Meent, Eds., in Proceedings of Machine
457 Learning Research, vol. 206. PMLR, 2023, pp. 11608–11629. [Online]. Available: [https://proceedings.](https://proceedings.mlr.press/v206/wongso23a.html)
458 [mlr.press/v206/wongso23a.html](https://proceedings.mlr.press/v206/wongso23a.html)
- 459 [28] S. Wongso, R. Ghosh, and M. Motani, “Pointwise Sliced Mutual Information for Neural Network
460 Explainability,” in *2023 IEEE International Symposium on Information Theory (ISIT)*, 2023, pp. 1776–
461 1781. doi: [10.1109/ISIT54713.2023.10207010](https://doi.org/10.1109/ISIT54713.2023.10207010).
- 462 [29] J. Dentan, D. Buscaldi, A. Shabou, and S. Vanier, “Predicting and analyzing memorization within fine-
463 tuned Large Language Models.” [Online]. Available: <https://arxiv.org/abs/2409.18858>
- 464 [30] S. Wongso, R. Ghosh, and M. Motani, “Sliced Information Plane for Analysis of Deep Neural Networks,”
465 Jan. 2025, doi: [10.36227/techrxiv.173833980.08812687/v1](https://doi.org/10.36227/techrxiv.173833980.08812687/v1).
- 466 [31] K. Nadjahi, K. Greenewald, R. B. Gabrielsson, and J. Solomon, “Slicing Mutual Information General-
467 ization Bounds for Neural Networks,” in *ICML 2023 Workshop Neural Compression: From Information*
468 *Theory to Applications*, 2023. [Online]. Available: <https://openreview.net/forum?id=cbLcwK3SZi>
- 469 [32] Z. Hu, S. Kang, Q. Zeng, K. Huang, and Y. Yang, “InfoNet: Neural Estimation of Mutual Information
470 without Test-Time Optimization,” in *Forty-first International Conference on Machine Learning*, 2024.
471 [Online]. Available: <https://openreview.net/forum?id=40hCy8n5XH>
- 472 [33] T. Nuradha and Z. Goldfeld, “Pufferfish Privacy: An Information-Theoretic Study,” *IEEE Trans. Inf.*
473 *Theor.*, vol. 69, no. 11, pp. 7336–7356, Nov. 2023, doi: [10.1109/TIT.2023.3296288](https://doi.org/10.1109/TIT.2023.3296288).
- 474 [34] I. Butakov *et al.*, “MUTINFO.” [Online]. Available: <https://github.com/VanessB/mutinfo>

- 475 [35] A. Kraskov, H. Stögbauer, and P. Grassberger, “Estimating mutual information,” *Phys. Rev. E*, vol. 69, no.
476 6, p. 66138, Jun. 2004, doi: [10.1103/PhysRevE.69.066138](https://doi.org/10.1103/PhysRevE.69.066138).
- 477 [36] F. Czyż Paweł and Grabowski, J. Vogt, N. Beerenwinkel, and A. Marx, “Beyond Normal: On the Evaluation
478 of Mutual Information Estimators,” in *Advances in Neural Information Processing Systems*, A. Oh, T.
479 Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., Curran Associates, Inc., 2023, pp.
480 16957–16990. [Online]. Available: [https://proceedings.neurips.cc/paper_files/paper/2023/file/36b80eae](https://proceedings.neurips.cc/paper_files/paper/2023/file/36b80eae70ff629d667f210e13497edf-Paper-Conference.pdf)
481 [70ff629d667f210e13497edf-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/36b80eae70ff629d667f210e13497edf-Paper-Conference.pdf)
- 482 [37] T. M. Cover and J. A. Thomas, *Elements of Information Theory (Wiley Series in Telecommunications and*
483 *Signal Processing)*. USA: Wiley-Interscience, 2006.
- 484 [38] M. D. Donsker and S. R. Varadhan, “Asymptotic evaluation of certain markov process expectations for
485 large time. IV,” *Communications on Pure and Applied Mathematics*, vol. 36, no. 2, pp. 183–212, Mar.
486 1983, doi: [10.1002/cpa.3160360204](https://doi.org/10.1002/cpa.3160360204).
- 487 [39] M. I. Belghazi *et al.*, “Mutual Information Neural Estimation,” in *Proceedings of the 35th International*
488 *Conference on Machine Learning*, J. Dy and A. Krause, Eds., in Proceedings of Machine Learning
489 Research, vol. 80. PMLR, 2018, pp. 531–540. [Online]. Available: [https://proceedings.mlr.press/v80/](https://proceedings.mlr.press/v80/belghazi18a.html)
490 [belghazi18a.html](https://proceedings.mlr.press/v80/belghazi18a.html)
- 491 [40] L. Deng, “The mnist database of handwritten digit images for machine learning research,” *IEEE Signal*
492 *Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- 493 [41] A. van den Oord, Y. Li, and O. Vinyals, “Representation Learning with Contrastive Predictive
494 Coding.” [Online]. Available: <https://arxiv.org/abs/1807.03748>
- 495 [42] A. Edelman and B. D. Sutton, “The beta-Jacobi matrix model, the CS decomposition, and generalized
496 singular value problems,” *Foundations of Computational Mathematics*, vol. 8, no. 2, pp. 259–285, 2008.
- 497 [43] A. McBride, “Special functions, by George E. Andrews, Richard Askey and Ranjan Roy. Pp. 664.£ 60.
498 1999. ISBN 0 521 62321 9 (Cambridge University Press.),” *The Mathematical Gazette*, vol. 83, no. 497,
499 pp. 355–357, 1999.
- 500 [44] N. Elezovic, C. Giordano, and J. Pecaric, “The best bounds in Gautschi’s inequality,” *Math. Inequal. Appl*,
501 vol. 3, no. 2, pp. 239–252, 2000.
- 502 [45] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization.” 2017.

503 A Supplementary theory

504 **Lemma A.1.** (Example 2.4 in [1]) $h(\mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \log((2\pi e)^d \det \Sigma)$.

505 **Corollary A.2.** For $(X, Y) \sim \mathcal{N}(\mu, \Sigma)$ with non-singular Σ

$$\begin{aligned} I(X; Y) &= \frac{1}{2} \log \det \Sigma_X + \frac{1}{2} \log \det \Sigma_Y - \frac{1}{2} \log \det \Sigma \\ &= -\frac{1}{2} \sum_{i=1}^d \log(1 - \rho_i^2), \end{aligned}$$

506 where Σ_X, Σ_Y are marginal covariances, Σ_{XY} is cross-covariance, $d = \min(d_x, d_y)$, and $\{\rho_i\}_{i=1}^d$
507 are singular values of $\Sigma_X^{-\frac{1}{2}} \Sigma_{XY} \Sigma_Y^{-\frac{1}{2}}$.

508 *Proof of Corollary A.2.* Combining Lemma A.1 and (2) yields the first result. Now note that

$$I(X; Y) = I\left(\Sigma_X^{-\frac{1}{2}} X; \Sigma_Y^{-\frac{1}{2}} Y\right) = I\left(U^\top \Sigma_X^{-\frac{1}{2}} X; V \Sigma_Y^{-\frac{1}{2}} Y\right),$$

509 where $U \text{diag}(\rho_i) V^\top$ is the SVD of $\Sigma_X^{-\frac{1}{2}} \Sigma_{XY} \Sigma_Y^{-\frac{1}{2}}$. However,

$$\left(U^\top \Sigma_X^{-\frac{1}{2}} X, V \Sigma_Y^{-\frac{1}{2}} Y\right) \sim \mathcal{N}\left(\mu', \begin{pmatrix} I & \text{diag}(\rho_i) \\ \text{diag}(\rho_i) & I \end{pmatrix}\right),$$

510 from which we arrive at the second expression. \square

511 **Lemma A.3.** Let $A \in \mathbb{R}^{n \times m}$ be full column-rank matrix and $\Theta \sim \mu_{\text{St}(n, k)}$. Then $\Theta^\top A$ is full-rank
512 with probability one.

513 *Proof of Lemma A.3.* Performing QR decomposition of A yields $\Theta^\top A = \Theta^\top Q R \stackrel{d}{=} \Theta^\top \begin{pmatrix} I_m \\ 0 \end{pmatrix} R$. Since
514 A is full-rank, R is invertible and $\text{rank } \Theta^\top A = \text{rank } \Theta^\top \begin{pmatrix} I_m \\ 0 \end{pmatrix}$. Therefore,

$$\mathbb{P}\{\Theta^\top A \text{ is full-rank}\} = 1 - \mathbb{P}\left\{\Theta^\top \begin{pmatrix} I_m \\ 0 \end{pmatrix} \text{ is not full-rank}\right\} = 1 - 0 = 1.$$

516 \square

517 **Lemma A.4.** (Theorem 1.5 in [42]) Let $W \sim \mu_{O(d)}$ and partition

$$W = \begin{pmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{pmatrix}.$$

518 with W_{11} of size k by k . Then the eigenvalues $\{\lambda_i\}_{i=1}^k$ of $W_{11} W_{11}^\top$ follow the Jacobi ensemble

$$p(\lambda) \propto \prod_{i < j} |\lambda_i - \lambda_j|^\beta \prod_{i=1}^k \lambda_i^{\frac{\beta}{2}(a+1)-1} (1 - \lambda_i)^{\frac{\beta}{2}(b+1)-1}$$

519 with parameters $a = 0, b = d - 2k$, and $\beta = 1$ (over \mathbb{R}).

520 *Proof of Lemma A.3.* Let $A_1 \in \mathbb{R}^{k \times d}$ and $A_2 \in \mathbb{R}^{(d-k) \times d}$ be independent matrices with i.i.d. entries
521 from $\mathcal{N}(0, 1)$. By stacking A_1 atop A_2 and then performing a block QR decomposition on the
522 resulting Gaussian matrix, the orthogonal invariance of the Gaussian law implies that the two
523 Q-blocks are independent of the upper-triangular factor R , with Q_1 and Q_2 uniformly distributed on
524 $O(k)$ and $\text{St}(k, d - k)$, respectively. Finally, computing the SVD of the block rows together with R
525 yields the generalized singular value decomposition (GSVD) of the pair (A_1, A_2) :

$$\begin{pmatrix} A_1 \\ A_2 \end{pmatrix} = \begin{pmatrix} Q_1 \\ Q_2 \end{pmatrix} R = \begin{pmatrix} U_1 & \\ & U_2 \end{pmatrix} \begin{pmatrix} \tilde{C} \\ 0 \\ -\tilde{S} \\ 0 \end{pmatrix} \tilde{V}^\top R,$$

where $U_1 \in O(k)$, $U_2 \in O(d-k)$, $\tilde{V} \in O(k)$, and $\tilde{C} = \text{diag}(c_i)$, $S = \text{diag}(s_i)$ with $c_i \geq 0$, $s_i \geq 0$, and $c_i^2 + s_i^2 = 1$ for all i . The diagonal entries of \tilde{C} are known as the generalized singular values of the pair (A_1, A_2) .

For a matrix $P = \text{diag}(p_1, \dots, p_k)$ with i.i.d. p_i sampled uniformly from $\{-1, 1\}$, we have $Q_1 P \stackrel{d}{=} W_{11}$. Let $W_{11} = UCV^\top$ be the SVD of W_{11} , then one has

$$U_1 \begin{pmatrix} \tilde{C} \\ 0 \end{pmatrix} \tilde{V}^\top P \stackrel{d}{=} UCV^\top.$$

Since U_1 , \tilde{V} , and U, V are uniformly distributed and independent of \tilde{C}, C , we have $\tilde{C} \stackrel{d}{=} C$ by the invariance of the Haar measure under orthogonal transformations. On the other hand, the generalized singular values \tilde{C} of a pair (A_1, A_2) follow the law of the Jacobi ensemble with parameters $a = 0$, $b = d - 2k$, and $\beta = 1$ (Proposition 1.2 in [42]). Therefore, the squared singular values of W_{11} follow the Jacobi ensemble with the same parameters. \square

Corollary A.5. The squared inner product $|\theta^\top \phi|^2$ between two independent random vectors $\theta, \phi \sim \mu_{\mathbb{S}^{d-1}}$ follows $\text{Beta}(\frac{1}{2}, \frac{d-1}{2})$. Moreover, the shifted inner product $(1 + \theta^\top \phi)/2$ is symmetrically distributed as $\text{Beta}(\frac{d-1}{2}, \frac{d-1}{2})$.

Proof of Corollary A.5. Setting Jacobi parameters $k = 1$, $a = 0$, $b = d - 2$ and $\beta = 1$, the density is proportional to $x^{-1/2}(1-x)^{(d-3)/2}$ on $[0, 1]$, which matches the $\text{Beta}(\frac{1}{2}, \frac{d-1}{2})$ distribution.

Next, observe that $\theta^\top \phi$ has a density proportional to $(1-t)^{\frac{d-3}{2}}$ for $t \in [-1, 1]$. Under the change of variables $\eta \sim \text{Beta}(\frac{d-1}{2}, \frac{d-1}{2})$.

\square

B Complete proofs

Proof of Lemma 4.1. One can acquire $I(X; Y) = -\frac{d}{2} \log(1 - \rho^2)$ from a general expression for MI of two jointly Gaussian random vectors (see Corollary A.2).

Recall that $(\theta^\top X, \phi^\top Y)$ is also Gaussian with cross-covariance $\rho \theta^\top \phi$. Therefore, by Corollary A.2 we have

$$SI(X; Y) = \mathbb{E}[I(\theta^\top X; \phi^\top Y) \mid \theta, \phi] = -\frac{1}{2} \mathbb{E}[\log(1 - \rho^2 |\theta^\top \phi|^2)].$$

From Corollary A.5, we note that $|\theta^\top \phi|^2 \sim \text{Beta}(\frac{1}{2}, \frac{d-1}{2})$, so

$$\begin{aligned} SI(X; Y) &= -\frac{1}{2B(\frac{1}{2}, \frac{d-1}{2})} \int_0^1 \log(1 - \rho^2 x) (1-x)^{\frac{d-3}{2}} x^{-\frac{1}{2}} dx \\ &= \frac{\rho^2}{2} \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d-1}{2})} \int_0^1 x^{\frac{1}{2}} (1-x)^{\frac{d-3}{2}} {}_2F_1(1, 1; 2; \rho^2 x) dx, \end{aligned} \tag{6}$$

where the last equality follows from the identity $\log(1-z) = -z {}_2F_1(1, 1; 2; z)$ with hypergeometric function ${}_2F_1$. Applying Euler's integral transform ([43], Eq. (2.2.3)) gives

$$\begin{aligned} SI(X; Y) &= \frac{\rho^2}{2d} \frac{\Gamma(\frac{d}{2} + 1)}{\Gamma(\frac{3}{2})\Gamma(\frac{d-2}{2})} \int_0^1 x^{\frac{3}{2}-1} (1-x)^{(\frac{d}{2}+1)-\frac{3}{2}-1} {}_2F_1(1, 1; 2; \rho^2 x) dx \\ &= \frac{\rho^2}{2d} {}_3F_2\left(1, 1, \frac{3}{2}; \frac{d}{2} + 1, 2; \rho^2\right). \end{aligned}$$

Here ${}_3F_2$ denotes the generalized hypergeometric function.

Finally, we calculate the limit of $SI(X; Y)$ as $\rho^2 \rightarrow 1$ using properties of beta-distribution. Denoting $\eta = (1 + \theta^\top \phi)/2 \sim \text{Beta}(\frac{d-1}{2}, \frac{d-1}{2})$ (see Corollary A.5), we get

$$\text{Sl}(X; Y) = -\log 2 - \mathbb{E} \log(1 - \eta) = -\log 2 - \mathbb{E} \log \eta = \psi(d-1) - \psi\left(\frac{d-1}{2}\right) - \log 2,$$

where ψ is the digamma function. Using the bounds on digamma function [44]

$$\log\left(x + \frac{1}{2}\right) - \frac{1}{x} \leq \psi(x) \leq \log\left(x + e^{\psi(1)}\right) - \frac{1}{x},$$

we derive an upper bound on this expression:

$$\psi(d-1) - \psi\left(\frac{d-1}{2}\right) - \log 2 \leq \frac{1}{d-1} + \log\left(1 + \frac{1 + e^{\psi(1)}}{d}\right)$$

To simplify the bound, one can note that $1 + e^{\psi(0)} < 2$, $\log(1+x) < x$ and $\frac{1}{d} < \frac{1}{d-1}$.

□

Proof of Proposition 4.2.

Let $Q_X, Q_Y \sim \mu_{\text{St}(k,d)}$. Then $[Q_X^\top X, Q_Y^\top Y] \sim \mathcal{N}(0, \Sigma)$, where Σ is a $2k \times 2k$ covariance matrix with the following block structure

$$\Sigma = \begin{pmatrix} I_k & \rho Q_X^\top Q_Y \\ \rho Q_Y^\top Q_X & I_k \end{pmatrix}.$$

Using the formula for the determinant of a block matrix Σ yields

$$\text{Sl}_k(X; Y) = -\frac{1}{2} \mathbb{E}[\log \det(\Sigma)] = -\frac{1}{2} \mathbb{E}\left[\log \det\left(I - \rho^2 (Q_X^\top Q_Y)(Q_X^\top Q_Y)^\top\right)\right].$$

By the invariance of the Haar measure under left and right multiplication, $Q_X^\top Q_Y \stackrel{d}{=} W_{11}$, where W_{11} is a k by k left upper block of the matrix $W \sim \mu_{O(d)}$. According to Lemma A.4, the eigenvalues of $W_{11} W_{11}^\top$ follow Jacobi ensemble with parameters $a = 0$, $b = d - 2k$ and $\beta = 1$:

$$p(\lambda) \propto \prod_{i < j} |\lambda_j - \lambda_i| \prod_{i=1}^k (1 - \lambda_i)^{\frac{d-2k-1}{2}}.$$

Thus, we get a general expression for k -SMI

$$\text{Sl}_k(X; Y) = -\frac{1}{2} \int_{[0,1]^k} \sum_{i=1}^k \log(1 - \rho^2 \lambda_i) p(\lambda) d\lambda.$$

□

Proof of Proposition 4.4. Using Lemma A.3 and $d_x, d_y < k$, we get that $\Theta^\top A$ and $\Phi^\top B$ are injective with probability one for independent Θ, Φ distributed uniformly on $\text{St}(d_x, k)$ and $\text{St}(d_y, k)$. Therefore, according to Theorem 3.1, $[l(\Theta^\top A X; \Phi^\top B Y) \mid \Theta, \Phi] = l(X; Y)$ almost sure. As a result, $\text{Sl}_k(A X; B Y) = l(\Theta^\top A X; \Phi^\top B Y \mid \Theta, \Phi) = l(X; Y)$. □

Proof of Proposition 4.7. Direct corollary of Corollary A.2. □

C Additional experiments

In this section, we conduct supplementary experiments to evaluate SMI under a broader range of setups. We begin by assessing k -SMI on the same set of benchmarks from Section 5. The results for $k = 1, 2, 3$ are presented in Figure 3, Figure 6, and Figure 7, respectively. Notably, saturation remains consistent even for $k = d - 1$ (i.e., when only one component is discarded).

Next, we examine a setup involving randomized distribution parameters, following the methodology of [34]. Among other adjustments, this includes randomizing per-component mutual information (e.g., assigning interactions unevenly in this experiment). In some cases (e.g., the log-gamma-

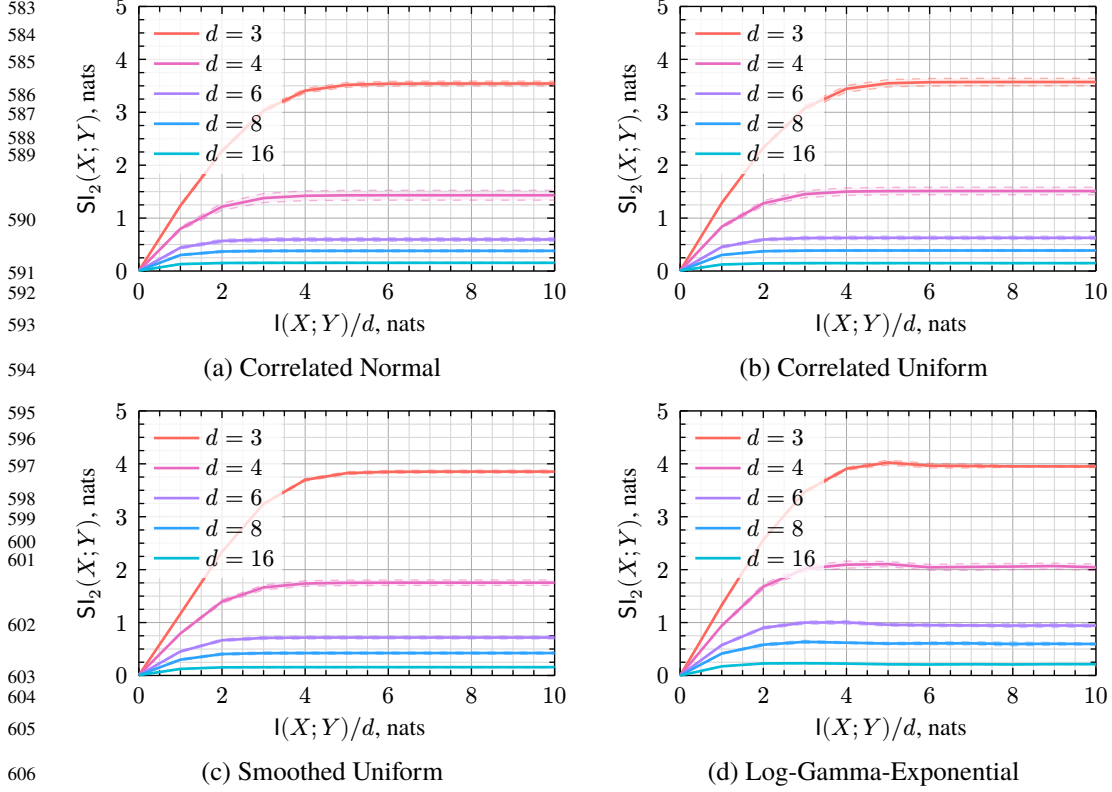


Figure 6: Results of synthetic experiments with different distributions for 2-SMI. We report mean values and standard deviations computed across 10 runs, with 10^4 samples used for MI estimation and 128 for averaging across projections.

exponential distribution), this increases linear redundancy, as component pairs with higher mutual information also exhibit higher variance in this particular scenario. Our results are displayed in Figure 8.

Due to numerical constraints, we do not track $l(X; Y)/d$, instead plotting the results against the total mutual information. While this makes saturation slightly less evident, the general trend of SMI decreasing with d remains observable. We also highlight the log-gamma-exponential distribution (Figure 8d), where SMI is less prone to saturation under parameter randomization due to the reasons mentioned earlier.

D Implementation details

D.1 Synthetic experiments

For the experiments from Section 5 and Section 6.1, we use implementation of Kraskov-Stoegbauer-Grassberger (KSG) [35] mutual information estimator and random slicing from [34]. The number of neighbors is set to $k_{\text{NN}} = 1$ for the KSG estimator. For each configuration, we conduct 10 independent runs with different random seeds to compute means and standard deviations. Our experiments use 10^4 samples for (X, Y) and 128 samples for (Θ, Φ) .

For the experiments from Section 5, we use independent components with equally distributed per-component MI. For the supplementary experiments from Figure 8, parameters of each distribution (e.g., covariance matrices) are randomized via the algorithm implemented in [34]. This includes randomization of per-component MI (which is done using a uniform distribution over a $(d - 1)$ -dimensional simplex).

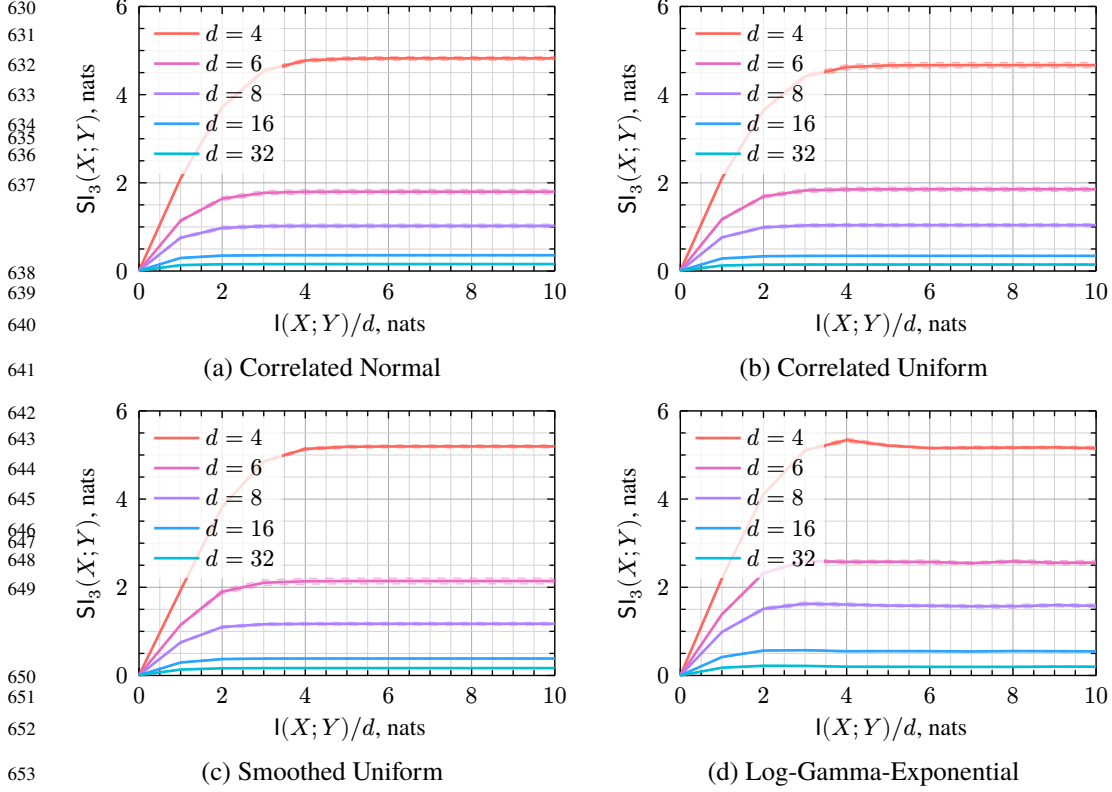


Figure 7: Results of synthetic experiments with different distributions for 3-SMI. We report mean values and standard deviations computed across 10 runs, with 10^4 samples used for MI estimation and 128 for averaging across projections.

For the experiments, we used AMD EPYC 7543 CPU, one core per distribution. Each experiment (fixed k , varying d) took no longer than 3 days to compute.

D.2 Representation learning experiments

For experiments on MNIST dataset, we use a simple ConvNet with three convolutional and two fully connected layers. A three-layer fully-connected perceptron serves as a critic network for the InfoNCE loss. We provide the details in Table 2. We use additive Gaussian noise with $\sigma = 0.2$ as an input augmentation. Training hyperparameters are as follows: batch size = 512, 2000 epochs, Adam optimizer [45] with learning rate 10^{-3} .

For the experiments, we used Nvidia A100 GPUs. Each experiment took no longer than 1 day to compute.

Table 2: The NN architectures used to conduct the tests on MNIST images in Section 6.2.

NN	Architecture
ConvNet,	$\times 1$: Conv2d(1, 32, ks=3), MaxPool2d(2), BatchNorm2d, LeakyReLU(0.01)
24×24	$\times 1$: Conv2d(32, 64, ks=3), MaxPool2d(2), BatchNorm2d, LeakyReLU(0.01)
images	$\times 1$: Conv2d(64, 128, ks=3), MaxPool2d(2), BatchNorm2d, LeakyReLU(0.01)
	$\times 1$: Dense(128, 128), LeakyReLU(0.01), Dense(128, dim)
Critic NN,	$\times 1$: Dense(dim + dim, 256), LeakyReLU(0.01)
pairs of vectors	$\times 1$: Dense(256, 256), LeakyReLU(0.01), Dense(256, 1)

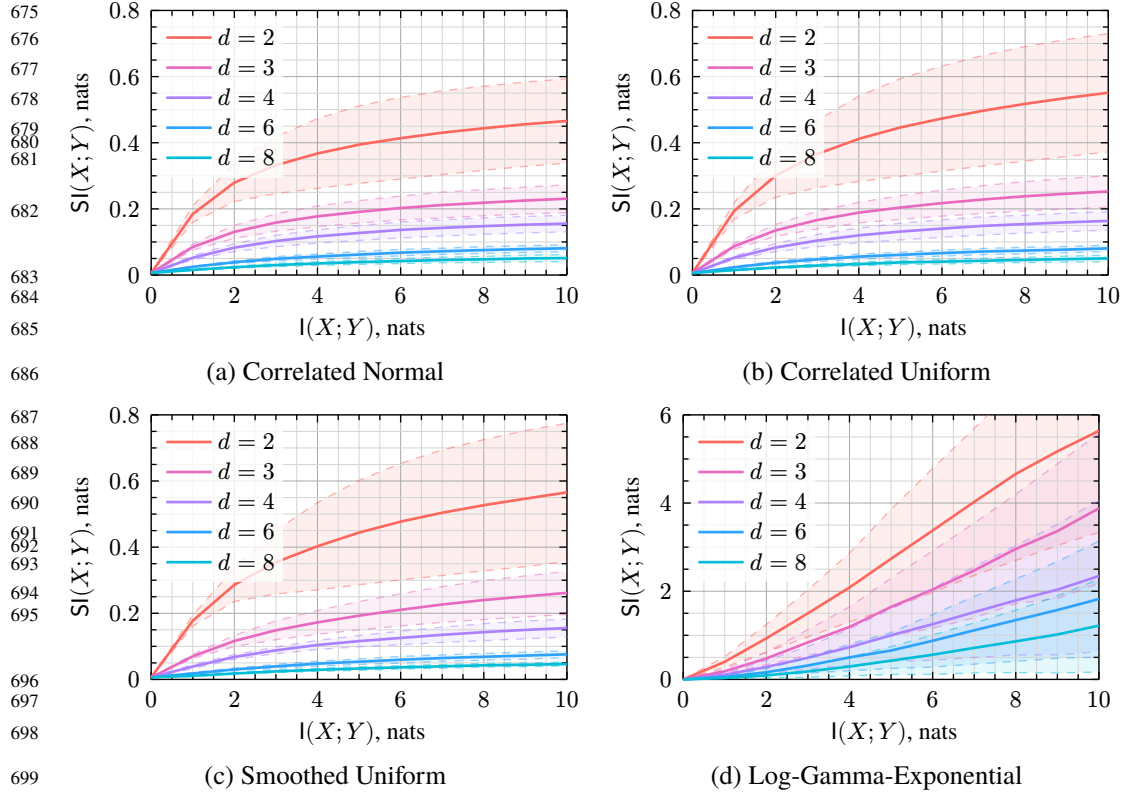


Figure 8: Results of synthetic experiments with different distributions. We report mean values and standard deviations computed across 10 runs, with 10^4 samples used for MI estimation and 128 for averaging across projections.

703 **NeurIPS Paper Checklist**

704 **1. Claims**

705 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's
706 contributions and scope?

707 Answer: [\[YES\]](#)

708 Justification: We state our claims clearly in the abstract and introduction. The claims are
709 supported by theoretical analysis and various experiments.

710 Guidelines:

- 711 • The answer NA means that the abstract and introduction do not include the claims made in
712 the paper.
- 713 • The abstract and/or introduction should clearly state the claims made, including the contri-
714 butions made in the paper and important assumptions and limitations. A No or NA answer
715 to this question will not be perceived well by the reviewers.
- 716 • The claims made should match theoretical and experimental results, and reflect how much
717 the results can be expected to generalize to other settings.
- 718 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are
719 not attained by the paper.

720 **2. Limitations**

721 Question: Does the paper discuss the limitations of the work performed by the authors?

722 Answer: [\[YES\]](#)

723 Justification: We discuss limitations in Section 7.

724 Guidelines:

- 725 • The answer NA means that the paper has no limitation while the answer No means that the
726 paper has limitations, but those are not discussed in the paper.
- 727 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 728 • The paper should point out any strong assumptions and how robust the results are to
729 violations of these assumptions (e.g., independence assumptions, noiseless settings, model
730 well-specification, asymptotic approximations only holding locally). The authors should
731 reflect on how these assumptions might be violated in practice and what the implications
732 would be.
- 733 • The authors should reflect on the scope of the claims made, e.g., if the approach was only
734 tested on a few datasets or with a few runs. In general, empirical results often depend on
735 implicit assumptions, which should be articulated.
- 736 • The authors should reflect on the factors that influence the performance of the approach. For
737 example, a facial recognition algorithm may perform poorly when image resolution is low
738 or images are taken in low lighting. Or a speech-to-text system might not be used reliably
739 to provide closed captions for online lectures because it fails to handle technical jargon.
- 740 • The authors should discuss the computational efficiency of the proposed algorithms and
741 how they scale with dataset size.
- 742 • If applicable, the authors should discuss possible limitations of their approach to address
743 problems of privacy and fairness.
- 744 • While the authors might fear that complete honesty about limitations might be used by
745 reviewers as grounds for rejection, a worse outcome might be that reviewers discover limi-
746 tations that aren't acknowledged in the paper. The authors should use their best judgment
747 and recognize that individual actions in favor of transparency play an important role in
748 developing norms that preserve the integrity of the community. Reviewers will be specifi-
749 cally instructed to not penalize honesty concerning limitations.

750 **3. Theory Assumptions and Proofs**

751 Question: For each theoretical result, does the paper provide the full set of assumptions and a
752 complete (and correct) proof?

753 Answer: [\[YES\]](#)

754 Justification: We provide comprehensive statements for theorems and lemmas. We also provide
755 complete proofs in Section B.

756 Guidelines:

- 757 • The answer NA means that the paper does not include theoretical results.
- 758 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
759 referenced.
- 760 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 761 • The proofs can either appear in the main paper or the supplemental material, but if they
762 appear in the supplemental material, the authors are encouraged to provide a short proof
763 sketch to provide intuition.
- 764 • Inversely, any informal proof provided in the core of the paper should be complemented by
765 formal proofs provided in appendix or supplemental material.
- 766 • Theorems and Lemmas that the proof relies upon should be properly referenced.

767 **4. Experimental Result Reproducibility**

768 Question: Does the paper fully disclose all the information needed to reproduce the main exper-
769 imental results of the paper to the extent that it affects the main claims and/or conclusions of the
770 paper (regardless of whether the code and data are provided or not)?

771 Answer: [\[YES\]](#)

772 Justification: We provide complete setup descriptions for the experiments in corresponding
773 sections.

774 Guidelines:

- 775 • The answer NA means that the paper does not include experiments.
- 776 • If the paper includes experiments, a No answer to this question will not be perceived well
777 by the reviewers: Making the paper reproducible is important, regardless of whether the
778 code and data are provided or not.
- 779 • If the contribution is a dataset and/or model, the authors should describe the steps taken to
780 make their results reproducible or verifiable.
- 781 • Depending on the contribution, reproducibility can be accomplished in various ways. For
782 example, if the contribution is a novel architecture, describing the architecture fully might
783 suffice, or if the contribution is a specific model and empirical evaluation, it may be
784 necessary to either make it possible for others to replicate the model with the same dataset,
785 or provide access to the model. In general, releasing code and data is often one good way
786 to accomplish this, but reproducibility can also be provided via detailed instructions for
787 how to replicate the results, access to a hosted model (e.g., in the case of a large language
788 model), releasing of a model checkpoint, or other means that are appropriate to the research
789 performed.
- 790 • While NeurIPS does not require releasing code, the conference does require all submissions
791 to provide some reasonable avenue for reproducibility, which may depend on the nature of
792 the contribution. For example
 - 793 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
794 reproduce that algorithm.
 - 795 (b) If the contribution is primarily a new model architecture, the paper should describe the
796 architecture clearly and fully.

- 797 (c) If the contribution is a new model (e.g., a large language model), then there should
798 either be a way to access this model for reproducing the results or a way to reproduce
799 the model (e.g., with an open-source dataset or instructions for how to construct the
800 dataset).
- 801 (d) We recognize that reproducibility may be tricky in some cases, in which case authors
802 are welcome to describe the particular way they provide for reproducibility. In the case
803 of closed-source models, it may be that access to the model is limited in some way (e.g.,
804 to registered users), but it should be possible for other researchers to have some path to
805 reproducing or verifying the results.

806 5. Open access to data and code

807 Question: Does the paper provide open access to the data and code, with sufficient instructions
808 to faithfully reproduce the main experimental results, as described in supplemental material?

809 Answer: [YES]

810 Justification: We use openly accessible data only for our experiments. Community-provided code
811 has been used for our experiments, and we reference its origin. Finally, we include additional
812 source code in the submission.

813 Guidelines:

- 814 • The answer NA means that paper does not include experiments requiring code.
- 815 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/public/
816 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 817 • While we encourage the release of code and data, we understand that this might not be
818 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including
819 code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 820 • The instructions should contain the exact command and environment needed to run to
821 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
822 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 823 • The authors should provide instructions on data access and preparation, including how to
824 access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 825 • The authors should provide scripts to reproduce all experimental results for the new pro-
826 posed method and baselines. If only a subset of experiments are reproducible, they should
827 state which ones are omitted from the script and why.
- 828 • At submission time, to preserve anonymity, the authors should release anonymized versions
829 (if applicable).
- 830 • Providing as much information as possible in supplemental material (appended to the paper)
831 is recommended, but including URLs to data and code is permitted.

832 6. Experimental Setting/Details

833 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpara-
834 meters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

835 Answer: [YES]

836 Justification: We provide all the necessary details in the corresponding sections, or in Appendix.

837 Guidelines:

- 838 • The answer NA means that the paper does not include experiments.
- 839 • The experimental setting should be presented in the core of the paper to a level of detail
840 that is necessary to appreciate the results and make sense of them.
- 841 • The full details can be provided either with the code, in appendix, or as supplemental
842 material.

843 7. Experiment Statistical Significance

844 Question: Does the paper report error bars suitably and correctly defined or other appropriate
845 information about the statistical significance of the experiments?

846 Answer: [YES]

847 Justification: We report mean and standard deviation.

848 Guidelines:

- 849 • The answer NA means that the paper does not include experiments.
- 850 • The authors should answer “Yes” if the results are accompanied by error bars, confidence
851 intervals, or statistical significance tests, at least for the experiments that support the main
852 claims of the paper.
- 853 • The factors of variability that the error bars are capturing should be clearly stated (for
854 example, train/test split, initialization, random drawing of some parameter, or overall run
855 with given experimental conditions).
- 856 • The method for calculating the error bars should be explained (closed form formula, call to
857 a library function, bootstrap, etc.)
- 858 • The assumptions made should be given (e.g., Normally distributed errors).
- 859 • It should be clear whether the error bar is the standard deviation or the standard error of
860 the mean.
- 861 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably
862 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality
863 of errors is not verified.
- 864 • For asymmetric distributions, the authors should be careful not to show in tables or figures
865 symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 866 • If error bars are reported in tables or plots, The authors should explain in the text how they
867 were calculated and reference the corresponding figures or tables in the text.

868 8. Experiments Compute Resources

869 Question: For each experiment, does the paper provide sufficient information on the computer
870 resources (type of compute workers, memory, time of execution) needed to reproduce the
871 experiments?

872 Answer: [YES]

873 Justification: We describe our setup and computational load of the experiments.

874 Guidelines:

- 875 • The answer NA means that the paper does not include experiments.
- 876 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or
877 cloud provider, including relevant memory and storage.
- 878 • The paper should provide the amount of compute required for each of the individual exper-
879 imental runs as well as estimate the total compute.
- 880 • The paper should disclose whether the full research project required more compute than the
881 experiments reported in the paper (e.g., preliminary or failed experiments that didn’t make
882 it into the paper).

883 9. Code Of Ethics

884 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS
885 Code of Ethics <https://neurips.cc/public/EthicsGuidelines>

886 Answer: [YES]

887 Justification:

888 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Although the paper may question methodology of some other works, there are no broader impacts of the research conducted.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.

- 937 • We recognize that providing effective safeguards is challenging, and many papers do not
938 require this, but we encourage authors to take this into account and make a best faith effort.

939 **12. Licenses for existing assets**

940 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the
941 paper, properly credited and are the license and terms of use explicitly mentioned and properly
942 respected?

943 Answer: [YES]

944 Justification: For all the assets, we use citations that the original authors provided.

945 Guidelines:

- 946 • The answer NA means that the paper does not use existing assets.
- 947 • The authors should cite the original paper that produced the code package or dataset.
- 948 • The authors should state which version of the asset is used and, if possible, include a URL.
- 949 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 950 • For scraped data from a particular source (e.g., website), the copyright and terms of service
951 of that source should be provided.
- 952 • If assets are released, the license, copyright information, and terms of use in the package
953 should be provided. For popular datasets, <https://paperswithcode.com/datasets> has
954 curated licenses for some datasets. Their licensing guide can help determine the license of
955 a dataset.
- 956 • For existing datasets that are re-packaged, both the original license and the license of the
957 derived asset (if it has changed) should be provided.
- 958 • If this information is not available online, the authors are encouraged to reach out to the
959 asset's creators.

960 **13. New Assets**

961 Question: Are new assets introduced in the paper well documented and is the documentation
962 provided alongside the assets?

963 Answer: [NA]

964 Justification:

965 Guidelines:

- 966 • The answer NA means that the paper does not release new assets.
- 967 • Researchers should communicate the details of the dataset/code/model as part of their
968 submissions via structured templates. This includes details about training, license, limita-
969 tions, etc.
- 970 • The paper should discuss whether and how consent was obtained from people whose asset
971 is used.
- 972 • At submission time, remember to anonymize your assets (if applicable). You can either
973 create an anonymized URL or include an anonymized zip file.

974 **14. Crowdsourcing and Research with Human Subjects**

975 Question: For crowdsourcing experiments and research with human subjects, does the paper
976 include the full text of instructions given to participants and screenshots, if applicable, as well
977 as details about compensation (if any)?

978 Answer: [NA]

979 Justification:

980 Guidelines:

981 • The answer NA means that the paper does not involve crowdsourcing nor research with
982 human subjects.

983 • Including this information in the supplemental material is fine, but if the main contribution
984 of the paper involves human subjects, then as much detail as possible should be included in
985 the main paper.

986 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or
987 other labor should be paid at least the minimum wage in the country of the data collector.

988 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
989 **Subjects**

990 Question: Does the paper describe potential risks incurred by study participants, whether such
991 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals
992 (or an equivalent approval/review based on the requirements of your country or institution) were
993 obtained?

994 Answer: [NA]

995 Justification:

996 Guidelines:

997 • The answer NA means that the paper does not involve crowdsourcing nor research with
998 human subjects.

999 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1000 may be required for any human subjects research. If you obtained IRB approval, you should
1001 clearly state this in the paper.

1002 • We recognize that the procedures for this may vary significantly between institutions and
1003 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines
1004 for their institution.

1005 • For initial submissions, do not include any information that would break anonymity (if
1006 applicable), such as the institution conducting the review.