

An Information Theoretic Interpretation to Deep Neural Networks

Shao-Lun Huang

DSIT Research Center

Tsinghua-Berkeley Shenzhen Institute

Shenzhen, China 518055

Email: shaolun.huang@sz.tsinghua.edu.cn

Xiangxiang Xu

Dept. of Electronic Engineering

Tsinghua University

Beijing, China 100084

Email: xuxx14@mails.tsinghua.edu.cn

Lizhong Zheng, Gregory W. Wornell

Dept. of Electrical & Computer Eng.

Massachusetts Institute of Technology

Cambridge, MA 02139-4307

Email: {lizhong, gww}@mit.edu

Abstract—It is commonly believed that the hidden layers of deep neural networks (DNNs) attempt to extract informative features for learning tasks. In this paper, we formalize this intuition by showing that the features extracted by DNN coincide with the result of an optimization problem, which we call the “universal feature selection” problem, in a local analysis regime. We interpret the weights training in DNN as the projection of feature functions between feature spaces, specified by the network structure. Our formulation has direct operational meaning in terms of the performance for inference tasks, and gives interpretations to the internal computation results of DNNs. Results of numerical experiments are provided to support the analysis.

I. INTRODUCTION

Due to the striking performance of deep learning in various fields, deep neural networks (DNNs) have gained great attentions in modern computer science. While it is a common understanding that the features extracted from the hidden layers of DNN are “informative” for learning tasks, the mathematical meaning of informative features in DNN is generally not clear. There have been numerous research efforts towards this direction [1]. For instance, the information bottleneck [2] employs the mutual information as the metric to quantify the informativeness of features in DNN, and other information metrics, such as the Kullback-Leibler (K-L) divergence [3] and Weissenstein distance [4] are also used in different problems. However, because of the complicated structure of DNNs, there is a disconnection between these information metrics and the performance objectives of the inference tasks that DNNs want to solve. Therefore, it is in general difficult to match the DNN learning with the optimization of a particular information metric.

In this paper, our first contribution is to propose a learning framework, called universal feature selection, which connects the information metric of features and the performance evaluation of inference problems. Specifically for a pair of data variables X and Y , the goal of universal feature selection is to select features from X to infer about a targeted attribute V of Y , where V is only assumed with a rotationally uniform prior over the attribute space of Y , but the precise statistical model between V and X is unknown. Thus, the selected features have to be good for solving multiple inference problems, and should be generally “informative” about Y . We show that in a local analysis regime, the averaged performance of inferring V by a

selected feature of X is measured via a linear projection of this feature, which leads to an information metric to features, and the optimal features can be computed from the singular value decomposition (SVD) of this linear projection.

More importantly, we show that in the local analysis regime, the optimal features selected in DNNs from log-loss optimization coincide with the solutions of universal feature selection. Therefore, the information metric developed in universal feature selection can be used to understand the operations in DNNs. As a result, we observe that the DNN weight updates in general can be interpreted as projecting features between the feature spaces of data and label for extracting the most correlated aspects between them, and the iterative projections can be viewed as computing the SVD of a linear projection between these feature spaces. Moreover, our results also give an explicit interpretation of the goal and the procedures of the BackProp/SGD operations in deep learning. Finally, the theoretic results are validated via numerical experiments.

Notations: Throughout this paper, we use X, \mathcal{X}, P_X , and x to represent a discrete random variable, the range, the probability distribution, and the value of X . In addition, for any function $s(X) \in \mathbb{R}^k$ of X , we use μ_s to denote the mean of $s(X)$, and “ $\tilde{\cdot}$ ” to denote the mean removed version of a variable; e.g., $\tilde{s}(X) = s(X) - \mu_s$. Finally, we use $\|\cdot\|$ and $\|\cdot\|_F$ to denote the ℓ_2 -norm and the Frobenius norm, respectively.

II. PRELIMINARY AND DEFINITION

Given a pair of discrete random variables X, Y with the joint distribution $P_{XY}(x, y)$, the $|\mathcal{Y}| \times |\mathcal{X}|$ matrix $\tilde{\mathbf{B}}$ is defined as

$$\tilde{\mathbf{B}}(y, x) \triangleq \frac{P_{XY}(x, y) - P_X(x)P_Y(y)}{\sqrt{P_X(x)P_Y(y)}}, \quad (1)$$

where $\tilde{\mathbf{B}}(y, x)$ is the (y, x) th entry of $\tilde{\mathbf{B}}$. The matrix $\tilde{\mathbf{B}}$ is referred to as the canonical dependence matrix (CDM). The SVD of $\tilde{\mathbf{B}}$ has the following properties [3].

Lemma 1. *The SVD of $\tilde{\mathbf{B}}$ can be written as $\tilde{\mathbf{B}} = \sum_{i=1}^K \sigma_i \psi_i^Y (\psi_i^X)^T$, where $K \triangleq \min\{|\mathcal{X}|, |\mathcal{Y}|\}$, and σ_i denotes the i th singular value with the ordering $1 \geq \sigma_1 \geq \dots \geq \sigma_K = 0$, and ψ_i^Y and ψ_i^X are the corresponding*

left and right singular vectors with $\psi_K^X(x) = \sqrt{P_X(x)}$ and $\psi_K^Y(y) = \sqrt{P_Y(y)}$.

This SVD decomposes the feature spaces of X, Y into maximally correlated features. To see that, consider the generalized canonical correlation analysis (CCA) problem:

$$\max_{\substack{\mathbb{E}[f_i(X)] = \mathbb{E}[g_i(Y)] = 0 \\ \mathbb{E}[f_i(X) f_j(X)] = \mathbb{E}[g_i(Y) g_j(Y)] = \mathbb{1}_{i=j}}} \sum_{i=1}^k \mathbb{E}[f_i(X) g_i(Y)].$$

It can be shown that for any $1 \leq k \leq K-1$, the optimal features are $f_i(x) = \psi_i^X(x)/\sqrt{P_X(x)}$, and $g_i(y) = \psi_i^Y(y)/\sqrt{P_Y(y)}$, for $i = 0, \dots, K-1$, where $\psi_i^X(x)$ and $\psi_i^Y(y)$ are the x th and y th entries of ψ^X and ψ^Y , respectively [3]. The special case $k = 1$ corresponds to the HGR maximal correlation [5]–[7], and the optimal features can be computed from the ACE algorithm [8].

Moreover, in this paper we focus on a particular analysis regime described as follows.

Definition 1 (ϵ -Neighborhood). Let \mathcal{P}^X denote the space of distributions on some finite alphabet \mathcal{X} , and let $\text{relint}(\mathcal{P}^X)$ denote the subset of strictly positive distributions. For a given $\epsilon > 0$, the ϵ -neighborhood of a distribution $P_X \in \text{relint}(\mathcal{P}^X)$ is defined by the χ^2 -divergence as

$$\mathcal{N}_\epsilon^X(P_X) \triangleq \left\{ P \in \mathcal{P}^X : \sum_{x \in \mathcal{X}} \frac{(P(x) - P_X(x))^2}{P_X(x)} \leq \epsilon^2 \right\}.$$

Definition 2 (ϵ -Dependence). The random variables X, Y is called ϵ -dependent if $P_{XY} \in \mathcal{N}_\epsilon^{X \times Y}(P_X P_Y)$.

Definition 3 (ϵ -Attribute). A random variable U is called an ϵ -attribute of X if $P_{X|U}(\cdot|u) \in \mathcal{N}_\epsilon^X(P_X)$, for all $u \in \mathcal{U}$.

Throughout this paper, we focus on the small ϵ regime, which we refer to as the local analysis regime. In addition, for any $P \in \mathcal{P}^X$, we define the information vector ϕ and feature function $L(x)$ corresponding to P , with respect to a reference distribution $P_X \in \text{relint}(\mathcal{P}^X)$, as

$$\phi(x) \triangleq \frac{P(x) - P_X(x)}{\sqrt{P_X(x)}}, \quad L(x) \triangleq \frac{\phi(x)}{\sqrt{P_X(x)}}. \quad (2)$$

This gives a three way correspondence $P \leftrightarrow \phi \leftrightarrow L$ for all distributions in $\mathcal{N}_\epsilon^X(P_X)$, which will be useful in our derivations.

III. UNIVERSAL FEATURE SELECTION

Suppose that given random variables X, Y with joint distribution P_{XY} , we want to infer about an attribute V of Y from observed i.i.d. samples x_1, \dots, x_n of X . When the statistical model $P_{X|V}$ is known, the optimal decision rule is the log-likelihood ratio test, where the log-likelihood function can be viewed as the optimal feature for inference. However, in many practical situations [3], it is hard to identify the model of the targeted attribute, and is necessary to select low-dimensional informative features of X for inference tasks before knowing the model. We call this universal feature selection problem. To formalize this problem, for an attribute V , we refer to $\mathcal{C}_y = \{ \mathcal{V}, \{P_V(v), v \in \mathcal{V}\}, \{\phi_v^{Y|V}, v \in \mathcal{V}\} \}$,

as the configuration of V , where $\phi_v^{Y|V} \leftrightarrow P_{Y|V}(\cdot|v)$ is the information vector specifying the corresponding conditional distribution $P_{Y|V}(\cdot|v)$. The configuration of V models the statistical correlation between V and Y . In the sequel, we focus on the local analysis regime, for which we assume that all the attributes V of our interests to detect are ϵ -attributes of Y . As a result, the corresponding configuration satisfies $\|\phi_v^{Y|V}\| \leq \epsilon$, for all $v \in \mathcal{V}$. We refer to this as the ϵ -configurations. The configuration of V is unknown in advance, but assumed to be generated from a rotational invariant ensemble (RIE).

Definition 4 (RIE). Two configurations \mathcal{C}_y and $\tilde{\mathcal{C}}_y$ defined as

$$\mathcal{C}_y = \{ \mathcal{V}, \{P_V(v), v \in \mathcal{V}\}, \{\phi_v^{Y|V}, v \in \mathcal{V}\} \},$$

$$\tilde{\mathcal{C}}_y \triangleq \{ \mathcal{V}, \{P_V(v), v \in \mathcal{V}\}, \{\tilde{\phi}_v^{Y|V}, v \in \mathcal{V}\} \}$$

are called rotationally equivalent, if there exists a unitary matrix \mathbf{Q} such that $\tilde{\phi}_v^{Y|V} = \mathbf{Q} \phi_v^{Y|V}$, for all $v \in \mathcal{V}$. Moreover, a probability measure defined on a set of configurations is called an RIE, if all rotationally equivalent configurations have the same measure.

The RIE can be interpreted as assigning a uniform measure to the attributes with the same level of distinguishability. To infer about the attribute V , we construct a k -dimensional feature vector $h^k = (h_1, \dots, h_k)$, for some $1 \leq k \leq K-1$, of the form $h_i = \frac{1}{n} \sum_{l=1}^n f_i(x_l)$, $i = 1, \dots, k$, for some choices of feature functions f_i . Our goal is to determine the f_i such that the optimal decision rule based on h^k achieves the smallest possible error probability, where the performance is averaged over the possible \mathcal{C}_y generated from an RIE. In turn, we denote $\xi_i^X \leftrightarrow f_i$ as the corresponding information vector, and define the matrix $\Xi^X \triangleq [\xi_1^X \ \dots \ \xi_k^X]$.

Theorem 1 (Universal Feature Selection). For $v, v' \in \mathcal{V}$, let $E_{h^k}(v, v')$ be the error exponent associated with the pairwise error probability distinguishing v and v' based on h^k , then the expectation of the error exponent over a given RIE defined on the set of ϵ -configuration is given by

$$\begin{aligned} & \mathbb{E} [E_{h^k}(v, v')] \\ &= \frac{\mathbb{E} [\|\phi_v^{Y|V} - \phi_{v'}^{Y|V}\|^2]}{8|\mathcal{Y}|} \left\| \tilde{\mathbf{B}} \Xi^X ((\Xi^X)^T \Xi^X)^{-\frac{1}{2}} \right\|_{\text{F}}^2 + o(\epsilon^2), \end{aligned} \quad (3)$$

where the expectations are taken over this RIE.

Proof. See Appendix A. \square

As a result of (3), designing the ξ_i^X as the singular vectors ψ_i^X of $\tilde{\mathbf{B}}$, for $i = 1, \dots, k$, optimizes (3) for all RIEs, pairs of (v, v') , and ϵ -configurations. Thus, the feature functions corresponding to ψ_i^X are *universally optimal* for inferring the unknown attribute V . Moreover, (3) naturally leads to an information metric $\|\tilde{\mathbf{B}} \Xi^X ((\Xi^X)^T \Xi^X)^{-\frac{1}{2}}\|_{\text{F}}^2$ for any feature Ξ^X of X , measured by projecting the normalized Ξ^X through a linear projection $\tilde{\mathbf{B}}$. This information metric quantifies how informative a feature of X is when solving inference problems with respect to Y , and is optimized when designing features by singular vectors of $\tilde{\mathbf{B}}$. Thus, we can interpret the universal

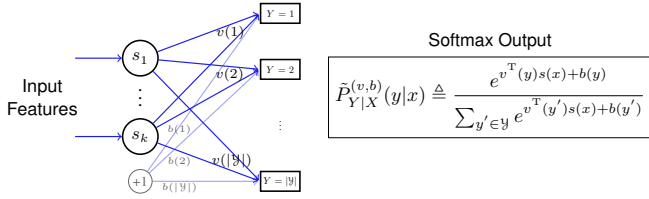


Fig. 1: A simple neural network with one layer of hidden nodes with softmax output.

feature selection as solving the most informative features for data inferences via the SVD of $\tilde{\mathbf{B}}$, which also coincides with the maximally correlated features in Section II. Later on we will show that the feature selections in DNN share the same information metric as universal feature selection in the local analysis regime.

IV. INTERPRETING SOFTMAX REGRESSION

To begin, recall that for a data vector X and label Y with labeled samples (x_i, y_i) , for $i = 1, \dots, N$, the softmax regression generally uses a discriminative model of the form

$$\tilde{P}_{Y|X}^{(v,b)}(y|x) \triangleq \frac{e^{v^T(y)s(x)+b(y)}}{\sum_{y' \in \mathcal{Y}} e^{v^T(y')s(x)+b(y')}} \quad (4)$$

to address the classification problems, where $s(x) \in \mathbb{R}^k$ is a k -dimensional representation of X used to predict the label, and $v(y) \in \mathbb{R}^k$ and $b(y) \in \mathbb{R}$ are the parameters required to be learned from

$$(v, b)^* = \arg \max_{(v,b)} \frac{1}{N} \sum_{i=1}^N \log \tilde{P}_{Y|X}^{(v,b)}(y_i|x_i). \quad (5)$$

As depicted in Fig. 1, the ordinary softmax regression corresponds to $s(x) = x$. More generally, $s(x)$ can be the output of the previous hidden layer of a neural network, i.e., the selected feature of x fed into the softmax regression. In the rest of this section, we will show that when X, Y are ϵ -dependent, the functions $s(x)$ and $v(y)$ coincide with the solutions of the universal feature selection.

First, we use P_{XY} to denote the joint empirical distribution of the labeled samples (x_i, y_i) , $i = 1, \dots, N$, and P_X, P_Y to denote the corresponding marginal distributions. Then, the objective function in the optimization problem (5) is precisely the empirical average of the log-likelihood, i.e., $\frac{1}{N} \sum_{i=1}^N \log \tilde{P}_{Y|X}^{(v,b)}(y_i|x_i) = \mathbb{E}_{P_{XY}} \left[\log \tilde{P}_{Y|X}^{(v,b)}(Y|X) \right]$. Therefore, maximizing this empirical average is equivalent as minimizing the K-L divergence:

$$(v, b)^* = \arg \min_{(v,b)} D(P_{XY} \| P_X \tilde{P}_{Y|X}^{(v,b)}). \quad (6)$$

This can be interpreted as finding the best fitting to empirical joint distribution P_{XY} by distributions of the form $P_X \tilde{P}_{Y|X}^{(v,b)}$. In our development, it is more convenient to denote the bias by $d(y) = b(y) - \log P_Y(y)$, for $y \in \mathcal{Y}$. Then, the following lemma illustrates the explicit constraint on the problem (6) in the local analysis regime.

Lemma 2. *If X, Y are ϵ -dependent, then the optimal v, d for (6) satisfy*

$$|\tilde{v}^T(y)s(x) + \tilde{d}(y)| = O(\epsilon), \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y}. \quad (7)$$

Proof. See Appendix B. \square

In turn, we take (7) as the constraint for solving the problem (6) in the local analysis regime. Moreover, we define the information vectors for zero-mean vectors \tilde{s}, \tilde{v} as $\xi^X(x) = \sqrt{P_X(x)} \tilde{s}(x)$, $\xi^Y(y) = \sqrt{P_Y(y)} \tilde{v}(y)$, and define matrices

$$\Xi^Y \triangleq [\xi^Y(1) \quad \dots \quad \xi^Y(|\mathcal{Y}|)]^T, \\ \Xi^X \triangleq [\xi^X(1) \quad \dots \quad \xi^X(|\mathcal{X}|)]^T.$$

Lemma 3. *The K-L divergence (6) in the local analysis regime (7) can be expressed as*

$$D(P_{XY} \| P_X \tilde{P}_{Y|X}^{(v,b)}) \\ = \frac{1}{2} \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_F^2 + \frac{1}{2} \eta^{(v,b)}(s) + o(\epsilon^2), \quad (8)$$

where $\eta^{(v,b)}(s) \triangleq \mathbb{E}_{P_Y} \left[(\mu_s^T \tilde{v}(Y) + \tilde{d}(Y))^2 \right]$.

Proof. See Appendix C. \square

Eq. (8) reveals key insights for feature selection in neural networks, which are illustrated by the following three learning problems, depending on if the weights, input feature, or both can be trained from data.

A. Forward Feature Projection

For the case that s is fixed, we can optimize (8) with Ξ^X fixed and get the following optimal weights:

Theorem 2. *For fixed Ξ^X and μ_s , the optimal Ξ^{Y*} to minimize (8) is given by*

$$\Xi^{Y*} = \tilde{\mathbf{B}} \Xi^X ((\Xi^X)^T \Xi^X)^{-1}, \quad (9)$$

and the optimal weights \tilde{v}^* and bias \tilde{d}^* are

$$\tilde{v}^*(y) = \mathbb{E}_{P_{X|Y}} \left[\Lambda_{\tilde{s}(X)}^{-1} \tilde{s}(X) \mid Y = y \right], \quad \tilde{d}^*(y) = -\mu_s^T \tilde{v}(Y). \quad (10)$$

where $\Lambda_{\tilde{s}(X)}$ denotes the covariance matrix of $\tilde{s}(X)$.

Proof. See Appendix D. \square

Eq. (9) can be viewed as a projection of the input feature $\tilde{s}(x)$, to a feature $v(y)$ computable from the value of y , which is the most correlated feature to $\tilde{s}(x)$. The solution is given by left multiplying the $\tilde{\mathbf{B}}$ matrix. We call this the ‘‘forward feature projection’’.

Remark 1. *While we assume the continuous input $s(x)$ is a function of a discrete variable X , we only need the labeled samples between s and Y to compute the weights and bias from the conditional expectation (10), and the correlation between X and s is irrelevant. Thus, our analysis for weights and bias can be applied to continuous input networks by just ignoring X and taking s as the real network input.*

B. Backward Feature Projection

It is also useful to consider the ‘‘backward problem’’, which attempts to find informative feature $s^*(X)$ to minimize the loss (8) with given weights and bias.

Theorem 3. For fixed \tilde{v} , Ξ^Y , and \tilde{d} , the optimal Ξ^{X^*} to minimize (8) is given by

$$\Xi^{X^*} = \tilde{\mathbf{B}}^T \Xi^Y ((\Xi^Y)^T \Xi^Y)^{-1}, \quad (11)$$

and the optimal feature function s^* , which are decomposed to \tilde{s}^* and μ_s^* , are given by

$$\begin{aligned} \tilde{s}^*(x) &= \mathbb{E}_{P_{Y|X}} \left[\Lambda_{\tilde{v}(Y)}^{-1} \tilde{v}(Y) \middle| X = x \right], \\ \mu_s^* &= -\Lambda_{\tilde{v}(Y)}^{-1} \mathbb{E}_{P_Y} \left[\tilde{v}(Y) \tilde{d}(Y) \right], \end{aligned} \quad (12)$$

where $\Lambda_{\tilde{v}(Y)}$ denotes the covariance matrix of $\tilde{v}(Y)$.

Proof. See Appendix D. \square

The solution of this backward feature projection is precisely symmetric to the forward one. Note we assumed here that the feature $s(X)$ can be selected as any desired function. This is only true in the ideal case where the previous hidden layers of the neural network have sufficient expressive power. That is, it can generate the desired feature function as given in (12). In general, however, the form of feature functions that can be generalized is often limited by the network structure. In the next section, we discuss such cases, where we do know the most desirable feature function as given in (12), and the question is how a network with limited expressive power approximates this optimal solution.

C. Universal Feature Selection

When both s and (v, b) (and hence Ξ^X, Ξ^Y, d) can be designed, the optimal (Ξ^Y, Ξ^X) corresponds to the low rank factorization of $\tilde{\mathbf{B}}$, and the solutions coincide with the universal feature selection.

Theorem 4. The optimal solutions for weights and bias to maximize (8) are given by $\tilde{d}(y) = -\mu_s^T \tilde{v}(y)$, and $(\Xi^Y, \Xi^X)^*$ chosen as the largest k left and right singular vectors of $\tilde{\mathbf{B}}$.

Proof. See Appendix E. \square

Therefore, we conclude that the softmax regression, when both s and (v, b) are designable, is to extract the most correlated aspects of the input data X and the label Y that are informative features for data inferences from universal feature selection.

In the learning process of DNN, the BackProp procedure alternatively chooses the weights of the softmax layer and those on the previous layer(s). In each step, the weights on the rest of the network are fixed. This is equivalent as alternating between the forward and the backward feature projections, i.e. it alternates between (9) and (11). This is in fact the power method to solve the SVD for $\tilde{\mathbf{B}}$ [9], which is also known as the Alternating Conditional Expectation (ACE) algorithm [8].

V. MULTI-LAYER NETWORK ANALYSIS

From the previous discussions, the performance of the softmax regression not only depends on the weight and bias

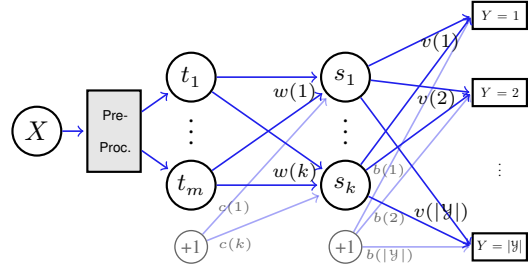


Fig. 2: A multi-layer network: all hidden layers previous to t are labeled as ‘‘pre-processing’’.

$(v(y), b(y))$, but the input feature $s(x)$ has to be informative. It turns out that the hidden layers of neural networks, which are known to have strong expressive power of features, are essentially extracting such informative features. For illustration, we consider the neural network with a hidden layer of k nodes, and a zero-mean continuous input $t = [t_1 \cdots t_m]^T \in \mathbb{R}^m$ to this hidden layer, where t is assumed to be a function $t(x)$ of some discrete variable X^1 . Our goal is to analyze the weights and bias in this layer with labeled samples $(t(x_i), y_i)$. Assume the activation function of the hidden layer is a generally smooth function $\sigma(\cdot)$, then the output $s_z(X)$ of the z -th hidden node is

$$s_z(x) = \sigma \left(w^T(z) t(x) + c(z) \right), \quad \text{for } z = 1, \dots, k, \quad x \in \mathcal{X}, \quad (13)$$

where $w(z) \in \mathbb{R}^m$ and $c(z) \in \mathbb{R}$ are the weights and bias from input layer to hidden layer as shown in Fig. 2. We denote $s = [s_1 \cdots s_k]^T$ as the input vector to the output softmax regression layer.

To interpret the feature selection in hidden layers, we fix $(v(y), b(y))$ at the output layer, and consider the problem of designing $(w(z), c(z))$ to minimize the loss function (6) of the softmax regression at the output layer. Ideally, we should have picked $w(z)$ and $c(z)$ to generate $s(x)$ to match $s^*(x)$ from (12), which minimizes the loss. However, here we have the constraint that $s(x)$ must take the form of (13), and intuitively the network should select $w(z), c(z)$ so that $s(x)$ is close to $s^*(x)$. Our goal is to quantify the notion of closeness in the local analysis regime.

To develop insights on feature selection in hidden layers, we again focus on the local analysis regime, where the weights and bias are assumed with the local constraint

$$|\tilde{v}^T(y) s(x) + \tilde{d}(y)| = O(\epsilon), \quad |w^T(z) \tilde{t}(x)| = O(\epsilon), \quad \forall x, y, z. \quad (14)$$

Then, since t is zero-mean, we can express (13) as

$$\begin{aligned} s_z(x) &= \sigma \left(w^T(z) t(x) + c(z) \right) \\ &= w^T(z) \tilde{t}(x) \cdot \sigma'(c(z)) + \sigma(c(z)) + o(\epsilon), \end{aligned} \quad (15)$$

¹As discussed in Remark 1, X is assumed only for the convenience of analysis, and the computation of weights and bias only needs t , but not X . Moreover, the input t to the hidden layer can be either directly from data or the output of previous hidden layers in a DNN, which we model as ‘‘pre-processing’’ as shown in Fig. 2.

Moreover, we define a matrix $\tilde{\mathbf{B}}_1$ with the (z, x) th entry $\tilde{\mathbf{B}}_1(z, x) = \frac{\sqrt{P_X(x)}}{\sigma'(c(z))} \tilde{s}_z^*(x)$, which can be interpreted as a generalized DTM for the hidden layer. Furthermore, we denote $\xi_1^X(x) = \sqrt{P_X(x)} \tilde{t}(x)$ as the information vector of $\tilde{t}(x)$ with the matrix Ξ_1^X defined as $\Xi_1^X \triangleq [\xi_1^X(1) \ \cdots \ \xi_1^X(|\mathcal{X}|)]^\top$, and we also define

$$\mathbf{W} \triangleq [w(1) \ \cdots \ w(k)]^\top$$

$$\mathbf{J} \triangleq \text{diag}\{\sigma'(c(1)), \sigma'(c(2)), \dots, \sigma'(c(k))\}.$$

The following theorem characterizes the loss (6).

Theorem 5. *Given the weights and bias (v, b) at the output layer, and for any input feature s , we denote $\mathcal{L}(s)$ as the loss (6) evaluated with respect to (v, b) and s . Then, with the constraints (14)*

$$\mathcal{L}(s) - \mathcal{L}(s^*)$$

$$= \frac{1}{2} \|\Theta \tilde{\mathbf{B}}_1 - \Theta \mathbf{W} (\Xi_1^X)^\top\|_{\text{F}}^2 + \frac{1}{2} \kappa^{(v,b)}(s, s^*) + o(\epsilon^2), \quad (16)$$

where $\Theta \triangleq ((\Xi^Y)^\top \Xi^Y)^{1/2} \mathbf{J}$, and the term $\kappa^{(v,b)}(s, s^*) = (\mu_s - \mu_{s^*})^\top \Lambda_{\tilde{v}(Y)} (\mu_s - \mu_{s^*})$.

Proof. See Appendix F. \square

Eq. (16) quantifies the closeness between s and s^* in terms of the loss (6). Then, our goal is to minimize (16), which can be separated to two optimization problems:

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \|\Theta \tilde{\mathbf{B}}_1 - \Theta \mathbf{W} (\Xi_1^X)^\top\|_{\text{F}}^2, \quad (17)$$

$$\mu_s^* = \arg \min_{\mu_s} \kappa^{(v,b)}(s, s^*). \quad (18)$$

First note that the optimization problem (17) is similar to the ordinary softmax regression depicted in Section IV, and the optimal solution is given by $\mathbf{W}^* = \tilde{\mathbf{B}}_1 \Xi_1^X ((\Xi_1^X)^\top \Xi_1^X)^{-1}$. Therefore, solving the optimal weights in the hidden layer can be interpreted as projecting $\tilde{s}^*(x)$ to the subspace of feature functions spanned by $t(x)$ to find the closest expressible function. Finally, the problem (18) is to choose μ_s (and hence the bias $c(z)$) to minimize the quadratic term similar to $\eta^{(v,b)}(s)$ in (8), and we refer to Appendix F for the optimal solution of (18).

Overall, we observe the correspondence between (9), (12), and (17), (18), and interpret both operations as feature projections. Our argument can be generalized to any intermediate layer in a multi-layer network, with all the previous layers viewed as the fixed pre-processing that specifies $t(x)$, and all the layers after determining s^* . Then the iterative procedure in back-propagation can be viewed as alternating projection finding the fixed-point solution over the entire network. This final fixed-point solution, even under the local assumption, might not be the SVD solution as in Theorem 4. This is because the limited expressive power of the network often makes it impossible to generate the desired feature function. In such cases, the concept of feature projection can be used to quantify this gap, and thus to measure the quality of the selected features.

VI. SCORING NEURAL NETWORKS

Given a learning problem, it is useful to tell whether or not some extracted features is informative [10]. Our previous development naturally gives rise to a performance metric.

Definition 5. *Given a feature $s(x) \in \mathbb{R}^k$ and weight $v(y) \in \mathbb{R}^k$ with the corresponding information matrices Ξ^X and Ξ^Y , the H-score $H(s, v)$ is defined as*

$$H(s, v) \triangleq \frac{1}{2} \|\tilde{\mathbf{B}}\|_{\text{F}}^2 - \frac{1}{2} \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^\top\|_{\text{F}}^2$$

$$= \mathbb{E}_{P_{XY}} \left[\tilde{s}^\top(X) \tilde{v}(Y) \right] - \frac{1}{2} \text{tr}(\Lambda_{\tilde{s}(X)} \Lambda_{\tilde{v}(Y)}). \quad (19)$$

In addition, we define the single-sided H-score $H(s)$ as

$$H(s) \triangleq \frac{1}{2} \|\tilde{\mathbf{B}} \Xi^X ((\Xi^X)^\top \Xi^X)^{-\frac{1}{2}}\|_{\text{F}}^2$$

$$= \frac{1}{2} \mathbb{E}_{P_Y} \left[\left\| \mathbb{E}_{P_{X|Y}} \left[\Lambda_{\tilde{s}(X)}^{-1/2} \tilde{s}(X) \mid Y \right] \right\|^2 \right]. \quad (20)$$

H-score can be used to measure the quality of features generated at any intermediate layer of the network. It is related to (16) when choosing the optimal bias and Θ as the identity matrix. This can be understood as taking the output of this layer $s(x)$ and directly feed it to a softmax output layer with $v(y)$ as the weights, and $H(s, v)$ measures the resulting performance. Note that $v(y)$ here can be an arbitrary function of Y . It is not necessarily the weights on the next layer computed by the network. When the optimal weights $v^*(y)$ is used, the resulting performance becomes the one-sided H-score $H(s)$, which measures the quality of $s(x)$, and coincides with the information metric (3).

In current practice the cross-entropy $\mathbb{E}[\log \tilde{P}_{Y|X}^{(v,b)}]$, is often used as the performance metric. One can in principle also use log-loss to measure the effectiveness of the selected feature at the output of an intermediate layer [10]. However, one problem of this metric is that for a given problem it is not clear what value of log-loss one should expect because the log-loss is generally unbounded. Moreover, the computation of the log-loss for optimal weights and bias with respect to a particular input feature requires solving a non-convex optimization problem with the issue of locking at local optimum.

In contrast, the H-score can be directly computed from the data samples, and has a clear upper bound from Lemma 1 that $H(s, v) \leq H(s) \leq (1/2) \sum_{i=1}^k \sigma_i^2 \leq k/2$. In this sequence of inequalities, the gap over the first " \leq " measures the optimality of the weights v ; the second gap is due to the difference between the chosen feature and the optimal solution, which is a useful measure of how restrictive (lack of expressive power) the network structure is; and the last one measures how good the dataset itself is. In Section VII, we validate this metric in real data.

VII. EXPERIMENTAL VALIDATION

We first validate the feature projection in Theorem 4. For this purpose, we construct the NN as shown in Fig. 1 with $k = 1$, $|\mathcal{X}| = 8$, and $|\mathcal{Y}| = 6$, and the input feature $s(X)$ is generated from a sigmoid layer with the one-hot encoded X as the input. Note that with proper weights in the sigmoid layer, $s(X)$ can

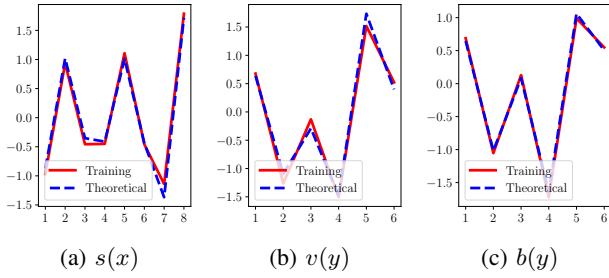


Fig. 3: The comparisons of the weights and bias in softmax regression.

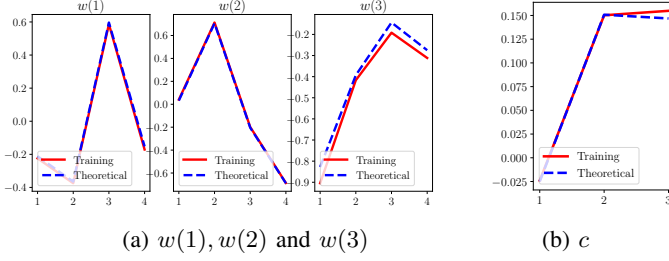


Fig. 4: The comparisons of the weights and bias in the hidden layer.

express any desired function, up to scaling and shifting. To compare the result trained by the neural network and that in Theorem 4, we first randomly generate a distribution P_{XY} , and then generate $n = 100,000$ samples of (X, Y) pair. Using these data to train the neural network, the corresponding results of $s(x)$, $v(y)$ and $b(y)$ are shown in Fig. 3 with a comparison to theoretical result, where the training results match our theory. In addition, we validate Theorem 5 by the NN depicted in Fig. 2, with the same setup of X, Y . The number of neurons in hidden layers are $m = 4$ and $k = 3$, and the input $t(X)$ is some randomly chosen function of X , and the activation $\sigma(\cdot)$ is the sigmoid function. We then fix the weights and bias at the output layer and train the weights $w(1), w(2), w(3)$, and bias c in the hidden layer to optimize the Log-Loss. Fig. 4 shows the matching between our results and the experiment.

Furthermore, we compare the performance of classification and the H-score evaluated from the features extracted from the last hidden layer of different DNNs. We use the ILSVRC2012 [11] as our validation data set, and train several state-of-art DNNs [12]–[17] to extract the features from the last hidden layers. The resulting H-scores are shown in TABLE I with the comparison to the classification accuracy, where H_{AIC} is the H-score with the correction of Akaike information criterion (AIC) [18] to reduce overfitting. In particular, $H_{\text{AIC}}(s)$ is given by

$$H_{\text{AIC}}(s) = H(s) - \frac{n_p}{n_s}, \quad (21)$$

where n_p is the number of parameters contained in the model, and $n_s = 1,300,000$ is the number of training samples in ImageNet. The corrected H-score is consistent with the accuracy, which validates the H-score.

Model	$H(s)$	$H_{\text{AIC}}(s)$	Accaracy
VGG16	148.3	41.9	0.642
VGG19	152.7	42.2	0.647
MobileNet	45.9	42.6	0.684
DenseNet121	59.5	53.3	0.714
DenseNet169	81.2	70.2	0.736
DenseNet201	89.1	73.5	0.744
Xception	179.8	162.2	0.775
InceptionV3	181.2	162.9	0.763
InceptionResNetV2	241.1	198.1	0.791

TABLE I: H_{AIC} is the H-score with AIC correction.

ACKNOWLEDGMENT

The research of Shao-Lun Huang was funded by the Natural Science Foundation of China 61807021, Shenzhen Science and Technology Research and Development Funds (JCYJ20170818094022586), and Innovation and entrepreneurship project for overseas high-level talents of Shenzhen (KQJSCX20180327144037831).

APPENDIX

A. Proof of Theorem 1

We commence with the characterization of the error exponent.

Lemma 4. *Given a reference distribution $P_X \in \text{relint}(\mathcal{P}^X)$, a constant $\epsilon > 0$ and integers n and k , let x_1, \dots, x_n denote i.i.d. samples from one of P_1 or P_2 , where $P_1, P_2 \in \mathcal{N}_\epsilon^X(P_X)$. To decide whether P_1 or P_2 is the generating distribution, a sequence of k -dimensional statistics $h^k = (h_1, \dots, h_k)$ is constructed as*

$$h_i = \frac{1}{n} \sum_{l=1}^n f_i(x_l), \quad i = 1, \dots, k, \quad (22)$$

where $(f_1(X), \dots, f_k(X))$ are zero mean, unit-variance, and uncorrelated with respect to P_X , i.e.,

$$\mathbb{E}_{P_X} [f_i(X)] = 0, \quad i \in \{1, \dots, k\} \quad (23a)$$

$$\mathbb{E}_{P_X} [f_i(X)f_j(X)] = \delta_{ij}, \quad i, j \in \{1, \dots, k\}. \quad (23b)$$

Then the error probability of the decision based on h^k decays exponentially in n as $n \rightarrow \infty$, with (Chernoff) exponent

$$\lim_{n \rightarrow \infty} \frac{-\log p_e}{n} \triangleq E_{h^k} = \sum_{i=1}^k E_{h_i}, \quad (24a)$$

where

$$E_{h_i} = \frac{1}{8} \langle \phi_1 - \phi_2, \xi_i \rangle^2 + o(\epsilon^2), \quad (24b)$$

and $\phi_1 \leftrightarrow P_1, \phi_2 \leftrightarrow P_2, \xi_i \leftrightarrow f_i(X), i \in \{1, \dots, k\}$ are the corresponding information vectors.

Proof. Since the rule is to decide based on comparing the projection

$$\sum_{i=1}^k h_i (\mathbb{E}_{P_1} [f_i(X)] - \mathbb{E}_{P_2} [f_i(X)])$$

to a threshold, via Cramér's theorem [19] the error exponent under P_j ($j = 1, 2$) is

$$E_j(\lambda) = \min_{P \in \mathcal{S}(\lambda)} D(P \| P_j), \quad (25)$$

where

$$\mathcal{S}(\lambda) \triangleq \left\{ P \in \mathcal{P}^{\mathcal{X}} : \right.$$

$$\left. \mathbb{E}_P[f^k(X)] = \lambda \mathbb{E}_{P_1}[f^k(X)] + (1-\lambda) \mathbb{E}_{P_2}[f^k(X)] \right\}. \quad (26)$$

Now since (23a) holds, we obtain

$$\begin{aligned} \mathbb{E}_{P_j}[f_i(X)] &= \sum_{x \in \mathcal{X}} P_j(x) f_i(x) \\ &= \sum_{x \in \mathcal{X}} P_X(x) f_i(x) + \sum_{x \in \mathcal{X}} (P_j(x) - P_X(x)) f_i(x) \\ &= \mathbb{E}_{P_X}[f_i(X)] + \sum_{x \in \mathcal{X}} \sqrt{P_X(x)} \phi_j(x) \cdot \frac{\xi_i(x)}{\sqrt{P_X(x)}} \\ &= \sum_{x \in \mathcal{X}} \phi_j(x) \xi_i(x) \\ &= \langle \phi_j, \xi_i \rangle, \quad j = 1, 2 \text{ and } i = 1, \dots, k, \end{aligned} \quad (27)$$

which we express compactly as

$$\mathbb{E}_{P_j}[f^k(X)] = \langle \phi_j, \xi^k \rangle, \quad j = 1, 2$$

with $\xi^k \triangleq (\xi_1, \dots, \xi_k)$.

Hence, the constraint (26) is expressed in information vectors as

$$\langle \phi, \xi_i \rangle = \langle \lambda \phi_1 + (1-\lambda) \phi_2, \xi_i \rangle, \quad i = 1, \dots, k,$$

i.e.,

$$\langle \phi, \xi^k \rangle = \langle \lambda \phi_1 + (1-\lambda) \phi_2, \xi^k \rangle. \quad (28)$$

In turn, the optimal P in (25), which we denoted by P^* , lies in the exponential family through P_j with natural statistic $f^k(x)$, i.e., the k -dimensional family whose members are of the form

$$\log \tilde{P}_{\theta^k}(x) = \sum_{i=1}^k \theta_i f_i(x) + \log P_j(x) - \alpha(\theta^k),$$

for which the associated information vector is

$$\tilde{\phi}_{\theta^k}(x) = \sum_{i=1}^k \theta_i \xi_i(x) + \phi_j(x) - \alpha(\theta^k) \sqrt{P_X(x)} + o(\epsilon), \quad (29)$$

where we have used the fact that

$$\begin{aligned} \log P(x) &= \log P_X(x) + \log \frac{P(x)}{P_X(x)} \\ &= \log P_X(x) + \log \left(1 + \frac{1}{\sqrt{P_X(x)}} \phi(x) \right) \\ &= \log P_X(x) + \frac{1}{\sqrt{P_X(x)}} \phi(x) + o(\epsilon) \end{aligned}$$

for all $P \in \mathcal{N}_\epsilon^{\mathcal{X}}(P_X)$ with the information vector $\phi \leftrightarrow P$. As a result,

$$\langle \tilde{\phi}_{\theta^k}, \xi_i \rangle = \theta_i + \langle \phi_j, \xi_i \rangle + o(\epsilon),$$

where we have used (23b). Hence, via (28) we obtain that the intersection with the linear family (26) is at $P^* = P_{\theta^{k*}}$ with

$$\theta_i^* = \langle \lambda \phi_1 + (1-\lambda) \phi_2 - \phi_j, \xi_i \rangle + o(\epsilon)$$

and thus

$$\begin{aligned} E_j(\lambda) &= D(P^* \| P_j) \\ &= \frac{1}{2} \|\tilde{\phi}_{\theta^{k*}} - \phi_j\|^2 + o(\epsilon^2) \end{aligned} \quad (30a)$$

$$= \frac{1}{2} \left\| \sum_{i=1}^k \theta_i^* \xi_i \right\|^2 + \frac{1}{2} \alpha(\theta^{k*})^2 + o(\epsilon^2) \quad (30b)$$

$$= \frac{1}{2} \sum_{i=1}^k (\theta_i^*)^2 + \frac{1}{2} \alpha(\theta^{k*})^2 + o(\epsilon^2) \quad (30c)$$

$$= \frac{1}{2} \sum_{i=1}^k \langle \lambda \phi_1 + (1-\lambda) \phi_2 - \phi_j, \xi_i \rangle^2 + o(\epsilon^2), \quad (30d)$$

where to obtain (30a) we have exploited the local approximation of K-L divergence [3], to obtain (30b) we have exploited (29), to obtain (30c) we have again exploited (23b), and to obtain (30d) we have used that

$$\alpha(\theta^{k*}) = o(\epsilon^2)$$

since $\theta^{k*} = O(\epsilon)$ and

$\alpha(0) = 0$, and $\nabla \alpha(0) = \mathbb{E}_{P_j}[f^k(X)] = \langle \phi_j, \xi^k \rangle = O(\epsilon)$. Finally, $E_1(\lambda) = E_2(\lambda)$ when $\lambda = 1/2$, so the overall error probability has exponent (24). \square

Then, the following lemma demonstrates a property of information vectors in a Markov chain.

Lemma 5. *Given the Markov relation $X \leftrightarrow Y \leftrightarrow V$ and any $v \in \mathcal{V}$, let $\phi_v^{X|V}$ and $\phi_v^{Y|V}$ denote the associated information vectors for $P_{X|V}(\cdot|v)$ and $P_{Y|V}(\cdot|v)$, then we have*

$$\phi_v^{X|V} = \tilde{\mathbf{B}}^T \phi_v^{Y|V}. \quad (31)$$

Proof. The Markov relation implies

$$\begin{aligned} P_X(x) &= \sum_{y \in \mathcal{Y}} P_{X|Y}(x|y) P_Y(y), \\ P_{X|V}(x|v) &= \sum_{y \in \mathcal{Y}} P_{X|Y,V}(x|y, v) P_{Y|V}(y|v) \\ &= \sum_{y \in \mathcal{Y}} P_{X|Y}(x|y) P_{Y|V}(y|v). \end{aligned}$$

As a result,

$$P_{X|V}(x|v) - P_X(x) = \sum_{y \in \mathcal{Y}} P_{X|Y}(x|y) [P_{Y|V}(y|v) - P_Y(y)],$$

and the corresponding information vectors satisfy

$$\begin{aligned} \phi_v^{X|V}(x) &= \frac{1}{\sqrt{P_X(x)}} \sum_{y \in \mathcal{Y}} P_{X|Y}(x|y) \sqrt{P_Y(y)} \phi_v^{Y|V}(y) \\ &= \sum_{y \in \mathcal{Y}} \left[\tilde{\mathbf{B}}(y, x) + \sqrt{P_X(x) P_Y(y)} \right] \phi_v^{Y|V}(y) \\ &= \sum_{y \in \mathcal{Y}} \tilde{\mathbf{B}}(y, x) \phi_v^{Y|V}(y), \end{aligned} \quad (32)$$

where the last equality follows from the fact that

$$\sum_{y \in \mathcal{Y}} \sqrt{P_Y(y)} \phi_v^{Y|V}(y) = \sum_{y \in \mathcal{Y}} [P_{Y|V}(y|v) - P_Y(y)] = 0.$$

Finally, express (32) in the matrix form and we obtain (31). \square

In addition, the following lemma is useful for dealing with the expectation over an RIE.

Lemma 6. *Let \mathbf{z} be a spherically symmetric random vector of dimension M , i.e., for any orthogonal \mathbf{Q} we have $\mathbf{z} \stackrel{d}{=} \mathbf{Q}\mathbf{z}$. If \mathbf{A} is a fixed matrix of compatible dimensions, then*

$$\mathbb{E} \left[\|\mathbf{z}^T \mathbf{A}\|^2 \right] = \frac{1}{M} \mathbb{E} \left[\|\mathbf{z}\|^2 \right] \|\mathbf{A}\|_{\text{F}}^2. \quad (33)$$

Proof. By definition we have $\Lambda_z = \mathbf{Q}\Lambda_z\mathbf{Q}^T$ for any orthogonal \mathbf{Q} , hence Λ_z is diagonal. Suppose $\Lambda_z = \lambda\mathbf{I}$, then from

$$\text{tr}(\Lambda_z) = \mathbb{E}[\|\mathbf{z}\|^2] = \lambda M$$

we obtain

$$\lambda = \frac{1}{M} \text{tr}(\Lambda_z).$$

As a result,

$$\begin{aligned} \mathbb{E}[\|\mathbf{z}^T \mathbf{A}\|^2] &= \text{tr}(\mathbf{A}^T \Lambda_z \mathbf{A}) = \lambda \text{tr}(\mathbf{A}^T \mathbf{A}) \\ &= \frac{1}{M} \mathbb{E}[\|\mathbf{z}\|^2] \|\mathbf{A}\|_{\text{F}}^2. \end{aligned} \quad (34)$$

□

We now have everything to prove Theorem 1.

Proof of Theorem 1. By definition of feature functions, we have $\mathbb{E}_{P_X}[f_i(X)] = 0, i = 1, \dots, k$. Suppose \mathbf{f} is the vector representation of f^k and denote by $\tilde{\mathbf{f}} \triangleq \Lambda_f^{-1/2} \mathbf{f}$ the normalized \mathbf{f} , with $\Lambda_f^{1/2}$ denoting any square root matrix of Λ_f , then the corresponding statistics $\tilde{f}^k = (\tilde{f}_1, \dots, \tilde{f}_k)$ satisfy the constraints (23). Further, construct the statistic $\tilde{h}^k = (\tilde{h}_1, \dots, \tilde{h}_k)$ as [cf. (22)]

$$\tilde{h}_i = \frac{1}{n} \sum_{l=1}^n \tilde{f}_i(x_l), \quad i = 1, \dots, k. \quad (35)$$

Then, from Lemma 4, the error exponent of distinguishing v and v' based on \tilde{h}^k is

$$\begin{aligned} E_{\tilde{h}^k}(v, v') &= \frac{1}{8} \sum_{i=1}^k \left[(\phi_v^{X|V} - \phi_{v'}^{X|V})^T \tilde{\xi}_i^X \right]^2 + o(\epsilon^2) \\ &= \frac{1}{8} \left\| (\phi_v^{X|V} - \phi_{v'}^{X|V})^T \tilde{\Xi}^X \right\|^2 + o(\epsilon^2), \end{aligned}$$

where $\phi_v^{X|V}$ denotes the associated information vector for $P_{X|V}(\cdot|v)$, $\tilde{\xi}_i^X$ denotes the information vectors of \tilde{f}_i , and $\tilde{\Xi}^X \triangleq [\tilde{\xi}_1^X, \dots, \tilde{\xi}_k^X]$. Since the optimal decision rule is linear, the error exponent is invariant with linear transformations of statistics, i.e.,

$$\begin{aligned} E_{\tilde{h}^k}(v, v') &= E_{\tilde{h}^k}(v, v') \\ &= \frac{1}{8} \left\| (\phi_v^{X|V} - \phi_{v'}^{X|V})^T \tilde{\Xi}^X \right\|^2 + o(\epsilon^2) \\ &= \frac{1}{8} \left\| (\phi_v^{Y|V} - \phi_{v'}^{Y|V})^T \tilde{\mathbf{B}} \tilde{\Xi}^X \right\|^2 + o(\epsilon^2), \end{aligned} \quad (36)$$

where the last equality follows from Lemma 5. Taking the expectation over a given RIE yields

$$\begin{aligned} \mathbb{E}[E_{\tilde{h}^k}(v, v')] &= \frac{1}{8} \mathbb{E} \left[\left\| (\phi_v^{Y|V} - \phi_{v'}^{Y|V})^T \tilde{\mathbf{B}} \tilde{\Xi}^X \right\|^2 \right] + o(\epsilon^2) \\ &= \frac{\mathbb{E} \left[\left\| \phi_v^{Y|V} - \phi_{v'}^{Y|V} \right\|^2 \right]}{8|\mathcal{Y}|} \left\| \tilde{\mathbf{B}} \tilde{\Xi}^X \right\|_{\text{F}}^2 + o(\epsilon^2), \end{aligned}$$

where we have exploited Lemma 6. Finally, the error exponent (3) can be obtained via noting

$$\tilde{\Xi}^X = \Xi^X ((\Xi^X)^T \Xi^X)^{-\frac{1}{2}}$$

from the definition of \tilde{f}^k .

□

B. Proof of Lemma 2

We first prove two useful lemmas.

Lemma 7. For distributions $P \in \text{relint}(\mathcal{P}^{\mathcal{X}})$, $Q, R \in \mathcal{P}^{\mathcal{X}}$, and sufficiently small ϵ , if $D(P\|Q) \leq \epsilon^2$ and $D(P\|R) \leq \epsilon^2$, then there exists a constant $C > 0$ independent of ϵ , such that $D(Q\|R) \leq C\epsilon^2$.

Proof. Denote by $\|\cdot\|_1$ the ℓ_1 -distance between distributions, i.e., $\|P - Q\|_1 \triangleq \sum_{x \in \mathcal{X}} |P(x) - Q(x)|$, then from Pinsker's inequality [20], we have

$$\|P - Q\|_1 \leq \sqrt{2D(P\|Q)} < \sqrt{2}\epsilon, \quad (37a)$$

$$\|P - R\|_1 \leq \sqrt{2D(P\|R)} < \sqrt{2}\epsilon, \quad (37b)$$

which implies

$$\|Q - R\|_1 \leq \|P - Q\|_1 + \|P - R\|_1 \leq 2\sqrt{2}\epsilon. \quad (38)$$

In addition, with the notation $p_{\min} \triangleq \min_{x \in \mathcal{X}} P(x)$, for all $x \in \mathcal{X}$ we have

$$R(x) > P(x) - |P(x) - R(x)| \quad (39a)$$

$$> \min_{x \in \mathcal{X}} P(x) - \sqrt{2}\epsilon \quad (39b)$$

$$= p_{\min} - \sqrt{2}\epsilon, \quad (39c)$$

where to obtain (39b) we have used (37b). Note that $p_{\min} > 0$ since $P \in \text{relint}(\mathcal{P}^{\mathcal{X}})$, thus $R(x) > p_{\min}/2$ for sufficiently small ϵ . As a result,

$$D(Q\|R) \leq \sum_{x \in \mathcal{X}} \frac{(Q(x) - R(x))^2}{R(x)} \quad (40a)$$

$$\leq \frac{2}{p_{\min}} \sum_{x \in \mathcal{X}} [Q(x) - R(x)]^2 \quad (40b)$$

$$\leq \frac{2\|Q - R\|_1^2}{p_{\min}} \quad (40c)$$

$$\leq \frac{16}{p_{\min}} \epsilon^2, \quad (40d)$$

where to obtain (40a) we have applied an upper bound of K-L divergence [21], and to obtain (40d) we have used (38). □

Lemma 8. For all $(x, y) \in \mathcal{X} \times \mathcal{Y}$, we have

$$\begin{aligned} D(P_X P_Y \| P_X \tilde{P}_{Y|X}^{(v,b)}) \\ \geq P_X(x) \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1 - P_Y(y)} \tau(x,y)} \right] \end{aligned}$$

where $\tilde{P}_{Y|X}^{(v,b)}$ is defined in (4) and $\tau(x, y)$ is defined as $\tau(x, y) \triangleq \tilde{v}^T(y)s(x) + \tilde{d}(y)$.

Proof. First, we can rewrite the conditional distribution $\tilde{P}_{Y|X}^{(v,b)}(y|x)$ as

$$\begin{aligned} \tilde{P}_{Y|X}^{(v,b)}(y|x) &= \frac{e^{v^T(y)s(x)+b(y)}}{\sum_{y' \in \mathcal{Y}} e^{v^T(y')s(x)+b(y')}} \\ &= \frac{P_Y(y) e^{v^T(y)s(x)+d(y)}}{\sum_{y' \in \mathcal{Y}} P_Y(y') e^{v^T(y')s(x)+d(y')}} \\ &= \frac{P_Y(y) e^{\tilde{v}^T(y)s(x)+\tilde{d}(y)}}{\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tilde{v}^T(y')s(x)+\tilde{d}(y')}} \\ &= \frac{P_Y(y) e^{\tau(x,y)}}{\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tau(x,y')}}. \end{aligned} \quad (41)$$

Then the K-L divergence $D(P_X P_Y \| P_X \tilde{P}_{Y|X}^{(v,b)})$ can be expressed as

$$\begin{aligned} & D(P_X P_Y \| P_X \tilde{P}_{Y|X}^{(v,b)}) \\ &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} P_X(x) P_Y(y) \log \frac{\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tau(x,y')}}{e^{\tau(x,y)}} \\ &= \sum_{x \in \mathcal{X}} P_X(x) \log \left[\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tau(x,y')} \right] - \mathbb{E}_{P_X P_Y} [\tau(X, Y)] \\ &= \sum_{x \in \mathcal{X}} P_X(x) \log \left[\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tau(x,y')} \right], \end{aligned} \quad (42)$$

where to obtain the last equality we have used the fact $\mathbb{E}_{P_X P_Y} [\tau(X, Y)] = 0$. As a result, we have

$$D(P_X P_Y \| P_X \tilde{P}_{Y|X}^{(v,b)}) \quad (43a)$$

$$\geq P_X(x) \log \left[\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tau(x,y')} \right] \quad (43b)$$

$$\geq P_X(x) \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)} \right], \quad (43c)$$

where (43c) follows from Jensen's inequality:

$$\begin{aligned} & \sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tau(x,y')} \\ &= P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) \sum_{y' \neq y} \frac{P_Y(y')}{1 - P_Y(y)} e^{\tau(x,y')} \\ &\geq P_Y(y) e^{\tau(x,y)} \\ &\quad + (1 - P_Y(y)) \exp \left(\frac{1}{1 - P_Y(y)} \sum_{y' \neq y} P_Y(y') \tau(x,y') \right) \\ &= P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)}. \end{aligned}$$

□

With the above lemmas, Lemma 2 can be proved as follows.

Proof of Lemma 2. Note that when $v = d = 0$, we have $\tilde{P}_{Y|X}^{(v,b)} = P_Y$. As a result, the optimal v, d for (6) satisfy

$$\begin{aligned} & D(P_{XY} \| P_X \tilde{P}_{Y|X}^{(v,b)}) \\ &\leq D(P_{XY} \| P_X P_Y) \\ &\leq \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \frac{[P_{X,Y}(x,y) - P_X(x)P_Y(y)]^2}{P_X(x)P_Y(y)} \quad (44) \\ &\leq \epsilon^2, \end{aligned}$$

where the second inequality again follows from the upper bound for K-L divergence [21], and the last inequality follows from the definition of ϵ -dependency.

As $P_{XY} \in \text{relint}(\mathcal{P}^{\mathcal{X} \times \mathcal{Y}})$, from Lemma 7, there exist $C > 0$ and $\epsilon_1 > 0$ such that $D(P_X P_Y \| P_X \tilde{P}_{Y|X}^{(v,b)}) < C\epsilon^2$ for all $\epsilon < \epsilon_1$. Furthermore, from Lemma 8, for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$

and $\epsilon \in (0, \epsilon_1)$, we have

$$C\epsilon^2 \geq P_X(x) \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)} \right]. \quad (45)$$

Since

$$\begin{aligned} & \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)} \right] \\ &= \frac{P_Y(y)}{2(1 - P_Y(y))} \tau^2(x,y) + o(\tau^2(x,y)), \end{aligned}$$

there exists $\delta > 0$ independent of ϵ_1 , such that for all $|\tau(x,y)| \leq \delta$, we have

$$\begin{aligned} & \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)} \right] \\ &> \frac{P_Y(y)}{2} \tau^2(x,y). \end{aligned} \quad (46)$$

In turn, if $|\tau(x,y)| > \delta$, we have

$$\begin{aligned} & \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)} \right] \\ &\geq \min \left\{ \log \left[P_Y(y) e^\delta + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \delta} \right], \right. \\ &\quad \left. \log \left[P_Y(y) e^{-\delta} + (1 - P_Y(y)) e^{\frac{P_Y(y)}{1-P_Y(y)} \delta} \right] \right\}, \\ &\geq \frac{P_Y(y)}{2} \delta^2, \end{aligned}$$

where to obtain the second inequality we have exploited the monotonicity of function $P_Y(y) e^t + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} t}$, and to obtain the third inequality we have exploited (46).

As a result, we have

$$\begin{aligned} & \log \left[P_Y(y) e^{\tau(x,y)} + (1 - P_Y(y)) e^{-\frac{P_Y(y)}{1-P_Y(y)} \tau(x,y)} \right] \\ &> \frac{P_Y(y)}{2} \cdot \min\{\delta^2, \tau^2(x,y)\}, \end{aligned} \quad (47)$$

hence (45) becomes

$$C\epsilon^2 \geq \frac{P_X(x)P_Y(y)}{2} \cdot \min\{\delta^2, \tau^2(x,y)\}, \quad (48)$$

from which we obtain $\tau(x,y) = O(\epsilon)$. Indeed, let $\epsilon_2 \triangleq \frac{\delta}{\sqrt{2C}} \cdot \min_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \sqrt{P_X(x)P_Y(y)}$, $\epsilon_0 \triangleq \min\{\epsilon_1, \epsilon_2\}$, then for all $\epsilon < \epsilon_0$, we have

$$C\epsilon^2 < \frac{P_X(x)P_Y(y)}{2} \cdot \delta^2,$$

and (48) implies $|\tau(x,y)| < C'\epsilon$ with $C' = \sqrt{\frac{2C}{P_X(x)P_Y(y)}}$. □

C. Proof of Lemma 3

Proof. From Lemma 2, there exists $C' > 0$ such that for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$, we have

$$|\tilde{v}^T(y)s(x) + \tilde{d}(y)| < C'\epsilon, \quad (49)$$

which implies

$$|\mu_s^T \tilde{v}(y) + \tilde{d}(y)| < C\epsilon, \quad (50)$$

$$|\tilde{v}^T(y)\tilde{s}(x)| < 2C\epsilon, \quad (51)$$

with $C = \max\{C', 1\}$.

From (41), we can assume $\mathbb{E}_{P_Y} [v(Y)] = \mathbb{E}_{P_Y} [d(Y)] = 0$ without loss of generality. Then (4) can be rewritten as

$$\tilde{P}_{Y|X}^{(v,b)}(y|x) = \frac{P_Y(y) e^{\tilde{v}^T(y)s(x) + \tilde{d}(y)}}{\sum_{y' \in \mathcal{Y}} P_Y(y') e^{\tilde{v}^T(y')s(x) + \tilde{d}(y')}}, \quad (52)$$

and the numerator has the approximation

$$\begin{aligned} & P_Y(y)e^{\tilde{v}^T(y)s(x)+\tilde{d}(y)} \\ &= P_Y(y) \left(1 + \tilde{v}^T(y)s(x) + \tilde{d}(y) + o(\epsilon) \right) \\ &= P_Y(y) \left(1 + \tilde{v}^T(y)s(x) + \tilde{d}(y) \right) + o(\epsilon), \end{aligned}$$

where we have used (49). Similarly, from

$$\begin{aligned} & \sum_{y' \in \mathcal{Y}} P_Y(y')e^{\tilde{v}^T(y')s(x)+\tilde{d}(y')} \\ &= \sum_{y' \in \mathcal{Y}} P_Y(y') \left(1 + \tilde{v}^T(y')s(x) + \tilde{d}(y') \right) + o(\epsilon) \\ &= 1 + \mathbb{E}_{P_Y} \left[\tilde{v}^T(Y)s(x) \right] + \mathbb{E}_{P_Y} \left[\tilde{d}(Y) \right] + o(\epsilon) \\ &= 1 + o(\epsilon) \end{aligned}$$

we obtain

$$\frac{1}{\sum_{y' \in \mathcal{Y}} P_Y(y')e^{\tilde{v}^T(y')s(x)+\tilde{d}(y')}} = \frac{1}{1 + o(\epsilon)} = 1 + o(\epsilon).$$

As a result, (52) has the approximation

$$\begin{aligned} & \tilde{P}_{Y|X}^{(v,b)}(y|x) \\ &= \left[P_Y(y) \left(1 + \tilde{v}^T(y)s(x) + \tilde{d}(y) \right) + o(\epsilon) \right] [1 + o(\epsilon)] \quad (53) \\ &= P_Y(y) \left(1 + \tilde{v}^T(y)s(x) + \tilde{d}(y) \right) + o(\epsilon), \end{aligned}$$

which implies $P_X \tilde{P}_{Y|X}^{(v,b)} \in \mathcal{N}_{C_\epsilon}^{\mathcal{X} \times \mathcal{Y}}(P_X P_Y)$ for sufficiently small ϵ . Besides, the local assumption of distributions implies that $P_{X,Y} \in \mathcal{N}_{C_\epsilon}^{\mathcal{X} \times \mathcal{Y}}(P_X P_Y) \subset \mathcal{N}_{C_\epsilon}^{\mathcal{X} \times \mathcal{Y}}(P_X P_Y)$. Again, from the local approximation of K-L divergence [3]

$$D(P_1 \| P_2) = \frac{1}{2} \|\phi_1 - \phi_2\|^2 + o(\epsilon^2), \quad (54)$$

we have

$$\begin{aligned} & D(P_{Y,X} \| P_X \tilde{P}_{Y|X}^{(v,b)}) \\ &= \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{\left[P_{Y,X}(y, x) - \tilde{P}_{Y|X}^{(v,b)}(y|x)P_X(x) \right]^2}{P_Y(y)P_X(x)} + o(\epsilon^2) \\ &= \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left[\frac{P_{Y,X}(y, x)}{\sqrt{P_Y(y)P_X(x)}} - \sqrt{P_Y(y)P_X(x)} \right. \\ &\quad \left. - \sqrt{P_Y(y)P_X(x)} \left(\tilde{v}^T(y)s(x) + \tilde{d}(y) + o(\epsilon) \right) \right]^2 + o(\epsilon^2) \\ &= \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left[\tilde{\mathbf{B}}(y, x) - \sqrt{P_Y(y)P_X(x)} \tilde{v}^T(y) \tilde{s}(x) \right. \\ &\quad \left. - \sqrt{P_Y(y)P_X(x)} \left(\tilde{d}(y) + \mu_s^T \tilde{v}(y) \right) \right. \\ &\quad \left. - \sqrt{P_Y(y)P_X(x)} o(\epsilon) \right]^2 + o(\epsilon^2) \\ &\stackrel{(*)}{=} \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left[\tilde{\mathbf{B}}(y, x) - \sqrt{P_Y(y)P_X(x)} \tilde{v}^T(y) \tilde{s}(x) \right]^2 \\ &\quad + \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left[\sqrt{P_Y(y)P_X(x)} \left(\tilde{d}(y) + \mu_s^T \tilde{v}(y) \right) \right]^2 \\ &\quad + o(\epsilon^2) \\ &= \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \left[\tilde{\mathbf{B}}(y, x) - (\xi^Y(y))^T \xi^X(x) \right]^2 \end{aligned}$$

$$\begin{aligned} & + \frac{1}{2} \mathbb{E}_{P_Y} \left[\left(\tilde{d}(y) + \mu_s^T \tilde{v}(y) \right)^2 \right] + o(\epsilon^2) \\ &= \frac{1}{2} \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_{\mathbb{F}}^2 + \frac{1}{2} \eta^{(v,b)}(s) + o(\epsilon^2), \end{aligned}$$

where to obtain (*) we have used (50)-(51) together with the fact $|\tilde{\mathbf{B}}(y, x)| < \epsilon$, and

$$\begin{aligned} & \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \tilde{\mathbf{B}}(y, x) \sqrt{P_Y(y)P_X(x)} \left(\tilde{d}(y) + \mu_s^T \tilde{v}(y) \right) = 0, \\ & \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_Y(y)P_X(x) \tilde{v}^T(y) \tilde{s}(x) \left(\tilde{d}(y) + \mu_s^T \tilde{v}(y) \right) = 0, \end{aligned}$$

since $\mathbb{E} \left[\tilde{d}(Y) \right] = 0, \mathbb{E} \left[\tilde{s}(X) \right] = \mathbb{E} \left[\tilde{v}(Y) \right] = 0$.

□

D. Proofs of Theorem 2 and Theorem 3

Theorem 2 and Theorem 3 can be proved based on Lemma 3.

Proofs of Theorem 2 and Theorem 3. Note that the value of $d(\cdot)$ only affects the second term of the K-L divergence, hence we can always choose $d(\cdot)$ such that $\tilde{d}(y) + \mu_s^T \tilde{v}(y) = 0$. Then the (Ξ^Y, Ξ^X) pair should be chosen as

$$(\Xi^Y, \Xi^X)^* = \arg \min_{(\Xi^Y, \Xi^X)} \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_{\mathbb{F}}^2. \quad (55)$$

Set the derivative²

$$\frac{\partial}{\partial \Xi^Y} \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_{\mathbb{F}}^2 = 2(\Xi^Y (\Xi^X)^T \Xi^X - \tilde{\mathbf{B}} \Xi^X) \quad (56)$$

to zero, and the optimal Ξ^Y for fixed Ξ^X is³

$$\Xi^{Y*} = \tilde{\mathbf{B}} \Xi^X \left((\Xi^X)^T \Xi^X \right)^{-1}. \quad (57)$$

As $\mathbf{1}^T \sqrt{P_Y} \tilde{\mathbf{B}} = 0$, we have $\mathbf{1}^T \sqrt{P_Y} \Xi^{Y*} = 0$, which demonstrates that Ξ^{Y*} is a valid matrix for a zero-mean feature vector.

To express Ξ^{Y*} of (57) in the form of s and v , we can make use of the correspondence between feature and information vectors. Note that, for a zero-mean feature function $f(X)$ with corresponding information vector ϕ , we have the correspondence $\mathbb{E}_{P_{X|Y}} [f(X)|Y] \leftrightarrow \tilde{\mathbf{B}}\phi$ since the y -th element of information vector $\tilde{\mathbf{B}}\phi$ is given by

$$\begin{aligned} & \sum_{x \in \mathcal{X}} \tilde{\mathbf{B}}(y, x) \phi(x) \\ &= \sum_{x \in \mathcal{X}} \frac{P_{X,Y}(x, y) - P_X(x)P_Y(y)}{\sqrt{P_X(x)P_Y(y)}} f(x) \sqrt{P_X(x)} \\ &= \frac{1}{\sqrt{P_Y(y)}} \sum_{x \in \mathcal{X}} P_{X,Y}(x, y) f(x) \\ &= \frac{1}{\sqrt{P_Y(y)}} \mathbb{E}_{P_{X|Y}} [f(X)|Y = y]. \end{aligned}$$

Using similar methods, we can verify that $\Lambda_{\tilde{s}(X)} = (\Xi^X)^T \Xi^X$.

As a result, (57) is equivalent to

$$\tilde{v}^*(y) = \mathbb{E}_{P_{X|Y}} \left[\Lambda_{\tilde{s}(X)}^{-1} \tilde{s}(X) \mid Y = y \right]. \quad (58)$$

²In this paper, we use the denominator-layout notation of matrix calculus where the scalar-by-matrix derivative will have the same dimension as the matrix.

³Here we assume the matrix $(\Xi^X)^T \Xi^X$, i.e., $\Lambda_{\tilde{s}(X)}$ is invertible. For the case where $(\Xi^X)^T \Xi^X$ is singular, all conclusions are the same when we use Moore-Penrose inverse to replace ordinary matrix inverse.

□

By symmetry, the first two equations of Theorem 3 can be proved using the same method. To obtain the third equations of these two theorems, we need to minimize $\eta^{(v,b)}(s) = \mathbb{E}_{P_Y} \left[(\mu_s^T \tilde{v}(Y) + \tilde{d}(Y))^2 \right]$. When \tilde{v} and μ_s are fixed, the optimal \tilde{d} is

$$\tilde{d}^*(y) = -\mu_s^T \tilde{v}(Y), \quad (59)$$

and the corresponding $\eta^{(v,b)}(s) = 0$.

When \tilde{d} and \tilde{v} are fixed, we have

$$\begin{aligned} & \eta^{(v,b)}(s) \\ &= \mathbb{E}_{P_Y} \left[(\mu_s^T \tilde{v}(Y) + \tilde{d}(Y))^2 \right] \\ &= \mu_s^T \mathbf{\Lambda}_{\tilde{v}(Y)} \mu_s + 2\mu_s^T \mathbb{E}_{P_Y} [\tilde{v}(Y)\tilde{d}(Y)] + \text{var}(\tilde{d}(Y)). \end{aligned} \quad (60)$$

Set $\frac{\partial}{\partial \mu_s} \eta^{(v,b)}(s) = 0$ and we obtain

$$\mu_s^* = -\mathbf{\Lambda}_{\tilde{v}(Y)}^{-1} \mathbb{E}_{P_Y} [\tilde{v}(Y)\tilde{d}(Y)]. \quad (61)$$

□

E. Proof of Theorem 4

Proof. From Lemma 3, choosing the optimal (Ξ^Y, Ξ^X) is equivalent to solving the matrix factorization problem of $\tilde{\mathbf{B}}$. Since both Ξ^Y and Ξ^X have rank no greater than k , from the Eckart-Young-Mirsky theorem [22], the optimal choice of $\Xi^Y (\Xi^X)^T$ should be the truncated singular value decomposition of $\tilde{\mathbf{B}}$ with top k singular values. As a result, $(\Xi^Y, \Xi^X)^*$ are the left and right singular vectors of $\tilde{\mathbf{B}}$ corresponding to the largest k singular values.

The optimality of bias $\tilde{d}(y) = -\mu_s^T \tilde{v}(y)$ has already been shown in Appendix D. □

F. Proof of Theorem 5

The following lemma is useful to prove Theorem 5.

Lemma 9 (Pythagorean theorem). *Let Ξ^{X^*} be the optimal matrix for given Ξ^Y as defined in (11). Then,*

$$\begin{aligned} & \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2 - \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T\|_{\text{F}}^2 \\ &= \|\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2. \end{aligned} \quad (62)$$

Proof of Lemma 9. Denote by $\langle \mathbf{U}, \mathbf{V} \rangle$ the Frobenius inner product of matrices \mathbf{U} and \mathbf{V} , i.e., $\langle \mathbf{U}, \mathbf{V} \rangle \triangleq \text{tr}(\mathbf{U}^T \mathbf{V})$, and we have

$$\begin{aligned} & \langle \tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T, \Xi^Y (\Xi^X)^T \rangle \\ &= \text{tr} \left(\tilde{\mathbf{B}} \Xi^X (\Xi^Y)^T \right) - \text{tr} \left(\Xi^{X^*} (\Xi^Y)^T \Xi^Y (\Xi^X)^T \right) \\ &= \text{tr} \left(\tilde{\mathbf{B}} \Xi^X (\Xi^Y)^T \right) - \text{tr} \left(\tilde{\mathbf{B}}^T \Xi^Y (\Xi^X)^T \right) \\ &= 0. \end{aligned}$$

As a result,

$$\begin{aligned} & \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2 \\ &= \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T + (\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T)\|_{\text{F}}^2 \\ &= \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T\|_{\text{F}}^2 + \|\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2 \\ & \quad + 2 \langle \tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T, \Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T \rangle \\ &= \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T\|_{\text{F}}^2 + \|\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2. \end{aligned}$$

Proof of Theorem 5. From Lemma 9, we have

$$\begin{aligned} & \mathcal{L}(s) - \mathcal{L}(s^*) \\ &= \frac{1}{2} \left[\|\tilde{\mathbf{B}} - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2 - \|\tilde{\mathbf{B}} - \Xi^Y (\Xi^{X^*})^T\|_{\text{F}}^2 \right] \\ & \quad + \frac{1}{2} \left[\eta^{(v,b)}(s) - \eta^{(v,b)}(s^*) \right] + o(\epsilon^2) \\ &= \frac{1}{2} \|\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2 + \frac{1}{2} \kappa^{(v,b)}(s, s^*) + o(\epsilon^2), \end{aligned}$$

where $\kappa^{(v,b)}(s, s^*) \triangleq \eta^{(v,b)}(s) - \eta^{(v,b)}(s^*)$. We then optimize $\|\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2$ and $\kappa^{(v,b)}(s, s^*)$ separately.

For the first term, we need to express Ξ^X in terms of \mathbf{W} and Ξ_1^X . From (15) we obtain

$$\mathbb{E}[s_z(X)] = \sigma(c(z)) + o(\epsilon), \quad (63a)$$

$$\tilde{s}_z(x) = w^T(z) \tilde{t}(x) \cdot \sigma'(c(z)) + o(\epsilon), \quad (63b)$$

which can be expressed in information vectors as

$$\Xi^X = \Xi_1^X \mathbf{W}^T \mathbf{J} + o(\epsilon). \quad (64)$$

From Theorem 3, we have

$$\Xi^{X^*} = \tilde{\mathbf{B}}^T \Xi^Y ((\Xi^Y)^T \Xi^Y)^{-1}. \quad (65)$$

As a result, we have

$$\begin{aligned} & \|\Xi^Y (\Xi^{X^*})^T - \Xi^Y (\Xi^X)^T\|_{\text{F}}^2 \\ &= \|((\Xi^Y)^T \Xi^Y)^{1/2} ((\Xi^{X^*})^T - (\Xi^X)^T)\|_{\text{F}}^2 \\ &= \left\| ((\Xi^Y)^T \Xi^Y)^{1/2} \cdot \left((\Xi^{X^*})^T - \mathbf{J} \mathbf{W} (\Xi_1^X)^T - o(\epsilon) \right) \right\|_{\text{F}}^2 \\ &= \left\| ((\Xi^Y)^T \Xi^Y)^{1/2} \cdot \left((\Xi^{X^*})^T - \mathbf{J} \mathbf{W} (\Xi_1^X)^T \right) \right\|_{\text{F}}^2 + o(\epsilon^2) \\ &= \left\| ((\Xi^Y)^T \Xi^Y)^{1/2} \mathbf{J} \right. \\ & \quad \cdot \left. \left(\mathbf{J}^{-1} (\Xi^{X^*})^T - \mathbf{W} (\Xi_1^X)^T \right) \right\|_{\text{F}}^2 + o(\epsilon^2) \\ &= \|\Theta \tilde{\mathbf{B}}_1 - \Theta \mathbf{W} (\Xi_1^X)^T\|_{\text{F}}^2 + o(\epsilon^2), \end{aligned} \quad (66)$$

where the third equality follows from the fact that [cf. (51)] $\tilde{s}(x) = O(\epsilon)$ and $\tilde{v}(y) = O(1)$, and the last equality follows from the definitions $\tilde{\mathbf{B}}_1 \triangleq \mathbf{J}^{-1} (\Xi^{X^*})^T$ and $\Theta \triangleq ((\Xi^Y)^T \Xi^Y)^{1/2} \mathbf{J}$.

For the second term, from (60) and (61), we have

$$\begin{aligned} & \kappa^{(v,b)}(s, s^*) \\ &= [(\mu_s - \mu_{s^*}) + \mu_{s^*}]^T \mathbf{\Lambda}_{\tilde{v}(Y)} [(\mu_s - \mu_{s^*}) + \mu_{s^*}] \\ & \quad - \mu_{s^*}^T \mathbf{\Lambda}_{\tilde{v}(Y)} \mu_{s^*} + 2(\mu_s - \mu_{s^*})^T \mathbb{E}_{P_Y} [\tilde{v}(Y)\tilde{d}(Y)] \\ &= (\mu_s - \mu_{s^*})^T \mathbf{\Lambda}_{\tilde{v}(Y)} (\mu_s - \mu_{s^*}) \\ & \quad + 2(\mu_s - \mu_{s^*})^T \left(\mathbf{\Lambda}_{\tilde{v}(Y)} \mu_{s^*} + \mathbb{E}_{P_Y} [\tilde{v}(Y)\tilde{d}(Y)] \right) \\ &= (\mu_s - \mu_{s^*})^T \mathbf{\Lambda}_{\tilde{v}(Y)} (\mu_s - \mu_{s^*}). \end{aligned} \quad (67)$$

Combining (66) and (67) finishes the proof. □

The results of μ_s^* and w^* can be obtained via minimizing the loss function \mathcal{L} . Again, these two terms can be optimized separately. To obtain μ_s^* , consider the case where the hidden layer has used a bounded activation function, i.e., $\sigma_{\min} \preceq \mu_s \preceq$

σ_{\max} , such as sigmoid function $1/(1+e^{-x})$ or $\tanh(x)$. Then the optimal μ_s is the solution of

$$\begin{aligned} & \underset{\mu_s}{\text{minimize}} \quad (\mu_s - \mu_{s^*})^T \mathbf{\Lambda}_{\tilde{v}(Y)} (\mu_s - \mu_{s^*}) \\ & \text{subject to} \quad \sigma_{\min} \preceq \mu_s \preceq \sigma_{\max}. \end{aligned} \quad (68)$$

If μ_{s^*} satisfies the constraint of (68), then it is the optimal solution. Otherwise, some elements of μ_{s^*} will become either σ_{\min} or σ_{\max} , which is known as the saturation phenomenon.

Further, from (63a), the bias $c(z)$ of hidden layer is⁴

$$c(z) = \sigma^{-1}(\mu_{s^*}(z)) + o(\epsilon).$$

To obtain \mathbf{W}^* , let

$$\begin{aligned} \tilde{\mathbf{B}}_1' & \triangleq \Theta \tilde{\mathbf{B}}_1 = ((\Xi^Y)^T \Xi^Y)^{-1/2} (\Xi^Y)^T \tilde{\mathbf{B}}, \\ \mathbf{W}' & \triangleq \Theta \mathbf{W} = ((\Xi^Y)^T \Xi^Y)^{1/2} \mathbf{J} \mathbf{W}, \end{aligned} \quad (69)$$

then the optimal \mathbf{W}' is the solution of

$$\underset{\mathbf{W}'}{\text{minimize}} \quad \|\tilde{\mathbf{B}}_1' - \mathbf{W}' (\Xi_1^X)^T\|_F^2, \quad (70)$$

i.e.,

$$\mathbf{W}^* = \tilde{\mathbf{B}}_1' \Xi_1^X ((\Xi_1^X)^T \Xi_1^X)^{-1}. \quad (71)$$

Hence, \mathbf{W}^* is given by

$$\begin{aligned} \mathbf{W}^* & = \Theta^{-1} \mathbf{W}'^* = \Theta^{-1} \tilde{\mathbf{B}}_1' \Xi_1^X ((\Xi_1^X)^T \Xi_1^X)^{-1} \\ & = \tilde{\mathbf{B}}_1 \Xi_1^X ((\Xi_1^X)^T \Xi_1^X)^{-1} \\ & = \mathbf{J}^{-1} \cdot [\Xi^Y ((\Xi^Y)^T \Xi^Y)^{-1}]^T \tilde{\mathbf{B}} \Xi_1^X ((\Xi_1^X)^T \Xi_1^X)^{-1}, \end{aligned}$$

where the term $\tilde{\mathbf{B}} \Xi_1^X ((\Xi_1^X)^T \Xi_1^X)^{-1}$ corresponds to a feature projection of $\tilde{t}(X)$:

$$\tilde{\mathbf{B}} \Xi_1^X ((\Xi_1^X)^T \Xi_1^X)^{-1} \leftrightarrow \mathbb{E}_{P_{X|Y}} \left[\mathbf{\Lambda}_{\tilde{t}(X)}^{-1} \tilde{t}(X) \mid Y \right]. \quad (72)$$

As a consequence, this multi-layer neural network is conducting a generalized feature projection between features extracted from different layers. In practice problems, the projected feature $\mathbb{E}_{P_{\tilde{t}|Y}} \left[\mathbf{\Lambda}_{\tilde{t}}^{-1} \tilde{t} \mid Y \right]$ depends only on the distribution $P_{\tilde{t}|Y}$, and does not depend on the distribution $P_{X|Y}$. Therefore, the above computations can be accomplished without knowing the hidden random variable X and can be applied to general cases.

REFERENCES

- [1] D. J. C. MacKay, *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, ISBN 9780521642989, 2003.
- [2] N. Tishby and N. Zaslavsky, "Deep learning and the information bottleneck principle," in *Information Theory Workshop (ITW), 2015 IEEE*. IEEE, 2015, pp. 1–5.
- [3] S.-L. Huang, A. Makur, L. Zheng, and G. W. Wornell, "On universal features for high-dimensional learning and inference," *submitted to IEEE Trans. Inform. Theory*, 2019. Preprint.
- [4] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," *arXiv preprint arXiv:1701.07875*, 2017.
- [5] H. O. Hirschfeld, "A connection between correlation and contingency," *Proc. Cambridge Phil. Soc.*, vol. 31, pp. 520–524, 1935.
- [6] H. Gebelein, "Das statistische problem der korrelation als variations-und eigenwertproblem und sein zusammenhang mit der ausgleichungsrechnung," *Z. für angewandte Math., Mech.*, vol. 21, pp. 364–379, 1941.
- [7] A. Rényi, "On measures of dependence," *Acta Mathematica Academiae Scientiarum Hungarica*, vol. 10, no. 3–4, pp. 441–451, 1959.
- [8] L. Breiman and J. H. Friedman, "Estimating optimal transformations for multiple regression and correlation," *Journal of American Statistical Association*, vol. 80, no. 391, pp. 614–619, 1985.
- [9] J. Stoer and R. Bulirsch, *Introduction to numerical analysis*. Springer Science & Business Media, 2013, vol. 12.

⁴When $\mu_t \neq 0$, the formula should be modified as $c(z) = \sigma^{-1}(\mu_{s^*}(z)) - \frac{\mu_t}{\sigma} w + o(\epsilon)$.

- [10] G. Alain and Y. Bengio, "Understanding intermediate layers using linear classifier probes," *arXiv preprint arXiv:1610.01644*, 2016.
- [11] R. Olga, D. Jia, S. Hao, K. Jonathan, S. Sanjeev, M. Sean, H. Zhiheng, K. Andrej, K. Aditya, B. Michael, C. B. Alexander, and F.-F. Li, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [12] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [13] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI*, 2017, pp. 4278–4284.
- [14] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," *arXiv preprint arXiv:1610.02357*, 2016.
- [15] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.
- [16] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [17] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, vol. 1, no. 2, 2017, p. 3.
- [18] H. Akaike, "Information theory and an extension of the maximum likelihood principle," in *Selected Papers of Hirotugu Akaike*. Springer, 1998, pp. 199–213.
- [19] A. Dembo and O. Zeitouni, "Large deviations techniques and applications. corrected reprint of the second (1998) edition. stochastic modelling and applied probability, 38," 2010.
- [20] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.
- [21] Y. Polyanskiy and Y. Wu, "Lecture notes on information theory," *Lecture Notes for ECE563 (UIUC) and*, vol. 6, pp. 2012–2016, 2014.
- [22] C. Eckart and G. Young, "The approximation of one matrix by another of lower rank," *Psychometrika*, vol. 1, no. 3, pp. 211–218, 1936.