
A Probabilistic Representation for Deep Learning: Delving into The Information Bottleneck Principle

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The Information Bottleneck (IB) principle has recently attracted great attention to
2 explaining Deep Neural Networks (DNNs), and the key is to accurately estimate the
3 mutual information between a hidden layer and dataset. However, some unsettled
4 limitations weaken the validity of the IB explanation for DNNs. To address these
5 limitations and fully explain deep learning in an information theoretic fashion, we
6 propose a probabilistic representation for deep learning that allows the framework
7 to estimate the mutual information, more accurately than existing non-parametric
8 models, and also quantify how the components of a hidden layer affect the mutual
9 information. Leveraging the probabilistic representation, we take into account the
10 back-propagation training and derive two novel Markov chains to characterize the
11 information flow in DNNs. We show that different hidden layers achieve different
12 IB trade-offs depending on the architecture and the position of the layers in DNNs,
13 whereas a DNN satisfies the IB principle no matter the architecture of the DNN.

14 1 Introduction

15 Deep learning [21] has already achieved great success in numerous applications. Deep Neural
16 Networks (DNNs), however, are still commonly viewed as ‘black boxes’ [32]. Considerable efforts
17 have been devoted to explaining the internal mechanism of DNNs from various perspectives, such as
18 mathematics [5, 14], statistics [16, 23, 28], computer vision [43, 25], *etc.* Recently, the Information
19 Bottleneck (IB) principle has attracted attention in opening the ‘black boxes’ of DNNs [35, 38].

20 Given a joint distribution $P(X, Y)$, the IB principle posits a random variable $T = f(X)$ obeying the
21 Markov chain $Y \rightarrow X \rightarrow T$ and optimizes T by the IB Lagrangian [37, 36]

$$\min_{P(T|X)} I(X; T) - \beta I(Y; T), \quad (1)$$

22 where $f(\cdot)$ is an arbitrary function, $I(\cdot; \cdot)$ denotes mutual information, and the Lagrange multiplier
23 $\beta > 0$ controls the IB trade-off between compressing the input X and preserving the information
24 of the label Y . In the seminal work [35], Tishby *et al.* manifest the IB trade-off in every layer of
25 DNNs $= \{\mathbf{x}; \mathbf{t}_1; \dots; \mathbf{t}_L; \hat{\mathbf{y}}\}$ via studying $I(X; T_i)$ and $I(Y; T_i)$, where T_i is the random variable of
26 the i th hidden layer \mathbf{t}_i . Especially, the authors ascribe DNN generalization to the compression [34].

27 In the context of deterministic DNNs, recent works reveal some limitations of the IB principle for
28 explaining DNNs. Amjad *et al.* argue that the IB principle becomes an ill-posed optimization problem
29 due to $I(X; T_i) = \infty$ [1], and Kolchinsky *et al.* demonstrate that not every layer of DNNs satisfies a
30 strict IB trade-off, *i.e.*, different layers only differ in $I(X; T_i)$ but $I(Y; T_i)$ keeps consistent in all
31 layers [17]. In addition, Saxe *et al.* experimentally show that the compression does not occur in
32 DNNs with non-saturating activation functions, *e.g.*, the popular ReLU function [33], and Goldfeld
33 *et al.* doubt the causality between the generalization of DNNs and the compression [11, 7]. These
34 unsettled limitations greatly weakens the validity of the IB explanations for DNNs.

35 The key to examining the IB principle in DNNs is the accurate estimation of the mutual information.
 36 However, regarding DNNs as deterministic models hinders us from specifying the random variable
 37 T_i and the distribution $P(T_i)$, thus it is difficult to accurately estimate $I(X; T_i)$ and $I(Y; T_i)$. More
 38 specifically, in the absence of a clear definition of T_i , simply assuming the activations of \mathbf{t}_i as the *i.i.d.*
 39 samples of T_i induces T_i being a continuous random variable and $I(X; T_i) = \infty$ in deterministic
 40 DNNs (see Appendix C in [33]). The complicated architecture of DNNs makes it challenging to
 41 specify $P(T_i)$. Therefore, most previous works have to indirectly estimate $P(T_i)$ via non-parametric
 42 models [40], such as the empirical distribution [35], Kernel Density Estimation (KDE) [33], and
 43 Gaussian convolution [11]. However, we experimentally confirm that classical non-parametric models
 44 derives poor mutual information estimation [29, 26] in DNNs, and one reason is because activations
 45 do not satisfy the *i.i.d.* prerequisite of non-parametric models (see Appendix G). In summary, the
 46 limitations mainly stem from the lack of an explicit probabilistic representation for deep learning.

47 The IB principle only formulates the information flow in DNNs = $\{\mathbf{x}, \mathbf{t}_1, \dots, \mathbf{t}_I, \hat{\mathbf{y}}\}$ after training,
 48 and the corresponding Markov chain (see Fig. 1 in [35])

$$Y \rightarrow X \rightarrow T_1 \cdots \rightarrow T_I \rightarrow \hat{Y} \quad (2)$$

49 indicates that the information of Y transfers to T_i in the forward direction and T_i receives the
 50 information of Y only via X . However, training DNNs by the back-propagation [30] implies that the
 51 information of Y transfers to T_i in the backward direction during training and retains information
 52 in T_i after training. Notably, Zhang *et al.* show that a DNN can fit labels well even using Gaussian
 53 noise as input to train the DNN [44], which implies that T_i can directly receive the information of Y .
 54 Hence, the IB principle does not comprehensively characterize the information flow in DNNs.

55 To address the above limitations and comprehensively explain DNNs in an information theoretic
 56 fashion, we introduce the probability space $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$ [6] for the i th hidden layer \mathbf{t}_i in DNNs.
 57 Compared to previous works, the probability space $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$ enables us to: (i) accurately estimate
 58 $I(X; T_i)$ and $I(Y; T_i)$ via specifying T_i and $P(T_i)$, and (ii) quantify the effect of the architecture of
 59 \mathbf{t}_i and the back-propagation on $I(X; T_i)$ and $I(Y; T_i)$ via explicitly modeling all the ingredients of \mathbf{t}_i ,
 60 such as the activation function and the weights in a probabilistic way. To the best of our knowledge,
 61 this is the first time the probability space of a hidden layer in DNNs is as defined.

62 Leveraging $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$, we derive information theoretic explanations for DNNs as follows:

- 63 • Two Markov chains¹ characterize the information flow in DNNs = $\{\mathbf{x}, \mathbf{t}_1, \dots, \mathbf{t}_I, \hat{\mathbf{y}}\}$

$$\begin{aligned} \bar{X} &\rightarrow T_1 \rightarrow \dots \rightarrow T_I \rightarrow \hat{Y} \\ T_1 &\leftarrow \dots \leftarrow T_I \leftarrow \hat{Y} \leftarrow Y. \end{aligned} \quad (3)$$

- 64 • Different hidden layers manifest different IB trade-offs depending on the architecture and
 65 the position of hidden layers in DNNs.
- 66 • A DNN satisfies the IB principle no matter the architecture of the DNN.

67 **Preliminaries.** $P(X, Y) = P(X)P(Y|X)$ is an unknown joint distribution between X and Y . A
 68 dataset $\mathcal{D} = \{(\mathbf{x}^j, y^j) | \mathbf{x}^j \in \mathbb{R}^M, y^j \in \mathbb{Z}\}_{j=1}^J$ consists of J *i.i.d.* samples generated from $P(X, Y)$
 69 with finite L labels, *i.e.*, $y^j \in \{1, \dots, L\}$. In the context of supervised learning, we focus on
 70 feedforward fully connected DNNs = $\{\mathbf{x}, \mathbf{t}_1, \dots, \mathbf{t}_I, \hat{\mathbf{y}}\}$, *i.e.*, Multi-Layer Perceptions (MLPs) [8]
 71 for the image classification task. Without loss of generality, we use the MLP = $\{\mathbf{x}, \mathbf{t}_1, \mathbf{t}_2, \hat{\mathbf{y}}\}$ with
 72 the cross-entropy loss ℓ_{CE} for most theoretical derivations. In addition, $H(\cdot)$ denotes entropy.

73 In the MLP, \mathbf{t}_1 and \mathbf{t}_2 have N and K neurons, respectively, and $\mathbf{t}_1 = \{t_{1n} = \sigma_1[\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle]\}_{n=1}^N$,
 74 where $\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle = \sum_{m=1}^M \omega_{mn}^{(1)} \cdot x_m + b_{1n}$ is the n th dot-product given the weight $\omega_{mn}^{(1)}$ and the bias
 75 b_{1n} , and $\sigma_1(\cdot)$ denotes an activation function, *e.g.*, ReLU. Similarly, $\mathbf{t}_2 = \{t_{2k} = \sigma_2[\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle]\}_{k=1}^K$,
 76 where $\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle = \sum_{n=1}^N \omega_{nk}^{(2)} \cdot t_{1n} + b_{2k}$. The output layer $\hat{\mathbf{y}}$ is softmax with L nodes

$$\hat{\mathbf{y}} = \{\hat{y}_l = \frac{1}{Z_Y} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle] = \frac{1}{Z_Y} \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))]\}_{l=1}^L, \quad (4)$$

77 where $\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle = \sum_{k=1}^K \omega_{kl}^{(3)} \cdot t_{2k} + b_{yl}$ and $Z_Y = \sum_{l=1}^L \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle]$ is the partition function.

¹In which the virtual random variable \bar{X} has all the information of X except Y , namely $H(\bar{X}) = H(X|Y)$.

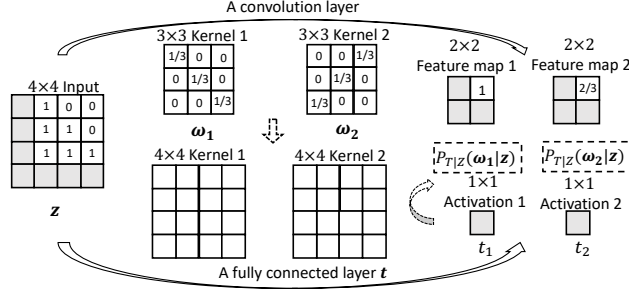


Figure 1: Given a 4×4 input z , a fully connected layer t is equivalent to a convolution layer with 4×4 convolution kernels. The definition of convolution (Chapter 9.1 in [12]) implies that the 4×4 weights ω_1 and ω_2 define two global features, and the two activations t_1, t_2 indicate the cross-correlation between ω_1, ω_2 and z , respectively. $P_{T|Z}(\omega_1|z)$ and $P_{T|Z}(\omega_2|z)$ measure the probability of ω_1 and ω_2 being recognized as the feature with the largest cross-correlation to z , respectively.

78 2 A probabilistic representation for deep learning

79 To accurately estimate $I(X; T_i)$ and $I(Y; T_i)$, in this section, we specify the probability space [6] for
80 a fully connected layer and derive the probabilistic explanations of the entire MLP.

81 It is known that a convolution kernel (namely the weights of convolution) defines a local feature,
82 and a convolution operation derives a feature map to measure the cross-correlation between the
83 local feature and input in a receptive field (Chapter 9.1 in [12]). Notably, a fully connected layer
84 is equivalent to a convolution layer with the kernel size having the same dimension as input. Thus
85 the weights of a neuron can be viewed as a global feature, and a fully connected layer with multiple
86 neurons derives activations to measure the cross-correlation between the multiple global features and
87 the input. The cross-correlation explanation for a fully connected layer is visualized in Figure 1.

88 Assuming that a fully connected layer t consists of N neurons $\{t_n = \sigma[\langle \omega_n, z \rangle]\}_{n=1}^N$, where $z \in \mathbb{R}^M$
89 is the input of t , $\langle \omega_n, z \rangle = \sum_{m=1}^M \omega_{mn} \cdot z_m + b_n$ is the dot-product between z and ω_n , and $\sigma(\cdot)$
90 is an activation function. Based on the cross-correlation explanation, the behavior of t is to measure
91 the cross-correlations between z and the N possible features defined by the the weights $\{\omega_n\}_{n=1}^N$.
92 In the context of pattern recognition [39], we define a virtual random process or ‘experiment’ as t
93 recognizing one of the patterns/features with the largest cross-correlation to z from the N possible
94 features. The experiment characterizes the behavior of t (*i.e.*, before recognizing the features with
95 the largest cross-correlation, t must measure the cross-correlations between z and all the N possible
96 features) while meets the requirement of the ‘experiment’ definition (*i.e.*, only one outcome will
97 occur on each trial of the experiment [6]). The probability space $(\Omega_T, \mathcal{F}, P_T)$ is defined as follows:

98 **Definition 1.** $(\Omega_T, \mathcal{F}, P_T)$ consists of three components: the sample space Ω_T has N possible
99 outcomes (features) $\{\omega_n = \{\omega_{mn}\}_{m=1}^M\}_{n=1}^N$ defined by the weights² of the N neurons; the event
100 space \mathcal{F} is the σ -algebra; and the probability measure P_T is a Gibbs distribution [22] to quantify the
101 probability of ω_n being recognized as the feature with the largest cross-correlation to z .

102 Taking into account the randomness of z , the conditional distribution $P_{T|Z}$ is formulated as

$$P_{T|Z}(\omega_n|z) = \frac{1}{Z_T} \exp(t_n) = \frac{1}{Z_T} \exp[\sigma(\langle \omega_n, z \rangle)], \quad (5)$$

103 where Z is the random variable of z and $Z_T = \sum_{n=1}^N \exp(t_n)$ is the partition function.

104 $(\Omega_T, \mathcal{F}, P_T)$ clearly explains all the ingredients of t in a probabilistic fashion. The n th neuron
105 defines a global feature by the weights w_n and the activation $t_n = \sigma(\langle \omega_n, z \rangle)$ measures the cross-
106 correlation between w_n and z . The Gibbs distribution $P_{T|Z}$ indicates that if w_n has the higher
107 activation, *i.e.*, the larger cross-correlation to z , it has the larger probability being recognized as
108 the feature with largest cross-correlation to z . For instance, if $z \in \mathbb{R}^{16}$ and t includes $N = 2$
109 neurons, then $\Omega_T = \{\omega_1, \omega_2\}$ defines two possible outcomes (features), where $\omega_n = \{\omega_{mn}\}_{m=1}^{16}$.
110 $\mathcal{F} = \{\emptyset, \{\omega_1\}, \{\omega_2\}, \{\omega_1, \omega_2\}\}$ means that neither, one, or both of the features are recognized
111 by t given z , respectively. $P_{T|Z}(\omega_1|z)$ and $P_{T|Z}(\omega_2|z)$ are the probability of ω_1 and ω_2 being
112 recognized as the feature with the largest cross-correlation to z , respectively.

²We do not take into account the scalar value b_n for defining Ω_T , as it not affects the feature defined by ω_n .

113 $(\Omega_T, \mathcal{F}, P_T)$ explains the representation ability of deep learning. Compared to Restricted Boltzmann
 114 Machines (RBMs) [31] simply using binary units to indicate features being recognized or not given
 115 input, the Gibbs distribution³ $P_{T|Z}(\omega_n|z)$ measures the probability of ω_n being recognized with
 116 the largest cross-correlation to z , *i.e.*, it characterizes the relation between features and input more
 117 accurately. Moreover, Equation 5 shows that $t_n = \sigma(\langle \omega_n, z \rangle)$ is the negative energy function [22] of
 118 the Gibbs distribution, thus $P_{T|Z}(\omega_n|z)$ can be derived as long as $\sigma(\langle \omega_n, z \rangle)$ are known because
 119 the energy function is the sufficient statistics [2] of the Gibbs distribution. That enables subsequent
 120 hidden layers to generate high-level features of input via directly processing the activations $\{t_n\}_{n=1}^N$,
 121 thus deep learning can form a hierarchical structure to represent much complex features.

122 $(\Omega_T, \mathcal{F}, P_T)$ answers a fundamental question: which component of a hidden layer contains the
 123 information of the layer? Since ω_n defines Ω_T , the weights contain all the information of a layer. In
 124 particular, since the activation $t_n = \sigma(\langle \omega_n, z \rangle)$ is a function of ω_n , the data processing inequality
 125 [4] indicates that the information of t_n is no more than the information of ω_n . Simulations in Section
 126 4.2 demonstrate that if activations do not correctly characterize the cross-correlation between weights
 127 and input, activations contain less information than weights do.

128 Based on $(\Omega_T, \mathcal{F}, P_T)$, we define the random variable T as follows:

129 **Definition 2.** Given the fully connected layer \mathbf{t} , we define the random variable $T : \Omega_T \rightarrow E_T$ as

$$T(\omega_n) \triangleq n, \quad (6)$$

130 where the measurable space $E_T = \{1, \dots, N\}$.

131 Since Ω_T is composed of finite N possible outcomes, T is a discrete random variable. Notably, the
 132 one-to-one correspondence between ω_n and n indicates

$$P_{T|Z}(\omega_n|z) = P_{T|Z}(n|z). \quad (7)$$

133 If not considering the back-propagation training, the weights (namely Ω_{T_i}) of each layer are fixed.
 134 Thus T_{i+1} entirely depends on T_i and the MLP = $\{\mathbf{x}; \mathbf{t}_1; \mathbf{t}_2; \hat{\mathbf{y}}\}$ forms a Markov chain

$$X \rightarrow T_1 \rightarrow T_2 \rightarrow \hat{Y}. \quad (8)$$

135 Based on the corresponding joint distribution $P(\hat{Y}, T_2, T_1|X) = P(T_1|X)P(T_2|T_1)P(\hat{Y}|T_2)$ and
 136 Definition 2, we derive a probabilistic explanation for the entire MLP, which is summarized in
 137 Theorem 1. The detailed derivation is presented in Appendix B.

138 **Theorem 1.** The MLP = $\{\mathbf{x}; \mathbf{t}_1; \mathbf{t}_2; \hat{\mathbf{y}}\}$ formulates a conditional Gibbs distribution

$$P_{\hat{Y}|X}(l|\mathbf{x}) = \sum_{k=1}^K \sum_{n=1}^N P(\hat{Y} = l, T_2 = k, T_1 = n|X = \mathbf{x}) = \frac{1}{Z_{\text{MLP}}(\mathbf{x})} \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))] \quad (9)$$

139 where $Z_{\text{MLP}}(\mathbf{x}) = \sum_{l=1}^L \sum_{k=1}^K \sum_{n=1}^N P_{\hat{Y}, T_2, T_1|X}(l, k, n|x)$ is the partition function.

140 Since $P_{\hat{Y}|X}(l|\mathbf{x})$ exactly equals the output \hat{y}_l of the MLP, namely Equation (4), we conclude that
 141 the entire architecture of the MLP forms a family of Gibbs distribution $P_{\hat{Y}|X}(l|\mathbf{x})$. In general, the
 142 back-propagation updates a weight ω based on the gradient of ℓ_{CE} with respect to ω ,

$$\omega(s+1) = \omega(s) - \alpha \cdot \frac{\partial \ell_{\text{CE}}}{\partial \omega(s)} = \omega(s) - \alpha \cdot \frac{\partial \text{KL}[P(Y|X)||P(\hat{Y}|X)]}{\partial \omega(s)}, \quad (10)$$

143 where s is the index of training iteration, α is the training rate, and $\text{KL}[\cdot||\cdot]$ is the KL-divergence.

144 Figure 2 summarizes the probabilistic explanation for deep learning based on the MLP. In general,
 145 a single learning iteration, an epoch, consists of two phases: training and inference (after training).
 146 During inference, the MLP bridges X and \hat{Y} via multiple intermediate features Ω_{T_1} , Ω_{T_2} , and $\Omega_{\hat{Y}}$
 147 defined by weights, and formulates the statistical relation between \hat{Y} and X as a family of conditional
 148 Gibbs distribution $P(\hat{Y}|X)$. During training, the back-propagation updates weights to learn optimal
 149 intermediate features for searching an optimal $P(\hat{Y}|X)$ to accurately approximate $P(Y|X)$.

³Recent works about Gibbs explanations for a hidden layer are discussed in Appendix A.

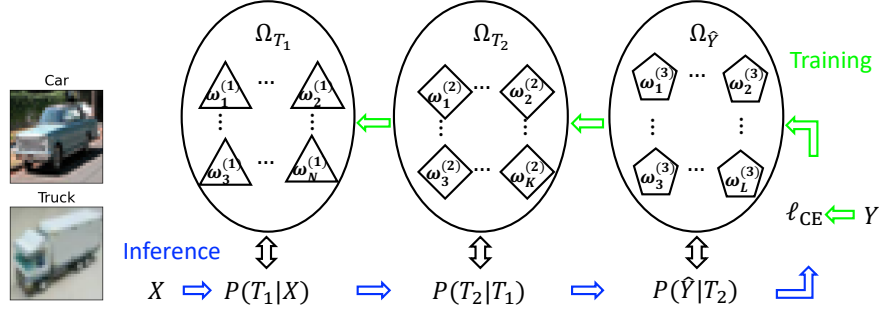


Figure 2: The visualization of the probabilistic explanation for deep learning based on the MLP.

150 3 The information theoretic explanations for deep learning

151 To address the limitations of existing IB explanations, this section proposes some novel information
 152 theoretic explanations for DNNs based on the proposed probabilistic representation.

153 **Proposition 1.** The mutual information between a fully connected layer and dataset is finite.

$$I(X; T) < \infty. \quad (11)$$

154 *Proof:* Definition 2 shows $E_T = \{1, \dots, N\}$. Thus T is a discrete random variable and $H(T) < \infty$,
 155 thereby $I(X; T) \leq H(T) < \infty$.

156 Proposition 1 circumvents the infinite mutual information problem. In the absence of a clear definition
 157 $T : \Omega_T \rightarrow E_T$, most previous works [33, 3, 1] simply viewing the activation t_n as the sample of T ,
 158 namely $t_n \in E_T = \mathbb{R}$, implies T being continuous and gives rise to the infinite mutual information
 159 problem in deterministic DNNs. However, $(\Omega_T, \mathcal{F}, P_T)$ indicates that t_n actually is a variable
 160 measuring the cross-correlation between w_n and z rather than the sample of T , namely $t_n \notin E_T$.

161 **Theorem 2.** The information of Y flows into the MLP in the backward direction during training

$$T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y. \quad (12)$$

162 *Proof:* First, since Ω_T is defined by ω in $(\Omega_T, \mathcal{F}, P_T)$ and Equation (10) shows that $\omega(s+1)$ is
 163 determined by all the previous gradients $\{\frac{\partial \ell_{CE}}{\partial \omega^{(s)}}\}_{s=1}^S$, and $\omega(0)$ is randomly initialized and α is a
 164 constant, we can derive that Ω_T is determined by $\frac{\partial \ell_{CE}}{\partial \omega}$. Second, based on the back-propagation, the
 165 relation between gradients in two adjacent layers in the MLP $= \{\mathbf{x}; \mathbf{t}_1; \mathbf{t}_2; \hat{\mathbf{y}}\}$ is formulated as

$$\begin{aligned} \frac{\partial \ell_{CE}^*}{\partial \omega_{kl}^{(3)}} &= [P_{\hat{Y}|X}(l|\mathbf{x}) - P_{Y|X}(l|\mathbf{x})] \cdot t_{2k}, \\ \frac{\partial \ell_{CE}^\circ}{\partial \omega_{nk}^{(2)}} &= \sum_{l=1}^L \frac{\partial \ell_{CE}^*}{\partial \omega_{kl}^{(3)}} \cdot \omega_{kl}^{(3)} \cdot \frac{\sigma_2'(\langle \omega_k^{(2)}, \mathbf{t}_1 \rangle)}{t_{2k}} \cdot t_{1n}, \quad \frac{\partial \ell_{CE}^\circ}{\partial \omega_{mn}^{(1)}} = \sum_{k=1}^K \frac{\partial \ell_{CE}^\circ}{\partial \omega_{nk}^{(2)}} \cdot \omega_{nk}^{(2)} \cdot \frac{\sigma_1'(\langle \omega_n^{(1)}, \mathbf{x} \rangle)}{t_{1n}} \cdot x_m. \end{aligned} \quad (13)$$

166 Equation 13 shows that $\frac{\partial \ell_{CE}^*}{\partial \omega^{(3)}}$ is a function of $P_{Y|X}(l|\mathbf{x})$ and $\frac{\partial \ell_{CE}^\circ}{\partial \omega^{(i)}}$ is a function of $\frac{\partial \ell_{CE}^*}{\partial \omega^{(i+1)}}$, where
 167 $\omega^{(3)}$ denotes the weight of $\hat{\mathbf{y}}$. The two points above enable us to derive that Ω_{T_i} is a function of $\Omega_{T_{i+1}}$
 168 and $\Omega_{\hat{Y}}$ is a function of $P(Y|X)$. Based on Definition 2, we can further derive that T_i is a function
 169 of T_{i+1} and \hat{Y} is a function of Y , i.e., $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$. (See the detailed proof in Appendix C).

170 Theorem 2 is consistent with the prevailing explanation for deep learning. LeCun *et al.* show that
 171 deep learning exploits the hierarchical property of signals [21], i.e., the layers farther from output
 172 learn lower-level features, such as edges, whereas the layers closer to output assemble lower-level
 173 features into the higher-level features corresponding to labels (see Figure 2 in [43]). Notably, since
 174 lower-level features commonly exist in signals with different labels (e.g., lower-level features, such
 175 as the edges of the vehicle frame and the circular contour of wheels, exist in both the car and the
 176 truck classes in the CIFAR-10 dataset [18] in Figure 2), lower-level features do not contain much
 177 information of labels. Therefore, the layers farther from output do not have much information of
 178 labels, which is consistent with the Markov chain $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$.

179 Since all the information of Y stems from X (i.e., $H(Y) = I(X; Y)$ proven in Appendix D),
 180 Theorem 2 implies that partial information of X flows into the MLP in the backward direction during
 181 training. Equation (2) shows the information of X flowing into the MLP in the forward direction
 182 during inference. Overall, the information of X flows in the backward and forward directions during
 183 training and inference, respectively. As a result, the Markov chain, Equation (2), proposed by recent
 184 works could not fully characterize the information flow of X in the MLP in each epoch. In other
 185 words, $I(X; T_i)$ is not necessarily greater than $I(X; T_{i+1})$ in the MLP in each epoch.

186 Equation (2) shows that T_i receives the information of Y via X during inference. Theorem 2 shows
 187 that T_i also directly receives information of Y during training, because the back-propagation updates
 188 weights (i.e., Ω_{T_i}) based on the label Y . Thus Equation (2) cannot fully characterize the information
 189 flow of Y in the MLP in each epoch, when we take into account the back-propagation training.

190 To fully characterize the information flow in the MLP in each epoch, we introduce Corollary 1.

191 **Corollary 1.** The information flow in the MLP can be characterized by two Markov chains as

$$\begin{aligned} \bar{X} &\rightarrow T_1 \rightarrow T_2 \rightarrow \hat{Y} \\ T_1 &\leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y. \end{aligned} \tag{14}$$

192 The virtual random variable \bar{X} contains all the information of X except Y , i.e., $H(\bar{X}) = H(X|Y)$.

193 *Proof of the first Markov chain:* Since \bar{X} does not have any information of Y , it can only flow into
 194 the MLP in the forward direction during inference. Again since \bar{X} does not have any information of
 195 Y , the information flow of Y during training will not affect the information flow of \bar{X} . Therefore,
 196 $\bar{X} \rightarrow T_1 \rightarrow T_2 \rightarrow \hat{Y}$ characterizes the information flow of \bar{X} in both training and inference phases.

197 *Proof of the second Markov chain:* Since the weights are fixed after training, the sample space and
 198 the distribution of hidden layers are fixed after training. Therefore, the information of Y transferred
 199 into hidden layers during training will retain there after training (i.e., during inference). In addition,
 200 Definition 1 indicates that a fully connected layer $\mathbf{t} = \{t_n = \sigma(\langle \omega_n, \mathbf{z} \rangle)\}_{n=1}^N$ measures the cross-
 201 correlation between ω_n and \mathbf{z} during inference, thus $\{\omega_n\}_{n=1}^N$ can be viewed as a representation of
 202 Z . As a result, even though Z has all the information of Y , the information of Y that \mathbf{t} can learn
 203 from Z is determined by how much information of Y the representation $\{\omega_n\}_{n=1}^N$ has. Overall, the
 204 information flow of Y during inference will be the same as that during training. Based on Theorem 2,
 205 we conclude that $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$ characterizes the information flow of Y in the MLP in both
 206 training and inference phases. Detailed derivations and explanations are presented in Appendix E.

207 To quantify how much information of X and Y is learned by the MLP, we introduce Corollary 2.

208 **Corollary 2.** The mutual information between dataset and the entire MLP can be expressed as

$$\begin{aligned} I(X; T_{\text{MLP}}) &= I(\bar{X}; T_1) + I(Y; \hat{Y}) \\ I(Y; T_{\text{MLP}}) &= I(Y; \hat{Y}) \end{aligned} \tag{15}$$

209 where T_{MLP} denotes a random variable corresponding to the entire architecture of the MLP.

210 *Proof:* Since $H(Y) = I(X; Y)$ (Appendix D), $H(X) = H(\bar{X}) + I(X; Y) = H(\bar{X}) + H(Y)$.
 211 Hence, Corollary 2 can be derived by Corollary 1 and the chain rule. The proof is in Appendix F.

212 4 Simulations

213 In this section, we propose a mutual information estimator based on $(\Omega_T, \mathcal{F}, P_T)$ and demonstrate the
 214 probabilistic representation and information theoretic explanations for deep learning on a synthetic
 215 dataset with known entropy. Additional experiments on benchmark datasets are in Appendix H.

216 4.1 Setup

217 **Mutual information estimator.** Based on the definition of mutual information, we have

$$I(X; T_i) = H(T_i) - H(T_i|X). \tag{16}$$

218 Previous works simply estimate $I(X; T_i) = H(T_i)$, because T_i is assumed to be entirely dependent
 219 on X in the Markov chain, Equation (2), thereby $H(T_i|X) = 0$. However, Corollary 1 shows that T_i
 220 depends on both X and Y if taking into account the training phase, thereby $H(T_i|X) \neq 0$.

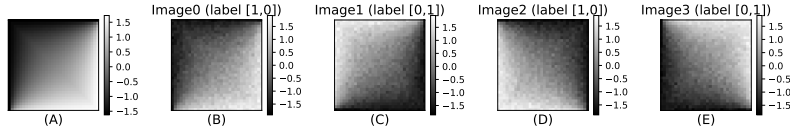


Figure 3: (A) the deterministic image $\hat{\mathbf{x}}$. Image0 is generated by adding $\mathcal{N}(\mu, \sigma^2)$ without rotation, Image1 is generated by rotating $\hat{\mathbf{x}}$ along the secondary diagonal direction and adding $\mathcal{N}(\mu, \sigma^2)$, Image2 and Image are generated by rotating $\hat{\mathbf{x}}$ along the vertical and horizontal directions, respectively, and adding $\mathcal{N}(\mu, \sigma^2)$.

Table 1: The number of neurons(nodes) and the activation function in the layers of the MLPs

	\mathbf{x}	\mathbf{t}_1	\mathbf{t}_2	$\hat{\mathbf{y}}$	$\sigma(\cdot)$
MLP1	1024 (32×32)	8	6	2	$\text{ReLU}(z) = \max(0, z)$
MLP2	1024 (32×32)	8	6	2	$\text{Tanh}(z) = (e^z - e^{-z}) / (e^z + e^{-z})$
MLP3	1024 (32×32)	2	6	2	ReLU

221 To accurately estimate $I(X; T_i)$, we need to specify $P(T_i|X)$ and $P(T_i)$. Based on $(\Omega_{T_i}, \mathcal{F}, P_{T_i})$,
 222 we formulate $P_{T_i|X}(n|\mathbf{x}^j)$ of the three fully connected layers in the MLP as

$$\begin{aligned}
 P_{T_1|X}(n|\mathbf{x}^j) &= \frac{1}{Z_{F_1}} \exp[\sigma_1(\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x}^j \rangle)], \quad P_{T_2|X}(k|\mathbf{x}^j) = \frac{1}{Z_{F_2}} \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1(\mathbf{x}^j) \rangle)], \\
 P_{\hat{Y}|X}(l|\mathbf{x}^j) &= \frac{1}{Z_{F_Y}} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2(\mathbf{t}_1(\mathbf{x}^j)) \rangle].
 \end{aligned} \tag{17}$$

223 To derive the marginal distribution $P(T_i)$, we sum the joint distribution $P(T_i, X)$ over $\mathbf{x} \in \mathcal{X}$,

$$P(T_i = n) = \sum_{\mathbf{x} \in \mathcal{X}} P_X(\mathbf{x}) P_{T_i|X}(n|\mathbf{x}) \approx \sum_{\mathbf{x}^j \in \mathcal{D}} P_X(\mathbf{x}^j) P_{T_i|X}(n|\mathbf{x}^j) = \frac{1}{J} \sum_{\mathbf{x}^j \in \mathcal{D}} P_{T_i|X}(n|\mathbf{x}^j), \tag{18}$$

224 where $P_X(\mathbf{x}^j)$ is estimated by the empirical distribution $1/J$ given \mathcal{D} . Finally, we can derive $I(X; T_i)$
 225 by Equation 16, 17, and 18. Similarly, based on the definition of mutual information, we have

$$I(Y; T_i) = H(T_i) - H(T_i|Y). \tag{19}$$

226 To estimate $H(T_i|Y)$, we reformulate $P(T_i|Y)$ as

$$P_{T_i|Y}(n|l) = \sum_{\mathbf{x} \in \mathcal{X}} P_{T_i|X}(n|\mathbf{x}) P_{X|Y}(\mathbf{x}|l) \approx \frac{1}{N(l)} \sum_{\mathbf{x}^j \in \mathcal{D}, y^j=l} P_{T_i|X}(n|\mathbf{x}^j), \tag{20}$$

227 where $P_{X|Y}(\mathbf{x}^j|l)$ is estimated by the empirical distribution $1/N(l)$ and $N(l)$ denotes the number of
 228 samples with the label l in \mathcal{D} . Finally, we can derive $I(Y; T_i)$ by Equation 18, 19, and 20.

229 **Synthetic dataset.** The dataset consists of 512 gray-scale 32×32 images, which are evenly generated
 230 by rotating a deterministic image $\hat{\mathbf{x}}$ in four different orientations and adding Gaussian noise with
 231 expectation $\mu = \mathbb{E}(\hat{\mathbf{x}})$ and variance $\sigma^2 = 1$, namely $\mathbf{x} = r(\hat{\mathbf{x}}) + \mathcal{N}(\mu, \sigma^2)$, where $r(\cdot)$ denotes
 232 the rotation method shown in Figure 3. The reason for adding Gaussian noise is to avoid DNNs
 233 directly memorizing the deterministic image. In addition, the binary labels [1,0] and [0,1] evenly
 234 divide the synthetic dataset into two classes. As a result, the synthetic dataset has (approximately)
 235 2 bits information and the labels have 1 bit information. Compared to popular benchmark dataset
 236 with unknown features and entropy, *e.g.*, MNIST [19] and Fashion-MNIST [41], the features and the
 237 entropy of the synthetic dataset are clear and known, which enables us to examine the probabilistic
 238 representation and the mutual information estimator.

239 **Neural Networks.** We train three MLPs, namely MLP1, MLP2 and MLP3, on the synthetic dataset
 240 by a variant of Stochastic Gradient Descent (SGD) method, namely Adam [15], over 1000 epochs
 241 with the learning rate $\alpha = 0.03$. Table 1 summarizes the architecture of the three MLPs.

242 4.2 Validating the probability space and the mutual information estimator

243 We demonstrate the sample space Ω_T by visualizing the weights⁴ of the eight neurons in \mathbf{t}_1 , *i.e.*,
 244 $\boldsymbol{\omega}_n^{(1)} = \{\omega_{mn}^{(1)}\}_{m=1}^{1024}$, in 5 different epochs (*i.e.*, 0,1,4,128,1000) in Figure 4 (Left). As training
 245 continues, we observe that $\boldsymbol{\omega}_n^{(1)}$ quickly learns all the spatial features of the synthetic dataset. For
 246 instance, $\boldsymbol{\omega}_2^{(1)}$ has low magnitude at top-left positions and high magnitude at bottom-right positions,
 247 which correctly characterizes the spatial feature of Image0. Similarly, $\boldsymbol{\omega}_3^{(1)}$, $\boldsymbol{\omega}_4^{(1)}$, and $\boldsymbol{\omega}_5^{(1)}$ correctly
 248 characterize the spatial feature of Image1, Image2, and Image3 in Figure 3, respectively.

⁴We only show the learned weights in MLP1 because we observe that the learned weights in MLP1 and MLP2 are very similar, though they use different activation functions.

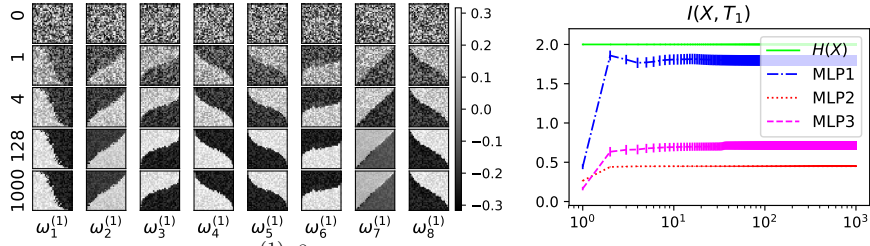


Figure 4: (Left) The eight features $\{\omega_n^{(1)}\}_{n=1}^8$ learned by the weights of the eight neurons in 5 different epochs (*i.e.*, 0,1,4,128,1000), where $\omega_n^{(1)} = \{\omega_{mn}^{(1)}\}_{m=1}^{1024}$ are reshaped into 32×32 to show the spatial structure. (Right) The variation of $I(X; T_1)$ in the MLP1, MLP2, and MLP3 during 1000 epochs.

Table 2: The Gibbs probability $P_{F_1|X}(\omega_n^{(1)}|\text{Image0})$ in MLP1 and MLP2 in the 1000 epoch

	$\omega_1^{(1)}$	$\omega_2^{(1)}$	$\omega_3^{(1)}$	$\omega_4^{(1)}$	$\omega_5^{(1)}$	$\omega_6^{(1)}$	$\omega_7^{(1)}$	$\omega_8^{(1)}$
$\langle \omega_n^{(1)}, \mathbf{x} \rangle$	-63.6	208.8	-181.6	45.1	-55.6	157.5	-210.0	-30.1
$f_{1n}^{\text{ReLU}}(\mathbf{x})$	0.0	208.8	0.0	45.1	0.0	157.5	0.0	0.0
$\exp[f_{1n}^{\text{ReLU}}(\mathbf{x})]$	1.0	4.79e+90	1.0	3.86e+19	1.0	2.51e+68	1.0	1.0
$P_{T_1 X}^{\text{ReLU}}$	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
$f_{1n}^{\text{Tanh}}(\mathbf{x})$	-1.0	1.0	-1.0	1.0	-1.0	1.0	-1.0	-1.0
$\exp[f_{1n}^{\text{Tanh}}(\mathbf{x})]$	0.36	2.71	0.36	2.71	0.36	2.71	0.36	0.36
$P_{T_1 X}^{\text{Tanh}}$	0.037	0.272	0.037	0.272	0.037	0.272	0.037	0.037

$f_{1n}^{\text{Tanh}}(\mathbf{x}) = \sigma^{\text{Tanh}}(\langle \omega_n^{(1)}, \mathbf{x} \rangle)$ and $f_{1n}^{\text{ReLU}}(\mathbf{x}) = \sigma^{\text{ReLU}}(\langle \omega_n^{(1)}, \mathbf{x} \rangle)$ are the activations given the same $\langle \omega_n^{(1)}, \mathbf{x} \rangle$.

249 We demonstrate that $P(T_1|X)$ correctly measures the probability of $\{\omega_n^{(1)}\}_{n=1}^8$ being recognized the
 250 feature with the largest cross-correlation to \mathbf{x} in Table 2. For instance, $\omega_2^{(1)}$ correctly characterizes
 251 the feature of Image0 and has the largest cross-correlation $\langle \omega_2^{(1)}, \mathbf{x} \rangle = 190.8$, thus it has the largest
 252 probability $P_{T_1|X}^{\text{ReLU}}(\omega_2^{(1)}|\text{Image0}) = 1.0$ being recognized as the feature with largest cross-correlation
 253 to Image0. In contrast, since $\omega_7^{(1)}$ incorrectly characterizes the feature of Image0 and has the lowest
 254 cross-correlation $\langle \omega_7^{(1)}, \mathbf{x} \rangle = -210.0$, so it has the lowest probability $P_{T_1|X}^{\text{ReLU}}(\omega_7^{(1)}|\text{Image0}) = 0.0$
 255 being recognized as the feature with largest cross-correlation to Image0.

256 We observe that an activation function (abbr. ACT) plays an important role in the distribution.
 257 Specifically, ReLU, a non-saturating (unbounded) ACT [10], preserves the positive cross-correlations
 258 while resets all the negative ones as zero. $P_{T_1|X}^{\text{ReLU}}(\omega_2^{(1)}|\text{Image0}) = 1.0$ shows that ReLU derives the
 259 correct probability of $\omega_2^{(1)}$ being recognized as the feature with largest cross-correlation. In contrast,
 260 though $\omega_2^{(1)}$ has stronger cross-correlation to Image0 than $\omega_4^{(1)}$, *i.e.*, $\langle \omega_2^{(1)}, \mathbf{x} \rangle > \langle \omega_4^{(1)}, \mathbf{x} \rangle$, Tanh, a
 261 saturating (bounded) ACT, derives $f_{12}^{\text{Tanh}}(\mathbf{x}) = f_{14}^{\text{Tanh}}(\mathbf{x}) = 1.0$, and makes $\omega_4^{(1)}$ to incorrectly have
 262 the same probability 0.272 to $\omega_2^{(1)}$ being recognized as the feature with the largest cross-correlation
 263 to Image0, *i.e.*, Tanh hinders t_1 from correctly recognizing the features of input. The simulations for
 264 validating the probability space based on other synthetic images are presented in Appendix G.

265 To validate the mutual information estimator, we follow recent works [35, 33] to train the three
 266 MLPs with 50 different random initialization and study the average mutual information. Figure 4
 267 (Right) shows that $I(X; T_1)$ quickly increases to 1.81 and keeps stable in the MLP1, *i.e.*, t_1 learns
 268 most information of the dataset as $H(X) = 2.0$. Notably, the result is consistent with the variation
 269 of the weights in Figure 4 (Left), which shows that the weights correctly characterize the features
 270 of the dataset and keeps stable after the fourth epoch. As a comparison, we observe that $I(X; T_1)$
 271 keeps stable at 0.44 in the MLP2, which confirms the statement that Tanh hinders t_1 from correctly
 272 recognizing the features of input. In addition, Figure 4 (Right) shows that $I(X; T_1) \approx 0.79$ in MLP3
 273 is smaller than $I(X; T_1) \approx 1.81$ in MLP1, which is consistent with Definition 1, *i.e.*, a layer with
 274 fewer neurons would represent fewer possible features, thus it contains less information.

275 In summary, we demonstrate the probability space $(\Omega_T, \mathcal{F}, P_T)$ and show that if an ACT cannot
 276 preserve the cross-correlation between weights(features) and input, it would distort the distribution
 277 of a layer, thereby affecting the mutual information between the layer and data/labels. In addition,
 278 we show that the proposed mutual information estimator outperforms the existing non-parametric
 279 models, *e.g.*, empirical distribution [35] and KDE [33], based on the synthetic dataset. Especially,
 280 activations do not satisfy the *i.i.d.* prerequisite of non-parametric models is an important reason for
 281 non-parametric models deriving inaccurate mutual information in DNNs. Due to limited space, the
 282 experimental comparison and study of non-parametric models are presented in Appendix G.

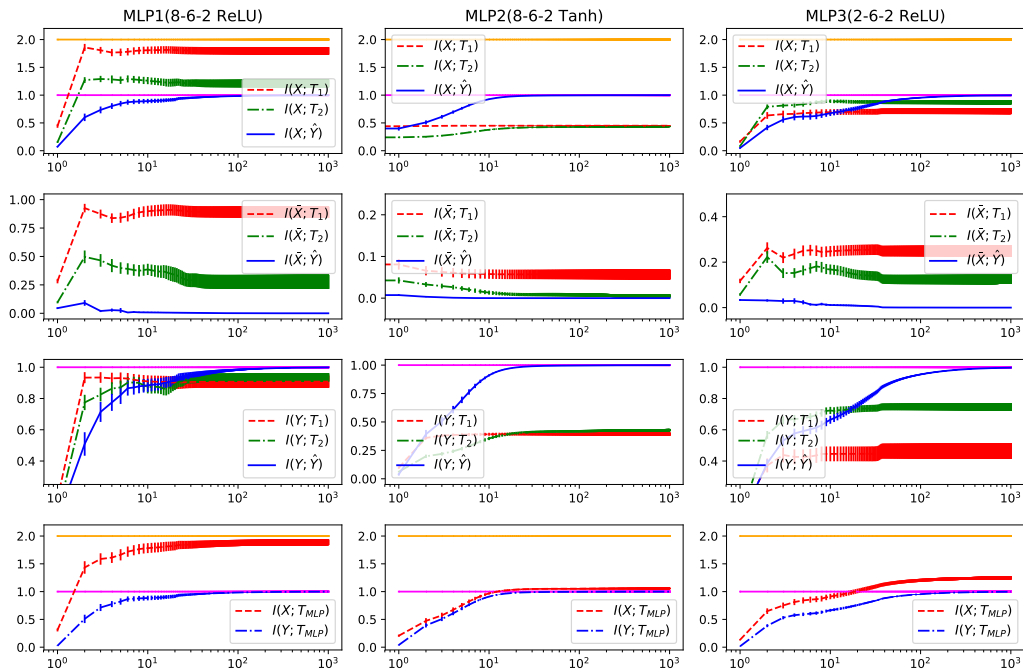


Figure 5: All the x-axis index training epochs. In each column, the first three figures show $I(X; T_i)$, $I(\bar{X}; T_i)$, and $I(Y; T_i)$ respectively. The fourth figure shows $I(X; T_{MLP})$ and $I(Y; T_{MLP})$ in a MLP. The pink line denotes $H(Y) = 1.0$ and the orange line denotes $H(X) = 2.0$.

283 4.3 Validating the information theoretic explanations for DNNs

284 In Figure 5, we observe $I(X; T_i) \leq I(X; \hat{Y})$ in MLP2 and MLP3, which confirms that the Markov
 285 chain proposed by previous works, Equation (2), cannot fully explain the information flow in MLPs,
 286 if taking into account the back-propagation training. As a comparison, the second and third row
 287 show $I(\bar{X}; T_1) \geq I(\bar{X}; T_2) \geq I(\bar{X}; \hat{Y})$ and $I(Y; T_1) \leq I(Y; T_2) \geq I(Y; \hat{Y})$ in all the three MLPs,
 288 which validates that Corollary 1, *i.e.*, Equation (14) characterizes the information flow in MLPs.

289 Figure 5 demonstrates that different hidden layers achieve different IB trade-offs depending on
 290 the architecture and the position of the layers in MLPs. In terms of architecture, $I(Y; T_1) > 0.8$
 291 and $I(\bar{X}; T_1) > 0.75$ in MLP1 indicate that t_1 , with ReLU, achieves a good prediction without
 292 much compression, whereas $I(Y; T_1) < 0.5$ and $I(\bar{X}; T_1) < 0.1$ in MLP2 show that t_1 , with Tanh,
 293 achieves a different IB trade-off. In addition, $I(Y; T_1) \approx 0.45$ and $I(\bar{X}; T_1) \approx 0.25$ in MLP3 show
 294 the effect of neuron numbers on the IB trade-off. In terms of position, $I(Y; \hat{Y}) = 1$ and $I(\bar{X}; \hat{Y}) = 0$
 295 in MLP1 means that \hat{y} has a different IB trade-off to t_1 in MLP1.

296 We demonstrate that a MLP satisfies the IB principle no matter what the architecture of the MLP
 297 is. Figure 5 visualizes $I(X; T_{MLP})$ and $I(Y; T_{MLP})$ based on Corollary 2. It shows that all of three
 298 MLPs satisfy the IB principle, namely $I(X; T_{MLP}) < H(X) = 2$ and $I(Y; T_{MLP}) = H(Y) = 1$,
 299 though they have different architectures. Importantly, in contrast to previous work [33] claiming that
 300 the compression not exists in DNNs with non-saturating ACT, such as ReLU, Figure 5 clearly shows
 301 that the compression exists in all the MLPs, no matter the activation function of MLPs.

302 We further demonstrate the information theoretic explanations for DNNs on the benchmark MNIST
 303 and Fashion-MNIST datasets. The experiments are presented in Appendix H.

304 5 Conclusion and future work

305 In this work, we (1) specify the probability space for a hidden layer for (2) accurately estimating the
 306 mutual information and (3) clearly explaining how the components of the layer affect the mutual
 307 information. We take into account the back-propagation training and derive two novel Markov chains
 308 to characterize the information flow in DNNs. Furthermore, we demonstrate that a DNN satisfies the
 309 IB principle no matter the architecture of the DNN. In contrast, different hidden layers show different
 310 IB trade-offs depending on the architecture and the position of the layers in DNNs. A potential
 311 direction is to study the generalization of DNNs based on the probabilistic representation.

References

- 312
- 313 [1] Rana Ali Amjad and Bernhard Claus Geiger. Learning representations for neural network-based classi-
314 fication using the information bottleneck principle. *IEEE transactions on pattern analysis and machine*
315 *intelligence*, 2019.
- 316 [2] George Casella and Roger L Berger. *Statistical inference*. Cengage Learning, 2021.
- 317 [3] Ivan Chelombiev, Conor Houghton, and Cian O’Donnell. Adaptive estimators show information compres-
318 sion in deep neural networks. In *International Conference on Learning Representations*, 2019.
- 319 [4] Thomas Cover and Joy Thomas. *Elements of Information Theory*. Wiley-Interscience, Hoboken, New
320 Jersey, 2006.
- 321 [5] Balázs Csanád Csáji et al. Approximation with artificial neural networks. *Faculty of Sciences, Eötvös Lornd*
322 *University, Hungary*, 24(48):7, 2001.
- 323 [6] Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- 324 [7] Marylou Gabrié, Andre Manoel, Clément Luneau, Nicolas Macris, Florent Krzakala, Lenka Zdeborová,
325 et al. Entropy and mutual information in models of deep neural networks. In *Advances in Neural*
326 *Information Processing Systems*, pages 1821–1831, 2018.
- 327 [8] Matt W Gardner and SR Dorling. Artificial neural networks (the multilayer perceptron)—a review of
328 applications in the atmospheric sciences. *Atmospheric environment*, 32(14-15):2627–2636, 1998.
- 329 [9] S. Geman and D. Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images.
330 *IEEE Transactions. on Pattern Analysis and Machine Intelligence*, pages 721–741, June 1984.
- 331 [10] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural
332 networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*,
333 pages 249–256, 2010.
- 334 [11] Ziv Goldfeld, Ewout Van Den Berg, Kristjan Greenewald, Igor Melnyk, Nam Nguyen, Brian Kingsbury,
335 and Yury Polyanskiy. Estimating information flow in deep neural networks. In *Proceedings of the 36th*
336 *International Conference on Machine Learning*, volume 97, pages 2299–2308, 2019.
- 337 [12] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT
338 press Cambridge, 2016.
- 339 [13] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computa-*
340 *tion*, 14:1771–1800, 2002.
- 341 [14] Patrick Kidger and Terry Lyons. Universal approximation with deep narrow networks. In *Conference on*
342 *Learning Theory*, pages 2306–2327. PMLR, 2020.
- 343 [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
344 *arXiv:1412.6980*, 2014.
- 345 [16] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
346 *International Conference on Machine Learning*, pages 1885–1894. PMLR, 2017.
- 347 [17] Artemy Kolchinsky, Brendan D Tracey, and Steven Van Kuyk. Caveats for information bottleneck in
348 deterministic scenarios. *arXiv preprint arXiv:1808.07593*, 2018.
- 349 [18] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technique Report*, 2009.
- 350 [19] Y. LeCun, BE. Boser, and JS. Denker. Handwritten digit recognition with a back-propagation network. In
351 *NeurIPS*, pages 396–494, 1990.
- 352 [20] Y. LeCun, L. Bottou, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings*
353 *of the IEEE*, 11:2278–2324, 1998.
- 354 [21] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- 355 [22] Yann LeCun, Sumit Chopra, Raia Hadsell, Marc’Aurelio Ranzato, and Fu Jie Huang. *A tutorial on*
356 *energy-based learning*. MIT Press, 2006.
- 357 [23] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha
358 Sohl-Dickstein. Deep neural networks as gaussian processes. In *ICLR*, 2018.

- 359 [24] Henry W Lin, Max Tegmark, and David Rolnick. Why does deep and cheap learning work so well?
360 *Journal of Statistical Physics*, 168(6):1223–1247, 2017.
- 361 [25] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them.
362 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5188–5196,
363 2015.
- 364 [26] David McAllester and Karl Stratos. Formal limitations on the measurement of mutual information. In
365 *International Conference on Artificial Intelligence and Statistics*, pages 875–884. PMLR, 2020.
- 366 [27] Pankaj Mehta and David J. Schwab. An exact mapping between the variational renormalization group and
367 deep learning. *arXiv preprint arXiv:1410.3831*, 2014.
- 368 [28] Julian D Olden and Donald A Jackson. Illuminating the “black box”: a randomization approach for
369 understanding variable contributions in artificial neural networks. *Ecological modelling*, 154(1-2):135–150,
370 2002.
- 371 [29] Liam Paninski. Estimation of entropy and mutual information. *Neural computation*, 15(6):1191–1253,
372 2003.
- 373 [30] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-
374 propagating errors. *Nature*, 323:533–536, October 1986.
- 375 [31] Ruslan Salakhutdinov and Geoffrey Hinton. Deep boltzmann machines. In *Artificial intelligence and
376 statistics*, pages 448–455. PMLR, 2009.
- 377 [32] Wojciech Samek, Thomas Wiegand, and Klaus-Robert Müller. Explainable artificial intelligence: Under-
378 standing, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296*, 2017.
- 379 [33] Andrew Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan Tracey, and
380 David Cox. On the information bottleneck theory of deep learning. In *International Conference on
381 Representation Learning*, 2018.
- 382 [34] Ohad Shamir, Sivan Sabato, and Naftali Tishby. Learning and generalization with the information
383 bottleneck. *Theoretical Computer Science*, 411(29-30):2696–2711, 2010.
- 384 [35] Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information.
385 *arXiv preprint arXiv:1703.00810*, 2017.
- 386 [36] Noam Slonim. *The information bottleneck: Theory and applications*. PhD thesis, Citeseer, 2002.
- 387 [37] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv
388 preprint physics/0004057*, 2000.
- 389 [38] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE
390 Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- 391 [39] Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*,
392 10(5):988–999, 1999.
- 393 [40] Larry Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.
- 394 [41] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
395 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 396 [42] Sho Yaida. Non-gaussian processes and neural networks at finite widths. *arXiv preprint arXiv:1910.00019*,
397 2019.
- 398 [43] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European
399 conference on computer vision*, pages 818–833. Springer, 2014.
- 400 [44] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
401 learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

402 **Checklist**

- 403 1. For all authors...
- 404 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-
405 tions and scope? [Yes]
- 406 (b) Did you describe the limitations of your work? [Yes] see Section 5
- 407 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 408 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 409 2. If you are including theoretical results...
- 410 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 411 (b) Did you include complete proofs of all theoretical results? [Yes]
- 412 3. If you ran experiments...
- 413 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
414 results (either in the supplemental material or as a URL)? [Yes] see the URL in Appendix G
- 415 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
416 [Yes] see Section 4.1, Appendix G, and Appendix H
- 417 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
418 multiple times)? [Yes] see Section 4
- 419 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
420 internal cluster, or cloud provider)? [Yes] see Appendix G
- 421 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 422 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 423 (b) Did you mention the license of the assets? [Yes]
- 424 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] see
425 Appendix H
- 426 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-
427 ing/curating? [N/A]
- 428 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-
429 tion or offensive content? [N/A]
- 430 5. If you used crowdsourcing or conducted research with human subjects...
- 431 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
432 [N/A]
- 433 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
434 approvals, if applicable? [N/A]
- 435 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
436 participant compensation? [N/A]

437 A Recent works about Gibbs explanations for a hidden layer

438 As a fundamental probabilistic graphic model, the Gibbs distribution (a.k.a., Boltzmann distribution, the energy
439 based model, or the renormalization group) formulates the dependence within X by associating an energy
440 $E(\mathbf{x}; \boldsymbol{\theta})$ to each dependence structure [9].

$$P(X; \boldsymbol{\theta}, \beta) = \frac{1}{Z(\boldsymbol{\theta}, \beta)} \exp[-\beta E(\mathbf{x}; \boldsymbol{\theta})], \quad (21)$$

441 where $E(\mathbf{x}; \boldsymbol{\theta})$ is the energy function, $\boldsymbol{\theta}$ denote the parameters of $E(\mathbf{x}; \boldsymbol{\theta})$, β is the inverse temperature constant.
442 Since β can be absorbed into $\boldsymbol{\theta}$, $P(X; \boldsymbol{\theta}, \beta)$ can be simplified as

$$P(X; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp[-E(\mathbf{x}; \boldsymbol{\theta})], \quad (22)$$

443 where the partition function⁵ is defined as

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \mathcal{X}} \exp[-E(\mathbf{x}; \boldsymbol{\theta})]. \quad (23)$$

444 The Gibbs distribution has three appealing properties. First, the deterministic energy function $E(\mathbf{x}; \boldsymbol{\theta})$ is a
445 sufficient statistics of $P(X; \boldsymbol{\theta})$. The property allows us to explain a deterministic function, e.g., a hidden layer,
446 in a probabilistic way. Second, a Gibbs distribution can be easily reformulated as various probabilistic models
447 via redefining $E(\mathbf{x}; \boldsymbol{\theta})$, which allows us to clarify the complicated architecture of a hidden layer. For example,
448 if the energy function is defined as the summation of multiple functions, namely $E(\mathbf{x}; \boldsymbol{\theta}) = -\sum_k f_k(\mathbf{x}; \boldsymbol{\theta}_k)$,
449 the Gibbs distribution would be the Product of Experts (PoE) model, *i.e.*, $P(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_k F_k$, where
450 $F_k = \exp[-f_k(\mathbf{x}; \boldsymbol{\theta}_k)]$ and $Z(\boldsymbol{\theta}) = \prod_k Z(\boldsymbol{\theta}_k)$ [13]. Third, the energy minimization is a typical optimization
451 for $\boldsymbol{\theta}$, namely $\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} E(\mathbf{x}; \boldsymbol{\theta})$ [22], which allows us to explain the back-propagation training, as the
452 energy minimization can be implemented by the gradient descent algorithm as long as $E(\mathbf{x}; \boldsymbol{\theta})$ is differentiable.

453 A well-known Gibbs distribution model in machine learning is the Restricted Boltzmann Machines (RBMs)
454 [31, 27]. Though Yaida indirectly proves the distribution of a fully connected layer as a Gibbs distribution [42],
455 and Lin *et al.* clarify certain advantages of DNNs based on the Gibbs distribution [24], there is few work to
456 extend the Gibbs explanation to complicated hidden layers, e.g., fully connected layers and convolutional layers.

457 B The marginal distribution of the MLP

458 Since the entire architecture of the MLP = $\{\mathbf{x}, \mathbf{t}_1, \mathbf{t}_2, \hat{\mathbf{y}}\}$ corresponds to a joint distribution

$$P(\hat{Y}, T_2, T_1 | X) = P(\hat{Y} | T_2) P(T_2 | T_1) P(T_1 | X), \quad (24)$$

459 the marginal distribution $P(\hat{Y} | X)$ can be formulated as

$$\begin{aligned} P_{\hat{Y}|X}(l|\mathbf{x}) &= \sum_{k=1}^K \sum_{n=1}^N P(\hat{Y} = l, T_2 = k, T_1 = n | X = \mathbf{x}) \\ &= \sum_{k=1}^K \sum_{n=1}^N P_{\hat{Y}|T_2}(l|k) P_{T_2|T_1}(k|n) P_{T_1|X}(n|\mathbf{x}). \end{aligned} \quad (25)$$

460 Based on the definition of the Gibbs probability measure (Equation 5), we have

$$P_{T_1|X}(n|\mathbf{x}) = \frac{1}{Z_{T_1}} \exp(t_{1n}) = \frac{1}{Z_{T_1}} \exp[\sigma_1(\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle)], \quad (26)$$

461 where $\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle = \sum_{m=1}^M \omega_{mn}^{(1)} \cdot x_m + b_{1n}$. Similarly, we have

$$P_{T_2|T_1}(k|n) = \frac{1}{Z_{T_2}} \exp(t_{2k}) = \frac{1}{Z_{T_2}} \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)], \quad (27)$$

462 where $\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle = \sum_{n=1}^N \omega_{nk}^{(2)} \cdot t_{1n} + b_{2k}$. Thus we have

$$\begin{aligned} &\sum_{n=1}^N P_{T_2|T_1}(k|n) P_{T_1|X}(n|\mathbf{x}) \\ &= \frac{1}{Z_{T_2}} \frac{1}{Z_{T_1}} \sum_{n=1}^N \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)] \exp[\sigma_1(\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle)]. \end{aligned} \quad (28)$$

⁵We only consider the discrete case in the paper.

463 Since $\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle = \sum_{n=1}^N \boldsymbol{\omega}_{nk}^{(2)} \cdot \mathbf{t}_{1n} + b_{2k}$ is a constant with respect to n , we have

$$\begin{aligned} & \sum_{n=1}^N P_{T_2|T_1}(k|n) P_{T_1|X}(n|\mathbf{x}) \\ &= \frac{1}{Z_{T_2}} \frac{1}{Z_{T_1}} \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)] \sum_{n=1}^N \exp[\sigma_1(\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle)]. \end{aligned} \quad (29)$$

464 In addition, $\sum_{n=1}^N \exp[\sigma_1(\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle)] = Z_{T_1}$, thus we have

$$\sum_{n=1}^N P_{T_2|T_1}(k|n) P_{T_1|X}(n|\mathbf{x}) = \frac{1}{Z_{T_2}} \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)]. \quad (30)$$

465 Therefore, we can simplify $P_{\hat{Y}|X}(l|\mathbf{x})$ as

$$\begin{aligned} P_{\hat{Y}|X}(l|\mathbf{x}) &= \sum_{k=1}^K P_{\hat{Y}|T_2}(l|k) \sum_{n=1}^N P_{T_2|T_1}(k|n) P_{T_1|X}(n|\mathbf{x}) \\ &= \sum_{k=1}^K P_{\hat{Y}|T_2}(l|k) \frac{1}{Z_{T_2}} \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)]. \end{aligned} \quad (31)$$

466 Since $P_{\hat{Y}|T_2}(l|k) = \frac{1}{Z_{\hat{Y}}} \exp[\sigma_3(\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle)]$ and $\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle = \sum_{k=1}^K \boldsymbol{\omega}_{lk}^{(3)} f_{2k} + b_{yl}$ is a constant with respect
467 to k , we can derive

$$P_{\hat{Y}|X}(l|\mathbf{x}) = P_{\hat{Y}|T_2}(l|k) \sum_{k=1}^K \frac{1}{Z_{T_2}} \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)]. \quad (32)$$

468 Since $Z_{T_2} = \sum_{k=1}^K \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)]$ is also constant to k ,

$$\begin{aligned} P_{\hat{Y}|X}(l|\mathbf{x}) &= P_{\hat{Y}|T_2}(l|k) \frac{1}{Z_{T_2}} \sum_{k=1}^K \exp[\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)]. \\ &= P_{\hat{Y}|T_2}(l|k) = \frac{1}{Z_{F_Y}} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle]. \end{aligned} \quad (33)$$

469 In addition, since $\mathbf{t}_2 = \{t_{2k}\}_{k=1}^K = \{\sigma_2(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle)\}_{k=1}^K$, we can extend $P_{\hat{Y}|X}(l|\mathbf{x})$ as

$$\begin{aligned} P_{\hat{Y}|X}(l|\mathbf{x}) &= P_{\hat{Y}|F_2}(l|k) = \frac{1}{Z_{F_Y}} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle] \\ &= \frac{1}{Z_{\hat{Y}}} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \begin{pmatrix} \sigma_2(\langle \boldsymbol{\omega}_1^{(2)}, \mathbf{t}_1 \rangle) \\ \vdots \\ \sigma_2(\langle \boldsymbol{\omega}_K^{(2)}, \mathbf{t}_1 \rangle) \end{pmatrix} \rangle]. \end{aligned} \quad (34)$$

470 Since $\mathbf{t}_1 = \{t_{1n}\}_{n=1}^N = \{\sigma_1(\langle \boldsymbol{\omega}_n^{(1)}, \mathbf{x} \rangle)\}_{n=1}^N$, we can further extend $P_{\hat{Y}|X}(l|\mathbf{x})$ as

$$\begin{aligned} P_{\hat{Y}|X}(l|\mathbf{x}) &= \frac{1}{Z_{\hat{Y}}} \exp[\langle \boldsymbol{\omega}_l^{(3)}, \begin{pmatrix} \sigma_2(\langle \boldsymbol{\omega}_1^{(2)}, \begin{pmatrix} \sigma_1(\langle \boldsymbol{\omega}_1^{(1)}, \mathbf{x} \rangle) \\ \vdots \\ \sigma_1(\langle \boldsymbol{\omega}_N^{(1)}, \mathbf{x} \rangle) \end{pmatrix} \rangle) \\ \vdots \\ \sigma_2(\langle \boldsymbol{\omega}_K^{(2)}, \begin{pmatrix} \sigma_1(\langle \boldsymbol{\omega}_1^{(1)}, \mathbf{x} \rangle) \\ \vdots \\ \sigma_1(\langle \boldsymbol{\omega}_N^{(1)}, \mathbf{x} \rangle) \end{pmatrix} \rangle) \end{pmatrix} \rangle] \\ &= \frac{1}{Z_{\text{MLP}}(\mathbf{x})} \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))]. \end{aligned} \quad (35)$$

471 Overall, we prove $P_{\hat{Y}|X}(l|\mathbf{x})$ as the Gibbs distribution expressed as

$$P_{\hat{Y}|X}(l|\mathbf{x}) = \frac{1}{Z_{\text{MLP}}(\mathbf{x})} \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))]. \quad (36)$$

472 where $E_l(\mathbf{x}) = -g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))$ is the energy function and the partition function

$$\begin{aligned} Z_{\text{MLP}}(\mathbf{x}) &= \sum_{l=1}^L \sum_{k=1}^K \sum_{t=1}^T P(\hat{Y}, T_2, T_1 | X = \mathbf{x}) \\ &= \sum_{l=1}^L \exp[g_l(\mathbf{t}_2(\mathbf{t}_1(\mathbf{x})))]. \end{aligned} \quad (37)$$

473 C The proof of Theorem 2

474 Based on the definition of the cross entropy, ℓ_{CE} can be formulated as

$$\ell_{\text{CE}} = - \sum_{l=1}^L P_{Y|X}(l|\mathbf{x}) \log P_{\hat{Y}|X}(l|\mathbf{x}). \quad (38)$$

475 where $P_{\hat{Y}|X}(l|\mathbf{x})$ is the output of the MLP, and $P_{Y|X}(l|\mathbf{x})$ is the one-hot probability of \mathbf{x} given the label y , i.e.,

$$P_{Y|X}(l|\mathbf{x}) = \begin{cases} 1 & \text{for } l = y \\ 0 & \text{for } l \neq y \end{cases} \quad (39)$$

476 The derivative of ℓ_{CE} with respect to $P_{\hat{Y}|X}(l|\mathbf{x})$ is

$$\frac{\partial \ell_{\text{CE}}}{\partial P_{\hat{Y}|X}(l|\mathbf{x})} = - \frac{P_{Y|X}(l|\mathbf{x})}{P_{\hat{Y}|X}(l|\mathbf{x})}. \quad (40)$$

477 Since $P_{\hat{Y}|X}(l|\mathbf{x})$ can be expressed as

$$P_{\hat{Y}|X}(l|\mathbf{x}) = \frac{1}{Z_{\text{MLP}}(\mathbf{x})} \exp[g_l(\mathbf{t}_2 \mathbf{t}_1(\mathbf{x}))], \quad (41)$$

478 the derivative of $P_{\hat{Y}|X}(z|\mathbf{x})$ with respect to $g_l(\mathbf{t}_2 \mathbf{t}_1(\mathbf{x}))$ is

$$\frac{\partial P_{\hat{Y}|X}(z|\mathbf{x})}{\partial g_l} = \frac{\frac{1}{Z_{\text{MLP}}} \exp(g_z)}{\frac{\partial g_l}{\partial g_l}} = \begin{cases} P_{\hat{Y}|X}(l|\mathbf{x}) \cdot [1 - P_{\hat{Y}|X}(l|\mathbf{x})] & \text{for } z = l \\ -P_{\hat{Y}|X}(l|\mathbf{x}) \cdot P_{\hat{Y}|X}(z|\mathbf{x}) & \text{for } z \neq l \end{cases}. \quad (42)$$

479 Overall, the derivative of ℓ_{CE} with respect to g_l can be expressed as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}}{\partial g_l} &= \sum_{z=1}^L \frac{\partial \ell_{\text{CE}}}{\partial P_{\hat{Y}|X}(z|\mathbf{x})} \frac{\partial P_{\hat{Y}|X}(z|\mathbf{x})}{\partial g_l} \\ &= -P_{Y|X}(l|\mathbf{x}) (1 - P_{\hat{Y}|X}(l|\mathbf{x})) + \sum_{z \neq l} P_{Y|X}(z|\mathbf{x}) P_{\hat{Y}|X}(l|\mathbf{x}) \\ &= P_{\hat{Y}|X}(l|\mathbf{x}) - P_{Y|X}(l|\mathbf{x}). \end{aligned} \quad (43)$$

480 Since $g_l = \langle \boldsymbol{\omega}_l^{(3)}, \mathbf{t}_2 \rangle = \sum_{k=1}^K \omega_{kl}^{(3)} \cdot t_{2k} + b_{yl}$, the derivative of ℓ_{CE} with respect to $\omega_{kl}^{(3)}$ can be expressed as

$$\frac{\partial \ell_{\text{CE}}}{\partial \omega_{kl}^{(3)}} = \frac{\partial \ell_{\text{CE}}}{\partial g_l} \frac{\partial g_l}{\partial \omega_{kl}^{(3)}} = [P_{\hat{Y}|X}(l|\mathbf{x}) - P_{Y|X}(l|\mathbf{x})] t_{2k}. \quad (44)$$

481 Similarly, the derivative of ℓ_{CE} with respect to $\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle$ can be expressed as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}}{\partial \langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle} &= \sum_{l=1}^L \frac{\partial \ell_{\text{CE}}}{\partial g_l} \frac{\partial g_l}{\partial t_{2k}} \frac{\partial t_{2k}}{\partial \langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle} \\ &= \sum_{l=1}^L [P_{\hat{Y}|X}(l|\mathbf{x}) - P_{Y|X}(l|\mathbf{x})] \omega_{kl}^{(3)} \sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle). \end{aligned} \quad (45)$$

482 Since $\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle = \sum_{n=1}^N \omega_{nk}^{(2)} \cdot t_{1n} + b_{2k}$, the derivative of ℓ with respect to $\omega_{nk}^{(2)}$ can be expressed as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}}{\partial \omega_{nk}^{(2)}} &= \frac{\partial \ell_{\text{CE}}}{\partial \langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle} \frac{\partial \langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle}{\partial \omega_{nk}^{(2)}} \\ &= \sum_{l=1}^L [P_{\hat{Y}|X}(l|\mathbf{x}) - P_{Y|X}(l|\mathbf{x})] \omega_{kl}^{(3)} \sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \mathbf{t}_1 \rangle) t_{1n} \end{aligned} \quad (46)$$

483 Similarly, the derivative of ℓ_{CE} with respect to $\langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle$ can be expressed as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}}{\partial \langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle} &= \sum_{k=1}^K \frac{\partial \ell_{\text{CE}}}{\partial \langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle} \frac{\partial \langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle}{\partial t_{1n}} \frac{\partial t_{1n}}{\partial \langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle} \\ &= \sum_{k=1}^K \sum_{l=1}^L [P_{\hat{Y}|X}(l|\boldsymbol{x}) - P_{Y|X}(l|\boldsymbol{x})] \omega_{kl}^{(3)} \sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle) \omega_{nk}^{(2)} \sigma_1'(\langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle). \end{aligned} \quad (47)$$

484 Since $\langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle = \sum_{m=1}^M \omega_{mn}^{(1)} \cdot x_m + b_{1n}$, the derivative of ℓ_{CE} with respect to $\omega_{mn}^{(1)}$ can be expressed as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}}{\partial \omega_{mn}^{(1)}} &= \frac{\partial \ell_{\text{CE}}}{\partial \langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle} \frac{\partial \langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle}{\partial \omega_{mn}^{(1)}} \\ &= \sum_{k=1}^K \sum_{l=1}^L [P_{\hat{Y}|X}(l|\boldsymbol{x}) - P_{Y|X}(l|\boldsymbol{x})] \omega_{kl}^{(3)} \sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle) \omega_{nk}^{(2)} \sigma_1'(\langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle) x_m. \end{aligned} \quad (48)$$

485 Overall, the derivative of ℓ_{CE} with respect to the weight in each layer is summarized as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}}{\partial \omega_{kl}^{(3)}} &= [P_{\hat{Y}|X}(l|\boldsymbol{x}) - P_{Y|X}(l|\boldsymbol{x})] t_{2k} \\ \frac{\partial \ell_{\text{CE}}}{\partial \omega_{nk}^{(2)}} &= \sum_{l=1}^L [P_{\hat{Y}|X}(l|\boldsymbol{x}) - P_{Y|X}(l|\boldsymbol{x})] \omega_{kl}^{(3)} \sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle) t_{1n} \\ \frac{\partial \ell_{\text{CE}}}{\partial \omega_{mn}^{(1)}} &= \sum_{k=1}^K \sum_{l=1}^L [P_{\hat{Y}|X}(l|\boldsymbol{x}) - P_{Y|X}(l|\boldsymbol{x})] \omega_{kl}^{(3)} \sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle) \omega_{nk}^{(2)} \sigma_1'(\langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle) x_m. \end{aligned} \quad (49)$$

486 Based on the above three equations, we can reformulate the derivatives as

$$\begin{aligned} \frac{\partial \ell_{\text{CE}}^*}{\partial \omega_{kl}^{(3)}} &= [P_{\hat{Y}|X}(l|\boldsymbol{x}) - P_{Y|X}(l|\boldsymbol{x})] \cdot t_{2k}, \\ \frac{\partial \ell_{\text{CE}}^\circ}{\partial \omega_{nk}^{(2)}} &= \sum_{l=1}^L \frac{\partial \ell_{\text{CE}}^*}{\partial \omega_{kl}^{(3)}} \cdot \omega_{kl}^{(3)} \cdot \frac{\sigma_2'(\langle \boldsymbol{\omega}_k^{(2)}, \boldsymbol{t}_1 \rangle)}{t_{2k}} \cdot t_{1n} \\ \frac{\partial \ell_{\text{CE}}^\circ}{\partial \omega_{mn}^{(1)}} &= \sum_{k=1}^K \frac{\partial \ell_{\text{CE}}^\circ}{\partial \omega_{nk}^{(2)}} \cdot \omega_{nk}^{(2)} \cdot \frac{\sigma_1'(\langle \boldsymbol{\omega}_n^{(1)}, \boldsymbol{x} \rangle)}{t_{1n}} \cdot x_m. \end{aligned} \quad (50)$$

487 The above three equations indicates that $\frac{\partial \ell_{\text{CE}}^*}{\partial \omega_{kl}^{(3)}}$ is a function of $P_{Y|X}(l|\boldsymbol{x})$, $\frac{\partial \ell_{\text{CE}}^\circ}{\partial \omega_{nk}^{(2)}}$ is a function of $\frac{\partial \ell_{\text{CE}}^*}{\partial \omega_{kl}^{(3)}}$, and

488 $\frac{\partial \ell_{\text{CE}}^\circ}{\partial \omega_{mn}^{(1)}}$ is a function of $\frac{\partial \ell_{\text{CE}}^\circ}{\partial \omega_{nk}^{(2)}}$. In addition, the back-propagation algorithm shows that

$$\begin{aligned} \omega_{mn}^{(1)}(s+1) &= \omega_{mn}^{(1)}(s) - \alpha \frac{\partial \ell_{\text{CE}}}{\partial \omega_{mn}^{(1)}(s)} \\ \omega_{nk}^{(2)}(s+1) &= \omega_{nk}^{(2)}(s) - \alpha \frac{\partial \ell_{\text{CE}}}{\partial \omega_{nk}^{(2)}(s)} \\ \omega_{kl}^{(3)}(s+1) &= \omega_{kl}^{(3)}(s) - \alpha \frac{\partial \ell_{\text{CE}}}{\partial \omega_{kl}^{(3)}(s)} \end{aligned} \quad (51)$$

489 where α is the learning rate and s denotes the index of the s th learning iteration. Therefore, $\omega(s+1)$ is
490 determined by all the previous gradients $\{\frac{\partial \ell_{\text{CE}}}{\partial \omega(s)}\}_{s=1}^S$ as $\omega(0)$ is randomly initialized and α is a constant.

491 Definition 1 indicates that the weights define the sample space Ω_{T_i} , thus we can derive that the gradients $\frac{\partial \ell_{\text{CE}}}{\partial \omega^{(i)}}$
492 determine Ω_{T_i} . As a result, Ω_{T_i} is a function of $\Omega_{T_{i+1}}$ and $\Omega_{\hat{Y}}$ is a function of $P(Y|X)$. Based on Definition
493 2, we can further derive that T_i is a function of T_{i+1} and \hat{Y} is a function of Y , i.e., $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$.

494 **D The proof of $H(Y) = I(X; Y)$**

495 Given a training sample \boldsymbol{x}^j and the corresponding label y^j , the target distribution $P_{Y|X}(y^j|\boldsymbol{x}^j)$ is commonly
496 formulated as the one-hot format, i.e.,

$$P_{Y|X}(l|\boldsymbol{x}^j) = \begin{cases} 1 & \text{for } l = y^j \\ 0 & \text{for } l \neq y^j \end{cases} \quad (52)$$

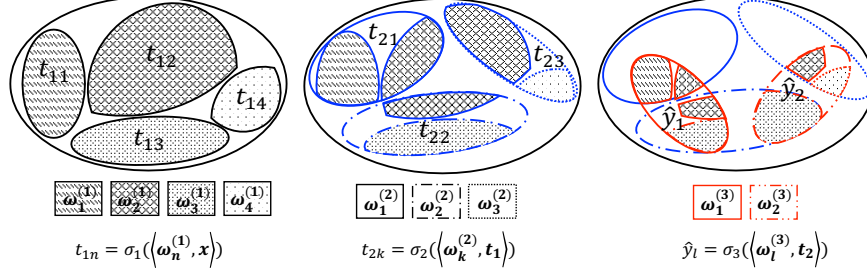


Figure 6: The graphical explanation for Corollary 1 based on the MLP = $\{\mathbf{x}, \mathbf{t}_1, \mathbf{t}_2, \hat{\mathbf{y}}\}$. The largest oval represents of the input \mathbf{x} , and each small shape indicates the representation capacity of a single feature. For example, if t_{12} is the largest activation, then the feature $\omega_2^{(1)}$ has the largest cross-correlation to \mathbf{x} , i.e., $\omega_2^{(1)}$ has largest representation capacity. Therefore, $\{\omega_n^{(1)}\}_{n=1}^4$ can be viewed as a representation of \mathbf{x} , and the representation capacity of $\{\omega_n^{(1)}\}_{n=1}^4$ is measured by $\{t_{1n}\}_{n=1}^4$, which is visualized by the left figure. The blue ovals indicates the representation capacity of the three features $\{\omega_k^{(2)}\}_{k=1}^3$ generated by combining the four features $\{\omega_n^{(1)}\}_{n=1}^4$. The two red ovals indicates the representation capacity of the two features $\{\omega_l^{(3)}\}_{l=1}^2$ generated by combining the three features $\{\omega_k^{(2)}\}_{k=1}^3$.

497 As a result, the conditional entropy $H(Y|X)$ can be formulated as

$$H(Y|X) = - \sum_{(\mathbf{x}^j, \mathbf{y}^j) \in \mathcal{D}} P_{X,Y}(\mathbf{x}^j, \mathbf{y}^j) \log P_{Y|X}(\mathbf{y}^j | \mathbf{x}^j) = 0. \quad (53)$$

498 Therefore, we can derive $H(Y) = I(X; Y)$ because $H(Y) = H(Y|X) + I(X; Y)$.

499 E The detailed derivations and explanations for Corollary 1

500 Definition 1 indicates that $\mathbf{t}_1 = \{t_{1n} = \sigma_1(\langle \omega_n^{(1)}, \mathbf{x} \rangle)\}_{n=1}^N$ defines N features of \mathbf{x} , namely $\{\omega_n^{(1)}\}_{n=1}^N$,
 501 thus $\{\omega_n^{(1)}\}_{n=1}^N$ can be viewed as a representation of \mathbf{x} . In addition, $\{t_{1n}\}_{n=1}^N$ measures the cross-correlation
 502 between $\{\omega_n^{(1)}\}_{n=1}^N$ and \mathbf{x} , (i.e., if $\omega_n^{(1)}$ describes \mathbf{x} more accurately and comprehensively, then t_{1n} is larger.),
 503 thus $\{t_{1n}\}_{n=1}^N$ quantifies the representation capacity of $\{\omega_n^{(1)}\}_{n=1}^N$. For example, in Figure 6 (Left), \mathbf{t}_1 defines
 504 4 features to describe \mathbf{x} and t_{12} is the largest activation, thus $\omega_2^{(1)}$ has the largest representation capacity of \mathbf{x} .

505 During inference, the second hidden layer \mathbf{t}_2 will process $\{t_{1n}\}_{n=1}^N$, and $t_{2k} = \sigma_2(\langle \omega_k^{(2)}, \mathbf{t}_1 \rangle)$ can be explained
 506 to generating a new feature via combining all the features $\{\omega_n^{(1)}\}_{n=1}^N$, i.e.,

$$\{\omega_{1k}^{(2)} \otimes \omega_1^{(1)}, \dots, \omega_{Nk}^{(2)} \otimes \omega_N^{(1)}\}. \quad (54)$$

507 Since the new feature is the linear combination of $\{\omega_n^{(1)}\}_{n=1}^N$, it can be simply noted as

$$\{\omega_{1k}^{(2)}, \dots, \omega_{Nk}^{(2)}\} = \omega_k^{(2)}, \quad (55)$$

508 and the representation capacity of the new feature $\omega_k^{(2)}$ is

$$t_{2k} = \omega_{1k}^{(2)} \cdot t_{11} + \dots + \omega_{Nk}^{(2)} \cdot t_{1N} \quad (56)$$

509 For example, if $N = 4$ and $K = 3$, the representation capacity of the three new features is visualized by Figure
 510 6 (Middle). Similarly, $\hat{\mathbf{y}}$ generates L new features via combining all the features $\{\omega_k^{(2)}\}_{k=1}^K$.

$$\{\omega_{1l}^{(3)} \otimes \omega_1^{(2)}, \dots, \omega_{Kl}^{(3)} \otimes \omega_K^{(2)}\}. \quad (57)$$

511 Since the new feature is the linear combination of $\{\omega_k^{(2)}\}_{k=1}^K$, it can be simply noted as

$$\{\omega_{1l}^{(3)}, \dots, \omega_{Kl}^{(3)}\} = \omega_l^{(3)}, \quad (58)$$

512 and the representation capacity of the new feature $\omega_l^{(3)}$ is

$$\hat{y}_l = \omega_{1l}^{(3)} \cdot t_{21} + \dots + \omega_{Kl}^{(3)} \cdot t_{2K} \quad (59)$$

513 For example, if $L = 2$, the representation capacity of the two new features is visualized by Figure 6 (Right).

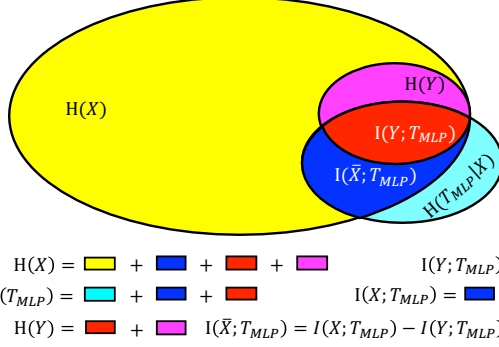


Figure 7: The Venn diagram of $H(X)$, $H(Y)$, and $I(X; T_{MLP})$.

514 Overall, the inference phase is a procedure of feature combination, *i.e.*, $\omega_l^{(3)}$ is a combination of $\{\omega_k^{(2)}\}_{k=1}^K$,
515 and $\omega_k^{(2)}$ is a combination of $\{\omega_n^{(1)}\}_{n=1}^N$. Theorem 2 proves that the layer closer to output has more information
516 of labels, *i.e.*, $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$, during training. Since the weights are fixed after training, the sample space
517 and the distribution of hidden layers are fixed after training. Therefore, the information of Y transferred into
518 hidden layers during training will retain there after training (*i.e.*, during inference), *i.e.*, $T_1 \leftarrow T_2 \leftarrow \hat{Y} \leftarrow Y$
519 characterizes the information flow of Y in the MLP in both training and inference phases.

520 For example, Figure 6 (Right) shows that the representation capacity of $\omega_1^{(3)}$ is the weighted combination of t_{11} ,
521 t_{12} , and t_{13} , and the representation capacity of $\omega_2^{(3)}$ is the weighted combination of t_{12} , t_{13} , and t_{14} . Therefore,
522 $\omega_2^{(1)}$ and $\omega_3^{(1)}$ exist in both classes, *i.e.*, the low-level features in \mathbf{t}_1 do not represent too much information of
523 the labels, though we combine low-level features to generate high-level features for representing labels.

524 F The proof of Corollary 2

525 Based on the property of mutual information, we have

$$\begin{aligned}
 H(X) &= H(X|Y) + I(X; Y) \\
 &= H(X|Y) + H(Y) \quad (\text{Appendix D}) \\
 &= H(\bar{X}) + H(Y)
 \end{aligned} \tag{60}$$

526 where \bar{X} is the virtual random variable containing all the information of X except Y , namely $H(\bar{X}) = H(X|Y)$.

527 Therefore, $I(X; T_{MLP})$ can be reformulated as

$$I(X; T_{MLP}) = I(\bar{X}; T_{MLP}) + I(Y; T_{MLP}). \tag{61}$$

528 The Venn diagram of $H(X)$, $H(Y)$, and $I(X; T_{MLP})$ are visualized in Figure 7. Corollary 1 indicates that all
529 the information of \bar{X} and Y learned by a MLP retains in T_1 and \hat{Y} , respectively. Therefore, we can derive

$$\begin{aligned}
 I(X; T_{MLP}) &= I(\bar{X}; T_1) + I(Y; \hat{Y}) \\
 I(Y; T_{MLP}) &= I(Y; \hat{Y})
 \end{aligned} \tag{62}$$

530 G Studying non-parametric models for mutual information estimation

531 In this section, we use the synthetic dataset to show that non-parametric models are sensitive to hyper-parameters
532 for mutual information estimation. In addition, we show that the proposed mutual information estimator derives
533 more accurate mutual information estimation than non-parametric models. Furthermore, we demonstrate that
534 one reason for non-parametric models deriving poor mutual information estimation is because activations do not
535 satisfy the *i.i.d.* prerequisite of non-parametric models. The experiment codes are available online⁶.

536 G.1 Non-parametric models are sensitive to hyper-parameters

537 To show non-parametric models being sensitive to hyper-parameters, we choose two commonly used non-
538 parametric models, namely the empirical distribution [35] and KDE [33], to measure the information flow in
539 MLP1 and MLP2 defined in Table 1 on the synthetic dataset.

⁶<https://github.com/Dlib-NeurIPS/Deep-Learning-Information-Theory>

Table 3: The hyper-parameters of empirical distributions and KDE

bs	0.001	0.01	0.1	1.0	2.0	4.0	6.0	8.0
σ_n^2	0.01	0.05	0.1	1.0	2.0	4.0	8.0	16.0

540 The empirical distribution is defined as

$$P(T = n) = \frac{1}{J} \mathbb{1}(\mathbf{t}, \mathbf{l}_n, \mathbf{r}_n) \quad (63)$$

541 where J is the number of samples, n denotes the n th bin, \mathbf{t} denotes an activation vector, \mathbf{l}_n and \mathbf{r}_n are the left
542 and right boundary vectors, respectively. The indicator function $\mathbb{1}(\mathbf{t}, \mathbf{l}_n, \mathbf{r}_n)$ is defined as

$$\mathbb{1}(\mathbf{t}, \mathbf{l}_n, \mathbf{r}_n) = \begin{cases} 1 & \text{for } \mathbf{l}_n \leq \mathbf{t} < \mathbf{r}_n \\ 0 & \text{otherwise} \end{cases} \quad (64)$$

543 Given a specific range, the hyper-parameter of the empirical distribution is the bin size, namely $bs = |\mathbf{r}_n - \mathbf{l}_n|$.
544 Based on the empirical distribution, Tishby *et al.* estimate $I(X; T_i)$ and $I(Y; T_i)$ (see Section 3.2 in [35]).

545 To estimate $I(X; T_i)$ and $I(Y; T_i)$ via KDE, Saxe *et al.* assume that the empirical distribution of input samples
546 is the true distribution and the distribution of a hidden layer is a mixture of Gaussian. In addition, Saxe *et al.*
547 regard a hidden layer as a deterministic function of input samples, thus the Gaussian noise $\mathcal{N}(0, \sigma_n^2)$ is added
548 into activations to avoid infinite mutual information, and $I(X; T_i)$ is estimated as

$$I(X; T_i) \leq -\frac{1}{J} \sum_j \log \frac{1}{J} \sum_{j'} \exp\left(-\frac{\|\mathbf{t}_j^{(i)} - \mathbf{t}_{j'}^{(i)}\|_2^2}{2\sigma_n^2}\right) \quad (65)$$

549 where J is the number of samples, $\mathbf{t}_j^{(i)}$ denote the activations vector of the i th hidden layer in response to the
550 input sample \mathbf{x}^j (see Appendix B.1 in [33]). Therefore, the hyper-parameter of KDE is the noise variance σ_n^2 .

551 Leveraging the same training method in Section 4.1, we achieves 100% training accuracy in MLP1 and MLP2
552 on the synthetic dataset. We specify 8 different values for each hyper-parameter, namely bs and σ_n^2 , in Table 3,
553 and use the empirical distribution and KDE to estimate $I(X; T_i)$ during training MLP1 and MLP2.

554 Figure 8 and 9 show that the empirical distribution is sensitive to the hyper-parameter, namely the bin size bs .
555 Figure 8 shows that $I(X; T_i)$ in different hidden layers of MLP1 converges to 1.5 as bs increases from 0.001 to
556 8.0. Figure 9 shows that $I(X; T_1)$, $I(X; T_2)$, and $I(X; \hat{Y})$ in MLP2 converge to 1.2, 0.8, and 0.7, respectively,
557 as bs increases from 0.001 to 8.0. Notably, since the synthetic dataset only has 2 bits information, $I(X; T_i)$
558 must be smaller than $H(X) = 2$ bits. However, we observe that if $bs < 1.0$, the empirical distribution derives
559 $I(X; T_i) > 2.0$ in both MLP1 and MLP2, thus the empirical distribution cannot correctly estimate $I(X; T_i)$ in
560 MLP1 and MLP2 on the synthetic dataset when $bs < 1.0$.

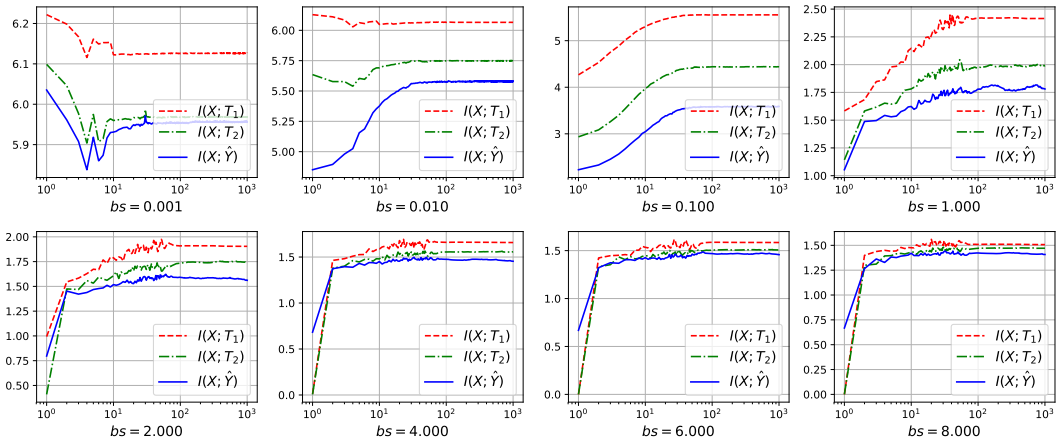


Figure 8: The estimation of $I(X; T_i)$ in MLP1 on the synthetic dataset via the empirical distribution with 8 different bs . All the x-axis index training epochs.

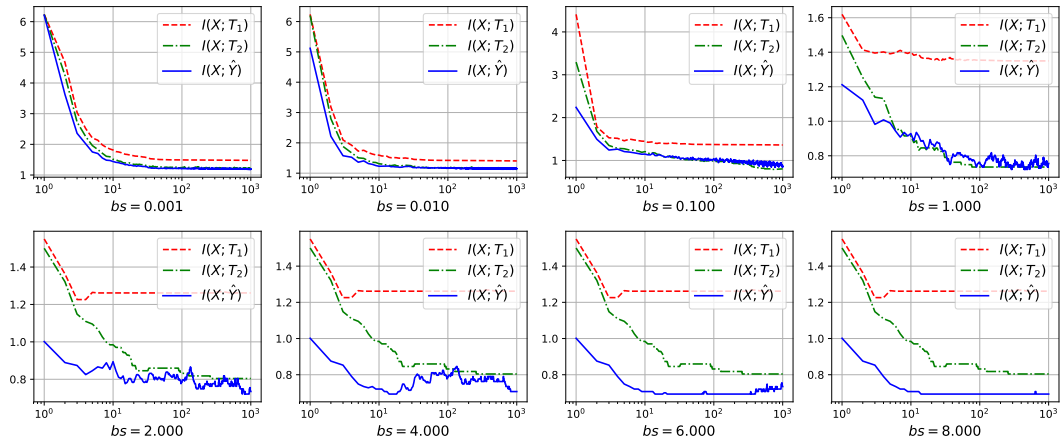


Figure 9: The estimation of $I(X; T_i)$ in MLP2 on the synthetic dataset via the empirical distribution with 8 different bs . All the x-axis index training epochs.

561 Similarly, Figure 10 and 11 show that KDE is also sensitive to the hyper-parameter, namely the noise variance
 562 σ_n^2 . Figure 10 shows that $I(X; T_i)$ in different hidden layers of MLP1 converges to 2.0 as σ_n^2 increases from
 563 0.01 to 16.0. Figure 11 shows that KDE derives different $I(X; T_1)$ and $I(X; T_2)$ in MLP2 given different
 564 σ_n^2 , except $I(X; \hat{Y})$ converges to 1.0, as bs increases from 0.01 to 16.0. Again, since the synthetic dataset
 565 only has 2 bits information, $I(X; T_i)$ must be smaller than $H(X) = 2$. However, we also observe that KDE
 566 derives $I(X; T_i) > 2.0$ when $\sigma_n^2 < 1.0$. Overall, different σ_n^2 make KDE to derive different mutual information
 567 estimations for $I(X; T_i)$ in MLP1 and MLP2 on the synthetic dataset, especially KDE does not correctly
 568 estimate $I(X; T_i)$ in MLP1 and MLP2 when $\sigma_n^2 < 1.0$.

569 In summary, the two non-parametric models are sensitive to hyper-parameters for mutual information estimation.
 570 Especially, since the entropy of the synthetic dataset is known, we can determine which hyper-parameter is
 571 appropriate to estimate the mutual information. However, if the entropy of dataset is unknown, it is very difficult
 572 to choose an appropriate hyper-parameter for non-parametric models to estimate the mutual information.

573 G.2 Comparison to non-parametric models on the synthetic dataset

574 In this section, we compare the proposed mutual information estimator to the empirical distribution and KDE
 575 in MLP1 and MLP2 on the synthetic dataset, and demonstrate that the proposed mutual information estimator
 576 derives more accurate mutual information estimation than non-parametric models. Based on Appendix G.1,
 577 we choose $bs = 2.0$ and $\sigma_n^2 = 2.0$ as the optimal hyper-parameters for the empirical distribution and KDE to
 578 estimate $I(X; T_i)$ and $I(Y; T_i)$. All the training methods are the same as Section 4.1.

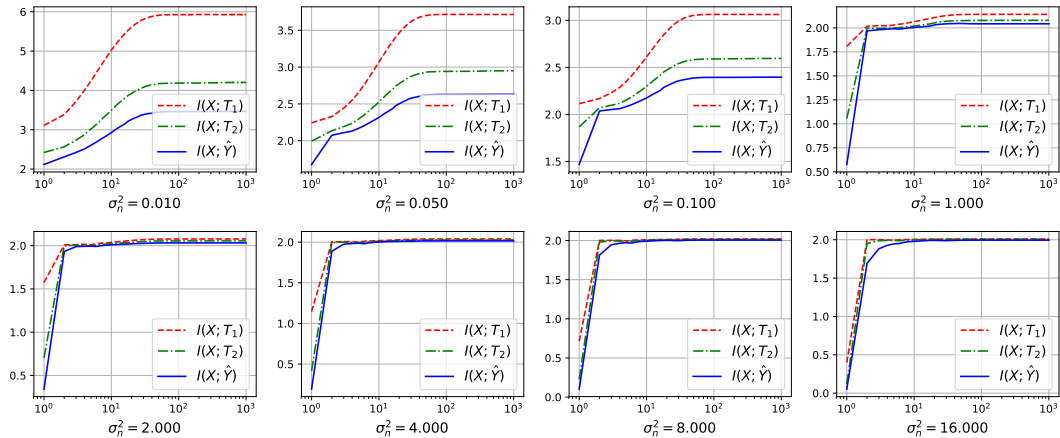


Figure 10: The estimation of $I(X; T_i)$ in MLP1 on the synthetic dataset by KDE with 8 different σ_n^2 . All the x-axis index training epochs.

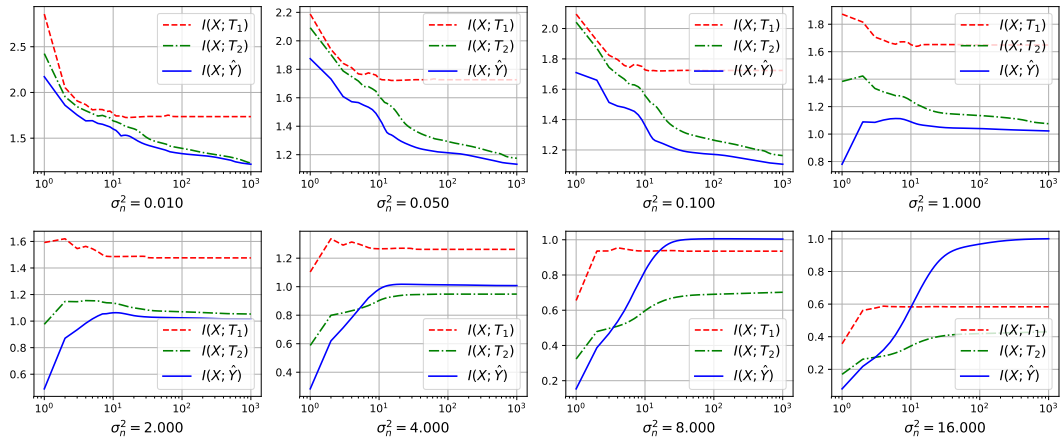


Figure 11: The estimation of $I(X; T_i)$ in MLP2 on the synthetic dataset by KDE with 8 different σ_n^2 . All the x-axis index training epochs.

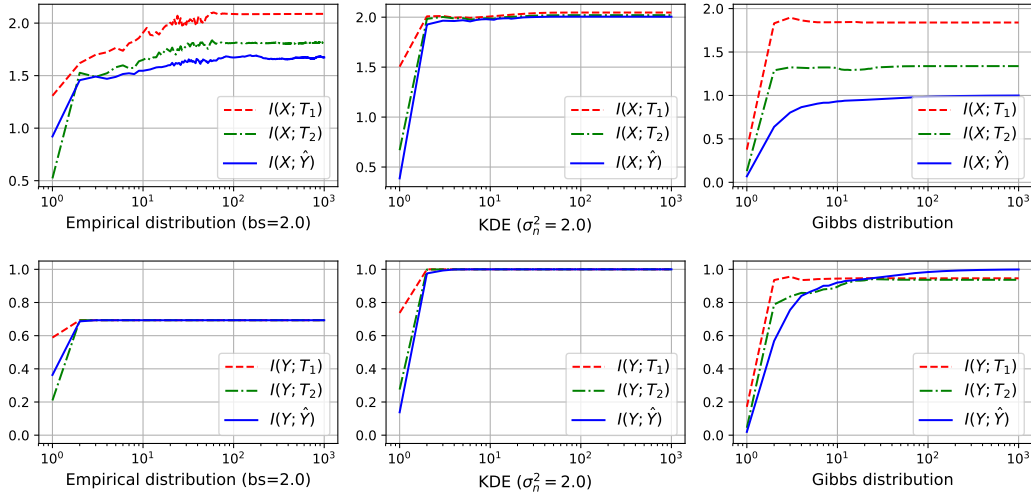


Figure 12: The estimation of $I(X; T_i)$ and $I(Y; T_i)$ in MLP1 based on the three mutual information estimators. All the x-axis index training epochs.

579 Figure 12 shows the estimation of $I(X; T_i)$ and $I(Y; T_i)$ in MLP1 derived by the three methods, namely the
580 empirical distribution ($bs = 2.0$), KDE ($\sigma_n^2 = 2.0$), and the Gibbs distribution. Since \hat{y} only has two nodes,
581 the maximal information of X that \hat{y} can have is 1 bit, i.e., $I(X; \hat{Y}) \leq 1$, based on Definition 1. However,
582 we observe that the empirical distribution derives $I(X; \hat{Y}) > 1.5$ and KDE derives $I(X; \hat{Y}) = 2.0$, thus the
583 empirical distribution and KDE do not accurately estimate $I(X; \hat{Y})$. In addition, since MLP1 correctly predicts
584 all the labels of synthetic images, it should have all the information of the labels. However, we observe that the
585 empirical distribution estimates $I(Y; T_i) = 0.7$ bits, which contradicts the fact. As a comparison, the proposed
586 method based on Gibbs distribution accurately estimate the information flow in MLP1.

587 Figure 13 shows the estimation of $I(X; T_i)$ and $I(Y; T_i)$ in MLP2 derived by the three methods. As shown in
588 Figure 4, MLP2 quickly learns all the features of the synthetic dataset, thus $I(X; T_i)$ should have an increasing
589 trend as training epochs increases. However, $I(X; T_i)$ estimated by the empirical distribution shows a decreasing
590 trend, which contradicts the variation of the weights shown in Figure 4. Therefore, the empirical distribution
591 does not accurately estimate $I(X; T_i)$ in MLP2. In addition, Section 4.2 shows that Tanh hinders t_1 from
592 correctly recognizing the features of input, thus t_1 in MLP2 does not contain too much information of X , i.e.,
593 $I(X; T_1)$ is small. However, KDE estimates $I(X; T_1) > 1.5$, i.e., t_1 in MLP2 has most information of X .
594 Therefore, KDE does not correctly measures the effect of activation functions on the mutual information.

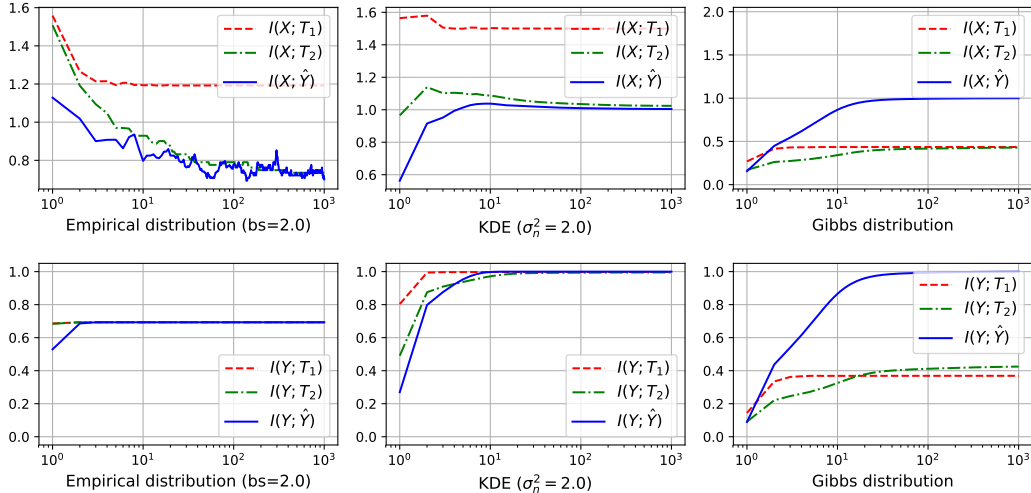


Figure 13: The estimation of $I(X; T_i)$ and $I(Y; T_i)$ in MLP2 based on the three mutual information estimators. All the x-axis index training epochs.

595 G.3 Activations do not satisfy the *i.i.d.* prerequisite of non-parametric models

596 In this section, we demonstrate that one reason for non-parametric models deriving poor mutual information
597 estimation is because activations do not satisfy the *i.i.d.* prerequisite of non-parametric models.

598 Given an input $\mathbf{x} \in \mathbb{R}^M$, we define the corresponding multivariate random variable as $X = [X_1, \dots, X_M]$,
599 where X_m is the scalar-valued random variable of x_m . In the context of frequentist probability, all the parameters
600 of MLPs are viewed as constants, thus the random variable of $\langle \omega_n^{(1)}, \mathbf{x} \rangle = \sum_{m=1}^M \omega_{mn}^{(1)} \cdot x_m + b_{1n}$ is defined
601 as $G_{1n} = \sum_{m=1}^M \omega_{mn}^{(1)} X_m + b_{1n}$, and the random variable of the activation $t_{1n} = \sigma_1(\langle \omega_n^{(1)}, \mathbf{x} \rangle)$ is defined
602 as $T_{1n} = \sigma_1(G_{1n})$. Therefore, the multivariate random variable of $\mathbf{t}_1 = [t_{11}, \dots, t_{1N}]$ can be defined as
603 $T_1 = [T_{11}, \dots, T_{1N}]$. Similarly, we define the multivariate random variable of \mathbf{t}_2 as $T_2 = [T_{21}, \dots, T_{2K}]$ and
604 the multivariate random variable of $\hat{\mathbf{y}}$ as $\hat{Y} = [\hat{Y}_1, \dots, \hat{Y}_L]$.

605 Samples being *i.i.d.* is the prerequisite of applying non-parametric models, e.g. the empirical distribution and
606 KDE, to model the true distribution of a random variable [40]. In the context of MLPs, most previous works
607 regard the activations of a layer as the samples of the random variable of the layer, and use non-parametric
608 models to simulate the distribution of the layer. As a result, activations must be *i.i.d.* samples.

609 Since the necessary condition for samples being *i.i.d.* is the samples being uncorrelated, we can use the sample
610 correlation to examine if activations being *i.i.d.*. More specifically, given two *i.i.d.* input samples \mathbf{x}^j and $\mathbf{x}^{j'}$,
611 the two activation vectors of the i th hidden layers are \mathbf{t}_i^j and $\mathbf{t}_i^{j'}$. If \mathbf{t}_i^j and $\mathbf{t}_i^{j'}$ are *i.i.d.* samples of T_i , the sample
612 correlation $R(\mathbf{t}_i^j, \mathbf{t}_i^{j'})$ must be zero, namely

$$R(\mathbf{t}_i^j, \mathbf{t}_i^{j'}) = \frac{\sum_{n=1}^N (t_{in}^j - \bar{t}_i^j)(t_{in}^{j'} - \bar{t}_i^{j'})}{\sqrt{\sum_{n=1}^N (t_{in}^j - \bar{t}_i^j)^2 \sum_{n=1}^N (t_{in}^{j'} - \bar{t}_i^{j'})^2}} = 0, \quad (66)$$

613 where $\bar{t}_i^j = \frac{1}{N} \sum_{n=1}^N t_{in}^j$, and N is the number of neurons in \mathbf{t}_i .

614 To study the sample correlation between activations given different samples, we use the Adam to train a MLP on
615 the MNIST dataset [20] over 200 epochs with the learning rate $\alpha = 0.0005$. Since the dimension of each image
616 is 28×28 , the number of the input nodes is $M = 784$. In addition, \mathbf{t}_1 , \mathbf{t}_2 , and $\hat{\mathbf{y}}$ have $N = 96$, $K = 32$, and
617 $L = 10$ neurons/nodes, respectively. All the activation functions are Tanh.

618 After training, we derive $R(\mathbf{t}_i^j, \mathbf{t}_i^{j'})$ on 5000 training samples $\{\mathbf{x}^j\}_{j=1}^{5000}$ and show the result in Figure 14. In
619 particular, we rearrange the order of $\{\mathbf{x}^j\}_{j=1}^{5000}$ such that images with the same label have consecutive index, i.e.,
620 images with the label l has the index $[l \times 500, (l + 1) \times 500)$, thus we can easily check the sample correlation
621 between activations with the same label. Figure 14 shows that the sample correlation between activations with
622 the same label becomes larger as the layer is closer to the output. In other words, activations are not *i.i.d.*.
623 Therefore, it is invalid to apply non-parametric models to model the true distribution of all the layers of the MLP,
624 because activations do not satisfy the *i.i.d.* prerequisite of non-parametric models.

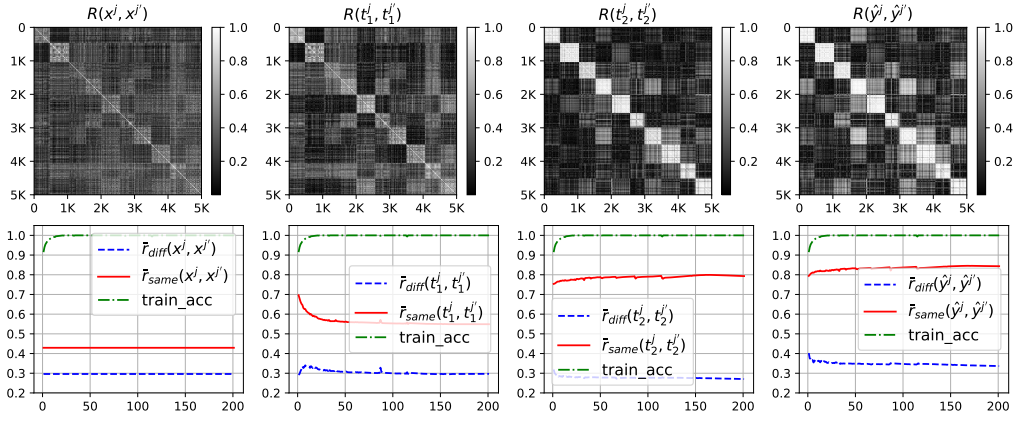


Figure 14: The first row shows that sample correlation between different samples/activations in each layer of the MLP after training. The second row shows the variation of the average sample correlation between different activations with different labels and with the same labels in each layer during training.

625 More specifically, Figure 14 shows that the sample correlation between each pair of training samples $\{\mathbf{x}^j\}_{j=1}^{5000}$
 626 is very small, thus *i.i.d.* can be viewed as a valid assumption for the input samples $\{\mathbf{x}^j\}_{j=1}^{5000}$. However, we
 627 observe an ascending trend for the sample correlation between different activations with the same label as the
 628 layer is closer to the output. For instance, the pixels at the top-left corner of $R(\mathbf{t}_i^j, \mathbf{t}_i^{j'})$ becomes lighter as the
 629 layer is closer to the output, i.e., the sample correlation between the activations with the label 0 becomes larger.

630 In addition, the second row of Figure 14 also show the ascending trend, i.e., $\bar{r}_{\text{same}}(\mathbf{t}_1^j, \mathbf{t}_1^{j'})$, $\bar{r}_{\text{same}}(\mathbf{t}_2^j, \mathbf{t}_2^{j'})$, and
 631 $\bar{r}_{\text{same}}(\hat{\mathbf{y}}^j, \hat{\mathbf{y}}^{j'})$ converge to 0.55, 0.79, and 0.84, respectively, where $\bar{r}_{\text{same}}(\mathbf{t}_i^j, \mathbf{t}_i^{j'})$ denotes the average sample
 632 correlation of $\{\mathbf{t}_i^j\}_{j=1}^{5000}$ with the same label in the i th hidden layer.

633 As a comparison, Figure 14 shows that the sample correlation of activations with different labels being relatively
 634 stable in different layers, because $\bar{r}_{\text{diff}}(\mathbf{t}_1^j, \mathbf{t}_1^{j'})$, $\bar{r}_{\text{diff}}(\mathbf{t}_2^j, \mathbf{t}_2^{j'})$, and $\bar{r}_{\text{diff}}(\hat{\mathbf{y}}^j, \hat{\mathbf{y}}^{j'})$ converge to 0.29, 0.27, and 0.33,
 635 respectively, where $\bar{r}_{\text{diff}}(\mathbf{t}_i^j, \mathbf{t}_i^{j'})$ denotes the average sample correlation of $\{\mathbf{t}_i^j\}_{j=1}^{5000}$ with different labels.

636 In summary, the sample correlation of activations with the same label becomes larger as the layer is closer to the
 637 output, thus activations being *i.i.d.* is not valid for all the layers of the MLP. As a result, non-parametric models,
 638 e.g., the empirical distribution and KDE, cannot correctly simulate the true distribution of all the layers, thus
 639 they are invalid for estimating the mutual information between each layer and dataset.

640 H Experiments on benchmark dataset

641 To further demonstrate the information theoretic explanations for DNNs, we design more complicated neural
 642 networks and conduct experiments on the benchmark MNIST and Fashion-MNIST (abbr. FMNIST) dataset. The
 643 experiment codes are also available online⁷.

644 H.1 Experiments on the MNIST dataset

645 We design three MLPs, namely MLP4, MLP5, and MLP6, and summarize the architectures of the three MLPs in
 646 Table 4. We train the three MLPs on the MNIST dataset by Adam [15] over 500 epochs with the learning rate
 647 $\alpha = 0.0005$. Based on the mutual information estimator proposed in Section 4.1, we measure the information
 648 flow in the three MLPs during 500 training epochs.

649 In Figure 15, we observe that the information flow of X in the three MLPs does not satisfy the Markov chain,
 650 namely Equation (2), proposed by previous works, i.e., we further confirm that Equation (2) does not fully
 651 characterize the information flow of X , especially when taking into account of the back-propagation training.

652 Moreover, the second and the third row of Figure 15 show $I(\bar{X}; T_1) \geq I(\bar{X}; T_2) \geq I(\bar{X}; \hat{Y})$ and $I(Y; T_1) \leq$
 653 $I(Y; T_2) \geq I(Y; \hat{Y})$ in all the three MLPs, which further validate that Corollary 1, i.e., Equation (14), correctly
 654 characterizes the information flow in MLPs.

655 The last row of Figure 15 shows that $I(X; T_{\text{MLP}}) > H(Y)$ and $I(Y; T_{\text{MLP}}) = H(Y)$ for most epochs in all the
 656 three MLPs. Though $H(X)$ is unknown for the MNIST dataset, we still can conclude that the three MLPs form
 657 three compressed representations of the data while preserve all the information of the labels. Hence, Figure 15
 658 further confirms that a MLP satisfies the IB principle no matter what the architecture of the MLP is.

⁷<https://github.com/Dlib-NeurIPS/Deep-Learning-Information-Theory>

Table 4: The number of neurons(nodes) and the activation function in MLP4 - MLP6

	x	t_1	t_2	\hat{y}	$\sigma(\cdot)$
MLP4	784 (28×28)	96	32	10	$\text{ReLU}(z) = \max(0, z)$
MLP5	784 (28×28)	96	32	10	$\text{Tanh}(z) = (e^z - e^{-z}) / (e^z + e^{-z})$
MLP6	784 (28×28)	32	96	10	ReLU

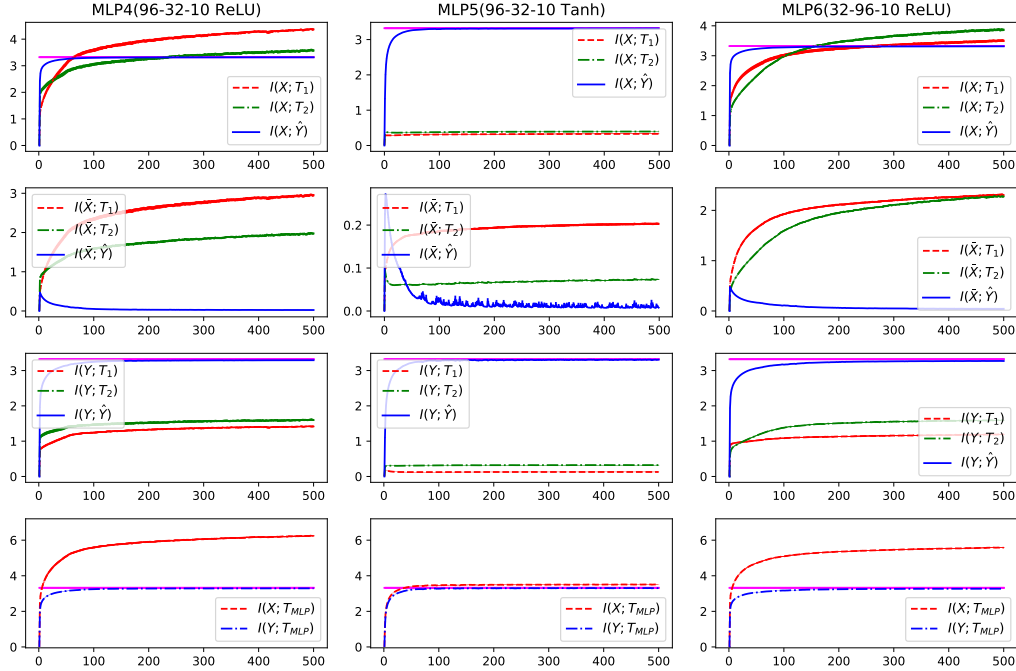


Figure 15: The information flow in MLP4, MLP5, and MLP6 on the MNIST dataset. All the x-axis index training epochs. In each column, the first three figures show $I(X; T_i)$, $I(\bar{X}; T_i)$, and $I(Y; T_i)$ respectively. The forth figure shows $I(X; T_{MLP})$ and $I(Y; T_{MLP})$ in a MLP. The pink line denotes $H(Y) = \log_2 10$.

659 **H.2 Experiments on the Fashion-MNIST dataset**

660 We design three MLPs, namely MLP7, MLP8, and MLP9, and summarize the architectures of the three MLPs in
 661 Table 5. Compared to the MLPs on the MNIST dataset, the three MLPs has one more hidden layer and each
 662 hidden layer has more neurons, *i.e.*, the MLPs are more complicated. Similarly, we train the three MLPs by
 663 Adam [15] over 500 epochs with the learning rate $\alpha = 0.0005$. Based on the mutual information estimator
 664 proposed in Section 4.1, we measure the information flow in the three MLPs during 500 training epochs.

665 Figure 16 shows similar results as Section 4.3 and Section H.1, thus it further confirms the information theoretic
 666 explanations for DNNs.

Table 5: The number of neurons(nodes) and the activation function in MLP7 - MLP9

	x	t_1	t_2	t_3	\hat{y}	$\sigma(\cdot)$
MLP7	784 (28×28)	256	128	96	10	$\text{ReLU}(z) = \max(0, z)$
MLP8	784 (28×28)	256	128	96	10	$\text{Tanh}(z) = (e^z - e^{-z}) / (e^z + e^{-z})$
MLP9	784 (28×28)	96	128	256	10	ReLU

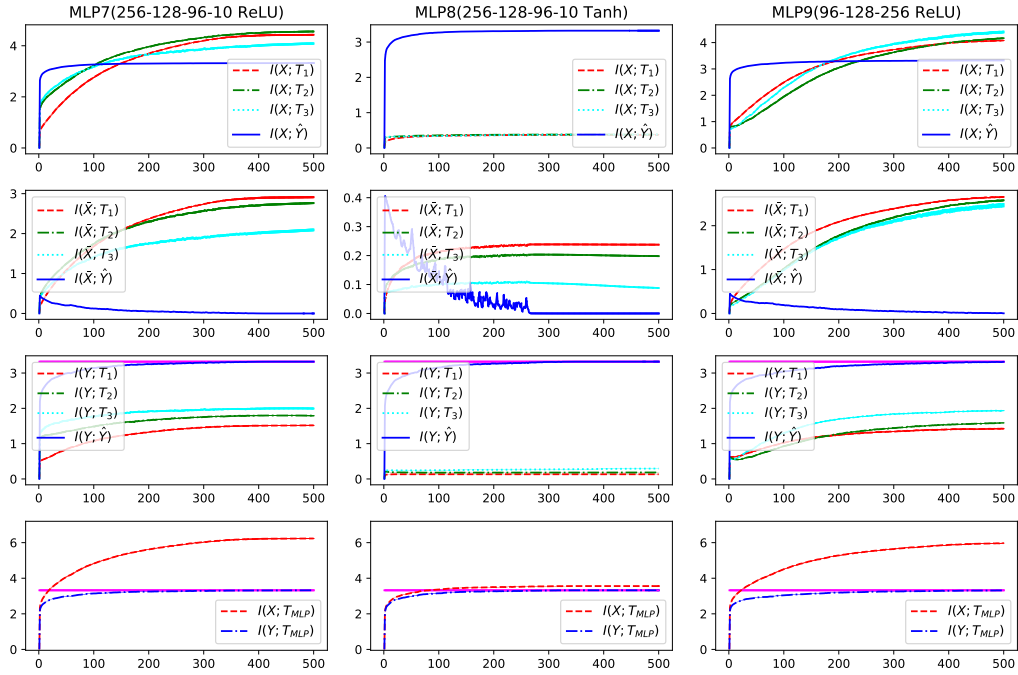


Figure 16: The information flow in MLP7, MLP8, and MLP9 on the MNIST dataset. All the x-axis index training epochs. In each column, the first three figures show $I(X; T_i)$, $I(\bar{X}; T_i)$, and $I(Y; T_i)$ respectively. The forth figure shows $I(X; T_{MLP})$ and $I(Y; T_{MLP})$ in a MLP. The pink line denotes $H(Y) = \log_2 10$.