# The Information Bottleneck Problem and Its Applications in Machine Learning

Ziv Goldfeld and Yury Polyanskiy

**Abstract**

Inference capabilities of machine learning (ML) systems skyrocketed in recent years, now playing a pivotal role in various aspect of society. The goal in statistical learning is to use data to obtain simple algorithms for predicting a random variable $Y$ from a correlated observation $X$. Since the dimension of $X$ is typically huge, computationally feasible solutions should summarize it into a lower-dimensional feature vector $T$, from which $Y$ is predicted. The algorithm will successfully make the prediction if $T$ is a good proxy of $Y$, despite the said dimensionality-reduction. A myriad of ML algorithms (mostly employing deep learning (DL)) for finding such representations $T$ based on real-world data are now available. While these methods are often effective in practice, their success is hindered by the lack of a comprehensive theory to explain it. The information bottleneck (IB) theory recently emerged as a bold information-theoretic paradigm for analyzing DL systems. Adopting mutual information as the figure of merit, it suggests that the best representation $T$ should be maximally informative about $Y$ while minimizing the mutual information with $X$. In this tutorial we survey the information-theoretic origins of this abstract principle, and its recent impact on DL. For the latter, we cover implications of the IB problem on DL theory, as well as practical algorithms inspired by it. Our goal is to provide a unified and cohesive description. A clear view of current knowledge is particularly important for further leveraging IB and other information-theoretic ideas to study DL models.

## I. INTRODUCTION

The past decade cemented the status of machine learning (ML) as a method of choice for a variety of inference tasks. The general learning problem considers an unknown probability distribution $P_{X,Y}$ that generates a target variable $Y$ and an observable correlated variable $X$. Given a dataset from $P_{X,Y}$, the goal is to learn a representation $T(X)$ of $X$ that is useful for inferring $Y$.[1] While these ideas date back to the late 1960's [2], [3], it was not until the 21st century that ML, and specifically, deep learning (DL), began to revolutionize data science practice [4]. Fueled by increasing computational power and data availability, deep neural networks (DNNs) are often unmatched in their effectiveness for classification, feature extraction and generative modeling.

This practical success, however, is not coupled with a comprehensive theory to explain how and why deep models work so well on real-world data. In recent years, information theory emerged as a popular lens through which to study fundamental properties of DNNs and their learning dynamics [5]–[15]. In particular, the information bottleneck

Z. Goldfeld is with the Electrical and Computer Engineering Department, Cornell University, Ithaca, NY, 14850, US (e-mail: goldfeld@cornell.edu). Y. Polyanskiy is with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, 02139, US (e-mail: yp@mit.edu).

[1]The reader is referred to [1] for background on statistical learning theory.

(IB) theory [5], [6] received significant attention. In essence, the theory extends an information-theoretic framework introduced in 1999 [16] to account for modern DL systems [5], [6]. It offers a novel paradigm for analyzing DNNs in an attempt to shed light on their layered structure, generalization capabilities and learning dynamics.

The IB theory for DL includes claims about the computational benefit of deep architectures, that generalization is driven by the amount of information compression the network attains, and more [5], [6]. This perspective inspired many followup works [9], [11]–[15], [17], some corroborating (and building on) and others refuting different aspects of the theory. This tutorial surveys the information-theoretic origins of the IB problem and its recent applications to and influences on DL. Our goal is to provide a comprehensive description that accounts for the multiple facets of this opinion-splitting topic.

### A. Origins in Information Theory

The IB problem was introduced in [16] as an information-theoretic framework for learning. It considers extracting information about a target signal $Y$ through a correlated observable $X$. The extracted information is quantified by a variable $T = T(X)$, which is (a possibly randomized) function of $X$, thus forming the Markov chain $Y \leftrightarrow X \leftrightarrow T$. The objective is to find a $T$ that minimizes the mutual information $I(X;T)$, while keeping $I(Y;T)$ above a certain threshold; the threshold determines how informative the representation $T$ is of $Y$. Specifically, the IB problem for a pair $(X, Y)$ (with a know joint probability law) is given by

$$\inf \ \ I(X;T)$$

$$\text{subject to:} \ \ I(Y;T) \geq \alpha,$$

where minimization is over all randomized mappings of $X$ to $T$. This formulation provides a natural *approximate* version of minimal sufficient statistic (MSS) [18] (cf. [19]). Notably, the same problem was introduced back in 1975 by Witsenhausen and Wyner [20] in a different context: as a tool for analyzing common information [21].

Beyond formulating the problem, [16] showed that, when alphabets are discrete, optimal $T$ can be found by iteratively solving a set of self-consistent equations. An algorithm (a generalization of Blahut-Arimoto [22], [23]) for solving the equations was also provided. Extensions of the IB problem and the aforementioned algorithm to the distributed setup were given in [24]. The only continuous-alphabet adaptation of this algorithm is known for jointly Gaussian $(X, Y)$ [25] (cf. [26] for the distributed Gaussian IB problem). In this case, the solution reduces to a canonical correlation analysis (CCA) projection with tunable rank. Solving the IB problem for complicated $P_{X,Y}$ distributions seems impossible (even numerically), even more so when only samples of $P_{X,Y}$ are available. The IB formulation, its solution, relations to MSSs, and the Gaussian case study are surveyed in Section II.

Optimized information measures under constraints often appear in information theory as solutions to operational coding problems. It is thus natural to ask whether there is an operational setup whose solution is the IB problem. Indeed, the classic framework of remote source coding (RSC) fits this description. RSC considers compressing a source $X^n$ so as to recover a correlated sequence $Y^n$ from that compression, subject to a distortion constraint [27], [28]. Choosing the distortion measure as the logarithmic loss and properly setting the compression threshold,

recovers the IB problem as the fundamental RSC rate. With that choice of loss function, RSC can roughly be seen as a multi-letter extension of the IB problem. The connection between the two problems is covered in Section III. The reader is referred to [29] for a recent survey covering additional connections between the IB problem and other information-theoretic setups.

### B. The Information Bottleneck Theory for Deep Learning

Although first applications of the IB problem to ML, e.g., for clustering [30], date two decades ago, it recently gained much traction in the context of DL. From the practical perspective, the IB principle was adopted as a design tool for DNN classifiers [8], [31] and generative models [32], with all three works published concurrently. Specifically, [31] optimized a DNN to solve the IB Lagrangian via gradient-based methods. Termed deep *variational IB* (VIB), the systems learns stochastic representation rules that were shown to generalize better and enjoy robustness to adversarial examples. The same objective was studied in [8], who argued it promotes minimality, sufficiency and disentanglement of representations. The disentanglement property was also employed in [32] for generative modeling purposes, where the $\beta$-variational autoencoder was developed. These applications of the IB framework are covered in Section IV-A.

The IB problem also had impact on DL theory. The main characterizing property of DNNs, as compared to general learning algorithms, is their layered structure. The IB theory suggests that deeper layers correspond to smaller $I(X;T)$ values, thus providing increasingly compressed sufficient statistics. In a classification task, the feature $X$ might contain information that is redundant for determining the label $Y$. It is therefore desirable to find representations $T$ of $X$ that shed redundancies, while retaining informativeness about $Y$. The argument of [5], [6] was that the IB problem precisely quantifies the fundamental tradeoff between informativeness (of $T$ for $Y$) and compression (of $X$ into $T$).

the IB theory for DL was first presented in [5], followed by the supporting empirical study [6]. The latter relied on a certain synthetic binary classification task as a running example. Beyond testing claims made in [5], the authors of [6] further evolved the IB theory. A main new argument was that classifiers trained with cross-entropy loss and stochastic gradient descent (SGD) inherently (try to) solve the IB optimization problem. As such, [6] posed the *information plane* (IP), i.e., the trajectory in $\mathbb{R}^2$ of the mutual information pair $\big(I(X;T), I(Y;T)\big)$ across training epochs, as a potent lens through which to analyze DNNs.

Based on this IP analysis, the IB theory proposed certain predictions about DNN learning dynamics. Namely, [6] argued that there is an inherent phase transition during DNN training, starting from a quick 'fitting' phase that aims to increase $I(Y;T)$ and followed by a long 'compression' phase that shrinks $I(X;T)$. This observation was explained as the network shedding redundant information and learning compressed representations, as described above. This is striking since the DNN has no explicit mechanism that encourages compressed representations. The authors of [6] further used the IP perspective to reason about computational benefits of deep architectures, phase transitions in SGD dynamics, and the network's generalization error. The IB theory for DL is delineated in section IV-B.

## C. Recent Results and Advances

The bold IB perspective on DL [5], [6] inspired followup research aiming to test and understand the observations therein. In [9], the empirical findings from [6] were revisited. The goal was to examine their validity in additional settings, e.g., across different nonlinearities, replacing SGD with batch gradient-descent, etc. The main conclusion of [9] was that the findings from [6] do not hold in general. Specifically, [9] provided experiments showing that IP dynamics undergo a compression phase only when double-sided saturating nonlinearities (e.g., $\tanh$, as used in [6]) are employed. Retraining the same network but with $\mathrm{ReLU}$ activations, produces a similarly-performing classifier whose IP trajectory does *not exhibit compression*. This refuted the fundamental role of compression in learning deep models, as posed in [6]. The results of [9], along with additional counterexamples they provided to claims from [6], are covered in Section V-A.

Theoretical aspect of the IB theory for DL were also reexamined. It was noted in [13], [15] that the mutual information measures of interest $\big(I(X;T), I(Y;T)\big)$ are vacuous in deterministic DNNs (i.e., networks that define a deterministic mapping for each fixed set of parameters). This happens because deterministic DNNs with strictly monotone activations (e.g., $\tanh$ or $\mathrm{sigmoid}$) can encode the entire input dataset into arbitrarily fine variations of $T$. Consequently, no information about $X$ is lost when traversing the network's layers and $I(X;T)$ is either the *constant* $H(X)$ (discrete $X$) or *infinite* (continuous $X$); a vacuous quantity, independent of the network parameters, either way [13].[2] Similar degeneracies occur in DNNs with any bi-Lipschitz nonlinearities, as well as for $\mathrm{ReLU}$s [15]. This implies that the (estimated) IP trajectories presented in [6], [9] for deterministic DNNs, cannot be fluctuating (e.g., undergoing fitting/compression phases) due to changes in mutual information. Indeed, [13] showed that these fluctuations are an artifact of the quantization-based method employed in [6], [9] to approximate the true mutual information. We describe the technical caveats in applying information-theoretic reasoning to deterministic DL systems and the misestimation of mutual information in Sections IV-B3 and V-B, respectively.

To circumvent these issues, [13] proposed an auxiliary framework, termed *noisy DNNs*, over which $I(X;T)$ and $I(Y;T)$ are meaningful, parameter-dependent quantities. In this framework, additive Gaussian noise is injected to the output of each of the network's neurons. [13] showed that noisy DNNs approximate deterministic ones both in how they generalize (just as well and sometimes better), and in terms of the learned representations. From a theoretical standpoint, the noise renders $X \mapsto T$ a stochastic parametrized channel, which makes $I(X;T)$ and $I(Y;T)$ functions of the network's parameters. Once the degeneracy of the considered mutual information terms was alleviated, [13] focused on accurately measuring their values across training epochs.

Building on recent results on differential entropy estimation under Gaussian smoothing [17], [33], [34], a rate-optimal estimator of $I(X;T)$ over noisy DNNs was developed [13]. This enabled tracking $I(X;T)$ across training of noisy DNN classifiers and empirically show that it also undergoes compression (just like the quantization-based estimate of mutual information over deterministic networks). To understand the relation between compression and the geometry of latent representations, [13] described an analogy between $I(X;T)$ and data transmission over additive

---

[2]A similar degeneracy occurs for $I(Y;T)$, which, e.g., equals $I(Y;X)$, whenever $X$ is discrete and the activations are injective. Note that, except for synthetic models, it is customary to replace the unknown joint distribution $P_{X,Y}$ with its empirical version, thus making $X$ discrete.

white Gaussian noise (AWGN) channels. The analogy gave rise to an argument that compression of $I(X;T)$ is driven by the progressive clustering of internal representations of equilabeled inputs.

Armed with the connection between compression and clustering, [13] traced back to deterministic DNNs and showed that, while compression of $I(X;T)$ is impossible therein, these networks nonetheless cluster equilabeled samples. They further demonstrated that the quantization-based proxy of mutual information used in [6], [9] in fact measures clustering in the latent spaces. This identified clustering of representations as the fundamental phenomenon of interest, while elucidating some of the machinery DNNs employ for learning. Section V-C elaborates on noisy DNNs, mutual information estimation, and the relation to clustering.

The IB problem remains an active area of research in ML and beyond. Its rich and opinion-splitting history inspired a myriad of followup works aiming to further explore and understand it. This tutorial surveys a non-exhaustive subset of these works as described above[3], that contributed to different facets of IB research, both historically and more recently. We aim to provide a balanced description that clarifies the current status of the IB problem applied to DL theory and practice. Doing so would help further leveraging the IB and other information-theoretic concepts for the study of DL systems.

## II. INFORMATION BOTTLENECK FORMULATION AND GAUSSIAN SOLUTION

The IB framework was proposed in [16] as a principled approach for extracting 'relevant' or 'useful' information from an observed signal about a target one. Consider a pair of correlated random variables $(X, Y)$, where $X$ is the observable and $Y$ is the object on interest. The goal is to compress $X$ into a representation $T$ that preserves as much information about $Y$ as possible.

The formulation of this idea uses mutual information. For a random variable pair $(A, B) \sim P_{A,B}$ with values in $\mathcal{A} \times \mathcal{B}$, set

$$I(A; B) := \int_{\mathcal{A} \times \mathcal{B}} \log \left( \frac{\mathrm{d}P_{A,B}}{\mathrm{d}P_A \otimes P_B}(a, b) \right) \mathrm{d}P_{A,B}(a, b),$$

as the mutual information between $A$ and $B$, where $\frac{\mathrm{d}P_{A,B}}{\mathrm{d}P_A \otimes P_B}$ is the Radon-Nikodym derivative of $P_{A,B}$ with respect to (w.r.t.) $P_A \otimes P_B$. Mutual information is a fundamental measure of dependence between random variables with many desirable properties. For instance, it nullifies if and only if $A$ and $B$ are independent, and is invariant to bijective transformations. In fact, mutual information can be obtained axiomatically as a unique (up to a multiplicative constant) functional satisfying several natural 'informativeness' conditions [37].

### A. The Information Bottleneck Framework

The IB framework for extracting the relevant information an $\mathcal{X}$-valued random variable $X$ contains about a $\mathcal{Y}$-valued $Y$ is described next. For a set $\mathcal{T}$, let $P_{T|X}$ be a transition kernel from $\mathcal{X}$ to $\mathcal{T}$.[4] The kernel $P_{T|X}$ can be viewed as transforming $X \sim P_X$ (e.g., via quantization) into a representation of $T \sim P_T(\cdot) := \int P_{T|X}(\cdot|x) \, \mathrm{d}P_X(x)$

---

[3]Others relevant papers include [7], [11], [12], [14], [32], [35], [36].

[4]Given two measurable spaces $(\mathcal{A}, \mathfrak{A})$ and $(\mathcal{B}, \mathfrak{B})$, $P_{A|B}(\cdot|\cdot) : \mathfrak{A} \times \mathcal{B} \to \mathbb{R}$ is a transition kernel from the former to the latter if $P_{A|B}(\cdot|b)$ is a probability measure on $(\mathcal{A}, \mathfrak{A})$, for all $b \in \mathcal{B}$, and $P_{A|B}(a|\cdot)$ is $\mathfrak{B}$-measurable for all $a \in \mathfrak{A}$.
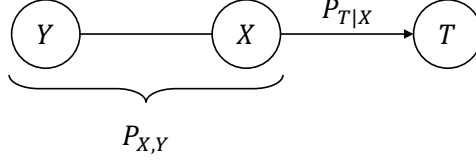
Fig. 1: Graphical representation of probabilistic relations in the IB framework. The triple $(X, Y, T)$ is jointly distributed according to $P_{X,Y} P_{T|X}$, thus forming a Markov chain $Y \leftrightarrow X \leftrightarrow T$. The goal is to find a compressed representation $T$ of $X$ (via the transition kernel $P_{T|X}$), i.e., minimize $I(X;T)$, while preserving at least $\alpha$ bits of information about $Y$, i.e., $I(Y;T) \geq \alpha$.

in the $\mathcal{T}$ space. The triple $Y \leftrightarrow X \leftrightarrow T$ forms a Markov chain in that order w.r.t. the joint probability measure $P_{X,Y,T} = P_{X,Y} P_{T|X}$. This joint measure specifies the mutual information terms $I(X;T)$ and $I(Y;T)$, that are interpreted as the amount of information $T$ conveys about $X$ and $Y$, respectively.

The IB framework concerns finding a $P_{T|X}$ that extracts information about $Y$, i.e., high $I(Y;T)$, while maximally compressing $X$, which is quantified as keeping $I(X;T)$ small. Since the Data Processing Inequality (DPI) implies $I(Y;T) \leq I(X;Y)$, the compressed representation $T$ cannot convey more information than the original signal. This gives rise to a tradeoff between compressed representations and preservation of meaningful information. The tradeoff is captured by a parameter $\alpha \in \mathbb{R}_{\geq 0}$ that specifies the lowest tolerable $I(Y;T)$ values. Accordingly, the IB problem is formulated through the constrained optimization

$$\inf_{P_{T|X}:\ I(Y;T) \geq \alpha} I(X;T), \tag{1}$$

where the underlying joint distribution is $P_{X,Y} P_{T|X}$. Thus, we pass the information that $X$ contains about $Y$ through a 'bottleneck' via the representation $T$ (see Fig. 1). An extension of the IB problem to the distributed setup was proposed in [24], [26]. In that problem, multiple sources $X_1, \ldots, X_K$ are compressed separately into representations $T_1, \ldots, T_K$, that, collectively, preserve as much information as possible about $Y$.

A slightly different form of the problem in (1), namely

$$\inf_{P_{T|X}:\ H(X|T) \geq \alpha} H(Y|T), \tag{2}$$

first appeared in a seminal paper of Witsenhausen and Wyner [20], who used it to simplify the proof of Gács and Körner's result on common information [21].

### B. Lagrange Dual Form and Information Bottleneck Curve

Note that the optimization problem (1) is not convex (in $P_{T|X}$), however it can be made into a convex rate-distortion problem (cf. (14)-(15)). In any case, (1) is commonly solved by introducing the Lagrange multiplier $\beta$ and considering the functional

$$\mathcal{L}_\beta(P_{T|X}) := I(X;T) - \beta I(T;Y). \tag{3}$$

With this definition, the IB problem can be recast as minimizing $\mathcal{L}_\beta(P_{T|X})$ over all possible $P_{T|X}$ kernels. Here $\beta$ controls the amount of compression in the representation $T$ of $X$: small $\beta$ implies more compression (sacrificing informativeness), while larger $\beta$ pushes towards finer representation granularity that favor informativeness. Varying $\beta \in [0, +\infty)$ regulates the tradeoff between informativeness and compression.

For any $\beta \in [0, +\infty)$, conditions for a stationary point of $\mathcal{L}_\beta(P_{T|X})$ can be expressed via the following self-consistent equations [16]. Namely, a stationary point $P_{T|X}^{(\beta)}$ must satisfy:

$$P_T^{(\beta)}(t) = \int_{\mathcal{X}} P_{T|X}^{(\beta)}(t|x)\, \mathrm{d}P_X(x) \tag{4a}$$

$$P_{Y|T}^{(\beta)}(y|t) = \frac{1}{P_T^{(\beta)}(t)} \int_{\mathcal{X}} P_{Y|X}(y|x) P_{T|X}^{(\beta)}(t|x)\, \mathrm{d}P_X(x) \tag{4b}$$

$$P_{T|X}^{(\beta)}(t|x) = \frac{P_T^{(\beta)}(t)}{Z_\beta(x)} e^{-\beta \mathsf{D}_{\mathsf{KL}}\left(P_{Y|X}(\cdot|x)\,\middle\|\,P_{Y|T}^{(\beta)}(\cdot|t)\right)}, \tag{4c}$$

where $Z_\beta(x)$ is the normalization constant (partition function). If $X$, $Y$ and $T$ take values in finite sets, and $P_{X,Y}$ is known, then alternating iterations of (4) locally converges to a solution, for any initial $P_{T|X}$ [16]. This is reminiscent of the Blahut-Arimoto algorithm for computing rate distortion functions or channel capacity [22], [23]. The discrete alternating iterations algorithms was later adapted to the Gaussian IB [25].[5] While the algorithm is infeasible for general continuous distributions, in the Gaussian case it translates into a parameter updating procedure.

The IB curve is obtained by solving $\inf_{P_{T|X}} \mathcal{L}_\beta(P_{T|X})$, for each $\beta \in [0, +\infty)$, and plotting the mutual information pair $\left(I_\beta(X;T), I_\beta(Y;T)\right)$ for an optimal $P_{T|X}^{(\beta)}$. In Section II-D we show the IB curve for jointly Gaussian $(X, Y)$ variables (Fig. 2). The two-dimensional plane in which the IB curve resides, was later coined as the *information plane* (IP) [5].

### C. Relation to Minimal Sufficient Statistics

A elementary concept that captures the notion of 'compressed relevance' is that of MSS [18], as defined next.

**Definition 1 (Minimal Sufficient Statistic)** *Let $(X, Y) \sim P_{X,Y}$. We call $T := t(X)$, where $t$ is a deterministic function, a sufficient statistic of $X$ for $Y$ if $Y \leftrightarrow T \leftrightarrow X$ forms a Markov chain. A sufficient statistic $T$ is minimal if for any other sufficient statistic $S$, there exists a function $f$, such that $T = f(S)$, $P_X$-almost surely (a.s.).*

The need to define minimally arises because $T = X$ is trivially sufficient, and one would like to avoid such degenerate choices. However, sufficient statistics themselves are rather restrictive, in the sense that their dimension always depends on the sample size, except when the data comes from an exponential family (cf., the PitmanKoopmanDarmois Theorem [38]). It is therefore useful to consider relaxations of the MSS framework.

Such a relaxation is given by the IB framework, which is evident by relating it to MSSs as follows. Allowing $T$ to be a stochastic (rather than deterministic) function of $X$, i.e., defined through a transition kernel $P_{T|X}$, we have

---

[5]see [24], [26] for extensions of both the discrete and the Gaussian cases to the distributed setup.

that $T$ is sufficient for $Y$ if and only if $I(X;Y) = I(X;T)$ [19]. Furthermore, set of MSSs coincides with the set of solutions of

$$\inf_{P_{T|X} \in \mathcal{F}} I(X;T), \tag{5}$$

where $\mathcal{F} := \big\{ P_{T|X} : \ I(Y;T) = \sup_{Q_{T'|X}} I(Y;T') \big\}$. Here, the mutual information terms $I(X;T)$ and $I(Y;T)$ are w.r.t. $P_{X,Y} P_{T|X}$, while $I(Y;T')$ is w.r.t. $P_{X,Y} Q_{T'|X}$.

The IB framework thus relaxes the notion of MSS in two ways. First, while a MSS is defined as a deterministic function of $X$, IB solutions are randomized mappings. Such mapping can attain strictly smaller values of $\mathcal{L}_\beta(P_{T|X})$ as compared to deterministic ones.[6] Second, in view of (5), the IB framework allows for $\beta$-approximate MSSs by regulating the amount of information $T$ retains about $Y$ (see also (1)).

### D. Gaussian Information Bottleneck

The Gaussian IB [25] has a closed form solution for (3), which is the focus of this section. Let $X \sim \mathcal{N}(\mathbf{0}, \Sigma_X)$ and $Y \sim \mathcal{N}(\mathbf{0}, \Sigma_Y)$ be centered multivariate jointly Gaussian (column) vectors of dimensions $d_x$ and $d_y$, respectively. Their cross-covariance matrix is $\Sigma_{XY} := \mathbb{E}\big[XY^\top\big] \in \mathbb{R}^{d_x \times d_y}$.[7]

*1) Analytic solution:* The first step towards analytically characterizing the optimal value of $\mathcal{L}_\beta(P_{T|X})$ is showing that a Gaussian $T$ is optimal. Using the entropy power inequality (EPI) in a similar vein to [39] shows that $\inf_{P_{T|X}} \mathcal{L}_\beta(P_{T|X})$ is achieved by $P_{T|X}^{(\beta)}$ for which $(X, Y, T)$ are jointly Gaussian. Since $Y \leftrightarrow X \leftrightarrow T$ forms a Markov chain, we may represent $T = \mathrm{A}X + Z$, where $Z \sim \mathcal{N}(\mathbf{0}, \Sigma_Z)$ is independent of $(X, Y)$. Consequently, the IB optimization problem reduces to

$$\inf_{P_{T|X}} \mathcal{L}_\beta(P_{T|X}) = \inf_{\mathrm{A}, \Sigma_Z} I(X; \mathrm{A}X + Z) - \beta I(\mathrm{A}X + Z; Y)$$

$$= \frac{1}{2} \inf_{\mathrm{A}, \Sigma_Z} \underbrace{(1 - \beta) \log\left( \frac{|\mathrm{A}\Sigma_X \mathrm{A}^\top + \Sigma_Z|}{|\Sigma_Z|} \right) + \beta \log\left( |\mathrm{A}\Sigma_{X|Y} \mathrm{A}^\top + \Sigma_Z| \right)}_{:= \mathcal{L}_\beta^{(\mathrm{G})}(\mathrm{A}, \Sigma_Z)}, \tag{6}$$

where $(X \ Y)^\top \sim \mathcal{N}\left( \mathbf{0}, \begin{bmatrix} \Sigma_X & \Sigma_{XY} \\ \Sigma_{XY}^\top & \Sigma_Y \end{bmatrix} \right)$ and $Z \sim \mathcal{N}(\mathbf{0}, \Sigma_Z)$ are independent. The second equality above follows by direct computation of the mutual information terms under the specified Gaussian law.

The optimal projection $T = \mathrm{A}X + Z$ (namely, explicit structure for A and $\Sigma_Z$) was characterized in [25, Theorem 3.1], and is restated in the sequel. Let $\Sigma_{X|Y} = \Sigma_X - \Sigma_{XY} \Sigma_Y^{-1} \Sigma_{YX}$ be the mean squared error (MSE) matrix for estimating $X$ from $Y$. Consider its normalized form $\Sigma_{X|Y} \Sigma^{-1}$ and let $(\mathbf{v}_i)_{i=1}^k$, $1 \le k \le d_x$ be the left eigenvectors of $\Sigma_{X|Y} \Sigma^{-1}$ with eigenvalues $(\lambda_i)_{i=1}^k$. We assume $(\lambda_i)_{i=1}^k$ are sorted in ascending order (with $(\mathbf{v}_i)_{i=1}^k$ ordered accordingly). For each $i = 1, \ldots, k$, define $\beta_i^\star = \frac{1}{1 - \lambda_i}$ and set $\alpha_i(\beta) := \frac{\beta(1 - \lambda_i) - 1}{\lambda_i r_i}$, where $r_i = \mathbf{v}_i^\top \Sigma_X \mathbf{v}_i$.

---

[6]Consider, e.g., $(X, Y)$ as correlated Bernoulli random variables.

[7]For simplicity, we assume all matrices are full rank.

The tradeoff parameter $\beta \in [0, \infty)$ defines the optimal A matrix through its relation to the critical $(\beta_i^\star)_{i=1}^k$ values. Fix $i = 1, \ldots, k+1$, and let $\beta \in [\beta_{i-1}^\star, \beta_i^\star)$, where $\beta_0^\star = 0$ and $\beta_{k+1}^\star = \infty$. For this $\beta$, define

$$A_i(\beta) = \left[ \alpha_1(\beta)\mathbf{v}_1, \ldots, \alpha_i(\beta)\mathbf{v}_i, \mathbf{0}, \ldots, \mathbf{0} \right]^\top \in \mathbb{R}^{d_x \times d_x}. \tag{7}$$

Thus, $A_i(\beta)$ has $\alpha_1(\beta)\mathbf{v}_1^\top, \ldots, \alpha_i(\beta)\mathbf{v}_i^\top$ as its first $i$ rows, followed by $d_x - i$ rows of all zeros.

**Theorem 1 (Gaussian IB [25])** *Fix $\beta \in [0, \infty)$. Choosing $\Sigma_Z = I_{d_x}$ and*

$$A(\beta) := \begin{cases} A_0(\beta), & \beta \in [0, \beta_1^\star) \\ A_1(\beta), & \beta \in [\beta_1^\star, \beta_2^\star) \\ \quad \vdots \\ A_k(\beta), & \beta \in [\beta_k^\star, \beta_{k+1}^\star) \end{cases}, \tag{8}$$

*where $A_0(\beta)$ is the all-zero matrix, achieves the infimum in* (6).

The proof of Theorem 1 is not difficult. First one shows that $\Sigma_Z = I_{d_x}$ is optimal by a standard whitening argument. This step assumes $\Sigma_Z$ is non-singular, which does not lose generality. Indeed, if $\Sigma_Z$ is low rank then $\mathcal{L}_\beta^{(G)}(A, \Sigma_Z) = \infty$. The result then follows by differentiating $\mathcal{L}_\beta^{(G)}(A, I_{d_x})$ with respect to A and solving analytically.

In light of Theorem 1, the optimal representation $T_\beta^\star = A(\beta)X + Z$, is a noisy linear projection to eigenvectors of the normalized correlation matrix $\Sigma_{X|Y}\Sigma_X^{-1}$. These eigenvectors coincide with the well-known CCA basis [40]. However, in the Gaussian IB, $\beta$ determines the dimensionality of the projection (i.e., how many CCA basis vectors are used). As $\beta$ grows in a continuous manner, the dimensionality of $T_\beta^\star$ increases discontinuously, with jumps at the critical $\beta_i^\star$ points. Larger $\beta$ implies higher-dimensional projections, while for $\beta \leq \beta_1^\star$, we have that $T_\beta^\star = I_{d_x}$ comprises only noise. The Gaussian IB thus serves as a complexity measure with $\beta$-tunable projection rank.

*2) Information bottleneck curve:* The analytic solution to the Gaussian IB problem enables plotting the corresponding IB curve. Adapted from [25, Fig. 3], Fig. 2 illustrates this curve. Before discussing it we explain how it is computed. This is done by substituting the optimal projection $A(\beta)$ and $\Sigma_Z = I_{d_x}$ into the mutual information pair of interest. Upon simplifying, for any $\beta \in [0, \infty)$, one obtains

$$I(X; T_\beta^\star) = \frac{1}{2} \sum_{i=1}^{n_\beta} \log \left( (\beta - 1)\frac{1 - \lambda_i}{\lambda_i} \right) \tag{9a}$$

$$I(T_\beta^\star; Y) = I(X; T_\beta) - \frac{1}{2} \sum_{i=1}^{n_\beta} \log \left( \beta(1 - \lambda_i) \right), \tag{9b}$$

where $n_\beta$ is the maximal index $i$ such that $\beta \geq \frac{1}{1-\lambda_i}$. Define the function

$$\mathsf{F}_{\mathsf{IB}}(x) := x - \frac{n_{\beta(x)}}{2} \log \left( \prod_{i=1}^{n_{\beta(x)}} (1 - \lambda_i)^{\frac{1}{n_{\beta(x)}}} + e^{\frac{2x}{n_{\beta(x)}}} \prod_{i=1}^{n_{\beta(x)}} \lambda_i^{\frac{1}{n_{\beta(x)}}} \right), \tag{10}$$

where $\beta(x)$ is given by isolating $\beta$ from (9a) as a function of $I(X; T_\beta^\star) = x$. Rearranging (9) one can show that $\mathsf{F}_{\mathsf{IB}}$ assigns each $I(X; T_\beta^\star)$ with its corresponding $I(T_\beta^\star; Y)$, i.e., $\mathsf{F}_{\mathsf{IB}}\left( I(X; T_\beta^\star) \right) = I(T_\beta^\star; Y)$, for all $\beta \in [0, \infty)$
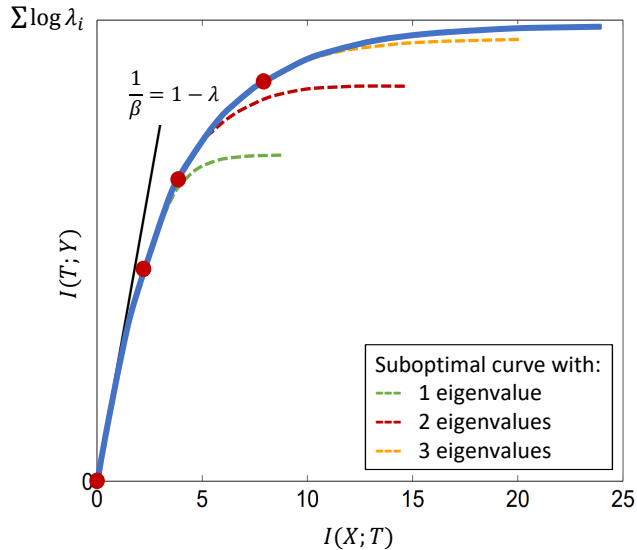
Fig. 2 [Adapted from Fig. 3 of [25]]: The Gaussian IB curve (computed with four eigenvalues $\lambda_1 = 0.1$, $\lambda_2 = 0.5$, $\lambda_3 = 0.7$, $\lambda_4 = 0.9$ and $\beta \in [0, 10^3)$) is shown in blue. Suboptimal curves that use only the first 1, 2, or 3 eigenvalues are illustrated by the dashed green, red and orange lines, respectively. Mutual information pairs at the critical values $\beta_i^\star = \frac{1}{1-\lambda_i}$, $i = 1, 2, 3, 4$, are marked by red circles. The slope of the curve at each point is the corresponding $1/\beta$. The tangent at zero, whose slope is $1 - \lambda_1$, is shown in black.

(see [25, Section 5]).

The Gaussian IB curve, as shown in blue in Fig. 2, is computed using (10). An interesting property of this curve is that while $\big(I(X;T_\beta^\star), I(T_\beta^\star;Y)\big)$ changes continuously with $\beta$, the dimensionality of $T_\beta^\star$ (or, equivalently, the number of eigenvectors used in $\mathrm{A}(\beta)$) increases discontinuously at the critical points $\beta_i^\star$, $i = 1, \ldots, k$, and is fixed between them. Restricting the dimension results in a suboptimal curve, that coincides with the optimal one up to a critical $\beta$ value and deviates from it afterwards. Some suboptimal curves are shown in Fig. 2 by the dashed, horizontally saturating segments. The critical values of $\big(I(X;T_\beta^\star), I(T_\beta^\star;Y)\big)$ after which the suboptimal curve deviate from the optimal one are marked with red circles. Notably, the optimal curve moves from one analytic segment to another in a smooth manner. Furthermore, one readily verifies that $\mathsf{F}_{\mathsf{IB}}$ is a concave function with slope tends to 0 as $\beta \to \infty$. This reflects the law of diminishing returns: encoding more information about $X$ in $T$ (higher $I(T;X)$) yields smaller increments in $I(T;Y)$.

## III. REMOTE SOURCE CODING

The IB framework is closely related to the classic information-theoretic problem of RSC. RSC dates back to the 1962 paper by Dobrushin and Tsybakov [27] (see also [28] for the additive noise case). As subsequently shown, the RSC rate-distortion function under logarithmic loss effectively coincides with 1. We start by setting up the operational problem.

## A. Operational Setup

Consider a source sequence $Y^n := (Y_1, \ldots, Y_n)$ of $n$ independent and identically distributed (i.i.d.) copies of a random variable $Y \sim P_Y \in \mathcal{P}(\mathcal{Y})$. An encoder observes the source $Y^n$ through a memoryless noisy channel $P_{X|Y}$. Namely, $P_{X|Y}$ stochastically maps each $Y_i$ to an $\mathcal{X}$-valued random variable $X_i$, where $i = 1, \ldots, n$. Denoting $P_{X,Y} = P_Y P_{X|Y}$, the pair $(X^n, Y^n)$ is distributed according to its $n$-fold product measure $P_{X,Y}^{\otimes n}$.

The sequence $X^n := (X_1, \ldots, X_n)$ is encoded through $f_n : \mathcal{X}^n \to \mathcal{M}_n$, where $|\mathcal{M}_n| < \infty$, producing the representation $M := f_n(X^n) \in \mathcal{M}_n$. The goal is to reproduce $Y^n$ from $M$ via a decoder $g_n : \mathcal{M}_n \to \hat{\mathcal{Y}}^n$, where $\hat{\mathcal{Y}}^n$ is called the reproduction space, subject to a distortion constraint.[8] The system is illustrated in Fig. 3

We adopt the logarithmic loss as the distortion measure. To set it up, let the reproduction alphabet be $\hat{\mathcal{Y}} = \mathcal{P}(\mathcal{Y})$, i.e., the set of probability measures on $\mathcal{Y}$. Thus, given a reproduction sequence $\hat{y}^n \in \hat{\mathcal{Y}}$, each $\hat{y}_i$, $i = 1, \ldots, n$, is a probability measure on $\mathcal{Y}$. This corresponds to a 'soft' estimates of the source sequence. The symbol-wise logarithmic loss distortion measure is

$$d(y, \hat{y}) := \log\left(\frac{1}{\hat{y}(y)}\right), \tag{11a}$$

which gives rise to

$$d(y^n, \hat{y}^n) := \frac{1}{n} \sum_{i=1}^{n} \log\left(\frac{1}{\hat{y}_i(y_i)}\right) \tag{11b}$$

as the distortion measure between the source $y^n \in \mathcal{Y}^n$ and its reproduction $\hat{y}^n \in \hat{\mathcal{Y}}$.

An RSC block code of length $n$ comprises a pair of encoder-decoder functions $(f_n, g_n)$. A rate-distortion pair $(R, D)$ is *achievable* if for every $\epsilon > 0$ there exists a large enough blocklength $n$ and a code $(f_n, g_n)$, such that

$$\frac{1}{n} \log |\mathcal{M}_n| < R + \epsilon \quad ; \quad \mathbb{E}\big[d(Y^n, \hat{Y}^n)\big] < D + \epsilon, \tag{12}$$

where $\hat{Y}^n = g_n\big(f_n(X^n)\big)$ and $(X^n, Y^n) \sim P_{X,Y}^{\otimes n}$. The *rate-distortion function* $R(D)$ is the infimum of rates $R$ such that $(R, D)$ is achievable for a given distortion $D$ [37, Section 10.2].

On the surface, RSC appears to be a new type of (lossy) compression problem, but it turns out to be a special case of it [41, Section 3.5]). Indeed, let us introduce another distortion metric $\tilde{d}(x, \hat{y}) := \mathbb{E}[d(Y, \hat{y})|X = x]$ and its $n$-fold extension $d(x^n, \hat{y}^n) = \frac{1}{n} \sum_{i=1}^{n} \tilde{d}(x_i, \hat{y}_i)$. One readily sees that $\mathbb{E}\big[d(Y^n, \hat{Y}^n)\big] = \mathbb{E}\big[\tilde{d}(X^n, \hat{Y}^n)\big]$. Thus, the standard rate-distortion theory implies that (subject to technical conditions, cf. [42, Chapter 26])

$$R(D) = \inf_{P_{\hat{Y}|X}: \; \mathbb{E}[\tilde{d}(X,\hat{Y})] \leq D} I(X; \hat{Y}). \tag{13}$$

## B. From Rate-Distortion Function to Information Bottleneck

Adapted to logarithmic loss, (13) can be further simplified. Indeed, consider an arbitrary distribution $P_{X,\hat{Y}}$ and extend it to $P_{Y,X,\hat{Y}} = P_{X,\hat{Y}} P_{Y|X}$, so that $Y \leftrightarrow X \leftrightarrow \hat{Y}$ forms a Markov chain. Consider any $\hat{y}$ and denote by $\hat{y}_0$ the following distribution:

$$\hat{y}_0(\cdot) := P_{Y|\hat{Y}=\hat{y}}.$$

---

[8]Rate distortion theory [37, Section 10] often considers $\hat{\mathcal{Y}} = \mathcal{Y}$, but this is not the case under logarithmic loss distortion, as described below.
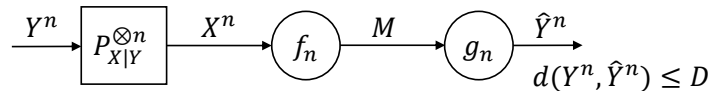
Fig. 3: Remote source coding operational setup.

From convexity, we can easily see that $\mathbb{E}\big[\tilde{d}(X,\hat{y})\big|Y = \hat{y}\big] \leq \mathbb{E}\big[\tilde{d}(X,\hat{y}_0)\big|Y = \hat{y}\big]$. Therefore, any distortion minimizing $P_{X,\hat{Y}}$ should satisfy the condition

$$\hat{y} = P_{Y|\hat{Y}=\hat{y}}, \qquad P_{\hat{Y}}\text{-a.s.}. \tag{14}$$

In other words, we have $\mathbb{E}\big[d(Y,\hat{Y})\big] = \mathbb{E}\big[\tilde{d}(X,\hat{Y})\big] \geq H(Y|\hat{Y})$ with equality holding whenever condition (14) is satisfied. Relabeling $\hat{Y}$ as $T$ we obtain

$$R(D) = \inf_{P_{T|X}:\ H(Y|T)\leq D} I(X;T), \tag{15}$$

where the underlying joint distribution is $P_{X,Y}P_{T|X}$, i.e., so that $Y \leftrightarrow X \leftrightarrow T$ forms a Markov chain.

Yet another way to express IB and $R(D)$ is the following. Let $\Pi_1$ be a random variable valued in $\mathcal{P}(\mathcal{Y})$ — the space of probability measures on $\mathcal{Y}$. Namely, $\Pi_1 = P_{Y|X=x_0}$ with probability $P_X(x_0)$ (and, more generally, $\Pi_1 = P_{Y|X=X}$ for a regular branch of conditional law $P_{Y|X}$). Then, we have

$$R(D) = \inf_{P_{\Pi_2|\Pi_1}} \left\{ I(\Pi_1;\Pi_2) : \mathbb{E}[\mathsf{D}_{\mathsf{KL}}(\Pi_1\|\Pi_2)] \leq D - H(Y|X) \right\}, \tag{16}$$

where $\mathsf{D}_{\mathsf{KL}}$ is the Kullback-Leibler (KL) divergence and the minimization is over all conditional distributions of $\Pi_2 \in \mathcal{P}(\mathcal{Y})$ given $\Pi_1$.

An interesting way to show the achievability of $R(D)$ in the RSC problem is via Berger-Tung multiterminal source coding (MSC) inner bound [43], [44]. RSC is a special case of a 2-user MSC obtained by nullifying the communication rate of one source and driving to infinity the distortion constraint on the other. Alternatively, one may employ an explicit coding scheme: first quantize (encode) the observation $X^n$ to a codeword $T^n$, such that the empirical distribution of $(X^n, T^n)$ is a good approximation of $P_X P_{T|X}$, and then communicate this quantization to the decoder via Slepian-Wolf coding [45]. The reader is referred to [46, Section IV-B] for a detailed treatment of MSC under logarithmic loss.

Setting $D = H(Y) - \alpha$ in (15) one recovers the IB problem (1). Thus, RSC with logarithmic loss can be viewed as the operational setup whose solution is given by the IB optimization problem. This provides an operational interpretation to the ad hoc definition of the IB framework as presented in Section II. Additional connections between the IB problem and other information-theoretic setups can be found in the recent survey [29].

## IV. INFORMATION BOTTLENECK IN MACHINE LEARNING

The IB framework had impact on both theory and practice of ML. While its first applications in the field was for clustering [30], more recently it is often explored as a learning objective for deep models. This was concurrently

done in [8], [31], [32] by optimizing the IB Lagrangian (3) via a variational bound compatible with gradient-based methods. Applications to classification and generative modeling were explored in [8], [31] and [32], respectively. A common theme in [8] and [32] was that the IB objective promotes disentanglement of representations.

From a theoretical standpoint, the IB gave rise to a bold information-theoretic DL paradigm [5], [6]. It was argued that, even when the IB objective in not explicitly optimized, DNNs trained with cross-entropy loss and SGD inherently solve the IB compression-prediction trade-off. This 'IB theory for DL' attracted significant attention, culminating in many follow-up works that tested the proclaimed narrative and its accompanying empirical observations. This section surveys the practical and theoretical roles of IB in DL, focusing on classification tasks.

### A. Information Bottleneck as Optimization Objective

*1) Variational approximation and efficient optimization:* In [31], a variational approximation to the IB objective (3) was proposed. The approximation parametrized the objective using a DNN and proposed an efficient training algorithm. As opposed to classic models, the obtained system is stochastic in the sense that it maps inputs to internal representations via randomized mappings. Empirical tests of the VIB system showed that it generalized better and was more robust to adversarial examples compared to competing methods.

**IB objective.** Given a DNN, regard its $\ell^{\text{th}}$ internal representation $T_\ell$, abbreviated as $T$, as a randomized mapping operating on the input feature $X$; the corresponding label is $Y$. Denote by $\mathcal{X}$, $\mathcal{Y}$ and $\mathcal{T}$ the sets in which $X$, $Y$ and $T$ take values. The encoding of $X$ into $T$ is defined through a conditional probability distribution, which the VIB parametrizes as $P_{T|X}^{(\theta)}$, $\theta \in \Theta$. Together with $P_{X,Y}$, $P_{T|X}^{(\theta)}$ defines the joint distribution of $\big(X, Y, T^{(\theta)}\big) \sim P_{X,Y,T}^{(\theta)} :=$ $P_{X,Y} P_{T|X}^{(\theta)}$. With respect to this distribution, one may consider the optimization objective

$$\mathcal{L}_\beta^{(\text{VIB})}(\theta) := \max_{\theta \in \Theta} I\big(T_\ell^{(\theta)}; Y\big) - \beta I\big(X; T_\ell^{(\theta)}\big). \tag{17}$$

In accordance to the original IB problem [16], the goal here is to learn representations that are maximally informative about the label $Y$, subject to a compression requirement on $X$. However, since the data distribution $P_{X,Y}$ is unknown and (even knowing it) direct optimization of (17) is intractable — a further approximation is needed.

**Variational approximation.** To overcome intractability of (17), the authors of [31] lower bound it by a form that is easily optimized via gradient-based methods. Using elementary properties of mutual information, entropy and KL divergence, (17) is lower bounded by

$$\mathbb{E}_{P_{Y,T}^{(\theta)}} \Big[\log Q_{Y|T}^{(\phi)}\big(Y\big|T^{(\theta)}\big)\Big] - \beta \mathsf{D}_{\mathsf{KL}}\big(P_{T|X}^{(\theta)} \big\| P_T^{(\theta)} \,\big| P_X\big), \tag{18}$$

where $Q_{Y|T}^{(\phi)}$ is a conditional distribution from $\mathcal{T}$ to $\mathcal{Y}$ parametrized by a NN with parameters $\phi \in \Phi$. This distribution plays the role of a variational approximation of the decoder $P_{Y|T}^{(\theta)}$. Another difficulty is that the marginal $P_T^{(\theta)}(\cdot) = \int P_{T|X}^{(\theta)}(\cdot|x)\, \mathrm{d}P_X(x)$, which is defined by $P_{T|X}^{(\theta)}$ and $P_X$, it is typically intractable. To circumvent this, one may replace it with some reference measure $R_T \in \mathcal{P}(\mathcal{T})$, thus further lower bounding (18). This is the strategy that was employed in [31].
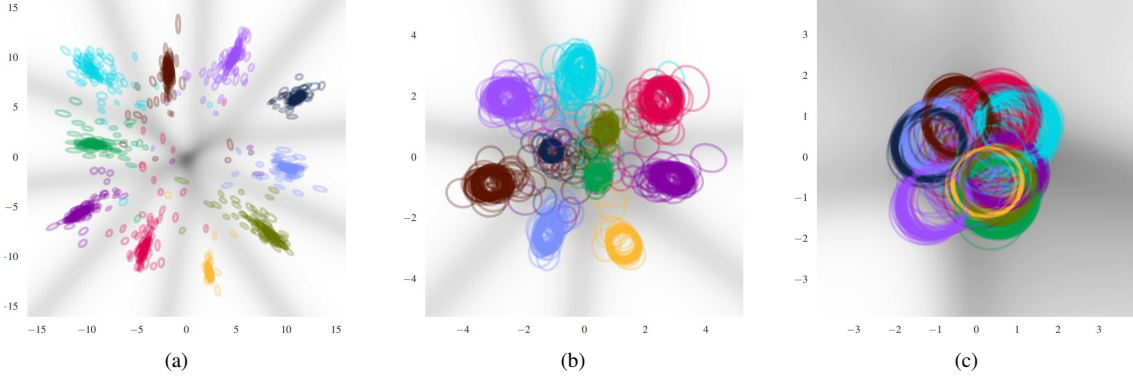
Fig. 4 [Fig. 2 from [31]]: 2-dimensional VIB embedding of $10^3$ MNIST images for: (a) $\beta = 10^{-3}$ ; (b) $\beta = 10^{-1}$; and (c) $\beta = 1$. The ellipses in each figure illustrate $95\%$ confidence intervals of the Gaussian encoding $P_{T|X}^{(\theta)}(\cdot|x) = \mathcal{N}\big(\mu_\theta(x), \Sigma_\theta(x)\big)$, for $10^3$ input MNIST images $x$. The larger $\beta$ is, the more compressed representations become. This is expressed in larger covariance matrices for the encoder. The background grayscale illustrates $H\left(Q_{Y|T}^{(\phi)}(\cdot|t)\right)$, for each $t \in \mathbb{R}^2$, which measures classification uncertainty at a given point.

Using the reparametrization trick from [47] and replacing $P_{X,Y}$ in (18) with its empirical proxy $P_n := \frac{1}{n}\sum_{i=1}^n \delta_{(x_i, y_i)}$, where $\delta_x$ is the Dirac measure centered at $x$ and $\mathcal{D}_n := \big\{(x_i, y_i)\big\}_{i=1}^n$ is the dataset, we arrive at the empirical loss function

$$\hat{\mathcal{L}}_\beta^{(\mathsf{VIB})}(\theta, \phi, \mathcal{D}_n) := \frac{1}{n}\sum_{i=1}^n \mathbb{E}\Big[-\log Q_{Y|T}^{(\phi)}\big(y_n\big|f(x_n, Z)\big)\Big] + \beta \mathsf{D}_{\mathsf{KL}}\Big(P_{T|X}^{(\theta)}(\cdot|x_n)\Big\|R_T(\cdot)\Big). \tag{19}$$

Here $Z$ is an auxiliary noise variable (see [47]) and the expectation is w.r.t. its law. The first term is the average cross-entropy loss, which is commonly used in DL. The second term serves as a regularizer that penalizes the dependence of $X$ on $T$, thus encouraging compression. The empirical estimator in (19) of the variational lower bound is differentiable and easily optimized via standard stochastic gradient-based methods. The obtained gradient is an unbiased estimate of the true gradient.

**Empirical study.** For their experiments, the authors of [31] set the encoder distribution to $P_{T|X}^{(\theta)}(\cdot|x) = \mathcal{N}\big(\mu_\theta(x), \Sigma_\theta(x)\big)$, with mean and convariance matrix parametrized by a DNN. The variational decoder $Q_{Y|T}^{(\phi)}$ was set to a logistic regression function, while $R_T$ was chosen as a (fixed) standard Gaussian distribution. The performance of VIB was tested on the MNIST and ImageNet datasets; we focus on the MNIST results herein. First, it was shown that a VIB classifier outperforms a multi-layer perceptron (MLP) fitted using (penalized) maximum likelihood estimation (MLE). Considered penalization techniques include dropout [48], confidence penalty, and label smoothing [49] (see [31, Table 1] for accuracy results).

To illustrate the operation of VIB, a 2-dimensional embedding of internal representations is examined for different $\beta$ values. The results are shown in Fig. 4 (reprinted from [31, Fig. 2]). The posteriors $P_{T|X}^{(\theta)}$ are represented as Gaussian ellipses (representing the $95\%$ confidence region) for $10^3$ images from the test set. Colors designate true class labels. The grayscale shade in the background corresponds to the entropy of the variational classifier $Q_{Y|T}^{(\phi)}$

at a given point, i.e., $H\left(Q_{Y|T}^{(\phi)}(\cdot|t)\right)$, for $t \in \mathbb{R}^2$. This entropy quantifies the uncertainty of the decoder regarding the class assignment to each point. Several interesting observations are: (i) as $\beta$ increases (corresponds to more compression in (17)), the Gaussian encoder covariances increase in relation to the distance between samples, and the classes start to overlap; (ii) beyond some critical $\beta$ value, the encoding 'collapses' essentially losing class information; and (iii) there is uncertainty in class predictions, as measured by $H\left(Q_{Y|T}^{(\phi)}(\cdot|t)\right)$, in the areas between the class embeddings. This illustrates the ability of VIB to learn meaningful and interpretable representations of data, while preserving good classification performance.

While $\beta = 10^{-1}$ (Fig. 4(b)) has relatively large covariance matrices, [31] showed that this this system has reasonable classification performance ($3.44\%$ test error). This implies there is significant in-class uncertainty about locations of internal representations, although the classes themselves are well-separated. Such encoding would make it hard to infer which input image corresponds to a given internal representation. This property leads to consider model robustness, which was demonstrated as another virtue of VIB.

To test model robustness, [31] employed the fast gradient sign (FGS) [50] and $L^2$ optimization [51] attacks for perturbing inputs to fool the classifier. Robustness was measured as classification accuracy of adversarial examples. The (stochastic) VIB classifier showed increasing robustness to these attacks, compared to competing deterministic models which misclassified all the perturbed inputs. This is an outcome of the randomized VIB mapping from $\mathcal{X}$ to $\mathcal{T}$, which suppresses the ability to tailor adversarial examples that flip the classification label. Robustness is also related to compression: indeed, larger $\beta$ values correspond to more robust systems, which can withstand adversarial attacks with large perturbation norms (see [31, Fig. 3]).

*2) IB objective for sufficiency, minimality and disentanglement:* Another perspective on the IB objective was proposed in [8], that was published concurrently with [31]. This work provided an information-theoretic formulation of some desired properties of internal representations, such as sufficiency, minimality and disentanglement, and showed that IB optimization favors models that posses them. To optimize VIB, [8] presented a method that closely resembles that of [31]. The method reparametrizes $T^{(\theta)}$ as a product of some function of $X$ and a noise variable $Z$. The noise distribution is for the system designer to choose. Setting $Z \sim \text{Bern}(p)$, recovers the popular dropout regularization, hence the authors of [8] called their method 'information dropout'. The noise distribution employed in [8] is log-normal with zero mean, and variance that is a parametrized function of $X$. We omit further details due to similarity to VIB optimization (Section IV-A1).

**Minimality and sufficiency.** That IB solutions are approximate MSSs is inherent to the problem formulation (see Section II-C). The authors of [8] discussed the relation between minimality of representation and invariance to nuisance factors. Such factors affect the data, but the label is invariant to them. By penalizing redundancy of representations (i.e., enforcing small $I(X; T^{(\theta)})$) it was heuristically argued that the model's sensitive to nuisances is mitigated. Some experiments to support this claim were provided, demonstrating the ability of VIB to classify hand-written digits from the cluttered MNIST dataset [52], or CIFAR-10 images occluded by MNIST digits.

**Disentanglement.** Another property considered in [8] is disentanglement, which refers to weak dependence

between elements of an internal representation vector. This idea was formalized using total correlation (TC):

$$\mathsf{D}_{\mathsf{KL}}\left(P_T^{(\theta)}\middle\|\prod_{j=1}^d Q_j\right) \tag{20}$$

where $\prod_{j=1}^d Q_j$ is some product measure on the elements of the $d$-dimensional representation $T$. Adding a TC regularizer to (18) will encourage the system to learn disentangled representation in the sense of small (20). If the TC regularization parameter is set equal to $\beta$ in (18), it trivially simplifies to

$$\mathbb{E}_{P_{Y,T}^{(\theta)}}\left[\log Q_{Y|T}^{(\phi)}\left(Y\middle|T^{(\theta)}\right)\right] - \beta\mathsf{D}_{\mathsf{KL}}\left(P_{T|X}^{(\theta)}\middle\|\prod_{j=1}^d Q_j\middle|P_X\right).$$

Thus, choosing $R_T \in \mathcal{P}(\mathcal{T})$, in the framework of [31], as a product measure is equivalent to regularizing for disentanglement. While the term 'disentanglement' was not used in [31], they do set $R_T$ as a product measure.

Altogether, Section IV-A demonstrates the practical usefulness of the IB as an optimization objective. It is easy to optimize under proper parameterization and learns representation with various desired properties. The work of [6] took these observation a step further. They claimed that DNNs trained with SGD and cross-entropy loss inherently solve the IB problem, even when there in no explicit reference to the IB problem in the system design. We elaborate on this theory next.

### B. Information Bottleneck Theory for Deep Learning

Recently, a information-theoretic paradigm for DL based on the IB framework was proposed [5], [6]. It claimed that DNN classifiers trained with cross-entropy loss and SGD inherently (try to) solve the IB optimization problem. Namely, the system aims to find internal representations that are maximally informative about the label $Y$, while compressing $X$, as captured by (3). Coupled with an empirical case study, this perspective was leveraged to reason about the optimization dynamics in the IP, properties of SGD training, and the computational benefit of deep architectures. This section reviews these ideas as originally presented in [5], [6], with the exception of Section IV-B3 that discusses the incompatibility of the IB framework to deterministic DNNs. Whether the proposed theory holds is general was challenged by several follow-up works, which are expounded upon in Section V.

*1) Setup and Preliminaries:* Consider a feedforward DNN with $L$ layers, operating on an input $x \in \mathbb{R}^{d_0}$ according to:

$$\phi_\ell(x) = \sigma(\mathrm{A}_\ell\phi_{\ell-1}(x) + b_\ell), \qquad \phi_0(x) = x, \qquad \ell = 1, \ldots, L, \tag{21}$$

where $\mathrm{A}_\ell \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$ and $b_\ell \in \mathbb{R}^{d_\ell}$ are, respectively, the $\ell^{\text{th}}$ weight matrix and bias vector, while $\sigma : \mathbb{R} \to \mathbb{R}$ is the activation function operating on vectors element-wise. Letting $(X, Y) \sim P_{X,Y}$ be the feature-label pair, we call $T_\ell \triangleq \phi_\ell(X)$, $\ell = 1, \ldots, L-1$, the $\ell^{\text{th}}$ internal representation. The output layer $T_L$ is a reproduction[9] of $Y$, sometimes denoted by $\hat{Y}$.

The goal of the DNN classifier is to learn a reproduction $\phi_L(X) = \hat{Y}$ that is a good approximation of the true label $Y$. Let $\theta$ represent the network parameters (i.e., weight matrices and bias vectors) and write $\phi_L^{(\theta)}$ for $\phi_L$ to

---

[9]Perhaps up to an additional soft-max operation.

stress the dependence of the DNN's output on $\theta$. Statistical learning theory measures the quality of the reproduction by the *population risk*

$$L_{P_{X,Y}}(\theta) := \mathbb{E}c\big(\phi_L^{(\theta)}(X), Y\big) = \int c\big(\phi_L^{(\theta)}(x), y\big)\, dP_{X,Y}(x, y),$$

where $c$ is the cost/loss function. Since $P_{X,Y}$ is unknown, a learning algorithm cannot directly compute $L_{P_{X,Y}}(\theta)$ for a given $\theta \in \Theta$. Instead, it can compute the empirical risk of $\theta$ on the dataset $\mathcal{D}_n := \big\{(x_i, y_i)\big\}_{i=1}^{n}$, which comprises $n$ i.i.d. samples from $P_{X,Y}$. The *empirical risk* is given by

$$L_{\mathcal{D}_n}(\theta) := \frac{1}{n}\sum_{i=1}^{n} c(\phi_L^{(\theta)}(x_i), y_i).$$

Minimizing $L_{\mathcal{D}_n}(\theta)$ over $\theta$ is practically feasible, but the end goal is to attain small population risk. The gap between them is captured by the *generalization error*

$$\mathsf{gen}(P_{X,Y}, \theta) := \mathbb{E}\big[L_{P_{X,Y}}(\theta) - L_{\mathcal{D}_n}(\theta)\big],$$

where the expectation is w.r.t $P_{X,Y}^{\otimes n}$. The *sample complexity* $n^\star(P_{X,Y}, \epsilon, \delta)$ is the least number of samples $n$ needed to ensure $L_{P_{X,Y}}(\theta) - L_{\mathcal{D}_n}(\theta) \leq \epsilon$ with probability at least $1 - \delta$. To reason about generalization error and sample complexity, the IB theory for DL views the learning dynamics through the so-called 'information plane'.

*2) The information plane:* Consider the joint distribution of the label, feature, internal representation and output random variables. By (21), they form a Markov chain $Y \leftrightarrow X \leftrightarrow T_1 \leftrightarrow T_2 \leftrightarrow \ldots \leftrightarrow T_L$, i.e., their joint law factors as $P_{X,Y,T_1,T_2,\ldots,T_L} = P_{X,Y}P_{T_1|X}P_{T_2|T_1}\cdots P_{T_L|T_{L-1}}$.[10] Based on this Markov relation, the data processing inequality (DPI) implies

$$I(X; Y) \geq I(T_1; Y) \geq I(T_2; Y) \geq \ldots \geq I(T_L; Y)$$

$$H(X) \geq I(X; T_1) \geq I(X; T_2) \geq \ldots \geq I(X; T_L). \tag{22}$$

illustrating how information inherently dissipates deeper in the network. The IB theory for DL centers around how each internal representation $T_\ell$ tries to balance its two associated mutual information term.

Fix $\ell = 1, \ldots, L-1$, and consider the conditional marginal distributions $P_{T_\ell|X}$ and $P_{T_L|T_\ell}$. Respectively, these distributions define an 'encoder' (of $X$ into $T_\ell$) and 'decoder' (of $T_\ell$ into $T_L = \hat{Y}$) for an IB problem associated with the $\ell^{\text{th}}$ hidden layer (see Fig. 5). The argument of [5] and [6] is that the IB framework captures the essence of learning to classify $Y$ from $X$: to reconstruct $Y$ from $X$, the latter has to go through the bottleneck $T_\ell$. This $T_\ell$ should shed information about $X$ that is redundant/irrelevant for determining $Y$, while staying maximally informative about $Y$ itself. The working assumption of [5], [6] is that by optimizing the DNN's layers $\{T_\ell\}_{\ell=1}^{L}$ for that task via standard SGD, the encoder-decoder pair for each hidden layer converges to its optimal IB solution.

With this perspective, [6] conducted an empirical case study of a DNN classifier trained on a synthetic task, reporting several striking findings. The study is centered around IP visualization of the learning dynamics, i.e.,

---

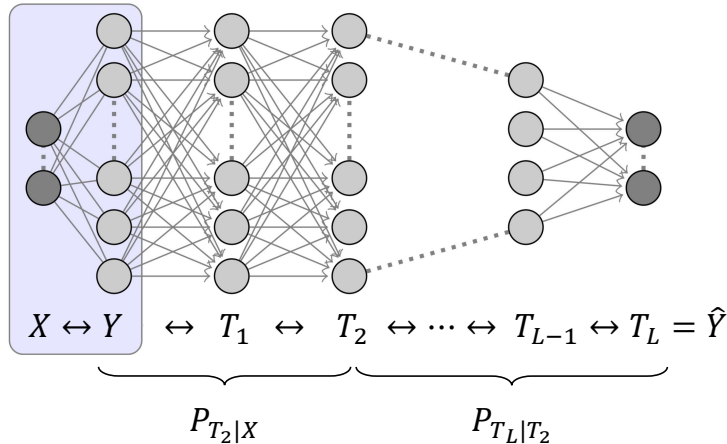[10]Marginal distributions of this law are designated by keeping only the relevant subscripts.

Fig. 5 [Adapted from Fig. 1 from [5]]: IB framework for a DNN classifier with $L$ layers. The label, feature, hidden representations and output form a Markov chain $Y \leftrightarrow X \leftrightarrow T_1 \leftrightarrow T_2 \leftrightarrow \ldots \leftrightarrow T_L$. An encoder $P_{T_\ell|X}$ and a decoder $T_{T_L|T_\ell}$ are associated with each hidden layer $T_\ell$, $\ell = 1, \ldots, L-1$.

tracking the evolution of $\big(I(X;T_\ell), I(T_\ell;Y)\big)$ across training epochs. First, they observed that the training process comprises two consecutive phases, termed 'fitting' and 'compression': after a quick initial fitting phase, almost the entire training process is dedicated to compressing $X$, while staying informative about $Y$. Second, it was observed that the compression phase starts when the mean and variance of the stochastic gradient undergoes a phase transition from high to low signal-to-noise (SNR). Third, [6] observed that the two training phases accelerate when more layers are added to the network, which led to an argument that the benefit of deeper architectures is computational.

Central to the empirical results of [6] is the ability to measure $I(X;T_\ell)$ and $I(T_\ell;Y)$ in a DNN with fixed parameters. We explain next that these information measures are ill-posed in deterministic networks (i.e., networks that define a deterministic mapping from input to output). Afterwards, we describe how [6] circumvented this issue by instead evaluating the mutual information terms after quantizing the internal representation vectors. We then turn to expound upon the empirical findings of [6] based on these quantized measurements.

*3) Vacuous mutual information in deterministic deep networks:* We address theoretical caveats in applying IB-based reasoning to deterministic DNNs. A deterministic DNN is one whose output, as well as every internal representation, is a deterministic function of its input $X$. The setup described in Section IV-B1, which is standard for DNNs and the one used in [5], [6], [9], adheres to the deterministic framework. In deterministic DNNs with continuous and strictly monotone nonlinearities (e.g., $\tanh$ or $\mathrm{sigmoid}$) or bi-Lipschitz (e.g., $\mathrm{leaky\text{-}ReLU}$), the mutual information $I(X;T_\ell)$ is provably either infinite (continuous $X$) or a constant that does not depend on the network parameters (discrete $X$). The behavior for $\mathrm{ReLU}$ or step activation functions is slightly different, though other issues arise, such as the IB functional being piecewise constant.

We start from continuous inputs. Elementary information-theoretic arguments show that $I(X;T_\ell) = \infty$ a.s.[11] if $X$ is a continuous $\mathbb{R}^d$-valued random variable and nonlinearities are continuous and strictly monotone. Indeed, in

---

[11]With respect to, e.g., the Lebesgue measure on the parameter space of $\big\{(A_i, b_i)\big\}_{i=1}^{\ell}$, or the entry-wise i.i.d. Gaussian measure.

this case the $d_\ell$-dimensional $T_\ell = \phi_\ell(X)$ is a.s. continuous whenever $d_\ell \leq d$, and so $I(X;T_\ell) \geq I(T_\ell;T_\ell) = \infty$ (cf. [42, Theorems 2.3 and 2.4]). A more general statement was given in [15, Theorem 1] and is restated next.

**Theorem 2 (Theorem 1 of [15])** *Let $X$ be a $d$-dimensional random variable, whose distribution has an absolutely continuous component with density function that is continuous on a compact subset of $\mathbb{R}^d$. Assume that the activation function $\sigma$ in (21) is either bi-Lipschitz or continuously differentiable with strictly positive derivative. Then, for every $\ell = 1, \ldots, L$ and almost all weight matrices $A_1, \ldots, A_\ell$, we have $I(X;T_\ell) = \infty$.*

The proof uses the notion of correlation dimension [53]. It shows that $X$ has a positive correlation dimension that remains positive throughout the DNN, from which the conclusion follows. This result broadens the conditions for which $I(X;T_\ell) = \infty$ beyond requiring that $T_\ell$ is continuous. This, for instance, accounts for cases when the number of neurons $d_\ell$ in $T_\ell$ exceeds the dimension of $X$. Theorem 2 implies that for continuous features $X$, the *true* mutual information $I(X;T_\ell) = \infty$, for any hidden layer $\ell = 1, \ldots, L$ in the $\tanh$ network from [6].

To avoid this issue, [6] model $X \sim \mathsf{Unif}(\mathcal{X}_n)$, where $\mathcal{X}_n = \{x_i\}_{i=1}^n$. While having a discrete distribution for $X$ ensures mutual information is finite, as $I(X;T_\ell) \leq H(X) = \log n$, a different problem arises. Specifically, whenever nonlinearities are injective (e.g., strictly monotone), the map from $\mathcal{X}_n$ to $\phi_\ell(\mathcal{X}_n) = \{\phi_\ell(x) : x \in \mathcal{X}_n\}$ (as a mapping between discrete sets) is a.s. injective. As such, we have that $I(X;T_\ell) = H(X) = \log n$ and $I(T_\ell;Y) = I(X;Y)$, which are constants independent of the network parameters.

Both continuous and discrete degeneracies are a consequence of the deterministic DNN's ability to encode information about $X$ in arbitrarily fine variations of $T_\ell$, essentially without loss, even if deeper layers have fewer neurons. Consequently, no information about $X$ is lost when traversing the network's layers, which renders $I(X;T_\ell)$ a vacuous quantity for almost all network parameters. In such cases approximating or estimating $I(X;T_\ell)$ and $I(T_\ell;Y)$ to study DNN learning dynamics is unwarranted. Indeed, the true value (e.g., infinity or $H(X)$ for $I(X;T_\ell)$) is known and does not depend on the network.

*4) Mutual information measurement via quantization:* The issues described above are circumvented in [6] by imposing a concrete model on $(X,Y)$ and quantizing internal representation vectors. Namely, $X$ is assumed to be uniformly distributed over the dataset $\mathcal{X}_n = \{x_i\}_{i=1}^n$, while $P_{Y|X}$ is defined through a logistic regression with respect to a certain spherically symmetric real-valued function of $X$. With this model, [6] quantize $T_\ell$ to compute $\big(I(X;T_\ell), I(T_\ell,Y)\big)$, as described nex.

Specifically, consider a DNN with bounded nonlinearities $\sigma : \mathbb{R} \to [a,b]$. Let $\mathsf{Q}_m[T_\ell]$ be the quantized version of the $d_\ell$-dimensional random vector $T_\ell$, which dissects its support (the hypercube $[a,b]^{d_\ell}$) into $m^{d_\ell}$ equal-sized cells ($m$ in each direction). The two mutual information terms are then approximated by

$$I(X;T_\ell) \approx I\big(X;\mathsf{Q}_m[T_\ell]\big) = H\big(\mathsf{Q}_m[T_\ell]\big) \tag{23a}$$

$$I(T_\ell;Y) \approx I\big(\mathsf{Q}_m[T_\ell];Y\big) = H\big(\mathsf{Q}_m[T_\ell]\big) - \sum_{y=0,1} P_Y(y) H\big(\mathsf{Q}_m[T_\ell]\big|Y=y\big) \tag{23b}$$

where the equality in (23a) is because $T_\ell$ (and thus $\mathsf{Q}_m[T_\ell]$) is a deterministic function of $X$, while $P_Y$ in (23a) is given by $P_Y(1) = \frac{1}{n}\sum_{i=1}^n P_{Y|X}(1|x_i)$. Computing $H\big(\mathsf{Q}_m[T_\ell]\big)$ amounts to counting how many inputs $x_i$,
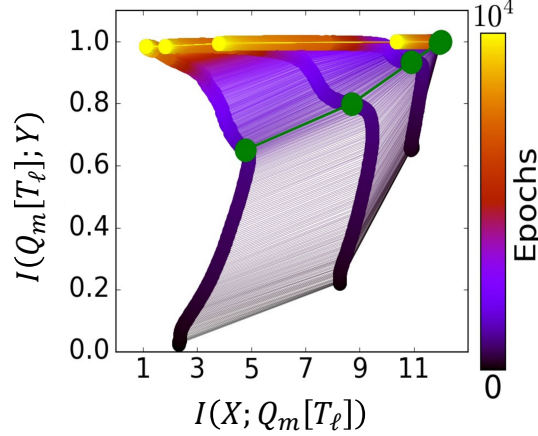
Fig. 6 [Fig. 3 from [6]]: IP dynamics for the fully connected 12–10–7–5–4–3–2 binary DNN classifier from [6]. The horizontal and vertical axis correspond to $I(X; Q_m[T_\ell])$ and $I(Q_m[T_\ell]; Y)$, respectively. The thick curves show $\big(I(X; Q_m[T_\ell]), I(Q_m[T_\ell]; Y)\big)$ values across training epoch (designated by the color map) for the different hidden layers. Curves for deeper layers appear more to the left. The thin lines between the curves connect $\big(I(X; Q_m[T_\ell]), I(Q_m[T_\ell]; Y)\big)$ values across layers at a given epoch. The two training phases of 'fitting' and 'compression' are clearly seen in this example, as the curve for each layer exhibits an elbow effect. The green marks correspond to the transition between the 'drift' and 'diffusion' phases of SGD (see Section IV-B6 and Fig. 7)

$i = 1, \ldots, n$, have their internal representation $\phi_\ell(x_i)$ fall inside each of the $m^{d_\ell}$ quantization cells under $Q_m$. For the conditional entropy, counting is restricted to $x_i$'s whose corresponding label is $y_i = y$.

The quantized mutual information proxies are motivated by the fact that for any random variable pair $(A, B)$,

$$I(A; B) = \lim_{j,k \to \infty} I\big([A]_j; [B]_k\big), \tag{24}$$

where $[A]_j$ and $[B]_k$ are any sequences of finite quantizations of $A$ and $B$, respectively, such that the quantization errors tend to zero as $j, k \to \infty$ (see [54, Section 2.3]). Thus, at the limit of infinitely-fine quantization, the approximations in (23) are exact. In practice, [6] fix the resolution $m$ in $Q_m$ and perform the computation w.r.t. this resolution. Doing so generally bounds the approximation away from the true value, resulting in a discrepancy between the computed values and the true DNN model (where activations are not quantized during training or inference). This discrepancy and its effect on the observations of [6] are discussed in detail in Section V. The remainder of the current section focuses on describing the findings of [6], temporarily overlooking these subtleties.

*5) Information plane dynamics and two training phases:* Using the above methodology for approximating $\big(I(X; T_\ell), I(T_\ell, Y)\big)$, [6] explores IP dynamics during training of a DNN classifier on a certain synthetic task. The task is binary classification of 12-dimensional inputs using a fully connected 12–10–7–5–4–3–2 architecture with $\tanh$ nonlinearities.[12] Training is performed via standard SGD and cross-entropy loss. The IP visualizations are produced by subsampling training epochs and computing $\big(I(X; Q_m[T_\ell]), I(Q_m[T_\ell]; Y)\big)$, with $m = 30$, w.r.t. the (fixed) DNN parameters at each epoch. To smooth out the curves, trajectories are averaged over 50 runs.

---

[12]For the full experimental setup see [6, Section 3.1].

Fig. 6 (reprinted from [6, Fig. 3]), demonstrates the IP dynamics for this experiment. The thick trajectories show $\big(I\big(X;\mathsf{Q}_m[T_\ell]\big), I\big(\mathsf{Q}_m[T_\ell];Y\big)\big)$ evolution across training epochs (which are designated by the color map) for the different hidden layers of the network. Deeper layers are bound from above (in the Pareto sense) by shallower ones, in accordance to the DPI (see Eq. (22)). The thin lines between the IP curves connect the mutual information pair values across layers at a given epoch. These IP dynamics reveal a remarkable trend, as the trajectories of $\big(I\big(X;\mathsf{Q}_m[T_\ell]\big), I\big(\mathsf{Q}_m[T_\ell];Y\big)\big)$ exhibit two distinct phases: an increase in both $I\big(X;\mathsf{Q}_m[T_\ell]\big)$ and $I\big(\mathsf{Q}_m[T_\ell];Y\big)$ at the beginning of training, followed by along-term decrease in $I\big(X;\mathsf{Q}_m[T_\ell]\big)$ that subsumes most of the training epochs. These two phases were termed, respectively, 'fitting' and 'compression'. While the increase in $I\big(\mathsf{Q}_m[T_\ell];Y\big)$ during the fitting phase is expected, there was no explicit regularization to encourage compression of representations.

Inspired by the classic IB problem (Section II), the compression phase was interpreted as the network 'shedding' information about $X$ that is 'irrelevant' for the classification task. The authors of [6] then argued that the observed two-phased dynamics are inherent to DNN classifiers trained with SGD, even when the optimization method/objective have no explicit reference to the IB principle. The compression phase was further claimed to be responsible for the outstanding generalization performance of deep networks, although no rigorous connection between $I\big(X;\mathsf{Q}_m[T_\ell]\big)$ and the generalization error is currently known.

*6) Connection to stochastic gradient descent dynamics:* To further understand the two IP phases of training, [6] compared them with SGD dynamics. For each layer $\ell = 1, \ldots, L$, define

$$\mu_\ell := \left\| \left\langle \frac{\partial c}{\partial \mathrm{A}_\ell} \right\rangle \right\|_{\mathsf{F}} \quad ; \quad \sigma_\ell := \left\| \mathsf{STD} \left( \frac{\partial c}{\partial \mathrm{A}_\ell} \right) \right\|_{\mathsf{F}},$$

where $\langle \cdot \rangle$ and $\mathsf{STD}(\cdot)$ denote, respectively, the mean and the element-wise standard deviation (std) across samples within a minibatch, and $\| \cdot \|_{\mathsf{F}}$ is the Frobenius norm. Thus, $\mu_\ell$ and $\sigma_\ell$ capture the gradient's mean and std w.r.t. to the $\ell^{\text{th}}$ weight matrix. Since weights tend to grow during training, $\mu_\ell$ and $\sigma_\ell$ were normalized by $\|\mathrm{A}_\ell\|_{\mathsf{F}}$. The evolution of the normalized mean and std during training is displayed in Fig. 7 (reprinted from [6, Fig. 4]).

The figure shows a clear phase transition around epoch 350, marked with the vertical grey line. At the first phase, termed 'drift', the gradient mean $\mu_\ell$ is much larger than its fluctuations, as measured by $\sigma_\ell$. Casting $\mu_\ell/\sigma_\ell$ as the gradient signal-to-noise ratio (SNR) at the $\ell^{\text{th}}$ layer, the drift phase is characterized by high SNR. This corresponds to SGD exploring the high-dimensional loss landscape, quickly converging from the random initialization to a near (locally) optimal region. In the second phase, termed 'diffusion', the gradient SNR abruptly drops. The low SNR regime, as explained in [6], is a consequence of empirical error saturating and SGD being dominated by its fluctuations. This observation corresponds to the earlier work of [55], [56], where two phases of gradient descent (convergence towards a near-optimal region and oscillation in that region) were also identified and described in greater generality . Rather than 'drift' and 'compression', these phases are sometimes termed 'transient' and 'stochastic' or 'search' and 'convergence'.

The correspondence between the two SGD phases and the IP trajectories from Fig. 6 was summarized in [6] as follows. First, the transition from 'fitting' to 'compression' in Fig. 6 happens roughly at the same epoch when SGD transitions from 'drift' to 'diffusion' (Fig. 7) — this is illustrated by the green marks in Fig. 6. The SGD drift
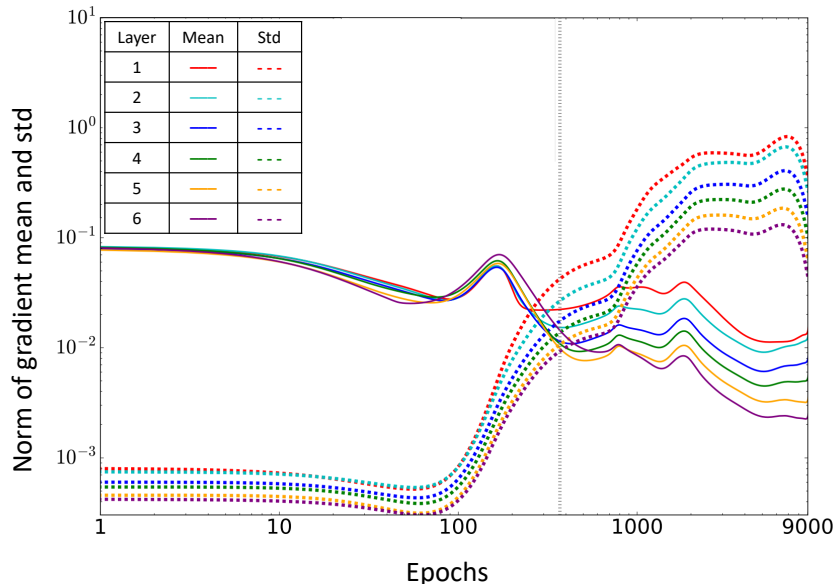
Fig. 7 [Fig. 4 from [6]]: Norms of gradient mean (solid lines) and standard deviation (dashed lines) for the different layers. For each layer, the values are normalized by the $L^2$ norm of the corresponding weight matrix. The grey vertical line marks the epoch when a phase transition between a high-SNR to a low-SNR occurs.

phase quickly reduces empirical error, thereby increasing $I(T_\ell; Y)$ (as captured by its approximation $I\big(\mathsf{Q}_m[T_\ell]; Y\big)$ in Fig. 6). The connection between the SGD diffusion phase and compression of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ in the IP is less clear. A heuristic argument given in [6] is that SGD diffusion mostly adds random noise to the weights, evolving them like Wiener processes. This diffusion-like behaviour inherently relies on the randomness in SGD (as opposed to, e.g., full gradient descent). Based on this hypothesis, [6] claimed that SGD diffusion can be described by the Fokker-Planck equation, subject to a small training error constraint. Together with the maximum entropy principle, this led to the conclusion that the diffusion phase maximizes the conditional entropy $H(X|T_\ell)$, or equivalently, minimizes $I(X; T_\ell)$ (approximated by $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ in Fig. 6). We note that [6] presented no rigorous derivations to support this explanation.

*7) Computation benefit of deep architectures:* It is known that neural networks gain in representation power with the addition of layers [57]–[60]. The argument of [6] is that deep architectures also result in a computational benefit. Specifically, adding more layers speeds up both the fitting and the compression phases in the IP dynamics.

Recall the synthetic binary classification task of 12-dimensional inputs described in Section IV-B5. Consider 6 neural network architectures of increasing depth (1 to 6 hidden layers) trained to solve that task. The first hidden layer has 12 neurons and each succeeding one has two neurons less. For example, the 3rd architecture is a fully connected 12–12–10–8–2 networks. Fig. 8 (reprinted from [6, Fig. 5]) shows the IP dynamics of these 6 architectures, each averaged over 50 runs.

The figure shows that additional layers speed up the IP dynamics. For instance, while last hidden layer of the deepest network (bottom-right subfigure) attains its maximal $I\big(\mathsf{Q}_m[T_6]; Y\big) \approx 0.7$ after about 400 epochs,
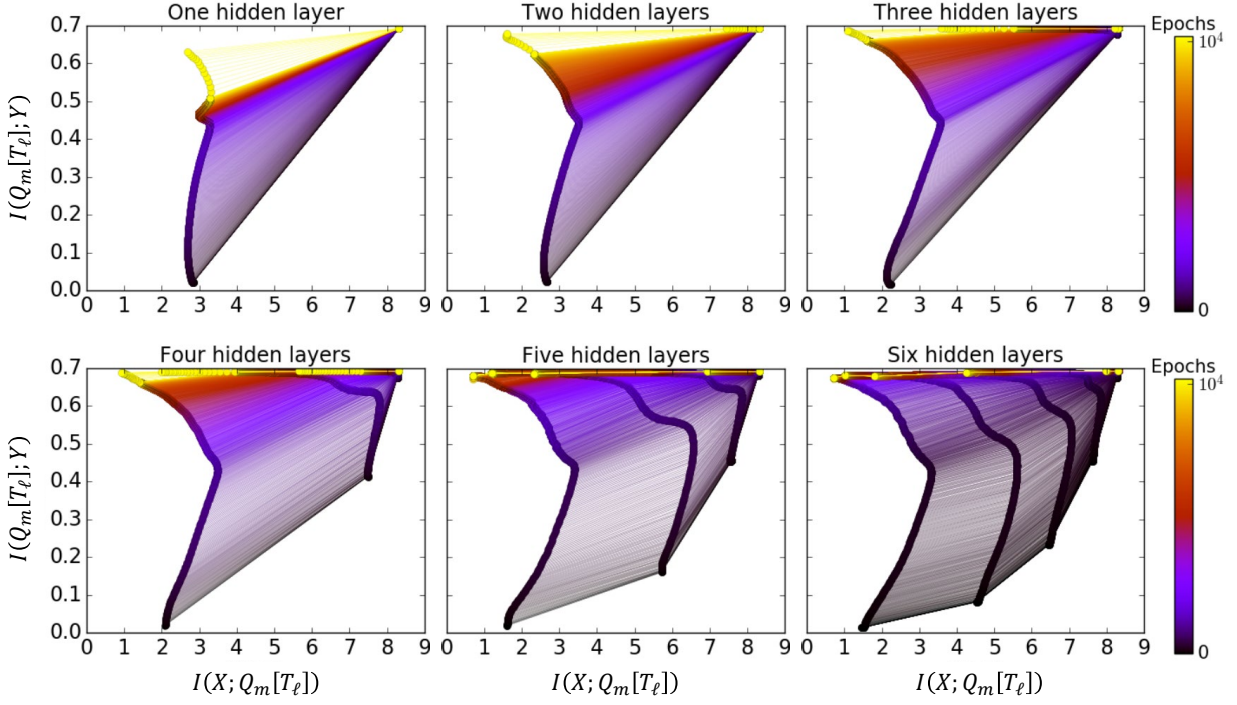
Fig. 8 [Fig. 5 from [6]]: IP dynamics for 6 neural networks with increasing depths. Each network has 12 input and 2 output units. The width of the hidden layers start from 12 and reduces by 2 with each added layer.

the shallowest network (top-left subfigure) does not reach that value throughout the entire training process. The compression of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$, for $\ell = 1, \ldots, L$, is also faster and more pronounced in deeper architectures. Thus, both IP phases accelerate as a result of more hidden layers. Furthermore, compression of a preceding $I\big(X; \mathsf{Q}_m[T_{\ell-1}]\big)$ seems to push $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ of the next layer to a smaller value.[13] We also note that IP values of shallow layers, even in deeper architectures, stay almost unchanged throughout training. These values are roughly $\big(H(X), I(X;Y)\big)$, which correspond to $\big(I(X;T_\ell), I(T_\ell;Y)\big)$ when $T_\ell = \phi_\ell(X)$ is a bijection (as a mapping from $\mathcal{X}_n$ to $\phi_\ell(\mathcal{X}_n) := \{\phi_\ell(x) : \ x \in \mathcal{X}_n\}$). This fact will be revisited and explained in the Section V.

The interpretation of Fig. 8 in [6] (Fig. 5 therein) synonymizes 'high $I\big(\mathsf{Q}_m[T_L];Y\big)$' and 'good generalization'. We refrain from this wording since there is no known rigorous connection between generalization error and $I\big(\mathsf{Q}_m[T_L];Y\big)$ (or $I(T_L;Y)$ for that matter).[14] If one adopts the IP values prescribed by the IB problem as the figure of merit, then indeed, deeper architectures enhance the corresponding dynamics. However, the implications of the IB theory on how depth affects generalization error or sample complexity remains unclear.

*8) Concluding remarks:* Several claims from [6] were not discussed here in detail. Our summary focused on the empirical observations of that work. Beyond these, [6] included heuristic arguments about how: (i) compression of representation is causally related to the outstanding generalization performance DNNs enjoy in practice; (ii) learned

---

[13]For the true (unquantized) $I(X;T_\ell)$ terms, $\ell = 1, \ldots, L$, their reduction with larger $\ell$ values is a consequence of the DPI. The DPI, however, does not hold for the quantized mutual information proxies.

[14]Some connections between other information measures and generalization error are known – see, e.g., [61].

hidden layers lie on (or very close to) the IB theoretical bound (3) for different $\beta$ values; (iii) the corresponding 'encoder' and 'decoder' maps satisfy the IB self-consistent equations (4); (iv) the effect of depth on the diffusion phase of SGD dynamics; and more. The reader is referred to the original paper for further details.

In sum, the IB theory for DL [5], [6] proposed a novel perspective on DNNs and their learning dynamics. At its core, the theory aims to summarize each hidden layer into the mutual information pair $\big(I(X;T_\ell), I(T_\ell;Y)\big)$, $\ell = 1, \ldots, L$, and study the systems through that lens. As the true mutual information terms degenerate in deterministic networks, [6] adopted the quantized versions $I\big(X;\mathsf{Q}_m[T_L]\big)$ and $I\big(\mathsf{Q}_m[T_L];Y\big)$ as figures of merit in their stead. While doing so created a gap between the empirical study of [6] and the theoretical IB problem, the evolution of the quantized terms throughout training revealed remarkable empirical trends. These observations were collected to a new information-theoretic paradigm to explain DL, which inspired multiple follow-up works. The next section describes some of these works and how they corroborate or challenge claims made in [5], [6].

## V. Revisiting the Information Bottleneck Theory for Deep Learning

Since [5], [6], the IB problem and the IP dynamics became sources of interest in DL research. Many works followed up both on the empirical findings and theoretical reasoning in [5], [6]; a nonexhastive list includes [7], [9], [11]–[15], [17], [36]. This section focuses on the empirical study conducted in [9], how quantization misestimates mutual information in deterministic networks (in light of the observations from [13], [15]), and the relation between compression and clustering revealed in [13].

### A. Revisiting the Empirical Study

In [9], an empirical study aiming to test the main claims of [6] was conducted. They focused on the two phases of training in the IP, the relation between compression of $I(X;T_\ell)$ and generalization, as well as the link between stochasticity in gradients and compression. Through a series of experiments, [9] produced counter examples to all these claims. We note that [9] employed methods similar to [6] for evaluating $I(X;T_\ell)$ and $I(T_\ell;Y)$. As subsequently explained in Sections V-B and V-C, these methods fail to capture the true mutual information values, which are vacuous in deterministic DNNs (i.e., when for fixed parameters, the DNN's output is a deterministic function of the input). Though the authors of [9] were aware of the this issue (see p. 5 and Appendix C therein), they adopted the methodology of [6] in favor of empirical comparability.[15]

*1) Information plane dynamics and the effect of activation function:* It was argued in [6] that the fitting and compression phases seen in the IP are inherent to DNN classifiers trained with SGD. On the contrary, [9] found that the IP profile of a DNN strongly depends on the employed nonlinear activation function. Specifically, double-sided saturating nonlinearities like $\tanh$ (used in [6]) or $\mathrm{sigmoid}$ yield a compression phase, but linear or single-sided saturating nonlinearities like $\mathrm{ReLU}(x) := \max\{0, x\}$ do not compress representations. The experiment showing

---

[15]Several methods for computing IP trajectories were employed in [9]. While they mostly used the quantization-based method of [6] (on which we focus herein), they also examined replacing $\mathsf{Q}_m[T_\ell]$ with $T_\ell + \mathcal{N}(\mathbf{0}, \sigma^2 \mathrm{I}_{d_\ell})$ (although no Gaussian noise was explicitly injected to the actual activations), as well as estimating mutual information from samples via $k$ nearest neighbor (kNN) [62] and kernel density estimation (KDE) [63] techniques.

this compared the IP dynamics of the same network once with $\mathrm{tanh}$ nonlinearities and then with $\mathrm{ReLU}$. Two tasks were examined: the synthetic experiment from [6] and MNIST classification. The architecture for the former was the same as in [6], while the latter used a 784–1024–20–20–20–10 fully connected architecture. The last layer in both $\mathrm{ReLU}$ networks employs $\mathrm{sigmoid}$ nonlinearities; all other neurons use $\mathrm{ReLU}$. Fig. 9 (reprinted from [9, Fig. 1]) shows the IP dynamics of all four models.

First note that top-right subfigure reproduces the IP dynamics from [6] (compare to Fig. 6 herein). The $\mathrm{ReLU}$ version of that network, however, does not exhibit compression, except in its last sigmoidal layer. Instead, the mutual information $I\big(X; \mathsf{Q}_m[t_\ell]\big)$ seems to monotonically increase in all $\mathrm{ReLU}$ layers. This stands in accordance with the argument of [9] that double-sided saturating nonlinearities can cause compression, while single-sided ones cannot. The same effect is observed for the MNIST network, by comparing its $\mathrm{tanh}$ version at the bottom-right with its $\mathrm{ReLU}$ version at bottom-left. Similar results were also observed for $\mathrm{soft\text{-}sign}$ (double-sided saturation) versus $\mathrm{soft\text{-}plus}$ (single-sided saturation) activations, given in Appendix B of [9]. They concluded that the choice of activation function has a significant effect on IP trajectories. In particular, the compression phase of training is caused by double-sided saturating activations, as opposed to being inherent to the learning dynamics.

*2) Relation between compression and generalization:* A main argument of [6] is that compression of representation is key for generalization. Specifically, by decreasing $I(X; T_\ell)$ (while keeping $I(T_\ell; Y)$ high), the DNN sheds information about $X$ that is irrelevant for learning $Y$, which in turn mitigates overfitting and promotes generalization.

To test this claim, [9] first considered deep linear networks [64] and leveraged recent results on generalization dynamics in the student-teacher setup [65], [66]. In this setup, one neural network ('student') learns to approximate the output of another ('teacher'). The linear student-teacher framework with Gaussian input allows exact computation of both the generalization error and the input mutual information[16] for a nontrivial task with interesting structure.

The exact setup is the following. Let $X \sim \mathcal{N}(0, \frac{1}{d}\mathrm{I}_d)$ be an isotropic $d$-dimensional Gaussian and set the teacher network output as $Y = B_0^\top X + N_0$, where $B_0 \in \mathbb{R}^d$ is the weight (column) vector and $N_0 \sim \mathcal{N}(0, \sigma_0^2)$ is an independent Gaussian noise. The teacher specifies a stochastic rule $P_{Y|X}$ that the student network needs to learn. Specifically, the student linear DNN is trained on a dataset generated by the teacher. Denoting the layers of the student network by $\{\mathsf{A}_\ell\}_{\ell=1}^L$, its reproduction of $Y$ is $\hat{Y} := \mathsf{A}_L \mathsf{A}_{L-1} \cdots \mathsf{A}_1 X$. Note that $X, Y, \hat{Y}$ and all the internal representation $T_\ell$ of the student network are jointly Gaussian. This allows analytic computation of generalization error and mutual information terms for any fixed parameters (i.e., at each epoch) – see Eq. (6)-(7) of [9].

Leveraging this fact, Fig. 3 of [9] (not shown here) compared the IP dynamics of the student linear network with a single hidden layer to the training and test errors. While the network generalized well, not compression of $I(X; T_1)$ was observed. Instead, the linear network qualitatively behaved like a $\mathrm{ReLU}$ network, presenting monotonically increasing mutual information trajectories. Building on the study of linear networks in [66], the authors then matched the size of the student network to the number of samples, causing it to severely overfit the data. Despite now having a large generalization gap, the IP trajectory of the network did not change, still showing monotonically increasing

---

[16]To be precise, the computed quantity is $I(X; T_\ell + Z)$, for some independent Gaussian noise $Z$. The addition of $Z$ is needed, as without it $I(X; T_\ell) = \infty$ because $X$ is Gaussian and $T_\ell$ is a linear deterministic function thereof.
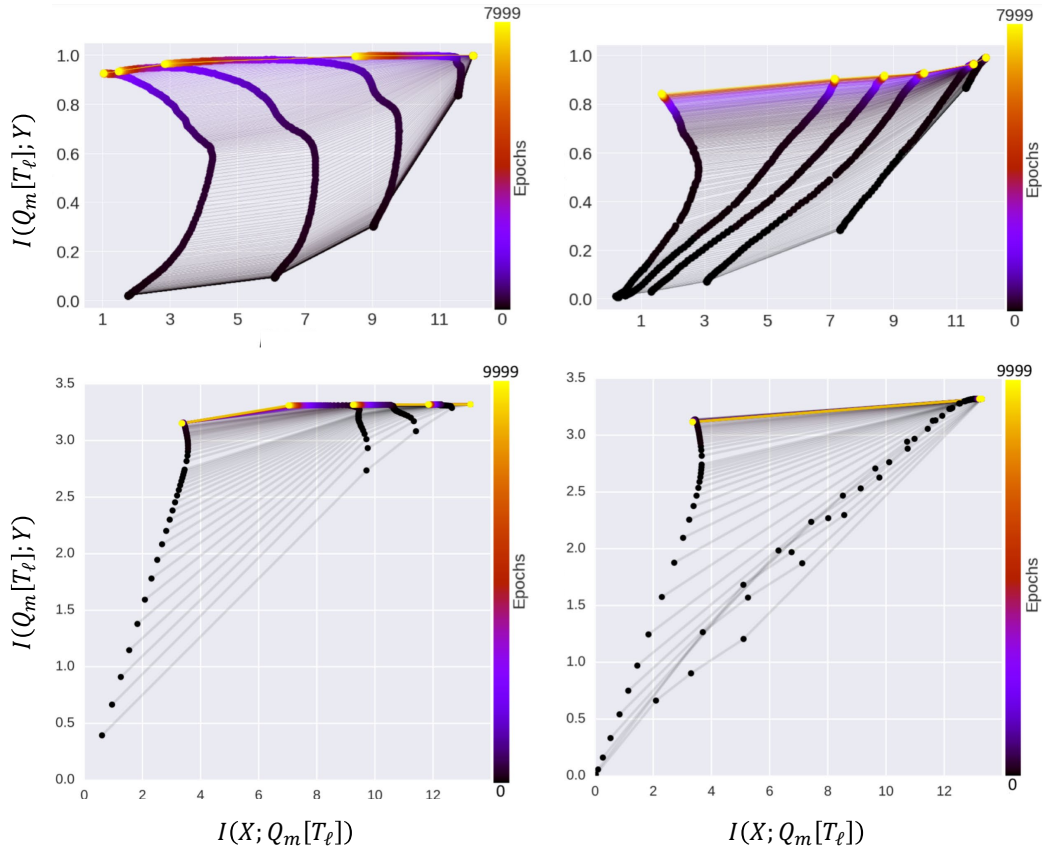
Fig. 9 [Fig. 1 from [9]]: IP dynamics for $\tanh$ and ReLU DNN classifiers (left and right, respectively) for the synthetic task from [6] and MNIST (top and bottom, respectively). The last layer in both ReLU networks (bottom row) has sigmoidal activations. While $\tanh$ networks exhibit compression of representation, no compression is observed in any of the ReLU layers.

$I(X; T_1)$ with epochs [9, Figs. 4A-4B]. This produced an example of two networks with the same IP profile (no compression) but widely different generalization performance.

Similar results were presented for nonlinear networks. The authors of [9] retrained the $\tanh$ network from [6] on the synthetic classification task therein, but using only $30\%$ of the data. This network significantly overfitted the data, resulting in a high generalization gap (left panel in Fig. 10, which is reprinted from [9, Fig. 4]). Nonetheless, its IP trajectories exhibit compression of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ (right panel of Fig. 10). A compression phase also occurs for the original network trained on $85\%$ of the data (Figs. 6 and 9), whose test performance is much better. The main difference between the IP profiles is that the overfitted network has lower $I\big(\mathsf{Q}_m[T_\ell]; Y\big)$ values, but $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ compressed in both cases.

Together, the linear and nonlinear examples dissociated compression of $I(X; T_\ell)$ (as approximated by $I\big(X; \mathsf{Q}_m[T_\ell]\big)$) and generalization: networks that compress may or may not generalize (nonlinear example), and the same applies for networks that do not compress (linear example). This suggest that the connection between compression and generalization, if exists, is not a simple causal relation. Furthermore, based on Fig. 10, a direct
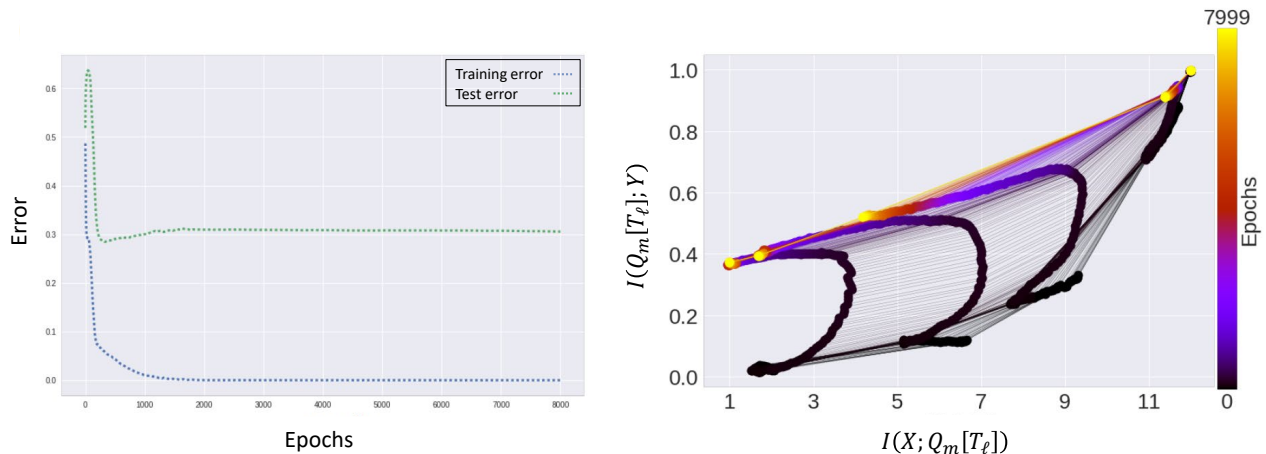
Fig. 10 [Fig. 4 from [9]]: Train/test errors and IP dynamics for a $\mathrm{tanh}$ DNN trained on $30\%$ of the synthetic dataset from [6]. The network has high generalization gap, yet exhibits compression of $I\big(X;\mathsf{Q}_m[T_\ell]\big)$ in the IP.

link between $I(X;T_\ell)$ and generalization that does not involve $I(T_\ell;Y)$, seems implausible.

*3) Stochastic gradients drive compression:* The third claim of [6] revisited by [9] is that the randomness in SGD causes its diffusion phase, which in turn drives compression (see Section IV-B6). According to this rationale, training a DNN with batch gradient decent (BGD), which updates weights using the gradient of the total error across all examples, should not induce diffusion nor result in compression.

The authors of [9] trained the $\mathrm{tanh}$ and $\mathrm{ReLU}$ networks on the synthetic task from [6] with SGD and BGD. Fig. 5[17] therein compares the obtained IP dynamics, showing no noticeable difference between the two training methodologies. Both SGD- and BGD-trained $\mathrm{tanh}$ networks present compression, while neither of the $\mathrm{ReLU}$ networks does. IP trajectories generated by SGD are shown in Fig. 9; the BGD trajectories, though not presented here, looks very much alike. Interestingly, [9] also examined the gradient's high-to-low SNR phase transition observed in [6] in multiple experiments. They found that it occurs every time, regardless of the employed training method, architecture or nonlinearity type, suggesting it is a general phenomenon inherent to DNN training, though not causally related to compression of representation.

## B. Mutual Information Misestimation in Deterministic Networks

As discussed in Section IV-B3, $I(X;T_\ell)$ and $I(T_\ell;Y)$ degenerate in deterministic DNNs with strictly monotone nonlinearities. The network from [6] (also studied in [9]), whose IP dynamics are shown in Figs. 6, 9 and 10, is deterministic with $\mathrm{tanh}$ nonlinearities. Therefore, $I(X;T_\ell) = H(X)$ and $I(T_\ell;Y) = I(X;Y)$ are constant, independent of the network parameters, under the $X \sim \mathsf{Unif}(\mathcal{X}_n)$ model used in these works (see Section IV-B3). As such, evaluating these information terms via quantization, noise injection, or estimation from samples is ill-advised. Yet, their estimates, as computed, e.g., in [6], [9], fluctuate with training epochs (that change parameter

---

[17]We believe that Figs. 5A and 5B in [9] should be switched to correspond to their captions; compare Fig. 5B and Fig. 1A.
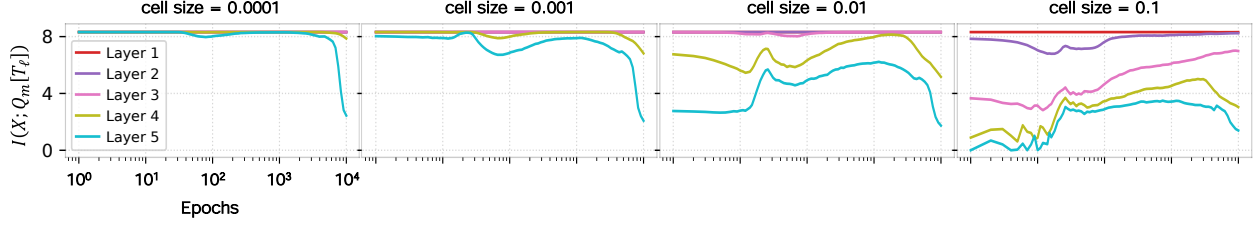
Fig. 11 [Fig. 1 from [13]]: $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ vs. epochs for different quantization cell (bin) sizes and the model in [6], where $X$ is uniformly distributed over a $2^{12}$-sized empirical dataset. The curves converge to $H(X) = \ln(2^{12}) \approx 8.3$ for small bins.

values) presenting IP dynamics. These fluctuations must therefore be a consequence of estimation errors rather than changes in mutual information. Indeed, quantization or noise injection are employed in [6], [9] as means for performing the measurements, but they are not part of the actual network. As explained next, this creates a mismatch between the measurement model and the system being analyzed.

To simplify discussion, we focus on $I(X; T_\ell)$ and its quantization-based approximation (similar reasoning applies for $I(T_\ell; Y)$ and to measurements via noise injection). Note that after quantization, the mapping from $X$ to $\mathsf{Q}_m[T_\ell]$ is no longer injective. This is since (distinct) representations $\phi_\ell(x)$ and $\phi_\ell(x')$, for $x, x' \in \mathcal{X}_n$, that lie sufficiently close to one another are mapped to the same value ('quantization cell' or 'bin') under $\mathsf{Q}_m$. The distances between representations are captured by the feature map $\phi_\ell(\mathcal{X}_n)$, which depends on the networks parameters. As the feature map and the quantization resolution $m$ determine the distribution of $\mathsf{Q}_m[T_\ell]$, the dependence of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ on the network parameters becomes clear. This results in a parameter-dependent estimate $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ of the parameter-independent quantity $I(X; T_\ell)$, which is undesirable. Strictly speaking, there is nothing here to estimate: since $X \sim \mathsf{Unif}(\mathcal{X}_n)$ (as assumed in [6], [9]) and $\phi_\ell$ is injective from $\mathcal{X}_n$ to $\phi_\ell(\mathcal{X}_n)$, the true $I(X; T_\ell)$ value is $\log n$.

Recalling (24), we expect that $I\big(X; \mathsf{Q}_m[T_\ell]\big) \to H(X) = \log n$ as $m \to 0$. For nonnegligible $m > 0$, the value of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ strongly depends on the quantization parameter. However, since no quantization is present in the actual network (see (21)), the value of $m$ is arbitrary and chosen by the user. Thus, the measured $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ reveal more about the estimator and the dataset than about the true mutual information, which is a global property of the underlying joint distribution. The dependence of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ on $m$ is illustrated in Fig. 11 (reprinted from [13, Fig. 1]), showing that widely different profiles can be obtained by changing $m$. The leftmost subfigure also shows how the quantized mutual information approaches $H(X) = \log n$ as $m$ shrinks. Recalling that the dataset size in [6] is $n = 2^{12}$ and taking logarithms to the base of $e$, we see that for $m = 10^{-4}$, $I\big(X; \mathsf{Q}_m[T_\ell]\big) \approx \log 2^{12} \approx 8.3$, for all hidden layers $\ell = 1, \ldots, 5$ and across (almost) all epochs of training.

In summary, while $I\big(X; \mathsf{Q}_m[T_\ell]\big)$ fails to capture the true (vacuous) mutual information value, it encodes nontrivial information about the feature maps. The compression phase observed in [6] is in fact a property of $I\big(X; \mathsf{Q}_m[T_\ell]\big)$, rather than $I(X; T_\ell)$, driven by the evolution of internal representations. Notably, for any of the quantization resolutions shown in Fig. 11, at least for the last hidden layer, $I\big(X; \mathsf{Q}_m[T_5]\big)$ always undergoes a decrease towards

the end of training. This raises the questions of what is the underlying phenomenon inside the internal representation spaces that causes this behavior. As was shown in [13], the answer turns out to be clustering. The next section focuses on the methodology and the results of [13] that led to this conclusion.

### C. Noisy Deep Networks and Relation to Clustering of Representations

While $I(X;T_\ell)$ is vacuous in deterministic DNNs, the compression phase that its estimate $I(X;\mathsf{Q}_m[T_\ell])$ undergoes during training seems meaningful. To study this phenomenon, [13] developed a rigorous framework for tracking the flow of information in DNNs. Specifically, they proposed an auxiliary 'noisy' DNN setup, under which the map $X \mapsto T_\ell$ is a *stochastic parameterized channel*, whose parameters are the network's weights and biases. This makes $I(X;T_\ell)$ over such networks, a meaningful system-dependent quantity. The authors of [13] then proposed a provably accurate estimator of $I(X;T_\ell)$ and studied its evolution during training. Although not covered in detail herein, additional ways to make the IB non-vacuous (beyond noising the activations) include: adding noise to the weights [7], [67], changing the objective [35], and changing the information measure [14], [36].

*1) Noisy DNNs:* The definition of a noisy DNN replaces $T_\ell = \phi_\ell(X)$ (see (21)) with $T_\ell^{(\sigma)} := T_\ell + Z_\ell^{(\sigma)}$, where $\{Z_\ell^{(\sigma)}\}_{\ell=1}^L$ are independent isotropic $d_\ell$-dimensional Gaussian vectors of parameter $\sigma > 0$, i.e., $Z_\ell^{(\sigma)} \sim \mathcal{N}_\sigma := \mathcal{N}(0,\sigma^2 \mathrm{I}_d)$. In other words, i.i.d. Gaussian noise is injected to the output of each hidden neuron. The noise here is intrinsic to the system, i.e, the network is trained with the noisy activation values. This stands in contrast to [6] and [9], where binning or noise injection were merely a part of the measurement (of mutual information) model. Intrinsic noise as in [13] ensures that the characteristics of *true* information measures over the network are tied to the network's dynamics and the representations it is learning. Furthermore, the isotropic noise model relates $I(X;T_\ell^{(\sigma)})$ to $I(X;\mathsf{Q}_m[T_\ell])$ when $\sigma$ is of the order of the quantization cell side length. This is important since it is the compression of the latter that was observed in preceding works.

To accredit the noisy DNN framework, [13] empirically showed that it forms a reasonable proxy of deterministic networks used in practice. Namely, when $\sigma > 0$ is relatively small (e.g., of the order of $10^{-2}$ for a $\mathtt{tanh}$ network), it was demonstrated that noisy DNNs not only perform similarly to deterministic ones, but also that the representations learned by both systems are closely related.

*2) Mutual information estimation:* Adopting $\left(I(X;T_\ell^{(\sigma)}), I(T_\ell^{(\sigma)};Y)\right)$ as the figure of merit in noisy DNNs, one still faces the task of evaluating these mutual information terms. Mutual information is a functional of the joint distribution of the involved random variables. While, $T_\ell = \phi_\ell(X) + Z_\ell^{(\sigma)}$ and $\phi_\ell$ is specified by the (known) DNN model, the data-label distribution $P_{X,Y}$ is unknown and we are given only the dataset $\mathcal{D}_n$. Statistical learning theory treats the elements of $\mathcal{D}_n$ as i.i.d. samples from $P_{X,Y}$. Under this paradigm, evaluating the mutual information pair of interest enters the realm of statistical estimation [17], [62], [68]–[74]. However, mutual information estimation from high-dimensional data is a notoriously difficult [75]. Corresponding error convergence rates (with $n$) in high-dimensional settings are too slow to be useful in practice.

Nevertheless, by exploiting the known distribution of the injected noise, [13] proposed a rate-optimal estimator of $I(X;T_\ell^{(\sigma)})$ that scales significantly better with dimension than generic methods (such as those mentioned above).

This was done by developing a forward-pass sampling technique that reduced the estimation of $I\big(X;T_\ell^{(\sigma)}\big)$ to estimating differential entropy under Gaussian noise as studied in [17] (see also [33], [34], [76]). Specifically, the latter considered estimating the differential entropy $h(S+Z) = h(P * \mathcal{N}_\sigma)$ based on 'clean' samples of $S \sim P$, where $P$ belongs to a nonparametric distribution class, and knowledge of the distribution of $Z \sim \mathcal{N}_\sigma$, which is independent of $S$. Here $(P*Q)(\mathcal{A}) := \int \int \mathbb{1}_\mathcal{A}(x+y)\,\mathrm{d}P(x)\,\mathrm{d}Q(y)$ is the convolution of two probability measures $P$ and $Q$, and $\mathbb{1}_\mathcal{A}$ is the indicator function of $\mathcal{A}$.

The reduction of mutual information estimation to estimating $h(S+Z)$ uses the decomposition

$$I\big(X;T_\ell^{(\sigma)}\big) = h(T_\ell^{(\sigma)}) - \int h\big(T_\ell^{(\sigma)}|X=x\big)\,\mathrm{d}P_X(x), \tag{25}$$

along with the fact that $T_\ell^{(\sigma)} = T_\ell + Z_\ell^{(\sigma)}$, where $T_\ell = \sigma(\mathrm{A}_\ell \phi_{\ell-1}(X) + b_\ell)$, is easily sampled via the forward-pass of the network (see [17, Section IV] and [13, Section III] for more details).[18] Building on [17], the employed estimator for $h(P * \mathcal{N}_\sigma)$ was $\hat{h}(S^n, \sigma) := h(\hat{P}_{S^n} * \mathcal{N}_\sigma)$, where $\hat{P}_{S^n} := \frac{1}{n}\sum_{i=1}^n \delta_{S_i}$ is the empirical distribution of the i.i.d. sample $S^n := (S_1, \ldots, S_n)$ of $S \sim P$. The estimation risk of $\hat{P}_{S^n}$ over the class of all compactly supported [17, Theorem 2] or sub-Gaussian [17, Theorem 3] $d$-dimensional distributions $P$ scales as $c^d n^{-\frac{1}{2}}$, for an explicitly characterized numerical constant $c$. Matching impossibility results showed that this rate is optimal both in $n$ and $d$. By composing multiple differential entropy estimators (for $h\big(T_\ell^{(\sigma)}\big)$ and each $h\big(T_\ell^{(\sigma)}|X=x\big)$, $x \in \mathcal{D}_n$), an estimator $\hat{I}(\mathcal{D}_n, \sigma)$ of $I\big(X;T_\ell^{(\sigma)}\big)$ was constructed. Its absolute-error estimation risk over, e.g., DNNs with $\tanh$ nonlinearities scales as follows.

**Proposition 1 (Mutual Information Estimator [17])** *For the described estimation setup, we have*

$$\sup_{P_X} \mathbb{E}\left|I(X;T) - \hat{I}_{\mathsf{Input}}\left(X^n, \hat{h}, \sigma\right)\right| \leq \frac{8c^{d_\ell} + d_\ell \log\left(1 + \frac{1}{\sigma^2}\right)}{4\sqrt{n}},$$

*where $c$ is a constant independent of $n$, $d$ and $P_X$, which is explicitly characterized in [17, Equation (61)].*

Notably, the right-hand side depends exponentially on dimensions, which limits the dimensionality of experiments for which the bound in non-vacuous. This limitation is inherent to the considered estimation problem, as [17] proved that the sample complexity of *any* good estimator depends exponentially on $d_\ell$.

*3) Empirical study and relation to clustering:* The developed toolkit enabled the authors of [13] to accurately track $I\big(X;T_\ell^{(\sigma)}\big)$ during training of (relatively small) noisy DNN classifiers. For the synthetic experiment from [6] with $\tanh$ activations, [13] empirically showed that $I\big(X;T_\ell^{(\sigma)}\big)$ indeed undergoes a long-term compression phase (see Fig. 12, which is reprinted from [13, Fig. 5(a)]). To reveal the geometric mechanism driving this phenomenon, they related $I\big(X;T_\ell^{(\sigma)}\big)$ to data transmission over AWGN channels. The mutual information $I\big(X;T_\ell^{(\sigma)}\big)$ can be viewed as the aggregate number of bits (reliably) transmittable over an AWGN channel using an input drawn from the latent representation space. As training progresses, the hidden representations of equilabeled inputs cluster together, becoming increasingly indistinguishable at the channel's output, thereby decreasing $I\big(X;T_\ell^{(\sigma)}\big)$. To test

---

[18]Samples for estimating the conditional differential entropy terms $h\big(T_\ell^{(\sigma)}|X=x\big)$, for $x \in \mathcal{X}$, are obtained by feeding the network with $x$ multiple times and reading $T_\ell^{(\sigma)}$ values corresponding to different noise realizations.
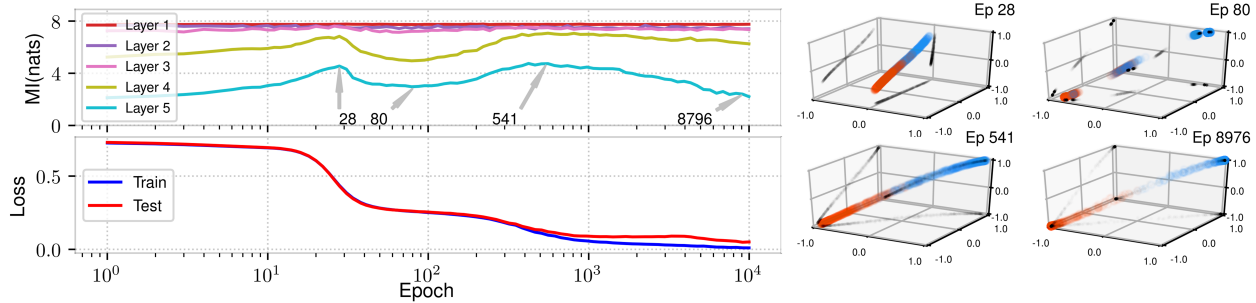
Fig. 12: $I\left(X; T_\ell^{(\sigma)}\right)$ evolution and training/test losses across training for a noisy version of the DNN from [6]. Scatter plots show the input samples as represented by the $5^{\text{th}}$ hidden layer (colors designate classes) at the arrow-marked epochs.

this empirically, [13] contrasted $I\left(X; T_\ell^{(\sigma)}\right)$ trajectories with scatter plots of the feature map $\phi_\ell(\mathcal{X}_n)$. Remarkably, compression of mutual information and clustering of latent representations clearly corresponded to one another. This can be seen in Fig. 12 by comparing the scatter plots on the right to the arrow-marked epochs in the information flow trajectory on the left.

Next, [13] accounted for the observations from [6], [9] of compression in deterministic DNNs. It was demonstrated that while the quantization-based estimator employed therein fails to capture the true (constant/infinite) mutual information, it does serve as a measure of clustering. Indeed, for a deterministic network we have

$$I\left(X; \mathsf{Q}_m[T_\ell]\right) = H\left(\mathsf{Q}_m[T_\ell]\right),$$

where $\mathsf{Q}_m$ partitions the dynamic range (e.g., $[-1,1]^{d_\ell}$ for a $\tanh$ layer) into $m^{d_\ell}$ cells. When hidden representations are spread out, many cells are non-empty, each having some small positive probability mass. Conversely, for clustered representations, the distribution is concentrated on a small number of cells, each with relatively high probability. Recalling that the uniform distribution maximizes Shannon entropy, we see that reduction in $H\left(\mathsf{Q}_m[T_\ell]\right)$ corresponds to tighter clusters in the latent representation space.

The results of [13] identified the geometric clustering of representations as the fundamental phenomenon of interest, while elucidating some of the machinery DNNs employ for learning. Leveraging the clustering perspective, [13] also provided evidence that compression and generalization may *not* be causally related. Specifically, they constructed DNN examples that generalize better when tight clustering is actively suppressed during training (compare Figs. 5(a) and 5(b) from [13]). This again showed that the relation between compression of mutual information and generalization is not a simple one, warranting further study. More generally, there seems to be a mismatch between the IB paradigm and common DL practice that mainly employs deterministic NNs, under which information measures of interest tend to degenerate. It remains unclear how to bridge this gap.

## VI. Summary and Concluding Remarks

This tutorial surveyed the IB problem, from its information-theoretic origins to the recent impact it had on ML research. After setting up the IB problem, we presented its relations to MSSs, discussed operational interpretations,

and covered the Gaussian IB setup. Together, these components provide background and context for the recent framing of IB as an objective/model for DNN-based classification. After describing successful applications of the IB framework as an objective for learning classifiers [8], [31] and generative models [32], we focused on the IB theory for DL [5], [6] and the active research area it inspired. The theory is rooted in the idea that DNN classifiers inherently aim to learn representations that are optimal in the IB sense. This novel perspective was combined in [6] with an empirical case-study to make claims about phase transitions in optimization dynamics, computational benefits of deep architectures, and relations between generalization and compressed representations. Backed by some striking empirical observations (though only for a synthetic classification task), the narrative from [6] ignited a series of followup works aiming to test its generality.

We focused here on works that contributed to different aspect of modern IB research. Our starting point is [9], that revisited the observations from [6] with a thorough empirical analysis. The experiments from [9] were designed to test central aspects of the IB theory for DL; their final conclusion was that the empirical findings from [6] do not hold in general. We then examined theoretical facets applying the (inherently stochastic) framework of IB to deterministic DL systems. Covering contributions from [15] and [13], caveats in measuring information flows in deterministic DNNs were explained. In a nutshell, key information measures degenerate over deterministic networks, becoming either constant (independent of the network's parameters) or infinite, depending on the modeling of the input feature. Either way, such quantities are vacuous over deterministic networks.

We then described the remedy proposed in [13] to the 'vacuous information measures' issue. Specifically, that work presented an auxiliary stochastic DNN framework over which the considered mutual information terms are meaningful and parameter-dependent. Using this auxiliary model, they demonstrated that compression of $I(X; T)$ over the course of training, is driven by clustering of equilabeled samples in the representation space of $T$. It was then shown that a similar clustering process occurs during training of deterministic DNN classifiers. Circling back to the original observations of compression [6] (see also [9]), the authors of [13] demonstrated that the measurement techniques employed therein in fact track clustering of samples. This clarified the geometric phenomena underlying the compression of mutual information during training. Still, many aspects of the IB theory for DL remain puzzling, awaiting further exploration.

REFERENCES

[1] T. Hastie, R. Tibshirani, and J. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.

[2] V. N. Vapnik and A. Y. Chervonenkis, "On the uniform convergence of relative frequencies of events to their probabilities," *Rep. Academy Sci. USSR*, no. 4, p. 181, 1968.

[3] M. L. Minsky and S. A. Papert, "Perceptrons," *MIT press*, p. 248, 1969.

[4] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.

[5] N. Tishby and N. Zaslavsky, "Deep learning and the information bottleneck principle," in *Proceedings of the Information Theory Workshop (ITW)*, Jerusalem, Israel, Apr.-May 2015, pp. 1–5.

[6] R. Shwartz-Ziv and N. Tishby, "Opening the black box of deep neural networks via information," 2017, arXiv:1703.00810 [cs.LG].

[7] A. Achille and S. Soatto, "On the emergence of invariance and disentangling in deep representations," *Journal of Machine Learning Research*, vol. 19, pp. 1–34, 2018.

[8] ——, "Information dropout: Learning optimal representations through noisy computation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 12, pp. 2897–2905, Jan. 2018.

[9] A. M. Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, B. D. Tracey, and D. D. Cox, "On the information bottleneck theory of deep learning," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.

[10] M. Gabrié, A. Manoel, C. Luneau, J. Barbier, N. Macris, F. Krzakala, and L. Zdeborová, "Entropy and mutual information in models of deep neural networks," *arXiv preprint arXiv:1805.09785*, 2018.

[11] S. Yu, R. Jenssen, and J. C. Principe, "Understanding convolutional neural network training with information theory," *arXiv preprint arXiv:1804.06537*, 2018.

[12] H. Cheng, D. Lian, S. Gao, and Y. Geng, "Evaluating capability of deep neural networks for image classification via information plane," in *European Conference on Computer Vision (ECCV-2018)*, Munich, Germany, September 2018, pp. 168–182.

[13] Z. Goldfeld, E. van den Berg, K. Greenewald, I. Melnyk, N. Nguyen, B. Kingsbury, and Y. Polyanskiy, "Estimating information flow in neural networks," in *Proceedings of the International Conference on Machine Learning (ICML-2019)*, vol. Long Beach, California, USA, Jun. 2019.

[14] K. Wickstrøm, S. Løkse, M. Kampffmeyer, S. Yu, J. C. Principe, and R. Jenssen, "Information plane analysis of deep neural networks via matrix-based renyi's entropy and tensor kernels," *arXiv preprint arXiv:1909.11396*, 2019.

[15] R. A. Amjad and B. C. Geiger, "Learning representations for neural network-based classification using the information bottleneck principle," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Apr. 2019.

[16] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," in *Proceedings of the Allerton Conference on Communication, Control and Computing*, Monticello, Illinois, US, Sep. 1999, pp. 368–377.

[17] Z. Goldfeld, K. Greenewald, Y. Polyanskiy, and J. Weed, "Convergence of smoothed empirical measures with applications to entropy estimation," *arXiv preprint arXiv:1905.13576*, 2019.

[18] E. L. Lehmann and H. H. Scheffé, "Completeness, similar regions, and unbiased estimation: Part i," *Sankhyā: The Indian Journal of Statistics*, vol. 10, no. 4, pp. 305–340, Nov. 1950.

[19] S. Kullback and R. A. Leibler, "On information and sufficiency," *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, Mar. 1951.

[20] H. Witsenhausen and A. D. Wyner, "A conditional entropy bound for a pair of discrete random variables," *IEEE Trans. Inf. Theory*, vol. 21, no. 5, pp. 493–501, Sep. 1975.

[21] P. Gács and J. Körner, "Common information is far less than mutual information," *Prob. Contr. Inf. Theory*, vol. 2, no. 2, pp. 149–162, 1973.

[22] R. E. Blahut, "Computation of channel capacity and rate-distortion functions," *IEEE Trans. Inf. Theory*, vol. 18, pp. 460–473, Jul. 1972.

[23] S. Arimoto, "An algorithm for computing the capacity of arbitrary discrete memoryless channels," *IEEE Trans. Inf. Theory*, vol. 18, pp. 14–20, Jan. 1972.

[24] I. Estella-Aguerri and A. Zaidi, "Distributed variational representation learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, July 2019.

[25] G. Chechik, A. Globerson, N. Tishby, and Y. Weiss, "Information bottleneck for Gaussian variables," *Journal of Machine Learning Research*, vol. 6, no. Jan., pp. 165–188, 2005.

[26] I. Estella-Aguerri and A. Zaidi, "Distributed information bottleneck method for discrete and gaussian sources," in *Proceedings of the International Zurich Seminar on Information and Communication (IZS-2018)*, Zurich, Switzerland, February 2018, pp. 35–39.

[27] R. Dobrushin and B. Tsybakov, "Information transmission with additional noise," *IRE Trans. Inf. Theory*, vol. 8, no. 5, pp. 293–304, Sep. 1962.

[28] J. Wolf and J. Ziv, "Transmission of noisy information to a noisy receiver with minimum distortion," *IEEE Trans. Inf. Theory*, vol. 16, no. 4, pp. 406–411, Jul. 1970.

[29] A. Zaidi, I. Estella-Aguerri, and S. Shamai, "On the information bottleneck problems: Models, connections, applications and information theoretic views," *Entropy*, vol. 22, no. 2, p. 151, February 2020.

[30] N. Slonim and N. Tishby, "Agglomerative information bottleneck," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS-2000)*, 2000, pp. 617–623.

[31] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," in *Proceedings of the International Conference on Learning Representations (ICLR-2017)*, Toulon, France, Apr. 2017.

[32] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "$\beta$-VAE: learning basic visual concepts with a constrained variational framework," in *Proceedings of the International Conference on Learning Representations (ICLR-2019)*, New Orleans, Louisiana, USA, May 2017.

[33] Z. Goldfeld, K. Greenewald, Y. Polyanskiy, and Y. Wu, "Differential entropy estimation under gaussian noise," in *IEEE International Conference on the Science of Electrical Engineering (ICSEE-2018)*, Eilat, Israel, December 2018.

[34] Z. Goldfeld, K. Greenewald, J. Weed, and Y. Polyanskiy, "Optimality of the plug-in estimator for differential entropy estimation under Gaussian convolutions," in *IEEE International Symposium on Information Theory (ISIT-2019)*, Paris, France, July 2019.

[35] D. J. Strouse and D. J. Schwab, "The deterministic information bottleneck," *Neural computation*, vol. 29, no. 6, pp. 1611–1630, Jun. 2017.

[36] M. Cvitkovic and G. Koliander, "Minimal achievable sufficient statistic learning," in *Proceedings of the International Conference on Machine Learning (ICML-2019)*, vol. Long Beach, California, USA, Jun. 2019, pp. 1465–1474.

[37] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. New-York: Wiley, 2006.

[38] B. O. Koopman, "On distributions admitting a sufficient statistic," *Transactions of the American Mathematical society*, vol. 39, no. 3, pp. 399–409, May 1936.

[39] T. Berger and R. Zamir, "A semi-continuous version of the Berger-Yeung problem," *IEEE Trans. Inf. Theory*, vol. 45, no. 5, pp. 1520–1526, Jul. 1999.

[40] H. Hotelling, "The most predictable criterion," *Journal of Educational Psychology*, vol. 26, no. 2, p. 139, Feb. 1935.

[41] T. Berger, *Rate-distortion theory: A mathematical basis for data compression*, ser. Information and System Sciences Series. Englewood Cliffs, NJ,USA: Prentice-Hall, 1971.

[42] Y. Polyanskiy and Y. Wu, "Lecture notes on information theory," *Lecture Notes 6.441, Department of Electrical Engineering and Computer Science, MIT, Cambridge, MA, USA*, 2012–2017.

[43] T. Berger, "Multiterminal source coding," in *Information Theory Approach to Communications*, G. Longo, Ed., vol. 229. CISM Cource and Lecture, 1978, pp. 171–231.

[44] S. Y. Tung, "Multiterminal source coding," Ph.D. dissertation, Cornell University, Ithaca, NY, USA, May 1978.

[45] D. Slepian and J. Wolf, "Noiseless coding of correlated information sources," *IEEE Trans. Inf. Theory*, vol. 19, no. 4, pp. 471–480, Jul. 1973.

[46] T. A. Courtade and T. Weissman, "Multiterminal source coding under logarithmic loss," *IEEE Trans. Inf. Theory*, vol. 60, no. 1, pp. 740–761, Nov. 2013.

[47] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *Proceedings of the International Conference on Learning Representations (ICLR-2014)*, Banff, Canada, Apr. 2014.

[48] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.

[49] G. Pereyra, G. Tucker, J. Chorowski, Ł. Kaiser, and G. Hinton, "Regularizing neural networks by penalizing confident output distributions," in *Proceedings of the International Conference on Learning Representations (ICLR-2017)*, Toulon, France, Apr. 2017.

[50] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *Proceedings of the International Conference on Learning Representations (ICLR-2015)*, San Diego, California, USA, May 2015.

[51] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *Proceedings of the IEEE Symposium on Security and Privacy (SP-2017)*, San Jose, California, USA, May 2017, pp. 39–57.

[52] V. Mnih, N. Heess, and A. Graves, "Recurrent models of visual attention," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS-2014)*, Montréal, Canada, Dec. 2014, pp. 2924–2932.

[53] I. Csiszár, "On the dimension and entropy of order $\alpha$ of the mixture of probability distributions," *Acta Mathematica Hungarica*, vol. 13, no. 3-4, pp. 245–255, Sep. 1962.

[54] A. El Gamal and Y.-H. Kim, *Network Information Theory*. Cambridge University Press, 2011.

[55] N. Murata, "A statistical study of on-line learning," *Online Learning and Neural Networks. Cambridge University Press, Cambridge, UK*, pp. 63–92, 1998.

[56] J. Chee and P. Toulis, "Convergence diagnostics for stochastic gradient descent with constant learning rate," in *Proceedings of the International Conference of Artificial Intelligence and Statistics (AISTATS-2018)*, Lanzarote, Canary Islands, April 2018, pp. 1476–1485.

[57] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio, "On the number of linear regions of deep neural networks," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS-2014)*, Montréal, Canada, Dec. 2014, pp. 2924–2932.

[58] R. Eldan and O. Shamir, "The power of depth for feedforward neural networks," in *Conference on Learning Theory (COLT-2016)*, New York City, New York, USA, Jun. 2016, pp. 907–940.

[59] M. Telgarsky, "Benefits of depth in neural networks," in *Conference on Learning Theory (COLT-2016)*, New York City, New York, USA, Jun. 2016, pp. 1517–1539.

[60] T. Poggio, H. Mhaskar, L. Rosasco, B. Miranda, and Q. Liao, "Why and when can deep-but not shallow-networks avoid the curse of dimensionality: a review," *International Journal of Automation and Computing*, vol. 14, no. 5, pp. 503–519, Oct. 2017.

[61] A. Xu and M. Raginsky, "Information-theoretic analysis of generalization capability of learning algorithms," in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS-2017)*, Long Beach, California, USA, Dec. 2017, pp. 2524–2533.

[62] H. S. A. Kraskov and P. Grassberger, "Estimating mutual information," *Phys. rev. E*, vol. 69, no. 6, p. 066138, June 2004.

[63] A. Kolchinsky and B. D. Tracey, "Estimating mixture entropy with pairwise distances," *Entropy*, vol. 19, no. 7, p. 361, Jul. 2017.

[64] P. Baldi and K. Hornik, "Neural networks and principal component analysis: Learning from examples without local minima," *Neural networks*, vol. 2, no. 1, pp. 53–58, Jan. 1989.

[65] H. S. Seung, H. Sompolinsky, and N. Tishby, "Statistical mechanics of learning from examples," *Physical review A*, vol. 45, no. 8, p. 6056, Apr. 1992.

[66] M. S. Advani and A. M. Saxe, "High-dimensional dynamics of generalization error in neural networks," *arXiv preprint arXiv:1710.03667*, 2017.

[67] A. Achille and S. Soatto, "Where is the information in a deep neural network?" *arXiv preprint arXiv:1905.12213*, May 2019.

[68] L. Paninski, "Estimating entropy on $m$ bins given fewer than $m$ samples," *IEEE Trans. Inf. Theory*, vol. 50, no. 9, pp. 2200–2203, Sep. 2004.

[69] G. Valiant and P. Valiant, "Estimating the unseen: improved estimators for entropy and other properties," in *Advances in Neural Information Processing Systems*, 2013, pp. 2157–2165.

[70] Y. Han, J. Jiao, and T. Weissman, "Adaptive estimation of Shannon entropy," in *IEEE International Symposium on Information Theory (ISIT-2016)*, Hong Kong, China, Jun. 2015, pp. 1372–1376.

[71] Y. Han, J. Jiao, T. Weissman, and Y. Wu, "Optimal rates of entropy estimation over Lipschitz balls," *arXiv preprint arXiv:1711.02141*, Nov. 2017.

[72] M. Noshad, Y. Zeng, and A. O. Hero, "Scalable mutual information estimation using dependence graphs," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP-2019)*, Brighton, UK, May 2019, pp. 2962–2966.

[73] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, D. Hjelm, and A. Courville, "Mutual information neural estimation," in *Proceedings of the International Conference on Machine Learning (ICML-2018)*, Stockholm, Sweden, Jul. 2018, pp. 530–539.

[74] C. Chung, A. Al-Bashabsheh, H. P. Huang, M. Lim, D. S. H. Tam, and C. Zhao, "Neural entropic estimation: A faster path to mutual information estimation," *arXiv preprint arXiv:1905.12957*, 2019.

[75] L. Paninski, "Estimation of entropy and mutual information," *Neural Computation*, vol. 15, pp. 1191–1253, Jun. 2003.

[76] Z. Goldfeld and K. Kato, "Limit distribution for smooth total variation and $\chi^2$-divergence in high dimensions," *arXiv preprint arXiv:2002.01013*, February 2020.